

Distortions and Learning

Richard Nock



Centre d'**E**tudes et de **R**echerches en **E**conomie,
Gestion, **M**athématiques et **I**nformatique **A**ppiquée

Université des Antilles-Guyane
Richard.Nock@martinique.univ-ag.fr

May 2007

Talk aims at showing **some ties** that (Computational) Geometries share with machine learning:

- 1 some are obvious (3% of the talk),
- 2 others are related to some of the most fundamental questions of machine learning:

What is learning ? How can an Algorithm “learn” ? (97%)

Geometries generally non Euclidean, basically non Riemannian.

Slides available at

<http://www.univ-ag.fr/~rnock/Slides/Geotopal07/>

- 1 (What is) *Supervised Learning* ?
- 2 Learning Strategies
- 3 **Boosting** algorithms
- 4 Distortions for Learning: Bregman Divergences
- 5 Axiomatization (of Supervised Learning)
- 6 Minimization Algorithms
- 7 Related Questions and Perspectives

- 1 What is Supervised learning ?
- 2 Key components, quantities ?
- 3 Formal models of Supervised learning ?

Supervised Learning > Example

Tic-tac-toe endgame configurations, X vs O, X has played first.

X		O
X	X	X
O		O

wins for X

X		X
O	O	O
	X	

,

X	O	X
X	O	O
O	X	X

do not win for X

Input

Output

Problem

Suppose you do not know the game.

*(How) can you infer an **accurate function** Input \rightarrow Output ?*

Supervised Learning > Key concepts

Example: (\mathbf{x}, y)

Ground information: each endgame configuration.

X		O
X	X	X
O		O

wins for X

is represented by: $(\underbrace{\text{XBOXXXOBO}}_{\text{observation}}, \underbrace{+1}_{\text{class}})$
example

Remarks:

- 1 Label = synonym of class ($+1$ = wins for X; -1 = does not win for X)
- 2 We assume two classes wlog

Supervised Learning > Key concepts (contd)

Domain: \mathcal{X}

Let \mathcal{X} denote the set of all board configurations (over $n = 9$ description variables). Some are endgames, some are not. Some are even not admissible for tic-tac-toe (xxxxxxxxxx).

Distributions: D

Example are drawn i.i.d., from some fixed and unknown distribution D over $\mathcal{X} \times \{-1, +1\}$. Let weight vector $\mathbf{w}_1 = \hat{D}$ over \mathcal{S} .

Learning Sample: \mathcal{S}

Suppose someone takes m endgame configurations (among the 958 distinct possible, according to D) and labels them. Let

$$\mathcal{S} \stackrel{\text{def}}{=} \{(\mathbf{x}_i, y_i^*) : \mathbf{x}_i \in \mathcal{X}, y_i^* \in \{-1, +1\}\}_{i=1}^m$$

be this set ($+1$ = wins for x, -1 = does not win for x).

Supervised Learning > Key concepts (contd)

Classifier: H

Any function $H : \mathcal{X} \rightarrow \mathbb{S} \subseteq \mathbb{R}$. For example $\mathbb{S} = \{-1, +1\}, \mathbb{R}, \dots$

$$\mathbf{x} \in \mathcal{X} \stackrel{\text{def}}{=} \left(\begin{array}{|c|c|c|} \hline v_1 & v_2 & v_3 \\ \hline v_4 & v_5 & v_6 \\ \hline v_7 & v_8 & v_9 \\ \hline \end{array} \right)$$

Example 1 Rule $H_1(\mathbf{x}) \stackrel{\text{def}}{=} \text{If } (v_1 = \mathbf{X}) \wedge (v_5 = \mathbf{X}) \wedge (v_9 = \mathbf{X}) \text{ Then } +1 \text{ Else } -1$

Example 2 Rule $H_2(\mathbf{x}) \stackrel{\text{def}}{=} \text{If } (v_2 = \mathbf{X}) \wedge (v_5 = \mathbf{O}) \wedge (v_9 = \mathbf{X}) \text{ Then } -1 \text{ Else } +1$

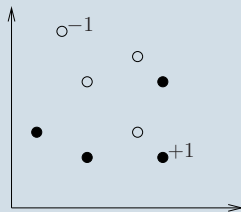
Example 3 Function $H_3(\mathbf{x}) \stackrel{\text{def}}{=} H_1(\mathbf{x})H_2(\mathbf{x})$

Example 4 Function $H_4(\mathbf{x}) \stackrel{\text{def}}{=} \text{sign}(3H_1(\mathbf{x}) - 2H_2(\mathbf{x}))$
($\text{sign}(z) = +1$ iff $z \geq 0$, and -1 otherwise)

\mathcal{H} = set of classifiers sharing the same model.

Supervised Learning > Another example

Geometric domain



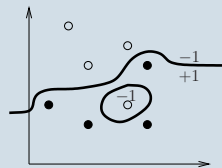
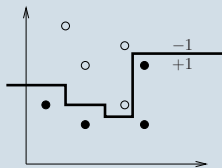
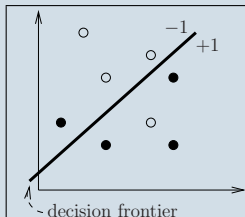
$$\mathcal{X} = \mathbb{R}^2$$

$$|\mathcal{S}| = 8$$

4 positive, 4 negative

Supervised Learning > Another example (contd)

Which classifier ?



decision frontier:

a line

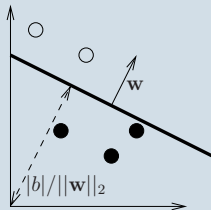
piecewise linear

curved,
regions not connex

Which **set of classifiers** \mathcal{H} is the best ? most popular ?

Supervised Learning > Key classifiers

Classifier: $(\mathcal{H} =) LS$



Linear Separators (linear models)

$$H(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b \in \mathbb{R}$$

Decision frontier: $\{\mathbf{x} : H(\mathbf{x}) = 0\}$,
an hyperplane

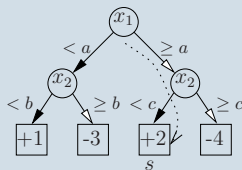
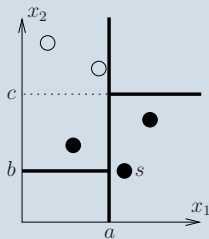
Needs $\mathcal{X} \subseteq \mathbb{R}^n$, but not a constraint: find **feature vector** \mathbf{f} s.t.

$$f_i : \mathcal{X} \rightarrow \mathbb{R} ,$$

and build $H(\mathbf{x}) = \langle \boldsymbol{\alpha}, \mathbf{f}(\mathbf{x}) \rangle + \beta$.

Supervised Learning > Key classifiers

Classifier: *DT*



Decision Trees

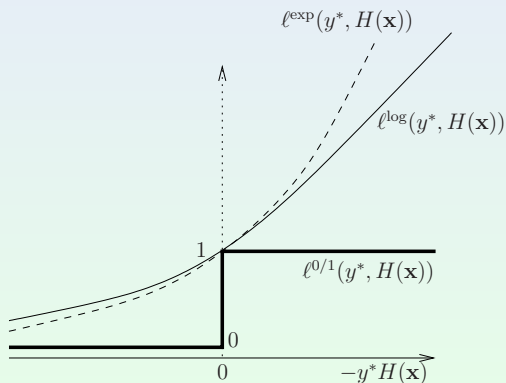
Decision frontier:
(portions of)
axis- \perp hyperplanes.

Does not require $\mathcal{X} \subseteq \mathbb{R}^n$. Works with any kind of description variable.

