

---

# How to Boost Any Loss Function

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

1        Boosting is a highly successful ML-born optimization setting in which one is  
2        required to computationally efficiently learn arbitrarily good models based on the  
3        access to a weak learner oracle, providing classifiers performing at least slightly dif-  
4        ferently from random guessing. A key difference with gradient-based optimization  
5        is that boosting’s original model does not require access to first order information  
6        about a loss, yet the decades long history of boosting has quickly evolved it into a  
7        first order optimization setting – sometimes even wrongfully *defining* it as such.  
8        Owing to recent progress extending gradient-based optimization to use only a  
9        loss’ zeroth ( $0^{th}$ ) order information to learn, this begs the question: what loss  
10      functions can be efficiently optimized with boosting and what is the information  
11      really needed for boosting to meet the *original* boosting blueprint’s requirements?  
12      We provide a constructive formal answer essentially showing that *any* loss function  
13      can be optimized with boosting and thus boosting can achieve a feat not yet  
14      known to be possible in the classical  $0^{th}$  order setting, since loss functions are not  
15      required to be convex, nor differentiable or Lipschitz – and in fact not required  
16      to be continuous either. Some tools we use are rooted in quantum calculus, the  
17      mathematical field – not to be confounded with quantum computation – that studies  
18      calculus without passing to the limit, and thus without using first order information.

19      

## 1 Introduction

20      In ML, zeroth order optimization has been devised as an alternative to techniques that would  
21      otherwise require access to  $\geq 1$ -order information about the loss to minimize, such as gradient  
22      descent (stochastic or not, constrained or not, etc., see Section 2). Such approaches replace the access  
23      to a so-called *oracle* providing derivatives for the loss at hand, operations that can be consuming  
24      or not available in exact form in the ML world, by the access to a cheaper function value oracle,  
25      providing loss values at queried points.

26      Zeroth order optimization has seen a considerable boost in ML over the past years, over many  
27      settings and algorithms, yet, there is one foundational ML setting and related algorithms that, to our  
28      knowledge, have not yet been the subject of investigations: boosting [33, 32]. Such a question is  
29      very relevant: boosting has quickly evolved as a technique requiring first-order information about  
30      the loss optimized [6, Section 10.3], [42, Section 7.2.2] [56]. It is also not uncommon to find  
31      boosting reduced to this first-order setting [10]. However, originally, the boosting model did not  
32      mandate the access to any first-order information about the loss, rather requiring access to a weak  
33      learner providing classifiers at least slightly different from random guessing [32]. In the context of  
34      zeroth-order optimization gaining traction in ML, it becomes crucial to understand not just whether  
35      differentiability is necessary for boosting, but more generally what are loss functions that can be  
36      boosted with a weak learner and *in fine* where boosting stands with respect to recent formal progress  
37      on lifting gradient descent to zeroth-order optimisation.

38 In this paper, we settle the question: we design a formal boosting algorithm for any loss function  
 39 whose set of discontinuities has zero Lebesgue measure. With traditional floating point encoding  
 40 (e.g. float64), any stored loss function would *de facto* meet this condition; mathematically speaking,  
 41 we encompass losses that are not necessarily convex, nor differentiable or Lipschitz. This is a key  
 42 difference with classical zeroth-order optimization results where the algorithms are zeroth-order *but*  
 43 their proof of convergence makes various assumptions about the loss at hand, such as convexity,  
 44 differentiability (once or twice), Lipschitzness, etc. . Our trick to avoid the use of derivatives in  
 45 boosting relies on using or extending tools from quantum calculus\*, some of which appear to be  
 46 standard in the analysis of zeroth-order optimization. To preserve readability and save space, all  
 47 proofs and additional information are postponed to an Appendix.

## 48 2 Related work

49 Over the past years, ML has seen a substantial push to get the cheapest optimisation routines, in  
 50 general batch [15], online [28], distributed [3], adversarial [21, 19] or bandits settings [2] or more  
 51 specific settings like projection-free [27, 29, 54] or saddle-point optimisation [26, 39]. We summarize  
 52 several dozen recent references in Table A1 in terms of assumptions for the analysis about the loss  
 53 optimized, provided in Appendix, Section I. Zeroth-order optimization reduces the information  
 54 available to the learner to the "cheapest" one which consists in (loss) function values, usually via  
 55 a so-called function value *oracle*. However, as Table A1 shows, the loss itself is always assumed  
 56 to have some form of "niceness" to study the algorithms' convergence, such as differentiability,  
 57 Lipschitzness, convexity, etc. . Another quite remarkable phenomenon is that throughout all their  
 58 diverse settings and frameworks, not a single one of them addresses boosting. Boosting is however  
 59 a natural candidate for such investigations, for two reasons. First, the most widely used boosting  
 60 algorithms are first-order information hungry [6, 42, 56]: they require access to derivatives to compute  
 61 examples' weights and classifiers' leveraging coefficients. Second and perhaps most importantly,  
 62 unlike other optimization techniques like gradient descent, the original boosting model *does not*  
 63 mandate the access to a first-order information oracle to learn, but rather to a weak learning oracle  
 64 which supplies classifiers performing slightly differently from random guessing [33, 32]. Only few  
 65 approaches exist to get to "cheaper" algorithms relying on less assumptions about the loss at hand,  
 66 and to our knowledge do not have boosting-compliant convergence proofs, as for example when  
 67 alleviating convexity [17, 48] or access to gradients of the loss [57]. Such questions are however  
 68 important given the early negative results on boosting convex potentials with first-order information  
 69 [38] and the role of the classifiers in the negative results [40].

70 Finally, we note that a rich literature has developed in mathematics as well for derivative-free  
 71 optimisation [35], yet methods would also often rely on assumptions included in the three above (e.g.  
 72 [43]). It must be noted however that derivative-free optimisation has been implemented in computers  
 73 for more than seven decades [25].

## 74 3 Definitions and notations

The following shorthands are used:  $[n] \doteq \{1, 2, \dots, n\}$  for  $n \in \mathbb{N}_*$ ,  $z \cdot [a, b] \doteq [min\{za, zb\}, max\{za, zb\}]$  for  $z \in \mathbb{R}, a \leq b \in \mathbb{R}$ . In the batch supervised learning setting, one is given a training set of  $m$  examples  $S \doteq \{(\mathbf{x}_i, y_i), i \in [m]\}$ , where  $\mathbf{x}_i \in \mathcal{X}$  is an observation ( $\mathcal{X}$  is called the domain: often,  $\mathcal{X} \subseteq \mathbb{R}^d$ ) and  $y_i \in \mathcal{Y} \doteq \{-1, 1\}$  is a label, or class. We study the empirical convergence of boosting, which requires fast convergence on training. We do not investigate the questions of generalization, which would entail specific design choices about the loss at hand (see e.g. [9]). The objective is to learn a *classifier*, i.e. a function  $h : \mathcal{X} \rightarrow \mathbb{R}$  which belongs to a given set  $\mathcal{H}$ . The goodness of fit of some  $h$  on  $S$  is evaluated from a given function  $F : \mathbb{R} \rightarrow \mathbb{R}$  called a loss function, whose expectation on training is sought to be minimized:

$$F(S, h) \doteq \mathbb{E}_{i \sim [m]}[F(y_i h(\mathbf{x}_i))]. \quad (1)$$

75 The set of most popular losses comprises convex functions: the exponential loss ( $F_{\text{EXP}}(z) \doteq \exp(-z)$ ),  
 76 the logistic loss ( $F_{\text{LOG}}(z) \doteq \log(1 + \exp(-z))$ ), the square loss ( $F_{\text{SQ}}(z) \doteq (1 - z)^2$ ), the Hinge loss  
 77 ( $F_{\text{H}}(z) \doteq \max\{0, 1 - z\}$ ). These losses have fundamental differences in terms of their relationships  
 78 to *proper losses* and *surrogate losses*, two dual views of losses related to the fact that Bayes rules  
 79 is an optimal predictor for the loss [46, 47, 51]. Our examples are surrogate losses because they all  
 80 define upperbounds of the 0/1-loss ( $F_{0/1}(z) \doteq 1_{z \leq 0}$ , "1" being the indicator variable), furthermore  
 81 *calibrated* because their derivative in  $z = 0$  is negative [8].

---

\*Calculus "without limits" [31] (thus without using derivatives), not to be confounded with calculus on quantum devices.

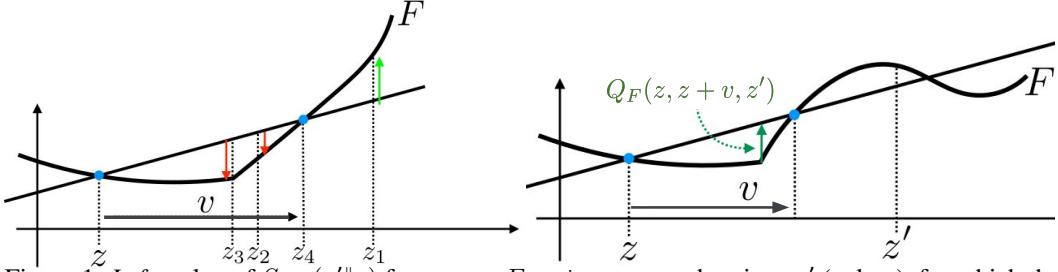


Figure 1: Left: value of  $S_{F|v}(z' \| z)$  for convex  $F$ ,  $v = z_4 - z$  and various  $z'$  (colors), for which the Bregman Secant distortion is positive ( $z' = z_1$ , green), negative ( $z' = z_2$ , red), minimal ( $z' = z_3$ ) or null ( $z' = z_4, z$ ). Right: depiction of  $Q_F(z, z + v, z')$  for non-convex  $F$  (Definition 4.6).

## 82 4 $v$ -derivatives, Bregman secant distortions

83 Unless otherwise stated, in this Section,  $F$  is a function defined over  $\mathbb{R}$ .

84 **Definition 4.1.** [31] For any  $z, v \in \mathbb{R}$ , we let  $\delta_v F(z) \doteq (F(z + v) - F(z))/v$  denote the  $v$ -derivative  
85 of  $F$  in  $z$ .

86 This expression, which gives the classical derivative when the offset  $v \rightarrow 0$ , is called the *h-derivative*  
87 in quantum calculus [31, Chapter 1]. We replaced the notation for the risk of confusion with classifiers.

88 Notice that the  $v$ -derivative is just the slope of the secant that passes through points  $(z, F(z))$  and  
89  $(z + v, F(z + v))$  (Figure 1). Higher order  $v$ -derivatives can be defined [31], though we shall need a  
90 more general definition that accommodates for variable offsets.

**Definition 4.2.** Let  $v_1, v_2, \dots, v_n \in \mathbb{R}$  and  $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$  and  $z \in \mathbb{R}$ . The  $\mathcal{V}$ -derivative  $\delta_{\mathcal{V}} F$  is:

$$\delta_{\mathcal{V}} F(z) \doteq \begin{cases} F(z) & \text{if } \mathcal{V} = \emptyset \\ \delta_{v_1} F(z) & \text{if } \mathcal{V} = \{v_1\} \\ \delta_{\{v_n\}}(\delta_{\mathcal{V} \setminus \{v_n\}} F)(z) & \text{otherwise} \end{cases}. \quad (2)$$

91 If  $v_i = v, \forall i \in [n]$  then we write  $\delta_v^{(n)} F(z) \doteq \delta_{\mathcal{V}} F(z)$ .

92 In the Appendix, Lemma A computes the unravelled expression of  $\delta_{\mathcal{V}} F(z)$ , showing that the order  
93 of the elements in  $\mathcal{V}$  does not matter.  $n$  is called the order of the  $\mathcal{V}$ -derivative.

94 We can now define a generalization of Bregman divergences called *Bregman Secant distortions*.

**Definition 4.3.** For any  $z, z', v \in \mathbb{R}$ , the Bregman Secant distortion  $S_{F|v}(z' \| z)$  with generator  $F$  and offset  $v$  is:

$$S_{F|v}(z' \| z) \doteq F(z') - F(z) - (z' - z)\delta_v F(z). \quad (3)$$

95 Even if  $F$  is convex, the distortion is not necessarily positive, though it is lowerbounded (Figure 1).

96 There is an intimate relationship between the Bregman Secant distortions and Bregman divergences.

97 We shall use a definition slightly more general than the original one when  $F$  is differentiable [12, eq.  
98 (1.4)], introduced in information geometry [5, Section 3.4] and recently reintroduced in ML [11].

99 **Definition 4.4.** The Bregman divergence with generator  $F$  (scalar, convex) between  $z'$  and  $z$  is  
100  $D_F(z' \| z) \doteq F(z') + F^*(z) - z'z$ , where  $F^*(z) \doteq \sup_t tz - F(t)$  is the convex conjugate of  $F$ .

101 We state the link between  $S_{F|v}$  and  $D_F$  (proof omitted).

102 **Lemma 4.5.** Suppose  $F$  strictly convex differentiable. Then  $\lim_{v \rightarrow 0} S_{F|v}(z' \| z) = D_F(z' \| F'(z))$ .

103 Relaxed forms of Bregman divergences have been introduced in information geometry [44].

**Definition 4.6.** For any  $a, b, \alpha \in \mathbb{R}$ , denote for short  $\mathbb{I}_{a,b} = [\min\{a, b\}, \max\{a, b\}]$  and  $(ab)_\alpha \doteq \alpha a + (1 - \alpha)b$ . The Optimal Bregman Information (OBI) of  $F$  defined by triple  $(a, b, c) \in \mathbb{R}^3$  is:

$$Q_F(a, b, c) \doteq \max_{\alpha: (ab)_\alpha \in \mathbb{I}_{a,c}} \{\alpha F(a) + (1 - \alpha)F(b) - F((ab)_\alpha)\}. \quad (4)$$

104 As represented in Figure 1 (right), the OBI is obtained by drawing the line passing through  $(a, F(a))$   
105 and  $(b, F(b))$  and then, in the interval  $\mathbb{I}_{a,c}$ , look for the maximal difference between the line and  
106  $F$ . We note that  $Q_F$  is non negative because  $a \in \mathbb{I}_{a,c}$  and for the choice  $\alpha = 1$ , the RHS in (4) is 0.  
107 We also note that when  $F$  is convex, the RHS is indeed the maximal Bregman information of two  
108 points in [7, Definition 2], where maximality is obtained over the probability measure. The following  
109 Lemma follows from the definition of the Bregman secant divergence and the OBI. An inspection of  
110 the functions in Figure 1 provides a graphical proof.

**Lemma 4.7.** For any  $F$ ,

$$\forall z, v, z' \in \mathbb{R}, S_{F|v}(z' \| z) \geq -Q_F(z, z + v, z'). \quad (5)$$

and if  $F$  is convex,

$$\forall z, v \in \mathbb{R}, \forall z' \notin \mathbb{I}_{z,z+v}, S_{F|v}(z' \| z) \geq 0, \quad (6)$$

$$\forall z, v, z' \in \mathbb{R}, S_{F|v}(z' \| z) \geq -Q_F(z, z + v, z + v). \quad (7)$$

We shall abbreviate the two possible forms of OBI in the RHS of (5), (7) as:

$$Q_F^*(z, z', v) \doteq \begin{cases} Q_F(z, z + v, z + v) & \text{if } F \text{ convex} \\ Q_F(z, z + v, z') & \text{otherwise} \end{cases}. \quad (8)$$

111

## 112 5 Boosting using only queries on the loss

113 We make the assumption that all training predictions of so-called "weak classifiers" are finite and  
114 non-zero.

115 **Assumption 5.1.**  $\forall t > 0, \forall i \in [m], |h_t(\mathbf{x}_i)| \in (0, +\infty)$  (we thus let  $M_t \doteq \max_i |h_t(\mathbf{x}_i)|$ ).

Excluding 0 ensures our algorithm does not make use of derivatives. If predictions can be zero, there is a simple tweak that still avoids the use of derivatives (Appendix, Section II.2). For short, we define two *edge* quantities for  $i \in [m]$  and  $t = 1, 2, \dots$ ,

$$e_{ti} \doteq \alpha_t \cdot y_i h_t(\mathbf{x}_i), \quad \tilde{e}_{ti} \doteq y_i H_t(\mathbf{x}_i), \quad (9)$$

where  $\alpha_t$  is a leveraging coefficient for the weak classifiers in an ensemble  $H_T(\cdot) \doteq \sum_{t \in [T]} \alpha_t h_t(\cdot)$ . We observe

$$\tilde{e}_{ti} = \tilde{e}_{(t-1)i} + e_{ti}. \quad (10)$$

116  
117

### 5.1 Algorithm: SECBOOST

#### 118 5.1.1 General steps

Without further ado, Algorithm SECBOOST presents our approach to boosting without using derivatives information. The key differences with traditional boosting algorithms are red color framed. We summarize its key steps.

**Step 1** This is the initialization step. Traditionally in boosting, one would pick  $h_0 = 0$ . Note that  $w$  is not necessarily positive.  $v_0$  is the initial offset (Section 4).

**Step 2.1** This step calls the weak learner, as in traditional boosting, using variable "weights" on examples (the absolute value of  $w$ ). The key difference with traditional boosting is that examples labels can switch between iterations as well.

**Step 2.3** This step computes the leveraging coefficient  $\alpha_t$  of the weak classifier  $h_t$ . It involves a quantity,  $\bar{W}_{2,t}$ , which we define as any strictly positive real satisfying

$$\mathbb{E}_{i \sim [m]} \left[ \delta_{\{e_{ti}, v_{(t-1)i}\}} F(\tilde{e}_{(t-1)i}) \cdot \left( \frac{h_t(\mathbf{x}_i)}{M_t} \right)^2 \right] \leq \bar{W}_{2,t}. \quad (14)$$

119 For boosting rate's sake, we should find  $\bar{W}_{2,t}$  as small as possible. We refer to (9) for the  $e, \tilde{e}$   
120 notations;  $v$  is the current (set of) offset(s) (Section 4 for their definition). The second-order  $\mathcal{V}$ -  
121 derivative in the LHS plays the same role as the second-order derivative in classical boosting rates,  
122 see for example [47, Appendix, Section 4]. As offsets  $\rightarrow 0$ , it converges to a second-order derivative;  
123 otherwise, they still share some properties, such as the sign for convex functions.