Supervised Learning > Key quantities (general)

Loss function

$\ell(y^*, H) : \mathbb{R}^2 \rightarrow \mathbb{R}_+$, computes to what extent the inputs match.



$$\ell^{0/1}(a, b) = 1_{[\text{sign}(a) \neq \text{sign}(b)]}$$

0/1 loss

($1_\pi = 1$ iff π **true**, 0 otherwise)

$$\ell^{\exp}(a, b) = \exp(-ab)$$

exponential loss

$$\ell^{\log}(a, b) = \log(1 + \exp(-ab))$$

logistic loss

(+ **many** others)

Supervised Learning > Key quantities (contd)

Empirical loss: $\varepsilon_S(H)$

Computes to what extent the labels match between S and H :

$$\varepsilon_{\mathbf{w}_1}(H) \stackrel{\text{def}}{=} \sum_{i=1}^m w_{1,i} \ell(y_i^*, H(\mathbf{x}_i)) = \mathbf{E}_{(\mathbf{x}, y^*) \sim \mathbf{w}_1} [\ell(y^*, H(\mathbf{x}))]$$

True loss: $\varepsilon_D(H)$

Computes the real matching between labels, by extending the prediction of H to D :

$$\varepsilon_D(H) \stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y^*) \sim D} [\ell(y^*, H(\mathbf{x}))]$$

$\ell(., .) : \mathbb{R}^2 \rightarrow \mathbb{R}_+ = \text{loss function.}$

Supervised Learning > Key quantities (contd)

Consider $\ell(.,.) = 0/1$ loss:

$$\ell^{0/1}(a, b) \stackrel{\text{def}}{=} 1_{[\text{sign}(a) \neq \text{sign}(b)]} ,$$

The losses specialize to the empirical and true **risks**:

$$\begin{aligned} \varepsilon_{\mathbf{w}_1}^{0/1}(H) &\stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y^*) \sim \mathbf{w}_1} [1_{\text{sign}(H(\mathbf{x})) \neq y^*}] , \\ \varepsilon_D^{0/1}(H) &\stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y^*) \sim D} [1_{\text{sign}(H(\mathbf{x})) \neq y^*}] \end{aligned}$$

Remark: if $\text{im}(H) \subseteq \mathbb{R}$, it gives **two** informations:

- 1 the class, $\text{sign}(H)$
- 2 a **confidence** in the class, $|H|$

Strong Learning

\mathcal{A} is a Strong learning algorithm for some class of classifiers \mathcal{H}
iff

Step A \mathcal{A} takes as input $0 < \epsilon, \delta \leq 1$

Step B In time $\text{poly}(1/\epsilon, 1/\delta, n, \dots)$, \mathcal{A} does the following:

Step B.1 \mathcal{A} requests a set \mathcal{S} sampled i.i.d. from D

Step B.2 \mathcal{A} returns $H \in \mathcal{H}$ such that:

$$\Pr_{\mathcal{S} \sim D^m} [\epsilon_D^{0/1}(H) \leq \epsilon] \geq 1 - \delta$$

Strong because requirements hold for *any* ϵ, δ within bounds.
 δ = **confidence** parameter.

Weak Learning

\cong Strong learning, *but* this time ϵ, δ are not user-fixed (Step A):

$$\epsilon \stackrel{\text{def}}{=} 1/2 - \gamma \text{ for some small } \gamma \text{ (e.g. cst, } 1/\text{poly}(n), \text{ etc...)}$$

$$\delta \stackrel{\text{def}}{=} \gamma$$

Thus, requirement in B.2 becomes:

$$\Pr_{S \sim D^m} [\underbrace{\epsilon_D^{0/1}(H) \leq 1/2 - \gamma}_{\text{close to } 1/2}] \geq \underbrace{\gamma}_{\text{close to 0}}$$

Weak because H is only required to carry little “information”:

$$\epsilon_D(\text{unbiased coin}) = 1/2, \forall D$$

Unbiased coin carries no information about the link Input \rightarrow Output.

The weakening on the confidence is superficial

$$\Pr_{S \sim D^m} [\underbrace{\varepsilon_D^{0/1}(H) \leq 1/2 - \gamma}_{\text{event } \mathbf{A}}] \geq \gamma$$

Run $T \geq (1/\gamma) \log(1/\delta)$ times Weak learning \mathbf{A} ($\forall 0 < \delta \leq 1$):

- the probability that event \mathbf{A} never occurs is
 $\leq (1 - \gamma)^T \leq \exp(-T\gamma) \leq \delta$
- thus, \mathbf{A} has occurred at least once with probability $\geq 1 - \delta$

..so we can consider that the main difference between Weak and Strong learning relies on the (empirical, true) risks.

To summarize, **Supervised learning**:

- 1 takes place in **polynomial time**,
- 2 means building a **classifier** that models a link between **inputs** (observations) and **outputs** (classes)
- 3 means **controlling** (up to the relevant extent) the **0/1 loss** of this classifier

The third point is the most important for our purpose.

Supervised Learning > Important problems

Strong learning

Is strong learning possible ?

Answers depend on \mathcal{H} ; *many* (early) complexity-theoretic negative results [Kearns and Vazirani(1994)].

Weak learning

Is weak learning possible ?

Some positive answers [Mannor and Meir(2000)]; recent complexity-theoretic negative results [Feldman(2006), Nock and Nielsen(2007b)].

Weak learning \Rightarrow Strong learning (**Boosting**)

Suppose \mathcal{A} has access to some Weak learning algorithm \mathcal{W} (Step B.1,5). (How) can \mathcal{A} meet the requirements of Strong learning ?

An early positive answer [Schapire(1990)].

- 1 Most frequent strategies for Strong, Weak learning, Boosting ?

Principle

Two-step strategy for some \mathcal{H} :

1 Find a **consistent** classifier H , i.e. with $\varepsilon_{w_1}^{0/1}(H) = 0, \forall S$
(in P-time).

2 Prove **structural properties**:

- 1 on the general \mathcal{H} ,
- 2 or on the subset of \mathcal{H} in which 1 is guaranteed to find H .

Step 2 ensure that we can lift at low (statistical, algorithmic) costs the consistency guarantee on S to the requirements of Strong learning over D .

Principle

Basically, **computationally tractable** search for a classifier with weak guarantees on the empirical risk.

Typically, two strategies relying on exhaustive search:

- 1 focused inside a “good” subset of \mathcal{H} if $|\mathcal{H}|$ too large [Mannor and Meir(2000)].
 - Algorithmic cost,
 - + Theoretically convenient: Weak learning guaranteed.
- 2 inside some \mathcal{H} of low poly-size. Main experimental approach.
 - Sometimes gives up with Weak learning,
 - + Empirically convenient: Very fast.

Boosting has to solve a trick:

$$\underbrace{\Pr_{\mathcal{S} \sim D^m} [\varepsilon_D^{0/1}(H) \leq 1/2 - \gamma] \geq \gamma}_{\text{Guaranteed on } \mathfrak{M}}$$



$$\underbrace{\Pr_{\mathcal{S} \sim D^m} [\varepsilon_D^{0/1}(H) \leq \epsilon] \geq 1 - \delta}_{\text{Required on } \mathfrak{A}}$$

Recall that boosting the confidence is not a problem. So, how can we “boost” the (empirical, true) risk: $1/2 - \gamma \rightarrow \epsilon$?

Learning Strategies > Boosting (contd)

(Rough, most popular) Principle, $\forall \mathcal{H}$

Two-step strategy for some \mathcal{H} :

\mathbb{M} Get many (T) classifiers from \mathbb{M} :

$$h_t \leftarrow \mathbb{M}(S, \mathbf{w}_t), t = 1, 2, \dots, T$$

(\mathbb{M} is trained over weight vector that may **differ** from \mathbf{w}_1)

H Combine the T classifiers to get some $H_T \in \mathcal{H}$.

Very appealing properties:

- for some Boosting algorithms, \mathbf{w}_t is repeatedly skewed towards the examples that have been hard to classify so far:

$$\ell^{0/1}(y_i^*, h_t(\mathbf{x}_i)) = 1 \Rightarrow \mathbf{w}_{t+1,i} > \mathbf{w}_{t,i}$$

Boosting Algorithms

- 1 Most popular algorithms ?
- 2 Common properties, differences ?

Boosting Algorithms > LS > AdaBoost

Suppose that \mathcal{W} returns $h_t : \mathcal{X} \rightarrow \{-1, +1\}$.