124 **Lemma 5.2.** Suppose  $F$  convex. For any  $a \in \mathbb{R}, b, c \in \mathbb{R}_*, \delta_{\{b,c\}} F(a) > 0$ .

(Proof in Appendix, Section II.3) We can also see a link with weights variation since, modulo a slight abuse of notation, we have  $\delta_{\{e_{ti}, v_{(t-1)i}\}} F(\tilde{e}_{(t-1)i}) = \delta_{e_{ti}} w_{ti}$ . A substantial difference with traditional boosting algorithms is that we have two ways to pick the leveraging coefficient  $\alpha_t$ ; the first one can be used when a convenient  $\bar{W}_{2,t}$  is directly accessible from the loss. Otherwise, there is a simple algorithm that provides parameters (including  $\bar{W}_{2,t}$ ) such that (14) is satisfied. Section 5.3

---

**Algorithm 1** SECBOOST( $\mathcal{S}, T$ )

---

**Input** sample  $\mathcal{S} = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, m\}$ , number of iterations  $T$ , initial  $(h_0, v_0)$  (constant classification and offset).

Step 1 : let  $H_0 \leftarrow 1 \cdot h_0$  and  $\mathbf{w}_1 = -\delta_{v_0} F(h_0) \cdot \mathbf{1}$ ; //  $h_0, v_0 \neq 0$  chosen s. t.  $\delta_{v_0} F(h_0) \neq 0$

Step 2 : **for**  $t = 1, 2, \dots, T$

Step 2.1 : let  $h_t \leftarrow \times(\mathcal{S}_t, |\mathbf{w}_t|)$  //weak learner call,  $\mathcal{S}_t \doteq \{(\mathbf{x}_i, y_i \cdot \text{sign}(w_{ti}))\}$

Step 2.2 : let  $\eta_t \leftarrow (1/m) \cdot \sum_i w_{ti} y_i h_t(\mathbf{x}_i)$  //unnormalized edge

Step 2.3 :	If bound on $\bar{W}_{2,t}$ available (Section 5.3) pick $\varepsilon_t > 0, \pi_t \in (0, 1)$ and $\alpha_t \in \frac{\eta_t}{2(1 + \varepsilon_t)M_t^2\bar{W}_{2,t}} \cdot [1 - \pi_t, 1 + \pi_t]; \quad (11)$	otherwise   general procedure $\alpha_t \leftarrow \text{SOLVE}_\alpha(\mathcal{S}, \mathbf{w}_t, h_t)$ // $\bar{W}_{2,t} > 0, \varepsilon_t > 0, \pi_t \in (0, 1)$ // Theorem 5.8
------------	--	---

Step 2.4 : let  $H_t \leftarrow H_{t-1} + \alpha_t \cdot h_t$  //classifier update

Step 2.5 : **if**  $\mathbb{I}_{ti}(\varepsilon_t \cdot \alpha_t^2 M_t^2 \bar{W}_{2,t}) \neq \emptyset, \forall i \in [m]$  **then** //new offsets

**for**  $i = 1, 2, \dots, m$ , let

$v_{ti} \leftarrow \text{OO}(t, i, \varepsilon_t \cdot \alpha_t^2 M_t^2 \bar{W}_{2,t})$  ;

**else return**  $H_t$ ;

Step 2.6 : **for**  $i = 1, 2, \dots, m$ , let //weight update

$w_{(t+1)i} \leftarrow -\delta_{v_{ti}} F(y_i H_t(\mathbf{x}_i))$  ;

Step 2.7 : **if**  $\mathbf{w}_{t+1} = \mathbf{0}$  **then break**;

**Return**  $H_T$ .

---

details those two possibilities and their implementation. In the more favorable case (the former one),  $\alpha_t$  can be chosen in an interval, furthermore defined by flexible parameters  $\varepsilon_t > 0, \pi_t \in (0, 1)$ . Note that fixing beforehand these parameters is not mandatory: we can also pick *any*

$$\alpha_t \in \eta_t \cdot (0, 1/(M_t^2 \bar{W}_{2,t})) , \quad (15)$$

and then compute choices for the corresponding  $\varepsilon_t$  and  $\pi_t$ .  $\varepsilon_t$  is important for the algorithm and both parameters are important for the analysis of the boosting rate. From the boosting standpoint, a smaller  $\varepsilon_t$  yields a larger  $\alpha_t$  and a smaller  $\pi_t$  reduces the interval of values in which we can pick  $\alpha_t$ ; both cases tend to favor better convergence rates as seen in Theorem 5.3.

Step 2.4 is just the crafting of the final model.

Step 2.5 is new to boosting, the use of a so-called offset oracle, detailed in Section 5.1.2.

Step 2.6 The weight update does not rely on a first-order oracle as in traditional boosting, but uses only loss values through  $v$ -derivatives. The finiteness of  $F$  implies the finiteness of weights.

Step 2.7 Early stopping happens if all weights are null. While this would never happen with traditional (*e.g.* strictly convex) losses, some losses that are unusual in the context of boosting can lead to early stopping. A discussion on early stopping and how to avoid it is in Section 6.

5.1.2 The offset oracle,  $\text{OO}$

Let us introduce notation

$$\mathbb{I}_{ti}(z) \doteq \{v : Q_F^*(\tilde{e}_{ti}, \tilde{e}_{(t-1)i}, v) \leq z\}, \forall i \in [m], \forall z > 0. \quad (16)$$

(see Figure 3 below to visualize  $\mathbb{I}_{ti}(z)$  for a non-convex  $F$ ) The offset oracle is used in Step 2.5, which is new to boosting. It requests the offsets to carry out weight update in (13) to an *offset oracle*, which achieves the following, for iteration # $t$ , example # $i$ , limit OBI  $z$ :

$$\text{OO}(t, i, z) \text{ returns some } v \in \mathbb{I}_{ti}(z) \quad (17)$$

Note that the offset oracle has the freedom to pick the offset in a whole set. Section 5.4 investigates implementations of the offset oracle, so let us make a few essentially graphical remarks here.  $\text{OO}$  does not need to build the whole  $\mathbb{I}_{ti}(z)$  to return some  $v \in \mathbb{I}_{ti}(z)$  for Step 2.5 in SECBOOST. In the construction steps of Figure 3, as soon as  $\mathcal{O} \neq \emptyset$ , one element of  $\mathcal{O}$  can be returned. Figure 4 presents more examples of  $\mathbb{I}_{ti}(z)$ . One can remark that the sign of the offset  $v_{ti}$  in Step 2.5 of SECBOOST is the same as the sign of  $\tilde{e}_{(t-1)i} - \tilde{e}_{ti} = -y_i \alpha_t h_t(\mathbf{x}_i)$ . Hence, unless  $F$  is derivable or all edges  $y_i h_t(\mathbf{x}_i)$  are of the same sign ( $\forall i$ ), the set of offsets returned in Step 2.5 always contain at least two different offsets, one non-negative and one non-positive (Figure 4, (a-b)).

145 **5.2 Convergence of SECBOOST**

146 The offset oracle has a technical importance for boosting:  $\mathbb{I}_{ti}(z)$  is the set of offsets that limit an  
 147 OBI for a training example (Definition 4.6). The importance for boosting comes from Lemma  
 148 4.7: upperbounding an OBI implies lowerbounding a Bregman Secant divergence, which will also  
 149 guarantee a sufficient slack between two successive boosting iterations. This is embedded in a  
 150 blueprint of a proof technique to show boosting-compliant convergence which is not new, see e.g.  
 151 [47]. We now detail this convergence.

Remark that the expected edge  $\eta_t$  in Step 2.2 of SECBOOST is not normalized. We define a normalized  
 version of this edge as:

$$[-1, 1] \ni \tilde{\eta}_t \doteq \sum_i \frac{|w_{ti}|}{W_t} \cdot \tilde{y}_{ti} \cdot \frac{h_t(\mathbf{x}_i)}{M_t}, \quad (18)$$

with  $\tilde{y}_{ti} \doteq y_i \cdot \text{sign}(w_{ti})$ ,  $W_t \doteq \sum_i |w_{ti}| = \sum_i |\delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})|$ . Remark that the labels are  
 corrected by the weight sign and thus may switch between iterations. In the particular case where the  
 loss is non-increasing (such as with traditional convex surrogates), the labels do not switch. We need  
 also a quantity which is, in absolute value, the expected weight:

$$\bar{W}_{1,t} \doteq |\mathbb{E}_{i \sim [m]} [\delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})]| \quad (19)$$

152 (we indeed observe  $\bar{W}_{1,t} = |\mathbb{E}_{i \sim [m]} [w_{ti}]|$ ) In classical boosting for convex decreasing losses<sup>†</sup>,  
 153 weights are non-negative and converge to a minimum (typically 0) as examples get the right class  
 154 with increasing confidence. Thus,  $\bar{W}_{1,t}$  can be an indicator of when classification becomes "good  
 155 enough" to stop boosting. In our more general setting, it shall be used in a similar indicator. We are  
 156 now in a position to show a first result about SECBOOST.

**Theorem 5.3.** *Suppose assumption 5.1 holds. Let  $F_0 \doteq F(S, h_0)$  in SECBOOST and  $z^*$  any real  
 such that  $F(z^*) \leq F_0$ . Then we are guaranteed that classifier  $H_T$  output by SECBOOST satisfies  
 $F(S, H_T) \leq F(z^*)$  when the number of boosting iterations  $T$  yields:*

$$\sum_{t=1}^T \frac{\bar{W}_{1,t}^2 (1 - \pi_t^2)}{\bar{W}_{2,t} (1 + \varepsilon_t)} \cdot \tilde{\eta}_t^2 \geq 4(F_0 - F(z^*)), \quad (20)$$

157 where parameters  $\varepsilon_t, \pi_t$  appear in Step 2.3 of SECBOOST.

158 (proof in Appendix, Section II.4) We observe the tradeoff between the freedom in picking parameters  
 159 and convergence guarantee as exposed by (20): to get more freedom in picking the leveraging  
 160 coefficient  $\alpha_t$ , we typically need  $\pi_t$  large (Step 2.3) and to get more freedom in picking the offset  
 161  $v_t \neq 0$ , we typically need  $\varepsilon_t$  large (Step 2.5). However, allowing more freedom in such ways  
 162 reduces the LHS and thus impairs the guarantee in (20). Therefore, there is a subtle balance between  
 163 "freedom" of choice and convergence. Figure 2 pictures notable regimes for  $\bar{W}_{1,t}$  and  $\bar{W}_{2,t}$ , leading  
 164 to varying contributions in (20).

**Boosting-compliant convergence** We characterize convergence in the boosting framework. We  
 define

$$\rho_t \doteq \bar{W}_{1,t}^2 / \bar{W}_{2,t}.$$

165 A small  $\bar{W}_{1,t}^2$  is an indicator as to whether SECBOOST is close to a minimum, as seen from the  
 166 left (blue) block in Figure 2. Analysis of convergence based on a minimal value of  $\bar{W}_{1,t}^2$  would be  
 167 standard with respect to classical non-convex optimisation. In our case, we construct a criterion  
 168 involving not just  $\bar{W}_{1,t}^2$ , but also  $\bar{W}_{2,t}$ : considering for the illustration a  $C^2$  loss function, when  
 169  $\bar{W}_{2,t}$  is small, it accounts for regions with "small second-order" variations, more "regularity", which  
 170 naturally offers more leeway for minimization, in particular if  $\bar{W}_{1,t}$  is large (see Figure 2). This  
 171 justifies the following assumption.

172 **Assumption 5.4.** ( $\rho_*$ -Convergence Regime,  $\rho_*$ -CR) We assume there exists  $\rho_* > 0$  such that  
 173  $\forall t \geq 1, \rho_t > \rho_*$ .

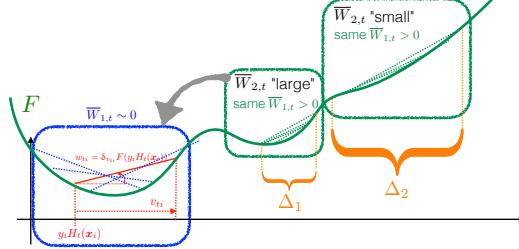


Figure 2: When SECBOOST approaches an optimum of the loss (blue + red), we can expect  $\bar{W}_{1,t}$  to be small. At a fixed  $\bar{W}_{1,t}$ , a larger  $\bar{W}_{2,t}$  indicates, at least locally near the secant points, a region of higher curvature (green blocks). Since the maximal difference of abscissae for a secant is proportional to the edge of the last weak classifier (Figure 1, right, (4) and (16)), the smaller  $\bar{W}_{2,t}$  indicates a greater potential for decrease of the loss (right green block,  $\Delta_2 > \Delta_1$ ) and leads to a greater contribution to the sum in (20). When  $\bar{W}_{1,t}^2$  is small enough compared to  $\bar{W}_{2,t}$ , we take it as an indicator of the proximity of a local optimum, which justifies fixing a lowerbound  $\rho_*$  to  $\rho_t$  to analyze convergence "outside" this regime (best viewed in color). Bold gray arrow is discussion in Section 6.

174 We also rely on boosting's traditional weak learning assumption.

175 **Assumption 5.5. ( $\gamma$ -Weak Learning Assumption,  $\gamma$ -WLA)** *We assume the following on the weak learner:  $\exists \gamma > 0$  such that  $\forall t > 0$ ,  $|\dot{\eta}_t| \geq \gamma$ .*

177 We are now in a position to state a simple corollary to Theorem 5.3.

**Corollary 5.6.** *Suppose assumptions 5.1, 5.4 and 5.5 hold. Let  $F_0 \doteq F(S, h_0)$  in SECBOOST and  $z^*$  any real such that  $F(z^*) \leq F_0$ . If SECBOOST is run for a number  $T$  of iterations satisfying*

$$T \geq \frac{4(F_0 - F(z^*))}{\gamma^2 \rho_*} \cdot \frac{1 + \max_{t \in [T]} \varepsilon_t}{1 - \max_{t \in [T]} \pi_t^2}, \quad (21)$$

178 then  $F(S, H_T) \leq F(z^*)$ .

179 We remark that the dependency in  $\gamma$  is optimal [4].

### 180 5.3 Finding $\bar{W}_{2,t}$

181 There is lots of freedom in the choice of  $\alpha_t$  in Step 2.3 of SECBOOST, and even more if we look  
182 at (15). This, however, requires access to some bound  $\bar{W}_{2,t}$ . In the general case, the quantity it  
183 upperbounds in (14) also depends on  $\alpha_t$  because  $e_{ti} \doteq \alpha_t \cdot y_i h_t(\mathbf{x}_i)$ . So unless we can obtain such a  
184 "simple"  $\bar{W}_{2,t}$  that does *not* depend on  $\alpha_t$ , (11) – and (15) – provide a *system* to solve for  $\alpha_t$ .

185  **$\bar{W}_{2,t}$  via properties of  $F$**  Classical assumptions on loss functions for zeroth-order optimization can  
186 provide simple expressions for  $\bar{W}_{2,t}$  (Table A1). Consider smoothness: we say that  $F$  is  $\beta$ -smooth if  
187 it is derivable and its derivative satisfies the Lipschitz condition  $|F'(z') - F'(z)| \leq \beta|z' - z|, \forall z, z'$   
188 [13]. Notice that this implies the condition on the  $v$ -derivative of the derivative:  $|\delta_v F'(z)| \leq \beta, \forall z, v$ .  
189 This also provides a straightforward useful expression for  $\bar{W}_{2,t}$ .

190 **Lemma 5.7.** *Suppose that the loss  $F$  is  $\beta$ -smooth. Then we can fix  $\bar{W}_{2,t} = 2\beta$ .*

(Proof in Appendix, Section II.5) What the Lemma shows is that a bound on the  $v$ -derivative of the derivative implies a bound on order-2  $\mathcal{V}$ -derivatives (in the quantity that  $\bar{W}_{2,t}$  bounds (14)). Such a condition on  $v$ -derivatives is thus weaker than a condition on derivatives, and it is strictly weaker if we impose a strictly positive lowerbound on the offset's absolute value, which would be sufficient to characterize the boosting convergence of SECBOOST.