(Discrete) AdaBoost

for $t = 1, 2, \dots, T$

1 $h_t \leftarrow \mathcal{W}(\mathcal{S}, \mathbf{w}_t);$

2 $\alpha_t \leftarrow \frac{1}{2} \log \frac{1 - \varepsilon_{\mathbf{w}_t}^{0/1}(h_t)}{\varepsilon_{\mathbf{w}_t}^{0/1}(h_t)};$

3 $\forall (\mathbf{x}_i, y_i) \in \mathcal{S}, w_{t+1,i} \leftarrow \frac{w_{t,i} \exp(-y_i \alpha_t h_t(\mathbf{x}_i))}{Z_t};$

($Z_t \stackrel{\text{def}}{=}$ normalization coefficient for \mathbf{w}_t)

return $H_T = \sum_{t=1}^T \alpha_t h_t;$

- + Straightforward to implement, “best off the shelf classifier in the world”;
- restricted outputs for $h_t, \{-1, +1\}$
(generalized in [Friedman et al.(2000)], [Nock and Nielsen(2007a)])

A **fundamental** and simple property of AdaBoost:

- 1 Compute explicitly the normalization coefficient:

$$Z_t = 2\sqrt{\varepsilon_{\mathbf{w}_t}^{0/1}(h_t)(1 - \varepsilon_{\mathbf{w}_t}^{0/1}(h_t))}$$

- 2 Unravel the update rule at the end of learning:

$$\mathbf{w}_{T+1,i} = \mathbf{w}_{1,i} \exp(-y_i^* H_T(\mathbf{x}_i)) / \prod_t Z_t$$

- 3 Recall that the 0/1 loss is \leq **exponential loss**:

$$\ell^{0/1}(y^*, H_T(\mathbf{x})) \leq \ell^{\exp}(y^*, H_T(\mathbf{x})) \stackrel{\text{def}}{=} \exp(-y^* H_T(\mathbf{x}))$$

Sum (2) over \mathcal{S} (this equals 1), multiply both sides by $\prod_t Z_t$, use (1) and get with (3):

$$\varepsilon_{\mathbf{w}_1}^{0/1}(H_T) \leq 2^T \prod_t \sqrt{\varepsilon_{\mathbf{w}_t}^{0/1}(h_t)(1 - \varepsilon_{\mathbf{w}_t}^{0/1}(h_t))}$$

Suppose that \mathcal{H} is a Weak learning algorithm, i.e.

$$\varepsilon_{\mathbf{w}_t}^{0/1}(h_t) \leq 1/2 - \gamma \quad (\text{whp}) .$$

We get:

$$\varepsilon_{\mathbf{w}_1}^{0/1}(H_T) \leq (1 - 4\gamma^2)^{T/2} \leq \exp(-2\gamma^2 T) .$$

Fixing $T = \Omega((1/\gamma^2) \log m)$ makes the rhs $< 1/m$, i.e. (when $\mathbf{w}_1 = \mathbf{u} \stackrel{\text{def}}{=} (1/m)\mathbf{1}$):

- ① H_T is consistent and AdaBoost is P-time;
- ② The final LS meets the structural assumptions for Strong learning;

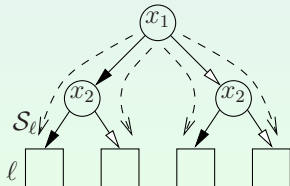
hence, AdaBoost is Strong learning, and thus Boosting.

Boosting Algorithms > DT

Property of a decision tree

The empirical risk can be decomposed at the leaves (each leaf sign is the majority class wlog):

$$\begin{aligned}\varepsilon_{\mathbf{w}_1}^{0/1}(H_T) &= \sum_{\ell \text{ leaf in } H_T} w_{1,\ell} \min\{w_{1,\ell}^+, 1 - w_{1,\ell}^+\} \\ &= \mathbf{E}_{\ell \sim \mathbf{w}_1} \left[\underbrace{\min\{w_{1,\ell}^+, 1 - w_{1,\ell}^+\}}_{\text{minority class \% in } \ell} \right]\end{aligned}$$



partition of \mathcal{S} in 4 subsets

\mathcal{S}_ℓ = subset of \mathcal{S} that reaches leaf ℓ

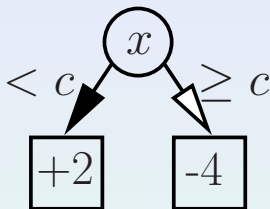
$w_{1,\ell}$ = total weight of \mathcal{S}_ℓ wrt \mathbf{w}_1

$w_{1,\ell}^+$ = total weight of class +1 in \mathcal{S}_ℓ wrt \mathbf{w}_1 , divided by $w_{1,\ell}$

e.g. $w_{1,\ell}^+ = \sum_{(x_i,+1) \in \mathcal{S}_\ell} w_{1,i} / \sum_{(x_i,y_i^*) \in \mathcal{S}_\ell} w_{1,i}$

Boosting Algorithms > DT (contd)

Suppose that \mathcal{M} returns **stumps**, i.e. depth-1 DTs.



- 1 These classifiers are so simple that \mathcal{A} generally implements \mathcal{M} as well...
- 2 Most popular DT induction algorithms have two stages:
 - 1 build a large tree (C4.5, ID3, etc.),
 - 2 prune the tree.

Here, we (deliberately) reduce them to their first stage.

Boosting Algorithms > DT > \mathbb{T} DIDT

Top-Down Induction of Decision Trees (\mathbb{T} DIDT)

\mathbb{T} DIDT (including \mathbb{W})

Initialize H_0 to be a single leaf node (=root, $\mathcal{S}_0 = \mathcal{S}$)

For $t = 1, 2, \dots, T$

1 pick some leaf ℓ in H_{t-1} .

\mathbb{W} choose the best stump on \mathcal{S}_ℓ :

$$h_t^* = \arg \min_{h_t} \underbrace{\sum_{\ell \text{ leaf in } H_{t-1} \oplus_\ell h_t} w_{1,\ell} \times \phi(w_{1,\ell}^+)}_{\text{Splitting criterion} = \varepsilon_{w_1}(H_{t-1} \oplus h_t, \phi)}$$

2 $H_t \leftarrow H_{t-1} \oplus_\ell h_t^*$

return $H_T \in \text{DT}$;

$\oplus_\ell h_t$ is the operation that replaces leaf ℓ by h_t , resulting in a DT with one more leaf.

Rewrite the splitting criterion as:

$$\varepsilon_{\mathbf{w}_1}(H_T, \phi) = \mathbf{E}_{\ell \sim \mathbf{w}_1} \phi(\mathbf{w}_{1,\ell}^+)$$

The empirical risk satisfies $\epsilon_{\mathbf{w}_1}^{0/1}(H) = \varepsilon_{\mathbf{w}_1}(H_T, \phi_{\text{emp}})$ for

$$\phi_{\text{emp}}(z) \stackrel{\text{def}}{=} \min\{z, 1 - z\}.$$

The first \mathcal{C} DIDT scheme, proposed in [Breiman et al.(1984)] proposes a **different** ϕ -criterion:

$$\phi_{\text{CART}}(z) = 2z(1 - z)$$

$\varepsilon_{\mathbf{w}_1}(H_T, \phi_{\text{CART}})$ is the (expected) **Gini index**

Other popular C4.5/DIDT schemes make the difference only on the choice of ϕ :

C4.5 [Quinlan(1993)] has

$$\phi_{C4.5}(z) = -z \log(z) - (1 - z) \log(1 - z)$$

(log base 2)

$\varepsilon_{w_1}(H_T, \phi_{C4.5})$ is the the (expected) **binary entropy**;

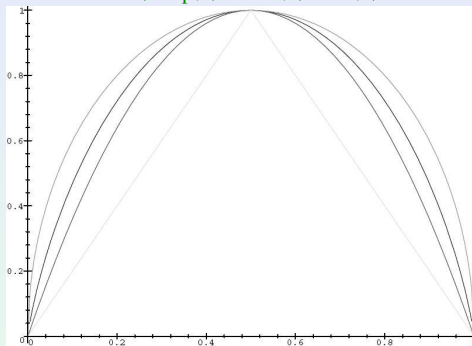
KM [Kearns and Mansour(1999)] has

$$\phi_{KM}(z) = 2\sqrt{z(1 - z)}$$

$\varepsilon_{w_1}(H_T, \phi_{KM})$ is the (expected) **Matsushita error**.

Boosting Algorithms > DT > Permissible functions

Functions ϕ_{emp} , ϕ_{CART} , $\phi_{\text{C4.5}}$, ϕ_{KM} have a crucial commonpoint:



Function ϕ is **permissible** iff:

- $\phi : [0, 1] \rightarrow [0, 1]$
- $\phi(0) = \phi(1) = 0$
- ϕ sym. wrt $z = 1/2$
- ϕ concave

ϕ **strictly permissible** iff:
permissible and strictly concave

From bottom to top: $2 \times \phi_{\text{emp}}$, ϕ_{CART} , $\phi_{\text{C4.5}}$, ϕ_{KM} .

All permissible, ϕ_{emp} not strictly permissible.

Normalization for $\phi(1/2) = 1$ not necessary (only for readability).

Boosting Algorithms > (CART, C4.5, KM) ∈ Boosting ?

Fundamental property, easily checked from the last picture:

$$\underbrace{\varepsilon_{\mathbf{w}_1}(H_T, \phi_{\text{emp}})}_{\text{empirical risk}} \leq 1/2 \times \begin{cases} \varepsilon_{\mathbf{w}_1}(H_T, \phi_{\text{CART}}) \\ \varepsilon_{\mathbf{w}_1}(H_T, \phi_{\text{C4.5}}) \\ \varepsilon_{\mathbf{w}_1}(H_T, \phi_{\text{KM}}) \end{cases}$$

Thus, minimizing any of the right criteria should amount to the minimization of the empirical risk.

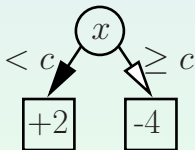
⇒ Why so many criteria ?

Boosting Algorithms > (Cart, 4.5, AM) ∈ Boosting ?

Consider the following assumption (**W**weak **L**earning **A**ssumption):

WLA $\exists \gamma > 0$ s.t. for any set S' , any distribution w' over S' , there exists a stump h for which:

$$\varepsilon_{w'}^{0/1}(h) \leq 1/2 - \gamma$$



We basically assume that step ((1)+~~W~~) is a Weak learning algorithm

Boosting Algorithms > $\mathbb{M} \in$ Boosting

Under assumption WLA, the following holds for \mathbb{M} and its output H_T :

$$T \geq \left(\frac{1}{\epsilon}\right)^{\frac{c}{\gamma^2}} \Rightarrow \varepsilon_{\mathbf{w}_1}^{0/1}(H_T) \leq \epsilon$$

c is a constant [Kearns and Mansour(1999), Henry et al.(2007)] (improved c). Fixing $\epsilon < 1/m$ ($\mathbf{w}_1 = \mathbf{u}$) makes:

- 1 H_T consistent, and \mathbb{M} is P-time;
- 2 The final DT meets the structural assumptions for Strong learning;

hence, \mathbb{M} is Strong learning, and thus Boosting.

- AdaBoost needs only \log calls to \mathbb{M} compared to \mathbb{M} , but it is a structural “difficulty” for DT
- the number of calls to \mathbb{M} is **optimal** for \mathbb{M} , from both the informational and complexity-theoretic standpoints [Kearns and Mansour(1999), Nock and Nielsen(2004)]

Boosting Algorithms > (CART, C4.5) ∈ Boosting ?

Under assumption WLA, the following holds for C4.5 and its output H_T :

$$T \geq \left(\frac{1}{\epsilon}\right)^{\frac{c \log(1/\epsilon)}{\gamma^2}} \Rightarrow \varepsilon_{\mathbf{w}_1}^{0/1}(H_T) \leq \epsilon$$

Under assumption WLA, the following holds for CART and its output H_T :

$$T \geq \left(\frac{1}{\epsilon}\right)^{\frac{c}{\gamma^2 \epsilon^2}} \Rightarrow \varepsilon_{\mathbf{w}_1}^{0/1}(H_T) \leq \epsilon$$

- 1 Fix $\epsilon < 1/m$ ($\mathbf{w}_1 = \mathbf{u}$): the consistency is met, but C4.5 is QP-time, while CART is EXP-time
- 2 Bounds above are lowerbounds. No tight bound is known, but experimental results seem to confirm the results
- 3 $\phi = \phi_{\text{emp}}$ would yield the poorest bounds of all (!)

Boosting Algorithms > General Observations

Observation 1 Most DT induction algorithms do not minimize **directly** the empirical risk, but (the expectation of) a *concave surrogate* (an upperbound: Gini index, entropy, Mastushita's error)

Observation 2 AdaBoost does not minimize directly the empirical risk, but (the expectation of) a *convex surrogate*, the **exponential loss**:

$$\begin{aligned}\varepsilon_{\mathbf{w}_1}^{\text{exp}}(H_T) &= \mathbf{E}_{(\mathbf{x}, y) \sim \mathbf{w}_1} [\exp(-y^* H_T(\mathbf{x}))] \\ &\geq \varepsilon_{\mathbf{w}_1}^{0/1}(H_T)\end{aligned}$$

Observation 3 On DT induction, the **more concave** the permissible function ϕ , the **better** the lowerbounds on T

Observation 4 On DT induction, the direct minimization of the empirical risk yields the **worst possible** lowerbound on T [Nock and Nielsen(2004)]

Boosting Algorithms > Question

Supervised learning roughly aims at minimizing empirical risk.

Why focusing on surrogates ?

Numerous well known surrogates:

Concave For DT (and related classes): Gini index, entropy, Matsushita's error

Convex For LS:

$$\begin{aligned}\varepsilon_{\mathbf{w}_1}^{\text{exp}}(H_T) &= \mathbf{E}_{(\mathbf{x}, y) \sim \mathbf{w}_1} [\exp(-y^* H_T(\mathbf{x}))] \\ &\quad \text{(Exponential loss: AdaBoost)} \\ \varepsilon_{\mathbf{w}_1}^{\text{log}}(H_T) &= \mathbf{E}_{(\mathbf{x}, y) \sim \mathbf{w}_1} [\log(1 + \exp(-2y^* H_T(\mathbf{x})))] \\ &\quad \text{(Logistic loss)} \\ \varepsilon_{\mathbf{w}_1}^{\text{squ}}(H_T) &= \mathbf{E}_{(\mathbf{x}, y) \sim \mathbf{w}_1} [(1 - y^* H_T(\mathbf{x}))^2] \\ &\quad \text{(Squared loss) ...}\end{aligned}$$

Convex surrogates **have the form** $F(y^* H_T(\mathbf{x}))$.

Why focusing on surrogates ?

Explanations so far in this talk:

- 1 Algorithmic: better convergence properties. Not satisfactory (lack of matching upperbounds).
- 2 Complexity-theoretic: empirical risk has more local minima (0/1 loss takes on 2 values, thus has less discrimination). Not satisfactory (can be hard for surrogates as well [Nock and Nielsen(2004)]).
- 3 Others (statistics).

No explanation drills down into the fundamental links between surrogates and classification.

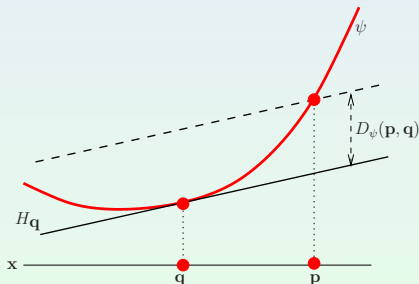
Bregman Divergences

- 1 Presentation of this class of distortion measures
- 2 An example of their widespread application in learning

Bregman Divergences (contd)

Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ strictly convex and differentiable, \mathcal{X} convex. The **Bregman divergence** with **generator** ψ is:

$$D_{\psi}(\mathbf{p}||\mathbf{q}) = \psi(\mathbf{p}) - \psi(\mathbf{q}) - \langle \mathbf{p} - \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle$$



In general, does not satisfy symmetry, triangular inequality.

Bregman Divergences > Example

Squared Euclidean distance

Generator $\psi(\mathbf{p}) = \|\mathbf{p}\|_2^2$: strictly convex and differentiable over \mathbb{R}^n

Divergence

$$\begin{aligned} D_\psi(\mathbf{p}||\mathbf{q}) &= \psi(\mathbf{p}) - \psi(\mathbf{q}) - \langle \mathbf{p} - \mathbf{q}, \nabla_\psi(\mathbf{q}) \rangle \\ &= \|\mathbf{p}\|_2^2 - \|\mathbf{q}\|_2^2 - \langle \mathbf{p} - \mathbf{q}, 2\mathbf{q} \rangle \\ &= \|\mathbf{p} - \mathbf{q}\|_2^2 \end{aligned}$$

Bregman Divergences > Example (contd)

Generalized I-Divergence

Generator $\psi(\mathbf{p}) = \sum_i p_i \log p_i$: strictly convex and differentiable over \mathbb{R}_+^n

Divergence

$$D_\psi(\mathbf{p}||\mathbf{q}) = \sum_i p_i \log \left(\frac{p_i}{q_i} \right) - p_i + q_i$$

If ψ restricted to the probability simplex, becomes Kullback-Leibler divergence.