**A general algorithm for  $\bar{W}_{2,t}$**  If we cannot make any assumption on  $F$ , there is a simple way to first obtain  $\alpha_t$  and then  $\bar{W}_{2,t}$ , from which all other parameters of Step 2.3 can be computed. We first need a few definitions. We first generalize the edge notation appearing in Step 2.2:

$$\eta(\mathbf{w}, h) \doteq \mathbb{E}_{i \sim [m]} [w_i y_i h(\mathbf{x}_i)],$$

---

<sup>†</sup>This is an important class of losses since it encompasses the convex surrogates of symmetric proper losses [45, 52]

---

**Algorithm 2** SOLVE $_{\alpha}(\mathcal{S}, \mathbf{w}, h)$ 


---

**Input** sample  $\mathcal{S} = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, m\}$ ,  $\mathbf{w} \in \mathbb{R}^m$ ,  $h : \mathcal{X} \rightarrow \mathbb{R}$ .  
Step 1 : find any  $a > 0$  such that

$$\frac{|\eta(\mathbf{w}, h) - \eta(\tilde{\mathbf{w}}(\text{sign}(\eta(\mathbf{w}, h)) \cdot a), h)|}{|\eta(\mathbf{w}, h)|} < 1. \quad (22)$$

**Return**  $\text{sign}(\eta(\mathbf{w}, h)) \cdot a$ .

---

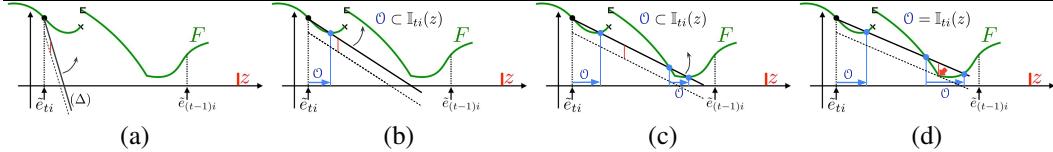


Figure 3: A simple way to build  $\mathbb{I}_{ti}(z)$  for a discontinuous loss  $F$  ( $\tilde{e}_{ti} < \tilde{e}_{(t-1)i}$  and  $z$  are represented),  $\mathcal{O}$  being the set of solutions as it is built. We rotate two half-lines, one passing through  $(\tilde{e}_{ti}, F(\tilde{e}_{ti}))$  (thick line,  $(\Delta)$ ) and a parallel one translated by  $-z$  (dashed line) (a). As soon as  $(\Delta)$  crosses  $F$  on any point  $(z', F(z'))$  with  $z \neq \tilde{e}_{ti}$  while the dashed line stays below  $F$ , we obtain a candidate offset  $v$  for  $\mathcal{O}$ , namely  $v = z' - \tilde{e}_{ti}$ . In (b), we obtain an interval of values. We keep on rotating  $(\Delta)$ , eventually making appear several intervals for the choice of  $v$  if  $F$  is not convex (c). Finally, when we reach an angle such that the maximal difference between  $(\Delta)$  and  $F$  in  $[\tilde{e}_{ti}, \tilde{e}_{(t-1)i}]$  is  $z$  ( $z$  can be located at an intersection between  $F$  and the dashed line), we stop and obtain the full  $\mathbb{I}_{ti}(z)$  (d).

so that  $\eta_t \doteq \eta(\mathbf{w}_t, h_t)$ . Remind the weight update,  $w_{ti} \doteq -\delta_{v_{(t-1)i}} F(y_i H_{t-1}(\mathbf{x}_i))$ . We define a "partial" weight update,

$$\tilde{w}_{ti}(\alpha) \doteq -\delta_{v_{(t-1)i}} F(\alpha y_i h_t(\mathbf{x}_i) + y_i H_{t-1}(\mathbf{x}_i)) \quad (23)$$

191 (if we were to replace  $v_{(t-1)i}$  by  $v_{ti}$  and let  $\alpha \doteq \alpha_t$ , then  $\tilde{w}_{ti}(\alpha)$  would be  $w_{(t+1)i}$ , hence the partial  
192 weight update). Algorithm 2 presents the simple procedure to find  $\alpha_t$ . Notice that we use  $\tilde{\mathbf{w}}$  with  
193 sole dependency on the prospective leveraging coefficient; we omit for clarity the dependences in the  
194 current ensemble ( $H_{\cdot}$ ), weak classifier ( $h_{\cdot}$ ) and offsets ( $v_{\cdot i}$ ) needed to compute (23).

**Theorem 5.8.** Suppose Assumptions 5.1 and 5.5 hold and  $F$  is continuous at all abscissae  $\{\tilde{e}_{(t-1)i} \doteq y_i H_{t-1}(\mathbf{x}_i), i \in [m]\}$ . Then there are always solutions to Step 1 of SOLVE $_{\alpha}$  and if we let  $\alpha_t \leftarrow \text{SOLVE}_{\alpha}(\mathcal{S}, \mathbf{w}_t, h_t)$  and then compute

$$\overline{W}_{2,t} \doteq \left| \mathbb{E}_{i \sim [m]} \left[ \frac{h_t^2(\mathbf{x}_i)}{M_t^2} \cdot \delta_{\{\alpha_t y_i h_t(\mathbf{x}_i), v_{(t-1)i}\}} F(\tilde{e}_{(t-1)i}) \right] \right|,$$

195 then  $\overline{W}_{2,t}$  satisfies (14) and  $\alpha_t$  satisfies (11) for some  $\varepsilon_t > 0, \pi_t \in (0, 1)$ .

196 The proof, in Section II.6, proceeds by reducing condition (15) to (22). The Weak Learning Assumption  
197 (5.5) is important for the denominator in the LHS of (22) to be non zero. The continuity  
198 assumption at all abscissae is important to have  $\lim_{a \rightarrow 0} \eta(\tilde{\mathbf{w}}_t(a), h_t) = \eta_t$ , which ensures the  
199 existence of solutions to (22), also easy to find, e.g. by a simple dichotomic search starting from an  
200 initial guess for  $a$ . Note the necessity of being continuous only at abscissae defined by the training  
201 sample, which is finite in size. Hence, if this condition is not satisfied but discontinuities of  $F$  are of  
202 Lebesgue measure 0, it is easy to add an infinitesimal constant to the current weak classifier, ensuring  
203 the conditions of Theorem 5.8 and keeping the boosting rates.

204 **5.4 Implementation of the offset oracle**

205 Figure 3 explains how to build graphically  $\mathbb{I}_{ti}(z)$  for a general  $F$ . While it is not hard to implement a  
206 general procedure following the blueprint (i.e. accepting the loss function as input), it would be far  
207 from achieving computational optimality: a much better choice consists in specializing it to the (set  
208 of) loss(es) at hand via hardcoding specific optimization features of the desired loss(es). This would  
209 not prevent "loss oddities" to get absolutely trivial oracles (see Appendix, Section II.7).

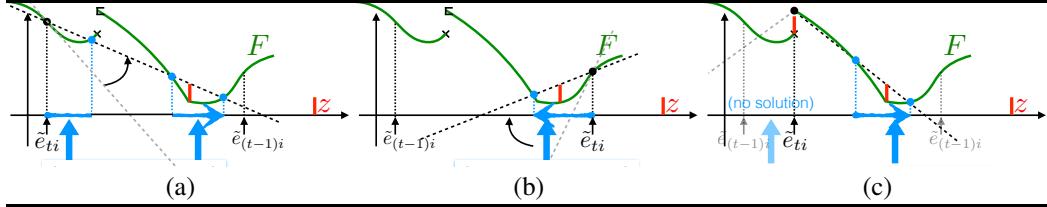


Figure 4: More examples of ensembles  $\mathbb{I}_{ti}(z)$  (in blue) for the  $F$  in Figure 3. (a):  $\mathbb{I}_{ti}(z)$  is the union of two intervals with all candidate offsets non negative. (b): it is a single interval with non-positive offsets. (c): at a discontinuity, if  $z$  is smaller than the discontinuity, we have no direct solution for  $\mathbb{I}_{ti}(z)$  for at least one positioning of the edges, but a simple trick bypasses the difficulty (see text).

## 210 6 Discussion

211 For an efficient implementation, boosting requires specific design choices to make sure the weak  
 212 learning assumption stands for as long as necessary; experimentally, it is thus a good idea to adapt  
 213 the weak learner to build more complex models as iterations increase (*e.g.* learning deeper trees),  
 214 keeping Assumption 5.5 valid with its advantage over random guessing parameter  $\gamma > 0$ . In our  
 215 more general setting, our algorithm SECBOOST pinpoints two more locations that can make use of  
 216 specific design choices to keep assumptions stand for a larger number of iterations.

217 The first is related to handling local minima. When Assumption 5.4 breaks, it means we are close to  
 218 a local optimum of the loss. One possible way of escaping those local minima is to adapt the offset  
 219 oracle to output larger offsets (Step 2.5) that get weights computed outside the domain of the local  
 220 minimum. As an illustration, in Figure 2, the middle green box shows a local minimum in which  
 221 SECBOOST can get trapped if offsets are small enough (slopes of mixed signs, just like in the blue  
 222 area;  $\bar{W}_{1,t}$  is small, Assumption 5.4 breaks). However, for larger offsets, many slopes will tip to  
 223 being positive with one of their intersection with  $F$  in the blue area, signalling a better optimum for  
 224 prediction in this blue area that SECBOOST can reach with its model update, schematized with a gray  
 225 arrow ( $\bar{W}_{1,t}$  is large, Assumption 5.4 does not break). Section 5.4 has presented a general blueprint  
 226 for the offset oracle but more specific implementation designs can be used; some are discussed in the  
 227 Appendix, Section II.7.

228 The second is related to handling losses that take on constant values over parts of their domain. To  
 229 prevent early stopping in Step 2.7 of SECBOOST, one needs  $w_{t+1} \neq 0$ . The update rule of  $w_t$   
 230 imposes that the loss must then have non-zero variation for some examples between two successive  
 231 edges (9). If the loss  $F$  is constant, then clearly the algorithm obviously stops without learning  
 232 anything. If  $F$  is piecewise-constant, this constrain the design of the weak learner to make sure that  
 233 some examples receive a different loss with the new model update  $H_+$ . As explained in Appendix,  
 234 Section II.11, this can be efficiently addressed by specific designs on  $SOLVE_\alpha$ .

235 In the same way as there is no "1 size fits all" weak learner for all domains in traditional boosting,  
 236 we expect specific design choices to be instrumental in better handling specific losses in our more  
 237 general setting. Our theory points two locations further work can focus on.

## 238 7 Conclusion

239 Boosting has rapidly moved to an optimization setting involving first-order information about the  
 240 loss optimized, rejoining, in terms of information needed, that of the hugely popular (stochastic)  
 241 gradient descent. But this was not a formal requirement of the initial setting and in this paper, we  
 242 show that essentially any loss function can be boosted without this requirement. From this standpoint,  
 243 our results put boosting in a slightly more favorable light than recent development on zeroth-order  
 244 optimization since, to get boosting-compliant convergence, we do not need the loss to meet any  
 245 of the assumptions that those analyses usually rely on. Of course, recent advances in zeroth-order  
 246 optimization have also achieved substantial design tricks for the implementation of such algorithms,  
 247 something that undoubtedly needs to be addressed in our case, such as for the efficient optimization of  
 248 the offset oracle. We leave this as an open problem but provide in Appendix some toy experiments  
 249 that a straightforward implementation achieves, hinting that SECBOOST can indeed optimize very  
 250 "exotic" losses.

251 **References**

- 252 [1] A. Akhavan, E. Chzhen, M. Pontil, and A.-B. Tsybakov. A gradient estimator via 11-  
253 randomization for online zero-order optimization with two point feedback. In *NeurIPS\*35*,  
254 2022.
- 255 [2] A. Akhavan, M. Pontil, and A.-B. Tsybakov. Exploiting higher order smoothness in derivative-  
256 free optimization and continuous bandits. In *NeurIPS\*33*, 2020.
- 257 [3] A. Akhavan, M. Pontil, and A.-B. Tsybakov. Distributed zero-order optimisation under adver-  
258 sarial noise. In *NeurIPS\*34*, 2021.
- 259 [4] N. Alon, A. Gonen, E. Hazan, and S. Moran. Boosting simple learners. In *STOC’21*, 2021.
- 260 [5] S.-I. Amari and H. Nagaoka. *Methods of Information Geometry*. Oxford University Press, 2000.
- 261 [6] F. Bach. *Learning Theory from First Principles*. Course notes, MIT press (to appear), 2023.
- 262 [7] A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh. Clustering with bregman divergences. In  
263 *Proc. of the 4<sup>th</sup> SIAM International Conference on Data Mining*, pages 234–245, 2004.
- 264 [8] P. Bartlett, M. Jordan, and J. D. McAuliffe. Convexity, classification, and risk bounds. *J. of the*  
265 *Am. Stat. Assoc.*, 101:138–156, 2006.
- 266 [9] P.-L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and  
267 structural results. *JMLR*, 3:463–482, 2002.
- 268 [10] G. Biau, B. Cadre, and L. Rouvière. Accelerated gradient boosting. *Mach. Learn.*, 108(6):971–  
269 992, 2019.
- 270 [11] M. Blondel, A.-F. T. Martins, and V. Niculae. Learning with Fenchel-Young losses. *J. Mach.*  
271 *Learn. Res.*, 21:35:1–35:69, 2020.
- 272 [12] L. M. Bregman. The relaxation method of finding the common point of convex sets and its  
273 application to the solution of problems in convex programming. *USSR Comp. Math. and Math.*  
274 *Phys.*, 7:200–217, 1967.
- 275 [13] S. Bubeck. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*,  
276 8(3-4):231–357, 2015.
- 277 [14] P.-S. Bullen. *Handbook of means and their inequalities*. Kluwer Academic Publishers, 2003.
- 278 [15] H. Cai, Y. Lou, D. McKenzie, and W. Yin. A zeroth-order block coordinate descent algorithm  
279 for huge-scale black-box optimization. In *38<sup>th</sup> ICML*, pages 1193–1203, 2021.
- 280 [16] C. Cartis and L. Roberts. Scalable subspace methods for derivative-free nonlinear least-squares  
281 optimization. *Math. Prog.*, 199:461–524, 2023.
- 282 [17] S. Cheamanunkul, E. Ettinger, and Y. Freund. Non-convex boosting overcomes random label  
283 noise. *CorR*, abs/1409.2905, 2014.
- 284 [18] L. Chen, J. Xu, and L. Luo. Faster gradient-free algorithms for nonsmooth nonconvex stochastic  
285 optimization. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023,*  
286 *Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages  
287 5219–5233. PMLR, 2023.
- 288 [19] X. Chen, S. Liu, K. Xu, X. Li, X. Lin, M. Hong, and D. Cox. ZO-AdaMM: Zeroth-order  
289 adaptive momentum method for black-box optimization. In *NeurIPS\*32*, 2019.
- 290 [20] X. Chen, Y. Tang, and N. Li. Improve single-point zeroth-order optimization using high-pass  
291 and low-pass filters. In *39<sup>th</sup> ICML*, volume 162 of *Proceedings of Machine Learning Research*,  
292 pages 3603–3620. PMLR, 2022.
- 293 [21] S. Cheng, G. Wu, and J. Zhu. On the convergence of prior-guided zeroth-order optimisation  
294 algorithms. In *NeurIPS\*34*, 2021.
- 295 [22] Z. Cranko and R. Nock. Boosted density estimation remastered. In *36<sup>th</sup> ICML*, pages 1416–  
296 1425, 2019.
- 297 [23] W. de Vazelhes, H. Zhang, H. Wu, X. Yuan, and B. Gu. Zeroth-order hard-thresholding:  
298 Gradient error vs. expansivity. In *NeurIPS\*35*, 2022.
- 299 [24] D. Dua and C. Graff. UCI machine learning repository, 2021.

- 300 [25] E. Fermi and N. Metropolis. Numerical solutions of a minimum problem. Technical Report TR  
 301 LA-1492, Los Alamos Scientific Laboratory of the University of California, 1952.
- 302 [26] L. Flokas, E.-V. Vlatakis-Gkaragkounis, and G. Piliouras. Efficiently avoiding saddle points  
 303 with zero order methods: No gradients required. In *NeurIPS\*32*, 2019.
- 304 [27] H. Gao and H. Huang. Can stochastic zeroth-order frank-wolfe method converge faster for  
 305 non-convex problems? In *37<sup>th</sup> ICML*, pages 3377–3386, 2020.
- 306 [28] A. Héliou, M. Martin, P. Mertikopoulos, and T. Rahier. Zeroth-order non-convex learning via  
 307 hierarchical dual averaging. In *38<sup>th</sup> ICML*, pages 4192–4202, 2021.
- 308 [29] F. Huang, L. Tao, and S. Chen. Accelerated stochastic gradient-free and projection-free methods.  
 309 In *37<sup>th</sup> ICML*, pages 4519–4530, 2020.
- 310 [30] B. Irwin, E. Haber, R. Gal, and A. Ziv. Neural network accelerated implicit filtering: Integrating  
 311 neural network surrogates with provably convergent derivative free optimization methods. In  
 312 *40<sup>th</sup> ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 14376–14389.  
 313 PMLR, 2023.
- 314 [31] V. Kac and P. Cheung. *Quantum calculus*. Springer, 2002.
- 315 [32] M. J. Kearns and U. V. Vazirani. *An Introduction to Computational Learning Theory*. M.I.T.  
 316 Press, 1994.
- 317 [33] M.J. Kearns. Thoughts on hypothesis boosting, 1988. ML class project.
- 318 [34] M.J. Kearns and Y. Mansour. On the boosting ability of top-down decision tree learning  
 319 algorithms. *J. Comp. Syst. Sc.*, 58:109–128, 1999.
- 320 [35] J. Larson, M. Menickelly, and S.-M. Wild. Derivative-free optimization methods. *Acta Numerica*,  
 321 pages 287–404, 2019.
- 322 [36] Z. Li, P.-Y. Chen, S. Liu, S. Lu, and Y. Xu. Zeroth-order optimization for composite problems  
 323 with functional constraints. In *AAAI’22*, pages 7453–7461. AAAI Press, 2022.
- 324 [37] T. Lin, Z. Zheng, and M.-I. Jordan. Gradient-free methods for deterministic and stochastic  
 325 nonsmooth nonconvex optimization. In *NeurIPS\*35*, 2022.
- 326 [38] P.-M. Long and R.-A. Servedio. Random classification noise defeats all convex potential  
 327 boosters. *MLJ*, 78(3):287–304, 2010.
- 328 [39] C. Maheshwari, C.-Y. Chiu, E. Mazumdar, S. Shankar Sastry, and L.-J. Ratliff. Zeroth-  
 329 order methods for convex-concave minmax problems: applications to decision-dependent risk  
 330 minimization. In *25<sup>th</sup> AISTATS*, 2022.
- 331 [40] Y. Mansour, R. Nock, and R.-C. Williamson. Random classification noise does not defeat all  
 332 convex potential boosters irrespective of model choice. In *40<sup>th</sup> ICML*, 2023.
- 333 [41] E. Mhanna and M. Assaad. Single point-based distributed zeroth-order optimization with a  
 334 non-convex stochastic objective function. In *40<sup>th</sup> ICML*, volume 202 of *Proceedings of Machine  
 335 Learning Research*, pages 24701–24719. PMLR, 2023.
- 336 [42] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. MIT Press,  
 337 2018.
- 338 [43] Y. Nesterov and V. Spokoiny. Random gradient-free optimization of convex functions. *Founda-  
 339 tions of Computational Mathematics*, 17:527–566, 2017.
- 340 [44] F. Nielsen and R. Nock. The Bregman chord divergence. In *Geometric Science of Information -  
 341 4th International Conference*, 2019, pages 299–308, 2019.
- 342 [45] R. Nock and A. K. Menon. Supervised learning: No loss no cry. In *37<sup>th</sup> ICML*, 2020.
- 343 [46] R. Nock and F. Nielsen. Bregman divergences and surrogates for learning. *IEEE Trans.PAMI*,  
 344 31:2048–2059, 2009.
- 345 [47] R. Nock and R.-C. Williamson. Lossless or quantized boosting with integer arithmetic. In *36<sup>th</sup>  
 346 ICML*, pages 4829–4838, 2019.
- 347 [48] N.-E. Pfetsch and Sebastian Pokutta. IPBoost - non-convex boosting via integer programming.  
 In *37<sup>th</sup> ICML*, volume 119, pages 7663–7672, 2020.
- 348 [49] Y. Qiu, U.-V. Shanbhag, and F. Yousefian. Zeroth-order methods for nondifferentiable, noncon-  
 349 vex and hierarchical federated optimization. In *NeurIPS\*36*, 2023.