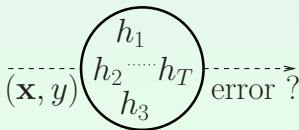
Bregman Divergences > On-line learning

The first Supervised learning setting in which they have been explicitly and extensively used.

On-line learning

a We are given a **fixed** set of experts $\{h_i : \mathcal{X} \rightarrow \{-1, +1\}\}_{i=1}^m$, a **stream** of examples. $\mathbf{w}_t \in \mathbb{R}_+^m$ is the current set of weights:

- 1 Receive example (\mathbf{x}_t, y_t^*)
- 2 Make prediction $H_m(\mathbf{x}) = \sum_i w_{t,i} h_i(\mathbf{x})$
- 3 Incur loss $\ell(y_t^*, H_m(\mathbf{x}_t))$
- 4 Modify the weights: $\mathbf{w}_{t+1} \leftarrow f(\mathbf{w}_t, \ell(y_t^*, H_m(\mathbf{x}_t)))$
- 5 Go to 1



Objective:
Minimize the number of mistakes

Bregman Divergences > On-line learning (contd)

On-line learning is a setting dual to Boosting (reverse the role of the examples and hypotheses in learning). Computation of \mathbf{w}_{t+1} involves the aggregation of two Bregman divergences:

$$\mathbf{w}_{t+1} \stackrel{\text{def}}{=} \arg \min_{\mathbf{w}} \left\{ \underbrace{D_{\psi'}(\mathbf{w} || \mathbf{w}_t)}_{\text{regularization}} + \underbrace{\eta D_{\psi} \left(\sum_{t=1}^T \mathbf{w}_t h_t(\mathbf{x}) || \underbrace{\nabla_{\psi}^{-1}(y)}_{\in \mathbb{R}} \right)}_{\text{matching loss}} \right\}$$

- 1 $y \stackrel{\text{def}}{=} (1 + y^*)/2 \in \{0, 1\}$ is the **Boolean** class.
- 2 η controls the tradeoff between the two losses.

- 1 What is the **true** loss $\ell(y^*, H(\mathbf{x}))$ that we (really) want to minimize on each example (\mathbf{x}, y^*) ?
(we have seen many losses so far: 0/1, convex/concave surrogates, Bregman divergences)
- 2 Can we find it based on its properties people usually **assume** ?
- 3 Links with conventional losses ?
- 4 New losses, families ?

Lifting Classification to **Estimation**

Early “ages” of supervised learning usually preferred the Boolean class:

$$y \stackrel{\text{def}}{=} (1 + y^*)/2 \in \{0, 1\}$$

y is thus a 0/1 estimator for $\mathbf{Pr}[y^* = +1|\mathbf{x}]$ (key ingredient for **Bayes rule**).

- 1 Wlog, assume H able to return an estimator

$$H(\mathbf{x}) \Leftarrow \hat{\mathbf{Pr}}_H[y^* = +1|\mathbf{x}]$$

If $\text{im}(H) \subseteq \mathbb{R}$, \rightarrow done by well-known transfo. (e.g. logistic)

If $\text{im}(H) \subseteq [0, 1]$, \leftarrow done by e.g. $\text{sign}(2H - 1) \in \{-1, +1\}$

- 2 We end up with the analysis of $\ell(.,.)$ with $\text{dom}(\ell) = [0, 1]^2$

Relies on three assumptions in Supervised learning:

- 1 On the **loss function**
- 2 On the **best possible rule**
- 3 On the **cost matrix** for learning

Non Negativity

People assume:

$$\ell(.,.) \geq 0$$

Axiomatization > Assumption 2

Suppose that **all** examples of \mathcal{S} share the **same** observation \mathbf{x}^* . What is the **best constant** prediction for \mathbf{x}^* in average:

$$c = \arg \min_{z \in [0,1]} \mathbf{E}_{(\mathbf{x}^*, y) \sim w_1} [\ell(y, z)] = ?$$

Bayes Optimality

People assume that the best prediction rule is Bayes rule:

$$\text{sign}(2\Pr[y^* = +1|\mathbf{x}] - 1)$$

Thus, the best constant prediction is the best estimator for $\Pr[y^* = +1|\mathbf{x}]$, i.e.:

$$c = \mathbf{E}_{(\mathbf{x}^*, y) \sim w_1} [y]$$

Axiomatization > Assumption 3

Fundamental (often implicit) input to supervised learning: the **cost matrix** $L \in \mathbb{R}_+^{2 \times 2}$.

		predicted class	
		1	0
true class	1	$\ell(1, 1)$	$\ell(1, 0)$
	0	$\ell(0, 1)$	$\ell(0, 0)$

The most general form for the empirical risk is:

$$\varepsilon_{w_1}^{0/1}(H) \stackrel{\text{def}}{=} \mathbf{E}_{(x,y) \sim w_1} [\ell(y, 1_{H(x) \geq 1/2})]$$

Remark that even right classifications may incur some $\neq 0$ cost.

Axiomatization > Assumption 3 (contd)

Symmetric Cost Matrix

People assume* that the **cost matrix** L satisfies the following symmetries:

Diagonal $\ell(1, 1) = \ell(0, 0)$ (= 0: no cost for right classifications)

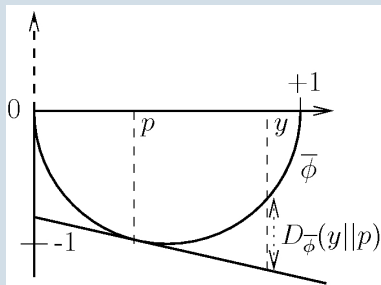
Outside $\ell(1, 0) = \ell(0, 1)$ (same cost for misclassifications)

(this simplifies the empirical risk to the one we have used since the beginning). We thus have:

$$\ell(y, z) = \ell(1 - y, 1 - z)$$

(*) Holds for domains that have no class-dependent misclassification costs. The others are much less formalized.

Theorem



(here, $\bar{\phi} = -\phi$)

Loss function $\ell(.,.) : [0, 1]^2 \rightarrow \mathbb{R}_+$
satisfies assumptions 1, 2, 3

if and only if

$\ell(y, z) = D_{-\phi}(y, z)$
with ϕ strictly permissible

(recall DT induction ?)

Strict subclass of Bregman divergences:

Strictly Permissible Bregman Loss Functions (BLF_{SP})

Axiomatization > BLF_{SP} (contd)

The loss we minimize has thus the general form ($\forall \phi$ strictly permissible):

$$\varepsilon_{\mathbf{w}_1}(H) = \varepsilon_{\mathbf{w}_1, \bar{\phi}}(H) \stackrel{\text{def}}{=} \underbrace{\mathbf{E}_{(\mathbf{x}, y) \sim \mathbf{w}_1} [D_{\bar{\phi}}(y || \hat{\mathbf{Pr}}_H[y = 1 | \mathbf{x}])]}_{\text{(expectation of) BLF}_{\text{SP}}}$$

and we can show:

$$\varepsilon_{\mathbf{w}_1}^{0/1}(H) \leq \varepsilon_{\mathbf{w}_1, \bar{\phi}}(H) / \phi(1/2)$$

Since $\phi(1/2) \neq 0$, minimizing any BLF_{SP} should amount to minimizing the empirical risk as well.

Thus, supervised classification aims at minimizing (the expectation of) some BLF_{SP} .

Up to a large extent, this means minimizing the empirical risk as well.

What else ?

- 1 links with surrogates ?
- 2 minimization algorithms ?

Definition

Suppose ψ strictly convex, differentiable over \mathbb{X} . Unique *Convex conjugate* function ψ^* obtained by the Legendre transformation:

$$\psi^*(\mathbf{q}) = \sup_{\mathbf{p} \in \mathbb{X}} \{\langle \mathbf{q}, \mathbf{p} \rangle - \psi(\mathbf{p})\}$$

Solve via $\nabla \psi^*(\mathbf{q}) = \nabla(\langle \mathbf{q}, \mathbf{p} \rangle - \psi(\mathbf{p})) = 0$, implying $\mathbf{q} = \nabla_{\psi}(\mathbf{p})$, $\mathbf{p} = \nabla_{\psi}^{-1}(\mathbf{q})$, and $\psi^*(\mathbf{q}) = \langle \mathbf{q}, \nabla_{\psi}^{-1}(\mathbf{q}) \rangle - \psi(\nabla_{\psi}^{-1}(\mathbf{q}))$.