- 351 [50] M. Rando, C. Molinari, L. Rosasco, and S. Villa. Structured zeroth-order for non-smooth  
352 optimization. In *NeurIPS\*36*, 2023.
- 353 [51] M.-D. Reid and R.-C. Williamson. Composite binary losses. *JMLR*, 11:2387–2422, 2010.
- 354 [52] M.-D. Reid and R.-C. Williamson. Information, divergence and risk for binary experiments.  
355 *JMLR*, 12:731–817, 2011.
- 356 [53] Z. Ren, Y. Tang, and N. Li. Escaping saddle points in zeroth-order optimization: the power of  
357 two-point estimators. In *40<sup>th</sup> ICML*, volume 202 of *Proceedings of Machine Learning Research*,  
358 pages 28914–28975. PMLR, 2023.
- 359 [54] A.-K. Sahu, M. Zaheer, and S. Kar. Towards gradient free and projection free stochastic  
360 optimization. In *22<sup>nd</sup> AISTATS*, pages 3468–3477, 2019.
- 361 [55] W. Shi, H. Gao, and B. Gu. Gradient-free method for heavily constrained nonconvex opti-  
362 mization. In *39<sup>th</sup> ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages  
363 19935–19955. PMLR, 2022.
- 364 [56] M.-K. Warmuth and S. V. N. Vishwanathan. Tutorial: Survey of boosting from an optimization  
365 perspective. In *26<sup>th</sup> ICML*, 2009.
- 366 [57] T. Werner and P. Ruckdeschel. The column measure and gradient-free gradient boosting, 2019.
- 367 [58] H. Zhang and B. Gu. Faster gradient-free methods for escaping saddle points. In *ICLR’23*.  
368 OpenReview.net, 2023.
- 369 [59] H. Zhang, H. Xiong, and B. Gu. Zeroth-order negative curvature finding: Escaping saddle  
370 points without gradients. In *NeurIPS\*35*, 2022.

# Supplementary Material

## Abstract

372 This is the Supplementary Material to Paper "How to Boost Any Loss Function"  
 373 submitted to NeurIPS'24.

374 This is the Appendix to paper "How to Boost Any Loss Function". To differentiate with the number-  
 375 ings in the main file, the numbering of Theorems, etc. is letter-based (A, B, ...).

376 **Table of contents**

377	<b>A quick summary of recent zeroth-order optimization approaches</b>	Pg 14
378		
379	<b>Supplementary material on proofs</b>	Pg 14
380		
381	→ Helper results	Pg 14
382	→ Removing the $\neq 0$ part in Assumption 5.1	Pg 15
383	→ Proof of Lemma 5.2	Pg 15
384	→ Proof of Theorem 5.3	Pg 16
385	→ Proof of Lemma 5.7	Pg 19
386	→ Proof of Theorem 5.8	Pg 19
387	→ Implementation of the offset oracle	Pg 20
388	→ Proof of Lemma E	Pg 21
389	→ Handling discontinuities in the offset oracle to prevent stopping in Step 2.5 of SECBOOST	Pg 23
390	→ A boosting pattern that can "survive" above differentiability	Pg 23
391	→ The case of piecewise constant losses for $SOLVE_\alpha$	Pg 25
392		
393	<b>Supplementary material on algorithms, implementation tricks and a toy experiment</b>	Pg 25
394		

reference	conv.	diff.	$F$	Lip.	smooth	Lb	$\nabla F$ diff.	main ML topic
[2]	✓	✓		✓	✓			online ML
[3]	✓			✓				distributed ML
[1]	✓			✓				online ML
[15]	✓	✓			✓		✓	alt. GD
[16]		✓			✓		✓	alt. GD
[19]		✓		✓				alt. GD
[18]				✓		✓		alt. GD
[20]	✓	✓		✓	✓			alt. GD
[21]	✓	✓			✓			alt. GD
[26]		✓		✓	✓			saddle pt opt
[27]		✓			✓			alt. FW
[29]		✓			✓			alt. FW
[23]	✓	✓						alt. GD
[28]				✓				online ML
[30]		✓			✓			deep ML
[36]	✓	✓			✓			alt. GD
[37]				✓				saddle pt opt
[39]	✓	✓		✓	✓			saddle pt opt
[41]		✓			✓		✓	distributed ML
[50]		✓			✓			alt. GD
[49]		✓		✓	✓			federated ML
[53]		✓		✓	✓			saddle pt opt
[54]		✓			✓			alt. FW
[55]		✓		✓	✓			alt. GD
[58]		✓		✓	✓			saddle pt opt
[59]		✓		✓	✓		✓	saddle pt opt

Table A1: Summary of formal assumptions about loss  $F$  used to prove algorithms' convergence in recent papers on zeroth order optimization, in different ML settings (see text for details). We use "smoothness" as a portmanteau for various conditions on the  $\geq 1$  order differentiability condition of  $F$ . "conv." = convex, "diff." = differentiable, "Lip." = Lipschitz, "Lb" = lower-bounded, "alt. GD" = general alternative to gradient descent (stochastic or not), "alt. FW" = idem for Frank-Wolfe. Our paper relies on no such assumptions.

## 395 I A quick summary of recent zeroth-order optimization approaches

396 Table A1 summarizes a few dozens of recent approaches that can be related to zeroth-order optimization  
397 in various topics of ML. Note that no such approaches focus on boosting.

## 398 II Supplementary material on proofs

### 399 II.1 Helper results

400 We now show that the order of the elements of  $\mathcal{V}$  does not matter to compute the  $\mathcal{V}$ -derivative as in  
401 Definition 4.2. For any  $\sigma \in \{0, 1\}^n$ , we let  $1_\sigma \doteq \sum_i \sigma_i$ .

**Lemma A.** For any  $z \in \mathbb{R}$ , any  $n \in \mathbb{N}_*$  and any  $\mathcal{V} \doteq \{v_1, v_2, \dots, v_n\} \subset \mathbb{R}$ ,

$$\delta_{\mathcal{V}} F(z) = \frac{\sum_{\sigma \in \{0, 1\}^n} (-1)^{n-1} \sigma F(z + \sum_{i=1}^n \sigma_i v_i)}{\prod_{i=1}^n v_i}. \quad (24)$$

402 Hence,  $\delta_{\mathcal{V}} F$  is invariant to permutations of the elements of  $\mathcal{V}$ .

*Proof.* We show the result by induction on the size of  $\mathcal{V}$ , first noting that

$$\delta_{\{v_1\}}F(z) = \delta_{v_1}F(z) \doteq \frac{F(z + v_1) - F(z)}{v_1} = \frac{1}{\prod_{i=1}^1 v_i} \cdot \sum_{\sigma \in \{0,1\}} (-1)^{1-1\sigma} F(z + \sigma v_1). \quad (25)$$

We then assume that (24) holds for  $\mathcal{V}_n \doteq \{v_1, v_2, \dots, v_n\}$  and show the result for  $\mathcal{V}_{n+1} \doteq \mathcal{V}_n \cup \{v_{n+1}\}$ , writing (induction hypothesis used in the second identity):

$$\begin{aligned} & \delta_{\mathcal{V}_{n+1}}F(z) \\ & \doteq \frac{\delta_{\mathcal{V}_n}F(z + v_{n+1}) - \delta_{\mathcal{V}_n}F(z)}{v_{n+1}} \\ & = \frac{\sum_{\sigma \in \{0,1\}^n} (-1)^{n-1\sigma} F(z + \sum_{i=1}^n \sigma_i v_i + v_{n+1}) - \sum_{\sigma \in \{0,1\}^n} (-1)^{n-1\sigma} F(z + \sum_{i=1}^n \sigma_i v_i)}{\prod_{i=1}^{n+1} v_i} \\ & = \frac{\sum_{\sigma \in \{0,1\}^n} (-1)^{n-1\sigma} F(z + \sum_{i=1}^n \sigma_i v_i + v_{n+1}) + \sum_{\sigma \in \{0,1\}^n} (-1)^{n-1\sigma+1} F(z + \sum_{i=1}^n \sigma_i v_i)}{\prod_{i=1}^{n+1} v_i} \\ & = \frac{\left\{ \begin{array}{l} \sum_{\sigma' \in \{0,1\}^{n+1}: \sigma'_{n+1}=1} (-1)^{n-(1\sigma'-1)} F(z + \sum_{i=1}^{n+1} \sigma'_i v_i) \\ + \sum_{\sigma' \in \{0,1\}^{n+1}: \sigma'_{n+1}=0} (-1)^{n+1-1\sigma'} F(z + \sum_{i=1}^{n+1} \sigma'_i v_i) \end{array} \right\}}{v^{n+1}} \\ & = \frac{\sum_{\sigma' \in \{0,1\}^{n+1}} (-1)^{n+1-1\sigma'} F(z + \sum_{i=1}^{n+1} \sigma'_i v_i)}{\prod_{i=1}^{n+1} v_i}, \end{aligned} \quad (26)$$

as claimed.  $\square$

We also have the following simple Lemma, which is a direct consequence of Lemma A.

**Lemma B.** For all  $z, v, z' \in \mathbb{R}_*$ , we have

$$\delta_v F(z + z') = \delta_v F(z) + z' \cdot \delta_{\{z', v\}} F(z). \quad (27)$$

*Proof.* It comes from Lemma A that  $\delta_{\{z', v\}} F(z) = \delta_{\{v, z'\}} F(z) = (\delta_v F(z + z') - \delta_v F(z))/z'$  (and we reorder terms).  $\square$

## II.2 Removing the $\neq 0$ part in Assumption 5.1

Because everything needs to be encoded, finiteness is not really an assumption. However, the non-zero assumption may be seen as limiting (unless we are happy to use first-order information about the loss (Section 5)). There is a simple trick to remove it. Suppose  $h_t$  zeroes on some training examples. The training sample being finite, there exists an open neighborhood  $\mathbb{I}$  in 0 such that  $h'_t \doteq h_t + \delta$  does not zero anymore on training examples, for any  $\delta \in \mathbb{I}$ . This changes the advantage  $\gamma$  in the WLA (Definition 5.5) to some  $\gamma'$  satisfying (we assume  $\delta > 0$  wlog)

$$\begin{aligned} \gamma' & \geq \frac{\gamma M_t}{M_t + \delta} - \frac{\delta}{M_t + \delta} \\ & \geq \gamma - \frac{\delta}{M_t} \cdot (1 + \gamma), \end{aligned}$$

from which it is enough to pick  $\delta \leq \varepsilon \gamma M_t / (1 + \gamma)$  to guarantee advantage  $\gamma' \geq (1 - \varepsilon)\gamma$ . If  $\varepsilon$  is a constant, this translates in a number of boosting iterations in Corollary 5.6 affected by a constant factor that we can choose as close to 1 as desired.

## II.3 Proof of Lemma 5.2

We reformulate

$$\delta_{\{b,c\}}F(a) = \frac{2}{b} \cdot \frac{1}{c} \cdot \left( \underbrace{\frac{F(a+b+c) + F(a)}{2} - \frac{F(a+b) + F(a+c)}{2}}_{\doteq \mu_2} \right). \quad (28)$$

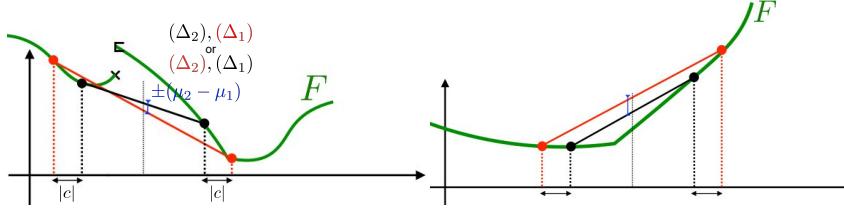


Figure II.1: *Left:* representation of the difference of averages in (28). Each of the secants  $(\Delta_1)$  and  $(\Delta_2)$  can take either the red or black segment. Which one is which depends on the signs of  $c$  and  $b$ , but the general configuration is always the same. Note that if  $F$  is convex, one necessarily sits above the other, which is the crux of the proof of Lemma 5.2. For the sake of illustration, suppose we can analytically have  $b, c \rightarrow 0$ . As  $c$  converges to 0 but  $b$  remains  $> 0$ ,  $\delta_{\{b,c\}}F(a)$  becomes proportional to the variation of the average secant midpoint; the then-convergence of  $b$  to 0 makes  $\delta_{\{b,c\}}F(a)$  converge to the second-order derivative of  $F$  at  $a$ . *Right:* in the special case where  $F$  is convex, one of the secants always sits above the other.

Both  $\mu_1$  and  $\mu_2$  are averages that can be computed from the midpoints of two secants (respectively):

$$\begin{aligned} (\Delta_1) &\doteq [(a+c, F(a+c)), (a+b, F(a+b))], \\ (\Delta_2) &\doteq [(a, F(a)), (a+b+c, F(a+b+c))]. \end{aligned}$$

412 Also, the midpoints of both secants have the same abscissa (and the ordinates are  $\mu_1$  and  $\mu_2$ ), so to  
413 study the sign of  $\delta_{\{b,c\}}F(a)$ , we can study the position of both secants with respect to each other.  $F$   
414 being convex, we show that the abscissae of one secant are included in the abscissae of the other, this  
415 being sufficient to give the position of both secants with respect to each other. We distinguish four  
416 cases.  
417

418 **Case 1:**  $c > 0, b > 0$ . We have  $a + b + c > \max\{a + b, a + c\}$  and  $a < \min\{a + b, a + c\}$ .  $F$  being  
419 convex,  $(\Delta_2)$  sits above  $(\Delta_1)$ . So,  $\mu_2 \geq \mu_1$  and finally  $\delta_{\{b,c\}}F(a) \geq 0$ .  
420

421 **Case 2:**  $c < 0, b < 0$ . We now have  $a + b + c < \min\{a + b, a + c\}$  while  $a > \max\{a + b, a + c\}$ ,  
422 so  $(\Delta_2)$  sits above  $(\Delta_1)$ . Again,  $\mu_2 \geq \mu_1$  and finally  $\delta_{\{b,c\}}F(a) \geq 0$ .  
423

424 **Case 3:**  $c > 0, b < 0$ . We have  $a + b < a$  and  $a + b < a + b + c$ . Also  $a + c > \max\{a + b + c, a\}$ ,  
425 so this time  $(\Delta_2)$  sits below  $(\Delta_1)$  but  $cb < 0$ , so  $\delta_{\{b,c\}}F(a) \geq 0$  again.  
426

427 **Case 4:**  $c < 0, b > 0$ . So  $a + c < a < a + b$  and  $a + c < a + b + c$ . So  $a + c < \min\{a, a + b + c\}$   
428 and  $a + b > \max\{a, a + c\}$ , so  $(\Delta_2)$  sits below  $(\Delta_1)$ . Since  $cb < 0$ , so  $\delta_{\{b,c\}}F(a) \geq 0$  again.  
429

#### 430 II.4 Proof of Theorem 5.3

Let us remind key simplified notations about edges,  $\forall t \geq 0$ :

$$\tilde{e}_{ti} = y_i \cdot H_t(\mathbf{x}_i), \quad (29)$$

$$e_{ti} = y_i \cdot \alpha_t h_t(\mathbf{x}_i) = \tilde{e}_{ti} - \tilde{e}_{(t-1)i}. \quad (30)$$

For short, we also let:

$$Q_{ti}^* \doteq Q_F^*(\tilde{e}_{ti}, \tilde{e}_{(t-1)i}, v_{i(t-1)}), \quad (31)$$

$$\Delta_{ti} \doteq \delta_{v_{i(t-1)}} F(\tilde{e}_{ti}) - \delta_{v_{i(t-1)}} F(\tilde{e}_{(t-1)i}), \quad (32)$$

where  $Q_{..}^*$  is defined in (8). We also split the computation of the leveraging coefficient  $\alpha_t$  in SECBOOST in two parts, the first computing a real  $a_t$  as:

$$a_t \in \frac{1}{2(1 + \varepsilon_t)M_t^2 \bar{W}_{2,t}} \cdot [1 - \pi_t, 1 + \pi_t], \quad (33)$$

and then using  $\alpha_t \leftarrow a_t \eta_t$ . We now use Lemma 4.7 (main file) and get

$$\mathbb{E}_{i \sim [m]} [S_{F|v_{ti}}(\tilde{e}_{ti} \| \tilde{e}_{(t+1)i})] \geq -\mathbb{E}_{i \sim D} [Q_{(t+1)i}^*], \forall t \geq 0. \quad (34)$$

If we reorganise (34) using the definition of  $S_{F| \cdot}(\cdot \| \cdot)$ , we get:

$$\begin{aligned} & \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{(t+1)i})] \\ & \leq \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - \mathbb{E}_{i \sim [m]} [(\tilde{e}_{ti} - \tilde{e}_{(t+1)i}) \cdot \delta_{v_{ti}} F(\tilde{e}_{(t+1)i})] + \mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*] \\ & = \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - \mathbb{E}_{i \sim [m]} [-e_{(t+1)i} \cdot \delta_{v_{ti}} F(\tilde{e}_{(t+1)i})] + \mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*] \end{aligned} \quad (35)$$

$$= \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] + \alpha_{t+1} \cdot \mathbb{E}_{i \sim [m]} [y_i h_{t+1}(\mathbf{x}_i) \cdot \delta_{v_{ti}} F(\tilde{e}_{(t+1)i})] + \mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*] \quad (36)$$

$$\begin{aligned} & = \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] + a_{t+1} \eta_{t+1} \cdot \mathbb{E}_{i \sim [m]} [y_i h_{t+1}(\mathbf{x}_i) \cdot \delta_{v_{ti}} F(\tilde{e}_{ti})] \\ & \quad + a_{t+1} \eta_{t+1} \cdot \mathbb{E}_{i \sim [m]} [y_i h_{t+1}(\mathbf{x}_i) \cdot \Delta_{(t+1)i}] + \mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*] \end{aligned} \quad (37)$$

$$\begin{aligned} & = \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - a_{t+1} \eta_{t+1} \cdot \underbrace{\mathbb{E}_{i \sim [m]} [w_{(t+1)i} y_i h_{t+1}(\mathbf{x}_i)]}_{=\eta_{t+1}} + a_{t+1} \eta_{t+1} \cdot \mathbb{E}_{i \sim [m]} [y_i h_{t+1}(\mathbf{x}_i) \cdot \Delta_{(t+1)i}] \\ & \quad + \mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*] \\ & = \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - a_{t+1} \eta_{t+1}^2 + a_{t+1} \eta_{t+1} \cdot \mathbb{E}_{i \sim [m]} [y_i h_{t+1}(\mathbf{x}_i) \cdot \Delta_{(t+1)i}] + \mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*]. \end{aligned} \quad (38)$$

431 (35) – (37) make use of definitions (30) (twice) and (32) as well as the decomposition of the leveraging  
432 coefficient in (33).

433 Looking at (38), we see that we can have a boosting-compliant decrease of the loss if the two  
434 quantities depending on  $\Delta_{(t+1)}$  and  $Q_{(t+1)}^*$  can be made small enough compared to  $a_{t+1} \eta_{t+1}^2$ . This  
435 is what we investigate.

436

**Bounding the term depending on  $\Delta_{(t+1)}$ .** – We use Lemma B with  $z \doteq \tilde{e}_{ti}$ ,  $z' \doteq e_{(t+1)i}$ ,  $v \doteq v_t$ , which yields (also using (30) and the assumption that  $h_{t+1}(\mathbf{x}_i) \neq 0$ ):

$$\begin{aligned} \Delta_{(t+1)i} & \doteq \delta_{v_{ti}} F(\tilde{e}_{(t+1)i}) - \delta_{v_{ti}} F(\tilde{e}_{ti}) \\ & = \delta_{v_{ti}} F(\tilde{e}_{ti} + e_{(t+1)i}) - \delta_{v_{ti}} F(\tilde{e}_{ti}) \\ & = e_{(t+1)i} \cdot \delta_{\{e_{(t+1)i}, v_{ti}\}} F(\tilde{e}_{ti}) \\ & = y_i \cdot \alpha_{t+1} h_{t+1}(\mathbf{x}_i) \cdot \delta_{\{e_{(t+1)i}, v_{ti}\}} F(\tilde{e}_{ti}), \end{aligned} \quad (39)$$

and so we get:

$$\begin{aligned} & a_{t+1} \eta_{t+1} \cdot \mathbb{E}_{i \sim [m]} [y_i h_{t+1}(\mathbf{x}_i) \cdot \Delta_{(t+1)i}] \\ & = a_{t+1} \eta_{t+1} \cdot \mathbb{E}_{i \sim [m]} [\alpha_{t+1} (y_i h_{t+1}(\mathbf{x}_i))^2 \cdot \delta_{\{e_{(t+1)i}, v_{ti}\}} F(\tilde{e}_{ti})] \\ & = a_{t+1}^2 \eta_{t+1}^2 \cdot \mathbb{E}_{i \sim [m]} [(h_{t+1}(\mathbf{x}_i))^2 \cdot \delta_{\{e_{(t+1)i}, v_{ti}\}} F(\tilde{e}_{ti})] \\ & \leq a_{t+1}^2 \eta_{t+1}^2 M_{t+1}^2 \cdot \overline{W}_{2,t+1}. \end{aligned} \quad (40)$$

**Bounding the term depending on  $Q_{(t+1)}^*$ .** – We immediately get from the value picked in argument of  $\mathbb{I}_{t+1}$  in step 2.5 of SECBOOST, the definition of  $\mathbb{I}_{ti}(\cdot)$  in (16) and our decomposition  $\alpha_t \leftarrow a_t \eta_t$  that  $Q_{(t+1)i}^* \leq \varepsilon_{t+1} \cdot a_{t+1}^2 \eta_{t+1}^2 M_{t+1}^2 \cdot \overline{W}_{2,t+1}$ ,  $\forall i \in [m]$ , so that:

$$\mathbb{E}_{i \sim [m]} [Q_{(t+1)i}^*] \leq \varepsilon_{t+1} \cdot a_{t+1}^2 \eta_{t+1}^2 M_{t+1}^2 \cdot \overline{W}_{2,t+1}. \quad (41)$$

**Finishing up with the proof** – Suppose that we choose  $\varepsilon_{t+1} > 0$ ,  $\pi_{t+1} \in (0, 1)$  and  $a_{t+1}$  as in (33). We then get from (38), (40), (41) that for any choice of  $v_{ti}$  in Step 2.5 of SECBOOST,

$$\begin{aligned} & \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{(t+1)i})] \\ & \leq \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - a_{t+1}\eta_{t+1}^2 + a_{t+1}^2\eta_{t+1}^2M_{t+1}^2 \cdot \bar{W}_{2,t+1} + \varepsilon_{t+1} \cdot a_{t+1}^2\eta_{t+1}^2M_{t+1}^2 \cdot \bar{W}_{2,t+1} \\ & = \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - a_{t+1}\eta_{t+1}^2 \cdot (1 - a_{t+1}(1 + \varepsilon_{t+1})M_{t+1}^2 \cdot \bar{W}_{2,t+1}) \\ & \leq \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{ti})] - \frac{\eta_{t+1}^2(1 - \pi_{t+1}^2)}{4(1 + \varepsilon_{t+1})M_{t+1}^2 \cdot \bar{W}_{2,t+1}}, \end{aligned} \quad (42)$$

where the last inequality is a consequence of (33). Suppose we pick  $H_0 \doteq h_0 \in \mathbb{R}$  a constant and  $v_0 > 0$  such that

$$\delta_{v_0} F(h_0) \neq 0. \quad (43)$$

The final classifier  $H_T$  of SECBOOST satisfies:

$$\mathbb{E}_{i \sim [m]} [F(y_i H_T(\mathbf{x}_i))] \leq F_0 - \frac{1}{4} \sum_{t=1}^T \frac{\eta_t^2(1 - \pi_t^2)}{(1 + \varepsilon_t)M_t^2 \bar{W}_{2,t}}, \quad (44)$$

with  $F_0 \doteq \mathbb{E}_{i \sim [m]} [F(\tilde{e}_{i0})] \doteq \mathbb{E}_{i \sim [m]} [F(y_i H_0)] = \mathbb{E}_{i \sim [m]} [F(y_i h_0)]$ . If we want  $\mathbb{E}_{i \sim [m]} [F(y_i H_T(\mathbf{x}_i))] \leq F(z^*)$ , assuming wlog  $F(z^*) \leq F_0$ , then it suffices to iterate until:

$$\sum_{t=1}^T \frac{1 - \pi_t^2}{\bar{W}_{2,t}(1 + \varepsilon_t)} \cdot \frac{\eta_t^2}{M_t^2} \geq 4(F_0 - F(z^*)). \quad (45)$$

Remind that the edge  $\eta_t$  is not normalised. We have defined a normalised edge,

$$[-1, 1] \ni \tilde{\eta}_t \doteq \sum_i \frac{|w_{ti}|}{W_t} \cdot \tilde{y}_{ti} \cdot \frac{h_t(\mathbf{x}_i)}{M_t}, \quad (46)$$

with  $\tilde{y}_{ti} \doteq y_i \cdot \text{sign}(w_{ti})$  and  $W_t \doteq \sum_i |w_{ti}| = \sum_i |\delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})|$ . We have the simple relationship between  $\eta_t$  and  $\tilde{\eta}_t$ :

$$\begin{aligned} \tilde{\eta}_t &= \sum_i \frac{|w_{ti}|}{W_t} \cdot (y_i \cdot \text{sign}(w_{ti})) \cdot \frac{h_t(\mathbf{x}_i)}{M_t} \\ &= \frac{1}{W_t M_t} \cdot \sum_i w_{ti} y_i h_t(\mathbf{x}_i) \\ &= \frac{m}{W_t M_t} \cdot \eta_t, \end{aligned} \quad (47)$$

resulting in ( $\forall t \geq 1$ ),

$$\begin{aligned} \frac{\eta_t^2}{M_t^2} &= \tilde{\eta}_t^2 \cdot \left( \frac{W_t}{m} \right)^2 \\ &= \tilde{\eta}_t^2 \cdot (\mathbb{E}_{i \sim [m]} [|\delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})|])^2 \\ &\geq \tilde{\eta}_t^2 \cdot (\mathbb{E}_{i \sim [m]} [|\delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})|])^2 \\ &= \tilde{\eta}_t^2 \cdot \bar{W}_{1,t}^2, \end{aligned} \quad (48)$$

recalling  $\bar{W}_{1,t} \doteq |\mathbb{E}_{i \sim D} [\delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})]|$ . It comes from (48) that a sufficient condition for (45) to hold is:

$$\sum_{t=1}^T \frac{\bar{W}_{1,t}^2(1 - \pi_t^2)}{\bar{W}_{2,t}(1 + \varepsilon_t)} \cdot \tilde{\eta}_t^2 \geq 4(F_0 - F(z^*)), \quad (49)$$

<sup>437</sup> which is the statement of Theorem 5.3.

We first observe that for any  $a \in \mathbb{R}, b, c \in \mathbb{R}_*$ ,

$$\begin{aligned} |\delta_{\{b,c\}} F(a)| &= \frac{1}{|bc|} \cdot \left| \begin{array}{c} F(a+b+c) - F(a+c) - bF'(a+c) \\ -(F(a+b) - F(a) - bF'(a)) \\ +b(F'(a+c) - F'(a)) \end{array} \right| \\ &\leq \frac{1}{|bc|} \cdot \left( \begin{array}{c} |F(a+b+c) - F(a+c) - bF'(a+c)| \\ +|(F(a+b) - F(a) - bF'(a))| \\ +|b(F'(a+c) - F'(a))| \end{array} \right) \\ &\leq \frac{1}{|bc|} \cdot \left( \frac{\beta}{2} \cdot b^2 + \frac{\beta}{2} \cdot b^2 + \beta|bc| \right) = \beta + \beta \cdot \frac{b^2}{|bc|}, \end{aligned} \quad (50)$$

where we used the  $\beta$ -smoothness of  $F$  and twice [13, Lemma 3.4]. We can also make a permutation in the expression of  $\delta_{\{b,c\}} F(a)$  and instead write

$$\begin{aligned} |\delta_{\{b,c\}} F(a)| &= \frac{1}{|bc|} \cdot \left| \begin{array}{c} F(a+b+c) - F(a+b) - cF'(a+b) \\ -(F(a+c) - F(a) - cF'(a)) \\ +c(F'(a+b) - F'(a)) \end{array} \right| \\ &\leq \frac{1}{|bc|} \cdot \left( \begin{array}{c} |F(a+b+c) - F(a+b) - cF'(a+b)| \\ +|(F(a+c) - F(a) - cF'(a))| \\ +|c(F'(a+b) - F'(a))| \end{array} \right) \\ &\leq \frac{1}{|bc|} \cdot \left( \frac{\beta}{2} \cdot c^2 + \frac{\beta}{2} \cdot c^2 + \beta|bc| \right) = \beta + \beta \cdot \frac{c^2}{|bc|}. \end{aligned} \quad (51)$$

We thus have

$$\begin{aligned} |\delta_{\{b,c\}} F(a)| &\leq \beta + \beta \cdot \left( \frac{\min\{|b|, |c|\}}{\sqrt{|bc|}} \right)^2 \\ &\leq 2\beta, \end{aligned} \quad (52)$$

by the power mean inequality [14, Chapter III, Theorem 2]. Since  $|h_t(\mathbf{x}_i)| \leq M_t$  by definition, we thus have

$$\left| \mathbb{E}_{i \sim [m]} \left[ \delta_{\{e_{ti}, v_{(t-1)i}\}} F(\tilde{e}_{(t-1)i}) \cdot \left( \frac{h_t(\mathbf{x}_i)}{M_t} \right)^2 \right] \right| \leq 2\beta, \quad (53)$$

439 which allows us to fix  $\bar{W}_{2,t} = 2\beta$  and completes the proof of Lemma 5.7.

440 **Remark C.** Our result is optimal in the sense that if we make one offset (say  $b$ ) go to zero, then  
441 the ratio in (52) goes to zero and we recover the condition on the  $v$ -derivative of the derivative,  
442  $|\delta_c F'(z)| \leq \beta$ .

We consider the upperbound::

$$\begin{aligned} \bar{W}_{2,t} &\doteq \left| \mathbb{E}_{i \sim [m]} \left[ \frac{h_t^2(\mathbf{x}_i)}{M_t^2} \cdot \delta_{\{e_{ti}, v_{(t-1)i}\}} F(\tilde{e}_{(t-1)i}) \right] \right| \\ &= \left| \mathbb{E}_{i \sim [m]} \left[ \frac{h_t^2(\mathbf{x}_i)}{M_t^2} \cdot \frac{1}{\tilde{e}_{ti}} \cdot \left( \frac{F(\tilde{e}_{ti} + v_{(t-1)i}) - F(\tilde{e}_{ti})}{v_{(t-1)i}} - \frac{F(\tilde{e}_{(t-1)i} + v_{(t-1)i}) - F(\tilde{e}_{(t-1)i})}{v_{(t-1)i}} \right) \right] \right| \\ &= \left| \frac{1}{\alpha_t} \cdot \mathbb{E}_{i \sim [m]} \left[ \frac{h_t(\mathbf{x}_i)}{y_i M_t^2} \cdot \left( \frac{F(\tilde{e}_{ti} + v_{(t-1)i}) - F(\tilde{e}_{ti})}{v_{(t-1)i}} - \frac{F(\tilde{e}_{(t-1)i} + v_{(t-1)i}) - F(\tilde{e}_{(t-1)i})}{v_{(t-1)i}} \right) \right] \right| \\ &= \left| \frac{1}{\alpha_t} \cdot \mathbb{E}_{i \sim [m]} \left[ \frac{y_i h_t(\mathbf{x}_i)}{M_t^2} \cdot (\delta_{v_{(t-1)i}} F(\tilde{e}_{ti}) - \delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i})) \right] \right| \end{aligned} \quad (54)$$

(The last identity uses the fact that  $y_i \in \{-1, 1\}$ ). Remark that we have extracted  $\alpha_t$  from the denominator but it is still present in the arguments  $\tilde{e}_{ti}$ . For any classifier  $h$ , we introduce notation

$$\eta(\mathbf{w}, h) \doteq \mathbb{E}_{i \sim [m]} [w_i y_i h(\mathbf{x}_i)],$$

and so  $\eta_t$  (Step 2.2 in SECBOOST) is also  $\eta(\mathbf{w}_t, h_t)$ , which is guaranteed to be non-zero by the Weak Learning Assumption (5.5). We want, for *some*  $\varepsilon_t > 0, \pi_t \in [0, 1]$ ,

$$\alpha_t \in \frac{\eta_t}{2(1 + \varepsilon_t)M_t^2\bar{W}_{2,t}} \cdot [1 - \pi_t, 1 + \pi_t]. \quad (55)$$

This says that the sign of  $\alpha_t$  is the same as the sign of  $\eta(\mathbf{w}_t, h_t) = \eta_t$ . Since we know its sign, let us look for its absolute value:

$$|\alpha_t| \in \frac{|\eta_t|}{2(1 + \varepsilon_t)M_t^2\bar{W}_{2,t}} \cdot [1 - \pi_t, 1 + \pi_t]. \quad (56)$$

From (15) (main file), we can in fact search  $\alpha_t$  in the union of all such intervals for  $\varepsilon_t > 0, \pi_t \in [0, 1]$ , which amounts to find first:

$$|\alpha_t| \in \left(0, \frac{|\eta_t|}{M_t^2\bar{W}_{2,t}}\right),$$

and then find any  $\varepsilon_t > 0, \pi_t \in [0, 1]$  such that (56) holds. Using (54) and simplifying the external dependency on  $\alpha_t$ , we then need

$$1 \in \left(0, \underbrace{\frac{|\eta_t|}{\left|\mathbb{E}_{i \sim [m]} [y_i h_t(\mathbf{x}_i) \cdot (\delta_{v_{(t-1)i}} F(\alpha_t y_i h_t(\mathbf{x}_i) + \tilde{e}_{(t-1)i}) - \delta_{v_{(t-1)i}} F(\tilde{e}_{(t-1)i}))]\right|}}_{\doteq B(\alpha_t)}\right), \quad (57)$$

under the constraint that the sign of  $\alpha_t$  be the same as that of  $\eta_t$ . But, using notation (23) (main file), we have

$$B(\alpha_t) = |\eta(\mathbf{w}_t, h_t) - \eta(\tilde{\mathbf{w}}_t(\alpha_t), h_t)|,$$

and so to get (57) satisfied, it is sufficient that

$$\frac{|\eta_t - \eta(\tilde{\mathbf{w}}_t(\alpha_t), h_t)|}{|\eta_t|} < 1, \quad (58)$$

which is Step 1 in  $\text{SOLVE}_\alpha$ . The Weak Learning Assumption (5.5) guarantees that the denominator is  $\neq 0$  so this can always be evaluated. The continuity of  $F$  in all  $\tilde{e}_{(t-1)i}$  guarantees  $\lim_{\alpha_t \rightarrow 0} \eta(\tilde{\mathbf{w}}_t(\alpha_t), h_t) = \eta_t$ , and thus guarantees the existence of solutions to (58) for some  $|\alpha_t| > 0$ .