Dual Bregman divergence

Fundamental link between Bregman divergences:

$$D_{\psi}(\mathbf{p}||\mathbf{q}) = D_{\psi^*}(\nabla_{\psi}(\mathbf{q})||\nabla_{\psi}(\mathbf{p}))$$

Axiomatization > BLF_{SP} > ...and Supervised learning

Convex conjugates bring the **link** between $[0, 1]$ classification (y) and real-valued classification (y^*). In the BLF_{SP} $\varepsilon_{\mathbf{w}_1, \bar{\phi}}(H)$, we can write:

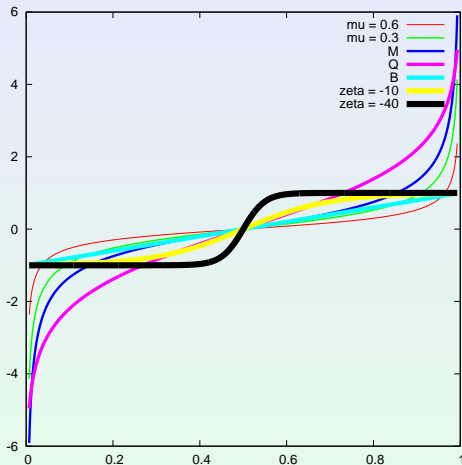
$$\underbrace{D_{\bar{\phi}}(y || \hat{\mathbf{P}}\mathbf{r}_H[y = 1 | \mathbf{x}])}_{\text{BLF}_{\text{SP}}: [0, 1] \text{ values}} = \underbrace{D_{\bar{\phi}^*}(\nabla_{\bar{\phi}}(\hat{\mathbf{P}}\mathbf{r}_H[y = 1 | \mathbf{x}]) || \nabla_{\bar{\phi}}(y))}_{\text{divergence of real values}}$$

Because ϕ strictly permissible, $\nabla_{\bar{\phi}}$ symmetric wrt $(1/2, 0)$:

- 1 $\nabla_{\bar{\phi}}(y)$, the \mathbb{R} real class, takes on two **opposite** values
- 2 Suppose $\text{im}(H) \subseteq \mathbb{R}$. We obtain the transformation rule for the $[0, 1]$ values:

$$\hat{\mathbf{P}}\mathbf{r}_H[y = 1 | \mathbf{x}] \stackrel{\text{def}}{=} \nabla_{\bar{\phi}}^{-1}(H(\mathbf{x}))$$

Axiomatization > BLF_{SP} > Example of plots for $\nabla_{\overline{\phi}}$



"mu" = *upper regime*

$$\phi_{M=KM}(z) = 2\sqrt{z(1-z)}$$

$$\phi_{Q=C4.5}(z) = -z \log(z) - (1-z) \log(1-z)$$

$$\phi_{B=CART}(z) = 4z(1-z)$$

"zeta" = *lower regime*

$\nabla_{\overline{\phi}}(z)$, for various strictly permissible ϕ , depending on its "concavity regime"

Let $\text{im}(H) \subseteq \mathbb{R}$.

Lemma

$$\underbrace{D_{\overline{\phi}}\left(y \parallel \nabla_{\overline{\phi}}^{-1}(H(\mathbf{x}))\right)}_{[0, 1] \text{ prediction}} = \underbrace{\overline{\phi}^*(-y^* H(\mathbf{x}))}_{\text{Real prediction}}$$

For any strictly permissible ϕ , $F_{\phi}(z) = \overline{\phi}^*(-z)/\phi(1/2)$ is called a **Permissible Convex Loss** (PCL).

Of course:

$$\varepsilon_{\mathbf{w}_1}^{0/1}(H) \leq \underbrace{\mathbf{E}_{(\mathbf{x}, y^*) \sim \mathbf{w}_1}[F_{\phi}(y^* H(\mathbf{x}))]}_{\text{(expectation of) PCL}}$$

Minimizing any PCL \Rightarrow minimizing the empirical risk.

PCL take on well known special expressions.

$\phi(z)$	$\text{im}H$ $= \text{im}(\nabla_{\bar{\phi}})$	$F_{\phi}(y^*H)$ $= \bar{\phi}^*(-y^*H)/\phi(1/2)$	$\hat{p}_H[y = y_1 \mathbf{o}]$ $= \nabla_{\bar{\phi}}^{-1}(H)$
$-z \log z$ $-(1-z) \log(1-z)$	\mathbb{R}	$\log(1 + \exp(-y^*H))^*$	$\frac{\exp(H)}{1+\exp(H)}^*$
$z(1-z)$	$[-1, 1]^{**}$	$(1 - y^*H)^{2***}$	$(1/2)(1 + H)$
$\sqrt{z(1-z)}$	\mathbb{R}	$-y^*H + \sqrt{1 + (y^*H)^2}$	$\frac{1}{2} \left(1 + \frac{H}{\sqrt{1+H^2}} \right)$

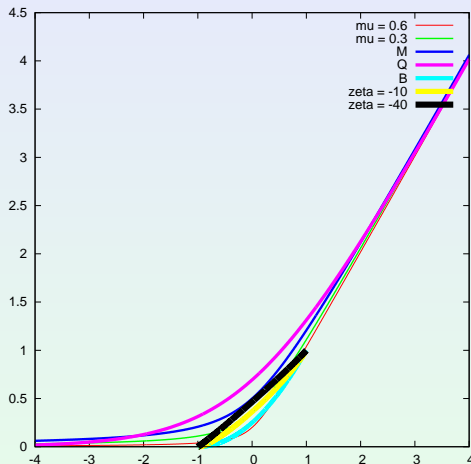
* Logistic loss and logistic transform

** Explains why problems when $H \in \text{LS}$.

*** Squared loss

Many other examples

Axiomatization > BLF_{sp} > Example of plots for PCL



"mu" = upper regime

$$\phi_{M=KM}(z) = 2\sqrt{z(1-z)}$$

$$\phi_{Q=C4.5}(z) = -z \log(z) - (1-z) \log(1-z)$$

$$\phi_{B=CART}(z) = 4z(1-z)$$

"zeta" = lower regime

$\overline{\phi}^*(z)$, for various strictly permissible ϕ , depending
on its "concavity regime"

Axiomatization > BLF_{SP} — PCL : The link

We have seen

Supervised Learning \Leftrightarrow min. BLF_{SP} \Leftrightarrow min. PCL

Different standpoints on Supervised classification:

- BLF_{SP}: $[0, 1]$ classification, H computes **probability estimates**
- PCL: \mathbb{R} Real classification, H computes **classes and confidences**

Lemma

AdaBoost's exponential loss is **not** a PCL:

$$\mathbf{E}_{(\mathbf{x}, y^*) \sim w_1} [\exp(-y^* H_T(\mathbf{x}))]$$

Minimization Algorithms

- 1 Algorithms that minimize some (any) BLF_{SP} , PCL ?
- 2 Link with existing algorithms ? New algorithms ?

Minimization Algorithms (contd)

Universal Minimization Algorithm

Let $\mathcal{A}(\mathcal{S}, \mathbf{w}_1, \phi)$ an algorithm that outputs classifiers from set \mathcal{H}

- 1 If, for any \mathcal{S} , \mathbf{w}_1 , for any strictly permissible ϕ , \mathcal{A} provably minimizes the corresponding PCL/BLF_{SP} (see below),
- 2 then \mathcal{A} is called a **Universal** Minimization Algorithm for \mathcal{H} .

$$\underbrace{\mathbf{E}_{(\mathbf{x}, y^*) \sim \mathbf{w}_1} [F_\phi(y^* H(\mathbf{x}))]}_{\text{PCL: } \text{im}(H) \subseteq \mathbb{R}} \quad \text{or} \quad \underbrace{\mathbf{E}_{(\mathbf{x}, y) \sim \mathbf{w}_1} [D_{-\phi}(y || \hat{\mathbf{Pr}}_H[y = 1 | \mathbf{x}])]}_{\text{BLF}_{\text{SP}}: \text{im}(H) = [0, 1]}$$

(No P-time complexity requirement)

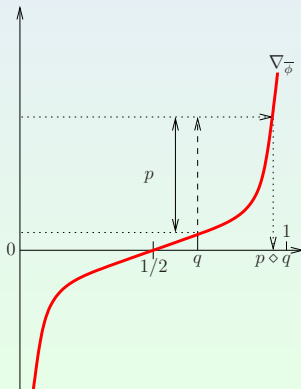
Minimization Algorithms > LS

Any BLF_{SP} is convex in its first argument.