To summarize, finding  $\alpha_t$  can be done in two steps, (i) solve

$$\frac{|\eta_t - \eta(\tilde{\mathbf{w}}_t(\text{sign}(\eta_t) \cdot a), h_t)|}{|\eta_t|} < 1$$

for some  $a > 0$  and (ii) let  $\alpha_t \doteq \text{sign}(\eta_t) \cdot a$ . This is the output of  $\text{SOLVE}_\alpha(\mathcal{S}, \mathbf{w}_t, h_t)$ , which ends the proof of Theorem 5.8.

## 449 II.7 Implementation of the offset oracle: particular cases

Consider the "spring loss" that we define, for  $[.]$  denoting the nearest integer, as:

$$F_{\text{SL}}(z) \doteq \log(1 + \exp(-z)) + 1 - \sqrt{1 - 4(z - [z])^2}. \quad (59)$$

Figure II.2 plots this loss, which composes the logistic loss with a "U"-shaped term. This loss would escape all optimization algorithms of Table A1 (Appendix), yet there is a trivial implementation of our offset oracle, as explained in Figure II.2:

- 453     1. if the interval  $\mathbb{I}$  defined by  $\tilde{e}_{(t-1)i}$  and  $\tilde{e}_{ti}$  contains at least one peak, compute the tangence  
 454       point  $(z_t)$  at the closest local "U" that passes through  $(\tilde{e}_{(t-1)i}, F(\tilde{e}_{(t-1)i}))$ ; then if  $z_t \in \mathbb{I}$   
 455       then  $v_{ti} \leftarrow z_t - \tilde{e}_{(t-1)i}$ , else  $v_{ti} \leftarrow \tilde{e}_{ti} - \tilde{e}_{(t-1)i}$ ;  
 456     2. otherwise  $F$  in  $\mathbb{I}$  is strictly convex and differentiable: a simple dichotomic search can retrieve  
 457       a feasible  $v_{ti}$  (see convex losses below);

458     Notice that one can alleviate the repetitive dichotomic search by pre-tabulating a feasible  $v$  for a set  
 459       of differences  $|a - b|$  ( $a, b$  belonging to the abscissae of the same "U") decreasing by a fixed factor,  
 460       choosing  $v_{ti} \leftarrow v$  of the largest tabulated  $|a - b|$  no larger than  $|\tilde{e}_{ti} - \tilde{e}_{(t-1)i}|$ .  
 461     **Discontinuities** discontinuities do not represent issues if the argument  $z$  of  $\mathbb{I}_{ti}(z)$  is large enough, as  
 462       shown from the following simple Lemma.

**Lemma D.** Define the discontinuity of  $F$  as:

$$\text{disc}(F) \doteq \max \left\{ \frac{\sup_z |F(z) - \lim_{z^-} F(z)|}{\sup_z |F(z) - \lim_{z^+} F(z)|} \right\}. \quad (60)$$

463     For any  $z \geq 0$ , if  $\text{disc}(F) \leq z$  then  $\mathbb{I}_{ti}(z) \neq \emptyset, \forall t \geq 1, \forall i \in [m]$ .

464     Figure 4 (c) shows a case where the discontinuity is larger than  $z$ . In this case, an issue eventually  
 465       happens for computing the next weight happens, only when the current edge is at the discontinuity.  
 466       We note that as iterations increase and the weak learner finds it eventually more difficult to return  
 467       weak hypotheses with  $\eta$  large enough, the discontinuities may become an issue for SECBOOST to  
 468       not stop at Step 2.5. Or one can always use a simple trick to avoid stopping and which relies on the  
 469       leveraging coefficient  $\alpha_t$ : this is described in the Appendix, Section II.9.

470     **The case of convex losses** If  $F$  is convex (not necessarily differentiable nor strictly convex), there is  
 471       a simple way to find a valid output for the offset oracle, which relies on the following Lemma.

**Lemma E.** Suppose  $F$  convex. Then for any  $z, z' \in \mathbb{R}, v \neq 0$ ,

$$\begin{aligned} & \{v > 0 : Q_F^*(z, z', v) = r\} \\ &= \left\{ v > 0 : D_F \left( z \left\| \frac{F(z+v) - F(z)}{v} \right\| \right) = r \right\}. \end{aligned} \quad (61)$$

472     (proof in Appendix, Section II.8) By definition,  $\mathbb{I}_{ti}(z') \subseteq \mathbb{I}_{ti}(z)$  for any  $z' \leq z$ , so a simple way to  
 473       implement the offset oracle's output  $\text{OO}(t, i, z)$  is, for some  $0 < r < z$ , to solve the Bregman identity  
 474       in the RHS of (61) and then return any relevant  $v$ . If  $F$  is strictly convex, there is just one choice.

475     If solving the Bregman identity is tedious but  $F$  is strictly convex, there a simple dichotomic search  
 476       that is guaranteed to find a feasible  $v$ . It exploits the fact that the abscissa maximizing the difference  
 477       between any secant of  $F$  and  $F$  has a simple closed form (see [22, Supplement, Figure 13]) and  
 478       so the OBI in (4) (Definition 4.6) has a closed form as well. In this case, it is enough, after taking  
 479       a first non-zero guess for  $v$  (either positive or negative), to divide it by a constant  $> 1$  until the  
 480       corresponding OBI is no larger than the  $z$  in the query  $\text{OO}(t, i, z)$ .

## 481     II.8 Proof of Lemma E

$F$  being convex, we first want to compute the set

$$\mathbb{I}_{z,v,r} \doteq \{v > 0 : Q_F(z, z+v, z+v) = r\}, \quad (62)$$

where  $r$  is supposed small enough for  $\mathbb{I}_{z,v,r}$  to be non-empty. There is a simple graphical solution to  
 this which, as Figure II.3 explains, consists in finding  $v$  solution of

$$\sup_t F(z+v) - \left( F(t) + \left( \frac{F(z+v) - F(z)}{v} \right) \cdot (z+v-t) \right) = r. \quad (63)$$

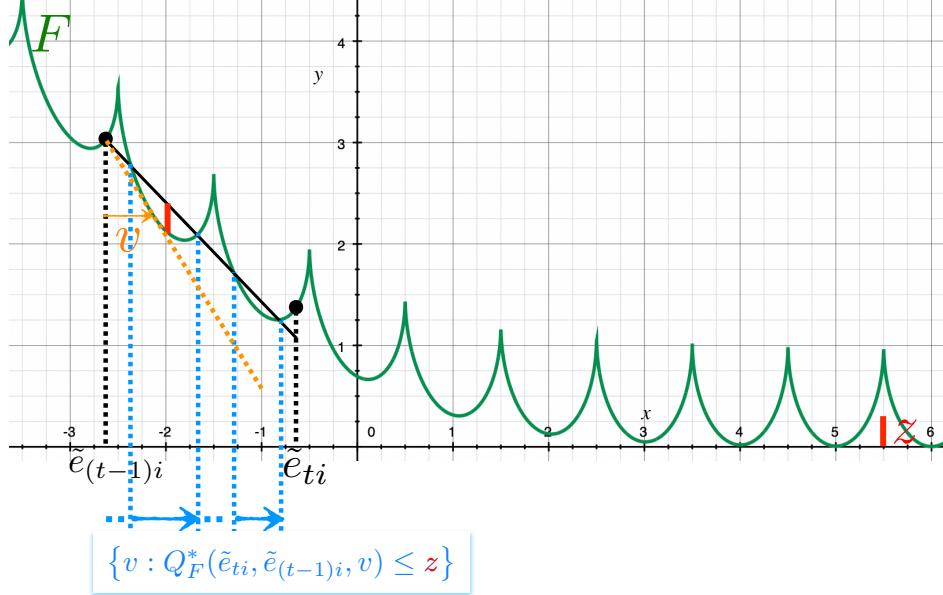


Figure II.2: The spring loss in (59) is neither convex, nor Lipschitz or differentiable and has an infinite number of local minima. Yet, an implementation of the offset oracle is trivial as an output for OO can be obtained from the computation of a single tangent point (here, the orange  $v$ , see text; best viewed in color).

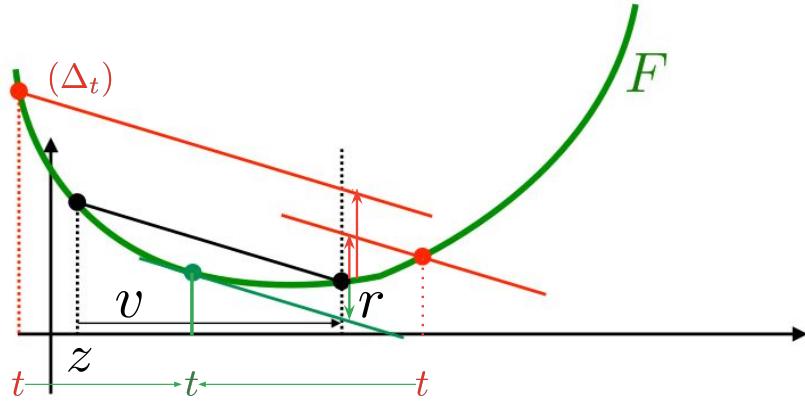


Figure II.3: Computing the OBI  $Q_F(z, z + v, z + v)$  for  $F$  convex,  $(z, v)$  being given and  $v > 0$ . We compute the line  $(\Delta_t)$  crossing  $F$  at any point  $t$ , with slope equal to the secant  $[(z, F(z)), (z + v, F(z + v))]$  and then the difference between  $F$  at  $z + v$  and this line at  $z + v$ . We move  $t$  so as to maximize this difference. The optimal  $t$  (in green) gives the corresponding OBI. In (64), we are interested in finding  $v$  given this difference,  $r$ . We also need to replicate this computation for  $v < 0$ .

The LHS simplifies:

$$\begin{aligned}
 & \sup_t F(z + v) - \left( F(t) + \left( \frac{F(z + v) - F(z)}{v} \right) \cdot (z + v - t) \right) \\
 &= \frac{(z + v)F(z) - zF(z + v)}{v} + \sup_t \left\{ t \cdot \frac{F(z + v) - F(z)}{v} - F(t) \right\} \\
 &= \frac{(z + v)F(z) - zF(z + v)}{v} + F^* \left( \frac{F(z + v) - F(z)}{v} \right) \\
 &= F(z) + F^* \left( \frac{F(z + v) - F(z)}{v} \right) - z \cdot \frac{F(z + v) - F(z)}{v} \\
 &= D_F \left( z \left| \frac{F(z + v) - F(z)}{v} \right. \right), 22
 \end{aligned}$$

so we end up with an equivalent but more readable definition for  $\mathbb{I}_{z,v,r}$ :

$$\mathbb{I}_{z,v,r} = \left\{ v > 0 : D_F \left( z \left\| \frac{F(z+v) - F(z)}{v} \right\| \right) = r \right\}, \quad (64)$$

482 which yields the statement of the Lemma.

483 **II.9 Handling discontinuities in the offset oracle to prevent stopping in Step 2.5**

484 Theorem 5.3 and Lemma 5.6 require to run SECBOOST for as many iterations are required. This  
 485 implies not early stopping in Step 2.5. Lemma D shows that early stopping can only be triggered by  
 486 too large local discontinuities at the edges. This is a weak requirement on running SECBOOST, but  
 487 there exists a weak assumption on the discontinuities of the loss itself that simply prevent any early  
 488 stopping and does not degrade the boosting rates. The result exploits the freedom in choosing  $\alpha_t$  in  
 489 Step 2.3.

490 **Lemma F.** *Suppose  $F$  is any function defined over  $\mathbb{R}$  discontinuities of zero Lebesgue measure. Then  
 491 Corollary 5.6 holds for boosting  $F$  with its inequality strict while never triggering early stopping in  
 492 Step 2.5 of SECBOOST.*

493 *Proof.* To show that we never trigger stopping in Step 2.5, it is sufficient to show that we can run  
 494 SECBOOST while ensuring  $F$  is continuous in an open neighborhood around all edges  $y_i H_t(\mathbf{x}_i), \forall i \in [m], \forall t \geq 0$  (by letting  $H_0 \doteq h_0$ ). Remind that  $\tilde{\epsilon}_{ti} \doteq \tilde{\epsilon}_{(t-1)i} + \alpha_t \cdot y_i h_t(\mathbf{x}_i)$ , so changing  $\alpha_t$  changes  
 495 all edges. We just have to show that either computing  $\alpha_t$  ensures such a continuity, or  $\alpha_t$  can be  
 496 slightly modified to do so. We have two ways to compute  $\alpha_t$ :

- 498 1. using a value for  $\bar{W}_{2,t}$  that represents an "absolute" upperbound in the sense of (14) (e.g.  
 499 Lemma 5.7) and then compute  $\alpha_t$  as in Step 2.3 of SECBOOST;
- 500 2. using algorithm  $\text{SOLVE}_\alpha$ .

501 Because of the assumption on  $F$ , we can always ensure that  $F$  is continuous in an open neighborhood  
 502 of all edges (the basis of the induction amounts to a straightforward choice for  $h_0$ ). This proves the  
 503 Lemma for [2.]

If we rely on [1.] and the  $\alpha_t$  computed leads to some discontinuities, then we have complete control  
 to change  $\alpha_t$ : any continuous change of  $\epsilon_t$  induces a continuous change in  $\alpha_t$  and thus a continuous  
 change of all edges as well. So, starting from the initial  $\epsilon_t$  chosen in Step 2.3, we increase it to a  
 value  $\epsilon_t^* > \epsilon_t$ , which we want to keep as small as possible. We can define for each  $i \in [m]$  an open  
 set  $(a_i, b_i)$  which is the interval spanned by the new  $\tilde{\epsilon}_{ti}(\epsilon'_t)$  using  $\epsilon'_t \in (\epsilon_t, \epsilon_t^*)$ . Since there are only  
 finitely many discontinuities on  $F$ , there exists a small  $\epsilon_t^* > \epsilon_t$  such that

$$\forall i \in [m], \forall z \in (a_i, b_i), F \text{ is continuous on } z.$$

504 This means that  $\forall \epsilon'_t \in (\epsilon_t, \epsilon_t^*)$ , we end up with a loss without any discontinuities on the new edges.  
 505 Now comes the reason why we want  $\epsilon_t^* - \epsilon_t$  small: we can check that there always exist a small  
 506 enough  $\epsilon_t^* > \epsilon_t$  such that for any  $\epsilon'_t$  we choose, the boosting rate in Corollary 5.6 is affected by at  
 507 most 1 additional iteration. Indeed, while we slightly change parameter  $\epsilon_t$  to land all new edges  
 508 outside of discontinuities of  $F$ , we also increase the contribution of the boosting iteration in the RHS  
 509 of (21) by a quantity  $\delta > 0$  which can be made as small as required — hence we can just replace the  
 510 inequality in (21) by a strict inequality. This proves the statement of the Lemma if we rely on [1.]  
 511 above.

512 This completes the proof of Lemma F. □

513 **II.10 A boosting pattern that can "survive" above differentiability**

514 Suppose  $F$  is strictly convex and strictly decreasing as for classical convex surrogates (e.g. logistic  
 515 loss). Assuming wlog all  $\alpha_t > 0$  and example  $i$  has both  $y_i h_t(\mathbf{x}_i) > 0$  and  $y_i h_{t-1}(\mathbf{x}_i) > 0$ , as long  
 516 as  $z$  is small enough, we are guaranteed that any choice  $v_{t-1} \in \mathbb{I}_{(t-1)i}(z)$  and  $v_t \in \mathbb{I}_{ti}(z)$  results in  
 517  $0 < w_{(t+1)i} < w_{ti}$ , which follows the classical boosting pattern that examples receiving the right  
 518 class by weak hypotheses have their weight decreased (See Figure II.4). If  $z = z'$  is large enough,  
 519 then this does not hold anymore as seen from Figure II.4.