Convex conjugates for BLF_{SP}

Let ϕ strictly permissible. $\forall p \in \mathbb{R}, \forall q \in [0, 1]$, the **Legendre dual** $p \diamond q$ of the ordered pair (p, q) is:

$$p \diamond q \stackrel{\text{def}}{=} \arg_{q' \in [0, 1]} \sup \{pq' - D_{\bar{\phi}}(q' || q)\} \quad (= \nabla_{\bar{\phi}}^{-1}(p + \nabla_{\bar{\phi}}(q)))$$



Legendre dual:

- 1- lifts q to $\text{im} \nabla_{\bar{\phi}}$
- 2- combines with p
- 3- maps back to $[0, 1]$

Minimization Algorithms > ULS

A Universal Minimization Algorithm for LS: ULS.

- 1 suppose that we already know the set $\{h_1, h_2, \dots, h_T\}$, for which $\text{im}(h_t) \subseteq \mathbb{R}$.
- 2 matrix $M \in \mathbb{R}^{m \times T}$ defined as:

$$m_{it} \stackrel{\text{def}}{=} -y_i^* h_t(\mathbf{x}_i), (\mathbf{x}_i, y_i^*) \in \mathcal{S}$$

- 3 vector notation $(M\alpha)_i \stackrel{\text{def}}{=} -y_i^* \underbrace{\sum_{t=1}^T \alpha_t h_t(\mathbf{x}_i)}_{H(\mathbf{x}_i)}, \forall \alpha \in \mathbb{R}^T$

- 4 Uniform distribution $\mathbf{w}_1 \stackrel{\text{def}}{=} \mathbf{u} = (1/m)\mathbf{1}$ wlog (duplicate examples)
- A None of the h_t has zero empirical risk (otherwise learning not necessary !)

AdaLS is mainly a two-step iterative algorithm:

- For $j = 1, 2, \dots, J$, do
 - 1 update the weights over the examples
 - 2 pick a subset $\mathcal{T}_j \subseteq \{1, 2, \dots, T\}$, update the leveraging coefficients of the classifiers $h_t, t \in \mathcal{T}_j$

“AdaBoosting flavor”. AdaLS specializes in different Boosting schemes:

- classical Boosting framework when $|\mathcal{T}_j| = 1$,
- totally corrective Boosting algorithm when $|\mathcal{T}_j| = \{1, 2, \dots, j\}$, etc.

ALS

Input: $M \in \mathbb{R}^{m \times T}$, strictly permissible ϕ ;

Initialize: $\alpha_1 \leftarrow \mathbf{0}$; (leveraging coefficient vector)

Initialize: $\mathbf{w}_0 \leftarrow (1/2)\mathbf{1}$; (uniform, **non-unit** weights)

For $j = 1, 2, \dots, J$:

① $\mathbf{w}_j \leftarrow (M\alpha_j) \diamond \mathbf{w}_0$; (Legendre dual componentwise)

② Pick $\mathcal{T}_j \subseteq \{1, 2, \dots, T\}$ and let $\delta_j \leftarrow \mathbf{0}$;

③ $\forall t \in \mathcal{T}_j$, find $\delta_{j,t}$ such that:

$$\sum_{i=1}^m m_{it}((M\delta_j) \diamond \mathbf{w}_j)_i = 0$$

④ $\alpha_{j+1} \leftarrow \alpha_j + \delta_j$;

Output: $H(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{t=1}^T \alpha_{J+1,t} h_t(\mathbf{x})$;

Property: (3) has always a solution under **A**.

Theorem

If \mathcal{T}_j is chosen as in classical Boosting, totally corrective Boosting (and others),

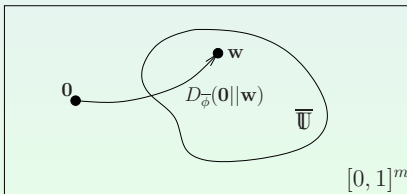
Then \mathcal{U} LS is a Universal Minimization Algorithm.

- 1 Full Theorem gives the necessary and sufficient conditions on the choice of \mathcal{T}_j for \mathcal{U} LS to remain Universal.
- 2 \mathcal{U} LS is the largest possible generalization to approaches in [Collins et al.(2002)] (generalizing more implies violating assumptions 1, 2 or 3)
- 3 Proof unveils the prominent role of “Bregman geometries”

(1): Shift from (P-time) Learning to (Computational) Geometry

$$\min_{\alpha \in \mathbb{R}^T} \underbrace{\sum_{i=1}^m F_{\phi}(y_i^* H(\mathbf{x}_i))}_{\text{PCL of the LS}} = \min_{\mathbf{w} \in \bar{\mathbb{U}}} \underbrace{\sum_{i=1}^m D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_i)}_{D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w})}$$

with $\bar{\mathbb{U}} \stackrel{\text{def}}{=} \{(M\alpha) \diamond \mathbf{w}_0 : \alpha \in \mathbb{R}^T\}$ (recall $\mathbf{w}_0 \stackrel{\text{def}}{=} (1/2)\mathbf{1}$)



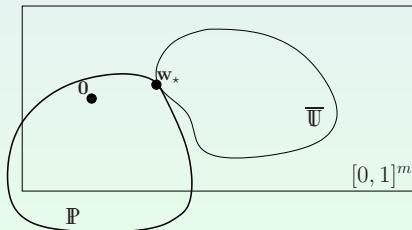
$\mathbf{0} \notin \bar{\mathbb{U}}$ under **A**

(2): Existence of a particular point in $\bar{\mathcal{U}}$

$$\forall \mathbf{w}_\star \in \mathbb{R}^m, \mathbf{w}_\star \in \mathbb{P} \cap \bar{\mathcal{U}} \Leftrightarrow \mathbf{w}_\star = \arg \min_{\mathbf{w} \in \bar{\mathcal{U}}} D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w})$$

with $\mathbb{P} \stackrel{\text{def}}{=} \{\mathbf{z} \in \mathbb{R}^m : M\mathbf{z} = \mathbf{0}\} = \text{Ker}M$

(recall $\mathcal{U} \stackrel{\text{def}}{=} \{(M\alpha) \diamond \mathbf{w}_0 : \alpha \in \mathbb{R}^T\}$ and $\mathbf{w}_0 \stackrel{\text{def}}{=} (1/2)\mathbf{1}$)



(\mathbf{w}_\star is unique)

Objective: find \mathbf{w}_\star

Minimization Algorithms > \mathcal{U} LS > Proof sketch

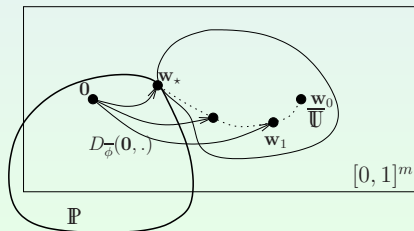
(3): \mathcal{U} LS is a constrained minimization algorithm

Recall that \mathcal{U} LS builds a sequence $\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_J \in \bar{\mathcal{U}}$.

Let an **auxiliary function** $u : [0, 1]^m \times [0, 1]^m \rightarrow \mathbb{R}$ for algorithm \mathcal{U} LS be a function that would satisfy:

$$D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_{j+1}) - D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_j) \leq u(\mathbf{w}_{j+1}, \mathbf{w}_j) \leq 0 \quad (\text{i})$$

$$u(\mathbf{w}_{j+1}, \mathbf{w}_j) = 0 \Rightarrow M\mathbf{w}_{j+1} = \mathbf{0} \quad (\text{ii})$$



If u exists,

(i) \mathcal{U} LS provably minimizes $D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w})$, and in $\bar{\mathcal{U}}$.

(ii) upon convergence, \mathcal{U} LS ends up with some $\mathbf{w}_J \in \mathbb{P}$

Hence, $\mathbf{w}_J \in \bar{\mathcal{U}} \cap \mathbb{P}$

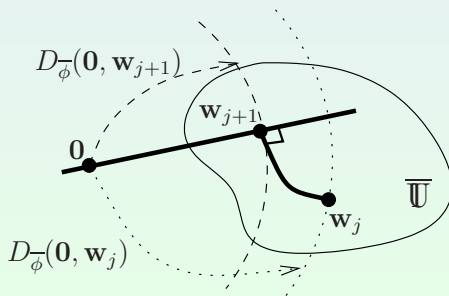
Hence, $\mathbf{w}_J = \mathbf{w}_* = \arg \min_{\mathbf{w}} D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w})$

(4): The auxiliary function for uLS

The computation of \mathbf{w}_j and δ_j in uLS yields

$$D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_{j+1}) - D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_j) = \underbrace{-D_{\bar{\phi}}(\mathbf{w}_{j+1} \parallel \mathbf{w}_j)}_{u(\mathbf{w}_{j+1}, \mathbf{w}_j)}$$

$u(\mathbf{w}_{j+1}, \mathbf{w}_j) \leq 0$, equality iff $\mathbf{w}_{j+1} = \mathbf{w}_j$ (prop of Bregman div.)



Generalized Pythagoras Theorem

$$D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_j) = D_{\bar{\phi}}(\mathbf{0} \parallel \mathbf{w}_{j+1}) + D_{\bar{\phi}}(\mathbf{w}_{j+1} \parallel \mathbf{w}_j)$$

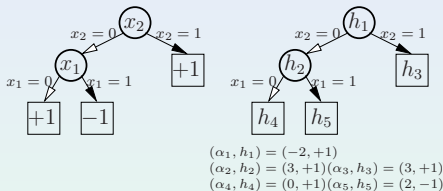
Summary:

- 1 ULS is a Universal Minimization Algorithm
- 2 Uses geometric properties on the weight vectors w to converge
- 3 Under some Weak Learning Assumption about the h_t , (loose) convergence rates

What about DT ? Any Universal minimization algorithm for DT ?

Minimization Algorithms > LDT

We do not have to think everything from scratch: use **Linearized** Decision Trees [Henry et al.(2007)]



In a LDT,

- 1 reals on every node (not just leaves)
- 2 sum the reals over a path to decide the class
- 3 each path is a **constant** LS

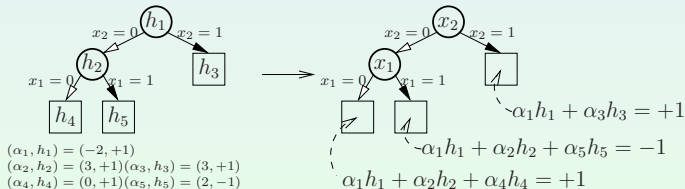
Minimization Algorithms > LDT (contd)

Twin DT

From any LDT, find the **twin** DT: for each path root — leaf,

- 1 computes the constant LS,
- 2 put the value at the leaf

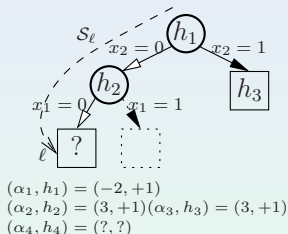
At the end, remove in all internal nodes any couple (α, h) . We obtain a DT equivalent to the LDT.



Minimization Algorithms > ADT

ADT recycles ALS on a strategy that minimizes the corresponding PCL

- To fit the internal couples (α, h) , use $\text{ALS}(\mathcal{S}_\ell, \mathbf{w}_{1,\ell}, \phi)$



- To find the splits, further minimize the global PCL over the choice of splits:

find the split replacing leaf ℓ which minimizes $\mathbf{E}_{(\mathbf{x}, y^*) \sim \mathbf{w}_1} [F_\phi(y^* H(\mathbf{x}))]$

Lemma

UDT is a Universal Minimization Algorithm

Is it known ?...

Theorem

Suppose we replace the LDT output by its twin DT.

\mathcal{A} DT($\mathcal{S}, \mathbf{w}_1, \phi$) simplifies **exactly** to the general \mathcal{A} DIDT scheme that minimizes $\varepsilon_{\mathbf{w}_1}(H_T, \phi)$ over the DT

Consequences (examples):

- 1 binds the most popular induction schemes for LS and DT as the **same** (master) algorithm, that uses the same geometric properties
- 2 AdaBoost and \mathcal{A} M are the same algorithm
- 3 the exact optimization of the logistic loss [Friedman et al.(2000)] is the **same** algorithm as \mathcal{A} 4.5

What about AdaBoost ?

Recall that the exponential loss of AdaBoost is not permissible... but:

Theorem

(a **very** slight) modification of AdaBoost is a Universal Minimization Algorithm

- 1 A single loss helps to minimize all !
- 2 Not surprising: minimizing any BLF_{SP} amounts to a Maximum Likelihood estimation to fit Bernoulli or Laplace priors.

1 Transfer positive results: we can use this master algorithm to obtain (new) formal Boosting algorithm for well known other classes \mathcal{H} :

- The algorithm fits local linear separators on a particular **decision graph**
- We can combine the same algorithm in a **recursive** fashion:
 - 1 we obtain \mathcal{A} ,
 - 2 we obtain \mathcal{W} **with the same algorithm**,
 - 3 and we can drill down even further...

- 2 **Transfer negative results:** we can translate back and forth bounds to get hint on the hardness of learning for particular \mathcal{H}
- Example: bounds of [Kearns and Mansour(1999)] on DT can be translated to LS
 - We get explicit bounds for the fact that the exact minimization of the logistic loss [Friedman et al.(2000)] may not be as efficient as AdaBoost
 - Optimizing the squared loss would be even less efficient
 - optimizing the empirical would be the less efficient of all criteria (!)

- 3 **No-Free lunch Theorems:** any algorithm has hard problems
- Find a good parameterization for ϕ : can we learn it while learning the data (self improving algorithms) ?

Thank you for your attention

Acknowledgments:

- work done in collaboration with
 - ① Frank Nielsen @ Sony CSL Tokyo
 - ② Claudia Henry (ANR/MESR PhD student)
- support from ANR, programme “Jeunes Chercheurs” JC 9009

Thank you for your attention

For more information:

- 1 see paper [Henry et al.(2007)] and longer version:
Boosting does not get Lost in Translation
(Nock, Henry, Nielsen), 45pp, submitted
- 2 see also:
On Permissible Surrogates for Classification
(Nock, Nielsen), 52pp, submitted
(available upon request)

Bibliography



Breiman, L., Freidman, J. H., Olshen, R. A., Stone, C. J., 1984. Classification and regression trees. Wadsworth.



Collins, M., Schapire, R., Singer, Y., 2002. Logistic regression, adaboost and Bregman distances. Machine Learning , 253–285.



Feldman, V., 2006. Optimal hardness results for maximizing agreements with monomials. In: Proc. of the 21st IEEE International Conference on Computational Complexity.



Friedman, J., Hastie, T., Tibshirani, R., 2000. Additive Logistic Regression : a Statistical View of Boosting. Ann. of Stat. 28, 337–374.



Henry, C., Nock, R., Nielsen, F., 2007. Real boosting *a la Carte* with an application to boosting Oblique Decision Trees. In: Proc. of the 21st International Joint Conference on Artificial Intelligence.



Kearns, M., Mansour, Y., 1999. On the boosting ability of top-down decision tree learning algorithms. Journal of Computer and System Sciences 58, 109–128.



Kearns, M. J., Vazirani, U. V., 1994. An Introduction to Computational Learning Theory. M.I.T. Press.



Mannor, S., Meir, R., 2000. Weak learners and improved rates of convergence in boosting. In: Advances in Neural Information Processing Systems 13.



Nock, R., Nielsen, F., 2004. On Domain-Partitioning Induction Criteria: Worst-case Bounds for the Worst-case Based. Theoretical Computer Science 321, 371–382.



Nock, R., Nielsen, F., 2007a. A Real Generalization of discrete AdaBoost. Artificial Intelligence 171, 25–41.



Nock, R., Nielsen, F., 2007b. Self-Improved gaps Almost Everywhere for the Agnostic Approximation of Monomials. Theoretical Computer Science 377, 139–150.



Quinlan, J. R., 1993. C4.5 : programs for machine learning. Morgan Kaufmann.