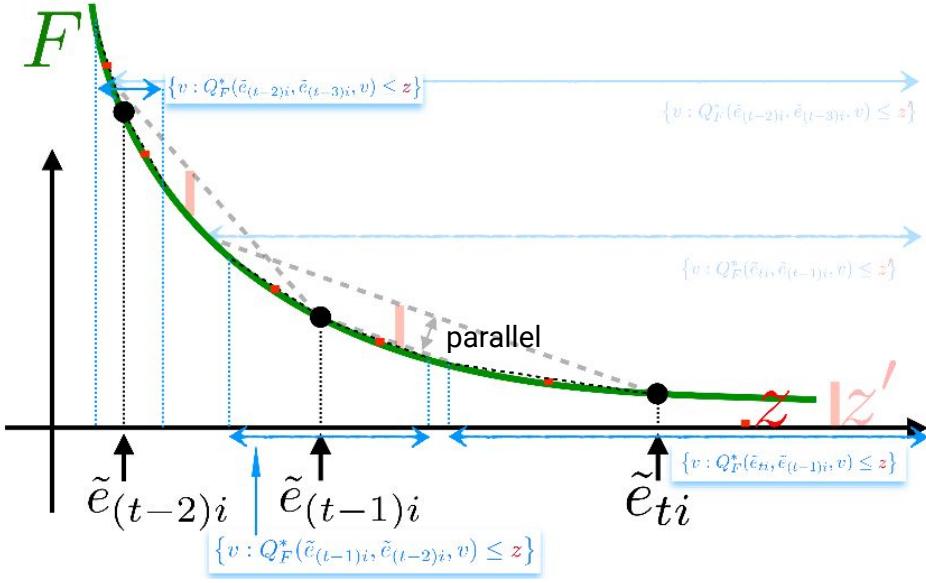


Figure II.4: Case F strictly convex, with two cases of limit OBI  $z$  and  $z'$  in  $\mathbb{I}_i(\cdot)$ . Example  $i$  has  $e_{ti} > 0$  and  $e_{(t-1)i} > 0$  (??) large enough (hence, edges with respect to weak classifiers  $h_t$  and  $h_{t-1}$  large enough) so that  $\mathbb{I}_{ti}(z) \cap \mathbb{I}_{(t-1)i}(z) = \mathbb{I}_{(t-1)i}(z) \cap \mathbb{I}_{(t-2)i}(z) = \mathbb{I}_{ti}(z) \cap \mathbb{I}_{(t-2)i}(z) = \emptyset$ . In this case, regardless of the offsets chosen by OO, we are guaranteed that its weights satisfy  $w_{(t+1)i} < w_{ti} < w_{(t-1)i}$ , which follows the boosting pattern that examples receiving the right classification by weak classifiers have their weights decreasing. If however the limit OBI changes from  $z$  to a larger  $z'$ , this is not guaranteed anymore: in this case, it may be the case that  $w_{(t+1)i} > w_{ti}$ .

520 **II.11 The case of piecewise constant losses for  $\text{SOLVE}_\alpha$**

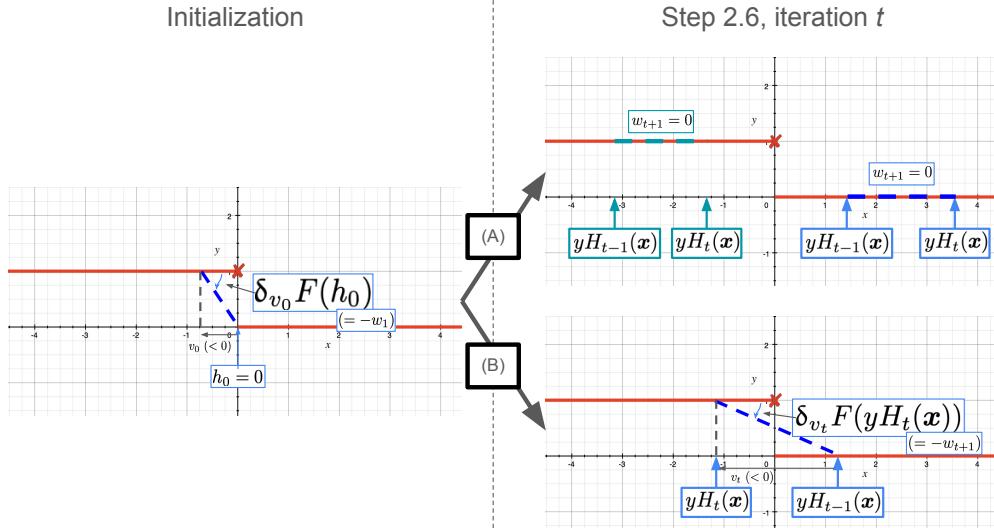


Figure II.5: How our algorithm works with the 0/1 loss (in red): at the initialization stage, assuming we pick  $h_0 = 0$  for simplicity and some  $v_0 < 0$ , all training examples get the same weight, given by negative the slope of the thick blue dashed line. All weights are thus  $> 0$ . At iteration  $t$  when we update the weights (Step 2.6), one of two cases can happen on some training example  $(x, y)$ . In (A), the edge of the strong model remains the same: either both are positive (blue) or both negative (olive green) (the ordering of edges is not important). In this case, regardless of the offset, the new weight will be 0. In (B), both edges have different sign (again, the ordering of edges is not important). In this case, the examples will keep non-zero weight over the next iteration. See text below for details.

521 Figure II.5 schematizes a run of our algorithm when training loss = 0/1 loss. At the initialization,  
 522 it is easy to get all examples to have non-zero weight. The weight update for example  $(x, y)$  of  
 523 our algorithm in Step 2.3 is (negative) the slope of a secant that crosses the loss in two points,  
 524 both being in between  $yH_{t-1}(x)$  and  $yH_t(x)$ . Hence, if the predicted label does not change  
 525 ( $\text{sign}(H_t(x)) = \text{sign}(H_{t-1}(x))$ ), then the next weight ( $w_{t+1}$ ) of the example *will be zero* (Figure  
 526 II.5, case (A)). However, if the predicted label does change ( $\text{sign}(H_t(x)) \neq \text{sign}(H_{t-1}(x))$ ) then  
 527 the example may get a non-zero weight depending on the offset chosen.

528 Hence, our generic implementation of Algorithms 3 and 4 may completely fail at providing non-zero  
 529 weights for the next iteration, which makes the algorithm stop in step 2.7. And even when not all  
 530 weights are zero, there may be just a too small subset of those, that would break the Weak Learning  
 531 Assumption for boosting compliance of the next iteration (Assumption 5.5).

532 **III Supplementary material on algorithms, implementation tricks and a toy  
 533 experiment**

534 **III.1 Algorithm and implementation of  $\text{SOLVE}_\alpha$  and how to find parameters from Theorem  
 535 5.8**

536 As Theorem 5.8 explains,  $\text{SOLVE}_\alpha$  can easily get to not just the leveraging coefficient  $\alpha_t$ , but also  
 537 other parameters that are necessary to implement SECBOOST:  $\bar{W}_{2,t}$  and  $\varepsilon_t$  (both used in Step 2.5).  
 538 We now provide a simple pseudo code on how to implement  $\text{SOLVE}_\alpha$  amnd get, on top of it, the two  
 539 other parameters. We do not seek  $\pi_t$  since it is useful only in the convergence analysis. Also, our  
 540 proposal implementation is optimized for complexity (because of the geometric updating of  $\delta, W$   
 541 in their respective loops) but much less so for accuracy. Algorithm  $\text{SOLVE\_extended}$  explains the  
 542 overall procedure.

---

**Algorithm 3** SOLVE\_extended( $\mathcal{S}, \mathbf{w}, h, M$ )

---

**Input** sample  $\mathcal{S} = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, m\}$ ,  $\mathbf{w} \in \mathbb{R}^m$ ,  $h : \mathcal{X} \rightarrow \mathbb{R}$ ,  $M \neq 0$ .  
 // in our case,  $\mathbf{w} \leftarrow \mathbf{w}_t$ ;  $h \leftarrow h_t$ ;  $M \leftarrow M_t$  (current weights, weak hypothesis and max confidence, see Step 2.3 in SECBOOST and Assumption 5.1)

Step 1 : // all initializations

$$\eta_{\text{init}} \leftarrow \eta(\mathbf{w}, h); \quad (65)$$

$$\delta \leftarrow 1.0; \quad (66)$$

$$W_{\text{init}} \leftarrow 1.0; \quad (67)$$

Step 2 : **do** // Step 2 computes the leveraging coefficient  $\alpha_t$

$$\alpha \leftarrow \delta \cdot \text{sign}(\eta_{\text{init}});$$

$$\eta_{\text{new}} \leftarrow \eta(\tilde{\mathbf{w}}(\alpha), h);$$

$$\text{if } |\eta_{\text{new}} - \eta_{\text{init}}| < |\eta_{\text{init}}| \text{ then } \text{found\_alpha} \leftarrow \text{true} \text{ else } \delta \leftarrow \delta/2;$$

**while**  $\text{found\_alpha} = \text{false}$ ;

Step 3 :  $W \leftarrow$  Left Hand Side of (14) (main file) // Step 3 computes  $\bar{W}_{2,t}$   
 // we can use (14) (main file) because we know  $\alpha$

**if**  $W =_{\text{machine}} 0$  **then**  
 // the LHS of (14) is (machine) 0: just need to find  $W$  such that (15) holds !

$$W \leftarrow W_{\text{init}};$$

$$\text{while } |\alpha| > |\eta_{\text{init}}|/(W \cdot M^2) \text{ do } W \leftarrow W/2;$$

**endif**

Step 4 :  $b_{\text{sup}} \leftarrow |\eta_{\text{init}}|/(W \cdot M^2);$  // Step 4 computes  $\varepsilon_t$   
 $\varepsilon \leftarrow (b_{\text{sup}}/\alpha) - 1;$

**Return**  $(\alpha, W, \varepsilon);$

---

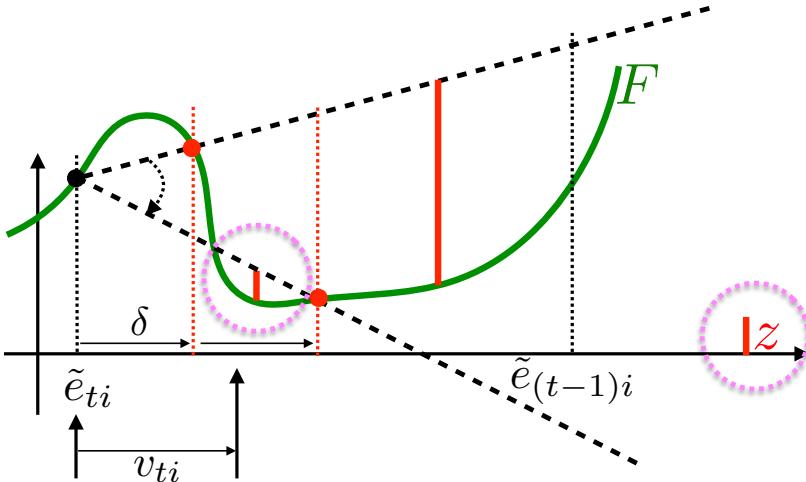


Figure III.1: How to find some  $v \in \mathbb{I}_{ti}(z)$ : parse the interval  $[\tilde{e}_{ti}, \tilde{e}_{(t-1)i}]$  with a regular step  $\delta$ , seek the secant with minimal slope (because  $\tilde{e}_{ti} < \tilde{e}_{(t-1)i}$ ; otherwise, we would seek the secant with maximal slope). It is necessarily the one minimizing the OBI among all regularly spaced choices. If the OBI is still too large, decrease the step  $\delta$  and start the search again.

543 **III.2 Algorithm and implementation of the offset oracle**

544 There exists a very simple trick to get some adequate offset  $v$  to satisfy (17) (main file), explained in  
 545 Figure III.1. In short, we seek the optimally bended secant and check that the OBI is no more than a  
 546 required  $z$ . This can be done via parsing the interval  $[\tilde{e}_{ti}, \tilde{e}_{(t-1)i}]$  using regularly spaced values. If  
 547 the OBI is too large, we can start again with a smaller step size. Algorithm OO\_simple details the key  
 548 part of the search.

---

**Algorithm 4** OO\_simple( $F, \tilde{e}_t, \tilde{e}_{t-1}, z, Z$ )

---

**Input** loss  $F$ , two last edges  $\tilde{e}_t, \tilde{e}_{t-1}$ , maximal OBI  $z$ , precision  $Z$ .  
 // in our case,  $\tilde{e}_t \leftarrow \tilde{e}_{ti}; \tilde{e}_{t-1} \leftarrow \tilde{e}_{(t-1)i}$ ; (for training example index  $i \in [m]$ )  
**Step 1 :** // all initializations

$$\delta \leftarrow \frac{\tilde{e}_{t-1} - \tilde{e}_t}{Z}; \quad (68)$$

$$z_c \leftarrow \tilde{e}_t + \delta; \quad (69)$$

$$i \leftarrow 0; \quad (70)$$

**Step 2 : do**

$$s_c \leftarrow \text{SLOPE}(F, \tilde{e}_t, z_c); \quad // \text{returns the slope of the secant passing through } (\tilde{e}_t, F(\tilde{e}_t)) \text{ and } (z_c, F(z_c))$$

$$\text{if } (i = 0) \vee ((\delta > 0) \wedge (s_c < s_*)) \vee ((\delta < 0) \wedge (s_c > s_*)) \text{ then } s_* \leftarrow s_c; z_* \leftarrow z_c$$

$$\text{endif}$$

$$z_c \leftarrow z_c + \delta;$$

$$i \leftarrow i + 1;$$

$$\text{while } (z_c - \tilde{e}_t) \cdot (z_c - \tilde{e}_{t-1}) < 0; \quad // \text{checks that } z_c \text{ is still in the interval}$$

**Return**  $z_* - \tilde{e}_t;$  // this is the offset  $v$

---

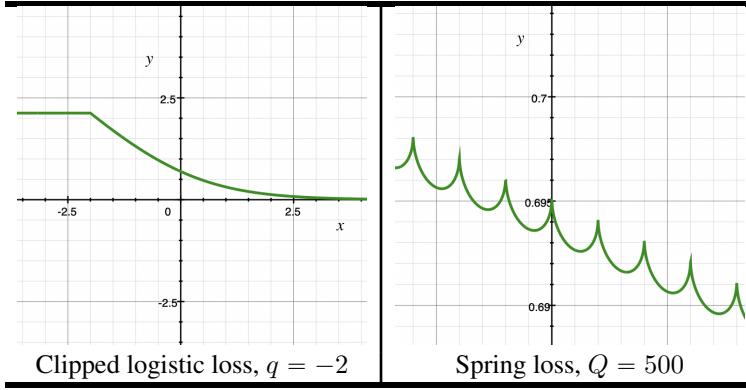


Figure III.2: Crops of the two losses whose optimization has been experimentally tested with SECBOOST, in addition to the logistic loss. See text for details.

549 **III.3 A toy experiments**

550 We provide here a few toy experiments using SECBOOST. These are just meant to display that a  
 551 simple implementation of the algorithm, following the blueprints given above, can indeed manage to  
 552 optimize various losses. These are not meant to explain how to pick the best hyperparameters (e.g.  
 553 (66)) nor how to choose the best loss given a domain, a problem that is far beyond the scope of our  
 554 paper.

555 In this implementation, the weak learner learns decision trees and we minimize Matusita's loss at  
 556 the leaves of decision trees to learn fixed size trees, see [34] for the criterion and induction scheme,  
 557 which is standard for decision trees. SECBOOST is implemented as is given in the paper, and so are  
 558 the implementation of  $\text{SOLVE}_\alpha$  and the offset oracle provided above. We have made no optimization  
 559 whatsoever, with one exception: when numerical approximation errors lead to an offset that is machine  
 560 0, we replace it by a small random value to prevent the use of derivatives in SECBOOST.

We have investigated three losses. The first is the well known logistic loss:

$$F_{\text{LOG}}(z) \doteq \log(1 + \exp(-z)). \quad (71)$$

The other two are tweaks of the logistic loss. We have investigated a clipped version of the logistic loss,

$$F_{\text{CL},q}(z) \doteq \min\{\log(1 + \exp(-z)), \log(1 + \exp(-q))\}, \quad (72)$$

with  $q \in \mathbb{R}$ , which clips the logistic loss above a certain value. This loss is non-convex and non-differentiable, but it is Lipschitz. We have also investigated a generalization of the spring loss (main file):

$$F_{\text{SL},Q}(z) \doteq \log(1 + \exp(-z)) + \frac{1 - \sqrt{1 - 4(z_Q - [z_Q])^2}}{Q}, \quad (73)$$

561 with  $z_Q \doteq Qz - 1/2$  ( $[.]$  is the closest integer), which adds to the logistic loss regularly spaced peaks  
562 of variable width. This loss is non-convex, non-differentiable, non-Lipschitz. Figure III.2 provides a  
563 crop of the clipped logistic loss and spring loss we have used in our test. Notice the “hardness” that  
564 the spring loss intuitively represents for ML.

565 We provide an experiment on public domain UCI `tictactoe` [24] (using a 10-fold stratified cross-  
566 validation to estimate test errors). In addition to the three losses, we have crossed them with several  
567 other variables: the size of the trees (either they have a single internal node = stumps or at most  
568 20 nodes) and, to give one example of how changing a (key) hyperparameter can change the result,  
569 we have tested for a scale of changes on the initial value of  $\delta$  in (66). Finally, we have crossed all  
570 these variables with the existence of symmetric label noise in the training data, following the setup  
571 of [38, 40]. We flip each label in the training sample with probability  $\eta$ . Table III.3 summarizes  
572 the results obtained. One can see that SECBOOST manages to optimize all losses in pretty much all  
573 settings, with an eventual early stopping required for the spring loss if  $\delta$  is too large. Note that the  
574 best initial value for  $\delta$  depends on the loss optimized in these experiments: for  $\delta = 0.1$ , test error  
575 from the spring loss decreases much faster than for the other losses, yet we remind that the spring  
576 loss is just the logistic loss plus regularly spaced peaks. This could signal interesting avenues for the  
577 best possible implementation of SECBOOST, or a further understanding of the best formal ways to fix  
578 those paramaters, all of which are out of the scope of this paper.

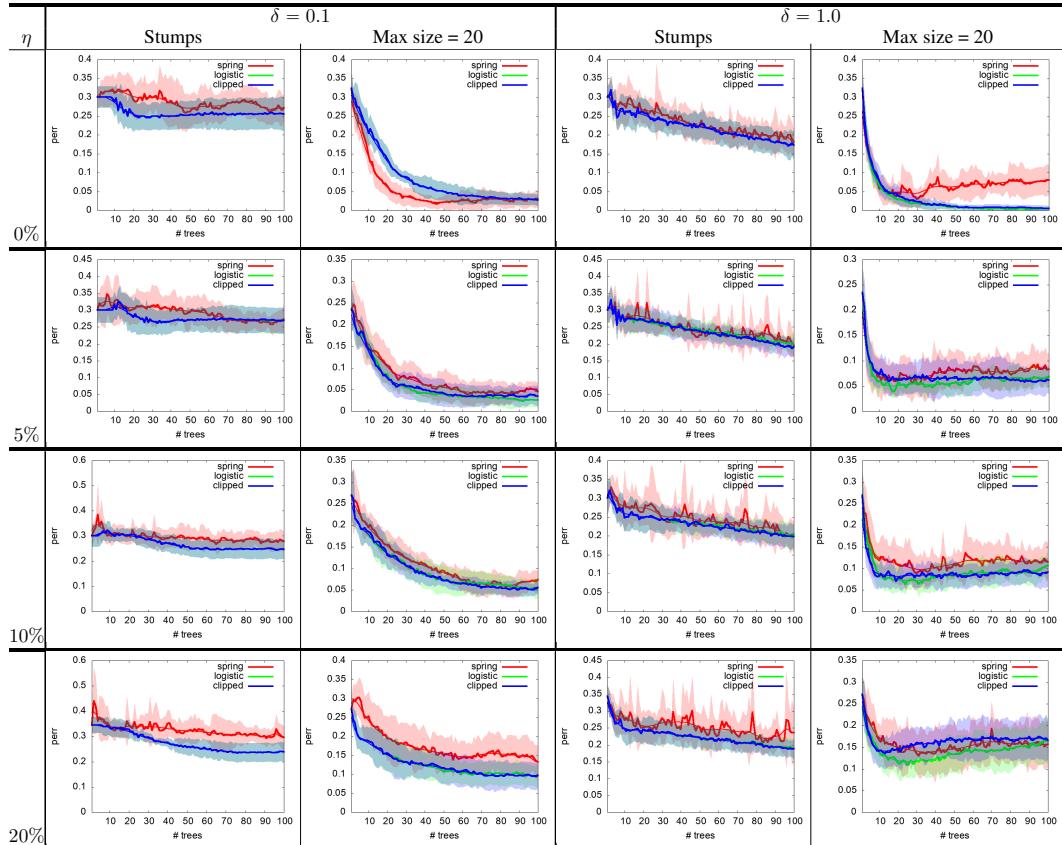


Figure III.3: Experiments on UCI tictactoe showing estimated test errors after minimizing each of the three losses we consider, with varying training noise level  $\eta$ , max tree size and initial hyperparameter  $\delta$  value in (66). See text.

579 **NeurIPS Paper Checklist**

580 **1. Claims**

581 Question: Do the main claims made in the abstract and introduction accurately reflect the  
582 paper's contributions and scope?

583 Answer: [Yes]

584 Justification: Our paper is a theory paper: all claims are properly formalized and used.

585 Guidelines:

- 586 • The answer NA means that the abstract and introduction do not include the claims  
587 made in the paper.
- 588 • The abstract and/or introduction should clearly state the claims made, including the  
589 contributions made in the paper and important assumptions and limitations. A No or  
590 NA answer to this question will not be perceived well by the reviewers.
- 591 • The claims made should match theoretical and experimental results, and reflect how  
592 much the results can be expected to generalize to other settings.
- 593 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
594 are not attained by the paper.

595 **2. Limitations**

596 Question: Does the paper discuss the limitations of the work performed by the authors?

597 Answer: [Yes]

598 Justification: The discussion section is devoted to limitations and improvement of our results

599 Guidelines:

- 600 • The answer NA means that the paper has no limitation while the answer No means that  
601 the paper has limitations, but those are not discussed in the paper.
- 602 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 603 • The paper should point out any strong assumptions and how robust the results are to  
604 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
605 model well-specification, asymptotic approximations only holding locally). The authors  
606 should reflect on how these assumptions might be violated in practice and what the  
607 implications would be.
- 608 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
609 only tested on a few datasets or with a few runs. In general, empirical results often  
610 depend on implicit assumptions, which should be articulated.
- 611 • The authors should reflect on the factors that influence the performance of the approach.  
612 For example, a facial recognition algorithm may perform poorly when image resolution  
613 is low or images are taken in low lighting. Or a speech-to-text system might not be  
614 used reliably to provide closed captions for online lectures because it fails to handle  
615 technical jargon.
- 616 • The authors should discuss the computational efficiency of the proposed algorithms  
617 and how they scale with dataset size.
- 618 • If applicable, the authors should discuss possible limitations of their approach to  
619 address problems of privacy and fairness.
- 620 • While the authors might fear that complete honesty about limitations might be used by  
621 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
622 limitations that aren't acknowledged in the paper. The authors should use their best  
623 judgment and recognize that individual actions in favor of transparency play an impor-  
624 tant role in developing norms that preserve the integrity of the community. Reviewers  
625 will be specifically instructed to not penalize honesty concerning limitations.

626 **3. Theory Assumptions and Proofs**

627 Question: For each theoretical result, does the paper provide the full set of assumptions and  
628 a complete (and correct) proof?

629 Answer: [Yes]

630 Justification: Our paper is a theory paper: all assumptions, statements and proofs provided.  
631

632 Guidelines:

- 633 • The answer NA means that the paper does not include theoretical results.
- 634 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 635 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 636 • The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- 637 • Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- 638 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 642 4. Experimental Result Reproducibility

643 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
644 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
645 of the paper (regardless of whether the code and data are provided or not)?

646 Answer: [Yes]

647 Justification: Though our paper is a theory paper, we have included in the supplement a  
648 detailed statement of all related algorithms and a toy experiment of a simple implementation  
649 of these algorithms showcasing a simple run on a public UCI domain.

650 Guidelines:

- 651 • The answer NA means that the paper does not include experiments.
- 652 • If the paper includes experiments, a No answer to this question will not be perceived  
653 well by the reviewers: Making the paper reproducible is important, regardless of  
654 whether the code and data are provided or not.
- 655 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
656 to make their results reproducible or verifiable.
- 657 • Depending on the contribution, reproducibility can be accomplished in various ways.  
658 For example, if the contribution is a novel architecture, describing the architecture fully  
659 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
660 be necessary to either make it possible for others to replicate the model with the same  
661 dataset, or provide access to the model. In general, releasing code and data is often  
662 one good way to accomplish this, but reproducibility can also be provided via detailed  
663 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
664 of a large language model), releasing of a model checkpoint, or other means that are  
665 appropriate to the research performed.
- 666 • While NeurIPS does not require releasing code, the conference does require all submis-  
667 sions to provide some reasonable avenue for reproducibility, which may depend on the  
668 nature of the contribution. For example
  - 669 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
670 to reproduce that algorithm.
  - 671 (b) If the contribution is primarily a new model architecture, the paper should describe  
672 the architecture clearly and fully.
  - 673 (c) If the contribution is a new model (e.g., a large language model), then there should  
674 either be a way to access this model for reproducing the results or a way to reproduce  
675 the model (e.g., with an open-source dataset or instructions for how to construct  
676 the dataset).
  - 677 (d) We recognize that reproducibility may be tricky in some cases, in which case  
678 authors are welcome to describe the particular way they provide for reproducibility.  
679 In the case of closed-source models, it may be that access to the model is limited in  
680 some way (e.g., to registered users), but it should be possible for other researchers  
681 to have some path to reproducing or verifying the results.

#### 682 5. Open access to data and code

683 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
684 tions to faithfully reproduce the main experimental results, as described in supplemental  
685 material?

686 Answer: [No]

687 Justification: Our paper is a theory paper. All algorithms we introduce are either in the main  
688 file or the appendix.

689 Guidelines:

- 690 • The answer NA means that paper does not include experiments requiring code.
- 691 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 692 • While we encourage the release of code and data, we understand that this might not be  
693 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
694 including code, unless this is central to the contribution (e.g., for a new open-source  
695 benchmark).
- 696 • The instructions should contain the exact command and environment needed to run to  
697 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 698 • The authors should provide instructions on data access and preparation, including how  
699 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 700 • The authors should provide scripts to reproduce all experimental results for the new  
701 proposed method and baselines. If only a subset of experiments are reproducible, they  
702 should state which ones are omitted from the script and why.
- 703 • At submission time, to preserve anonymity, the authors should release anonymized  
704 versions (if applicable).
- 705 • Providing as much information as possible in supplemental material (appended to the  
706 paper) is recommended, but including URLs to data and code is permitted.

## 707 6. Experimental Setting/Details

710 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
711 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
712 results?

713 Answer: [NA]

714 Justification: Our paper is a theory paper.

715 Guidelines:

- 716 • The answer NA means that the paper does not include experiments.
- 717 • The experimental setting should be presented in the core of the paper to a level of detail  
718 that is necessary to appreciate the results and make sense of them.
- 719 • The full details can be provided either with the code, in appendix, or as supplemental  
720 material.

## 721 7. Experiment Statistical Significance

722 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
723 information about the statistical significance of the experiments?

724 Answer: [NA]

725 Justification: Our paper is a theory paper.

726 Guidelines:

- 727 • The answer NA means that the paper does not include experiments.
- 728 • The authors should answer “Yes” if the results are accompanied by error bars, confi-  
729 dence intervals, or statistical significance tests, at least for the experiments that support  
730 the main claims of the paper.
- 731 • The factors of variability that the error bars are capturing should be clearly stated (for  
732 example, train/test split, initialization, random drawing of some parameter, or overall  
733 run with given experimental conditions).

- 734           • The method for calculating the error bars should be explained (closed form formula,  
 735           call to a library function, bootstrap, etc.)  
 736           • The assumptions made should be given (e.g., Normally distributed errors).  
 737           • It should be clear whether the error bar is the standard deviation or the standard error  
 738           of the mean.  
 739           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
 740           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
 741           of Normality of errors is not verified.  
 742           • For asymmetric distributions, the authors should be careful not to show in tables or  
 743           figures symmetric error bars that would yield results that are out of range (e.g. negative  
 744           error rates).  
 745           • If error bars are reported in tables or plots, The authors should explain in the text how  
 746           they were calculated and reference the corresponding figures or tables in the text.

747           **8. Experiments Compute Resources**

748           Question: For each experiment, does the paper provide sufficient information on the com-  
 749           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
 750           the experiments?

751           Answer: [NA].

752           Justification: Our paper is a theory paper.

753           Guidelines:

- 754           • The answer NA means that the paper does not include experiments.
- 755           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
 756           or cloud provider, including relevant memory and storage.
- 757           • The paper should provide the amount of compute required for each of the individual  
 758           experimental runs as well as estimate the total compute.
- 759           • The paper should disclose whether the full research project required more compute  
 760           than the experiments reported in the paper (e.g., preliminary or failed experiments that  
 761           didn't make it into the paper).

762           **9. Code Of Ethics**

763           Question: Does the research conducted in the paper conform, in every respect, with the  
 764           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

765           Answer: [Yes]

766           Justification: The research of the paper follows the code of ethics.

767           Guidelines:

- 768           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 769           • If the authors answer No, they should explain the special circumstances that require a  
 770           deviation from the Code of Ethics.
- 771           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
 772           eration due to laws or regulations in their jurisdiction).

773           **10. Broader Impacts**

774           Question: Does the paper discuss both potential positive societal impacts and negative  
 775           societal impacts of the work performed?

776           Answer: [NA]

777           Justification: Our paper is a theory paper.

778           Guidelines:

- 779           • The answer NA means that there is no societal impact of the work performed.
- 780           • If the authors answer NA or No, they should explain why their work has no societal  
 781           impact or why the paper does not address societal impact.
- 782           • Examples of negative societal impacts include potential malicious or unintended uses  
 783           (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
 784           (e.g., deployment of technologies that could make decisions that unfairly impact specific  
 785           groups), privacy considerations, and security considerations.

- 786           • The conference expects that many papers will be foundational research and not tied  
 787           to particular applications, let alone deployments. However, if there is a direct path to  
 788           any negative applications, the authors should point it out. For example, it is legitimate  
 789           to point out that an improvement in the quality of generative models could be used to  
 790           generate deepfakes for disinformation. On the other hand, it is not needed to point out  
 791           that a generic algorithm for optimizing neural networks could enable people to train  
 792           models that generate Deepfakes faster.
- 793           • The authors should consider possible harms that could arise when the technology is  
 794           being used as intended and functioning correctly, harms that could arise when the  
 795           technology is being used as intended but gives incorrect results, and harms following  
 796           from (intentional or unintentional) misuse of the technology.
- 797           • If there are negative societal impacts, the authors could also discuss possible mitigation  
 798           strategies (e.g., gated release of models, providing defenses in addition to attacks,  
 799           mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
 800           feedback over time, improving the efficiency and accessibility of ML).

## 801          11. Safeguards

802          Question: Does the paper describe safeguards that have been put in place for responsible  
 803          release of data or models that have a high risk for misuse (e.g., pretrained language models,  
 804          image generators, or scraped datasets)?

805          Answer: [NA]

806          Justification: No release of data or models.

807          Guidelines:

- 808           • The answer NA means that the paper poses no such risks.
- 809           • Released models that have a high risk for misuse or dual-use should be released with  
 810           necessary safeguards to allow for controlled use of the model, for example by requiring  
 811           that users adhere to usage guidelines or restrictions to access the model or implementing  
 812           safety filters.
- 813           • Datasets that have been scraped from the Internet could pose safety risks. The authors  
 814           should describe how they avoided releasing unsafe images.
- 815           • We recognize that providing effective safeguards is challenging, and many papers do  
 816           not require this, but we encourage authors to take this into account and make a best  
 817           faith effort.

## 818          12. Licenses for existing assets

819          Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
 820          the paper, properly credited and are the license and terms of use explicitly mentioned and  
 821          properly respected?

822          Answer: [NA]

823          Justification: no outside code, data or models used requiring licensing.

824          Guidelines:

- 825           • The answer NA means that the paper does not use existing assets.
- 826           • The authors should cite the original paper that produced the code package or dataset.
- 827           • The authors should state which version of the asset is used and, if possible, include a  
 828           URL.
- 829           • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 830           • For scraped data from a particular source (e.g., website), the copyright and terms of  
 831           service of that source should be provided.
- 832           • If assets are released, the license, copyright information, and terms of use in the  
 833           package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
 834           has curated licenses for some datasets. Their licensing guide can help determine the  
 835           license of a dataset.
- 836           • For existing datasets that are re-packaged, both the original license and the license of  
 837           the derived asset (if it has changed) should be provided.

- 838           • If this information is not available online, the authors are encouraged to reach out to  
839           the asset's creators.

840           **13. New Assets**

841           Question: Are new assets introduced in the paper well documented and is the documentation  
842           provided alongside the assets?

843           Answer: [NA]

844           Justification: No new assets provided.

845           Guidelines:

- 846           • The answer NA means that the paper does not release new assets.  
847           • Researchers should communicate the details of the dataset/code/model as part of their  
848           submissions via structured templates. This includes details about training, license,  
849           limitations, etc.  
850           • The paper should discuss whether and how consent was obtained from people whose  
851           asset is used.  
852           • At submission time, remember to anonymize your assets (if applicable). You can either  
853           create an anonymized URL or include an anonymized zip file.

854           **14. Crowdsourcing and Research with Human Subjects**

855           Question: For crowdsourcing experiments and research with human subjects, does the paper  
856           include the full text of instructions given to participants and screenshots, if applicable, as  
857           well as details about compensation (if any)?

858           Answer: [NA]

859           Justification: No crowdsourcing or research with human subjects.

860           Guidelines:

- 861           • The answer NA means that the paper does not involve crowdsourcing nor research with  
862           human subjects.  
863           • Including this information in the supplemental material is fine, but if the main contribu-  
864           tion of the paper involves human subjects, then as much detail as possible should be  
865           included in the main paper.  
866           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
867           or other labor should be paid at least the minimum wage in the country of the data  
868           collector.

869           **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
870           Subjects**

871           Question: Does the paper describe potential risks incurred by study participants, whether  
872           such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
873           approvals (or an equivalent approval/review based on the requirements of your country or  
874           institution) were obtained?

875           Answer: [NA]

876           Justification: No research with human subjects.

877           Guidelines:

- 878           • The answer NA means that the paper does not involve crowdsourcing nor research with  
879           human subjects.  
880           • Depending on the country in which research is conducted, IRB approval (or equivalent)  
881           may be required for any human subjects research. If you obtained IRB approval, you  
882           should clearly state this in the paper.  
883           • We recognize that the procedures for this may vary significantly between institutions  
884           and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
885           guidelines for their institution.  
886           • For initial submissions, do not include any information that would break anonymity (if  
887           applicable), such as the institution conducting the review.