

On Regularizing Rademacher Observation Losses

Authored by:

Richard Nock

Abstract

It has recently been shown that supervised learning linear classifiers with two of the most popular losses, the logistic and square loss, is equivalent to optimizing an equivalent loss over sufficient statistics about the class: Rademacher observations (rados). It has also been shown that learning over rados brings solutions to two prominent problems for which the state of the art of learning from examples can be comparatively inferior and in fact less convenient: protecting and learning from private examples, learning from distributed datasets without entity resolution. *Bis repetita placent*: the two proofs of equivalence are different and rely on specific properties of the corresponding losses, so whether these can be unified and generalized inevitably comes to mind. This is our first contribution: we show how they can be fit into the same theory for the equivalence between example and rado losses. As a second contribution, we show that the generalization unveils a surprising new connection to regularized learning, and in particular a sufficient condition under which regularizing the loss over examples is equivalent to regularizing the rados (i.e. the data) in the equivalent rado loss, in such a way that an efficient algorithm for one regularized rado loss may be as efficient when changing the regularizer. This is our third contribution: we give a formal boosting algorithm for the regularized exponential rado-loss which boost with any of the ridge, lasso, slope, L_{inf}, or elastic nets, using the same master routine for all. Because the regularized exponential rado-loss is the equivalent of the regularized logistic loss over examples we obtain the first efficient proxy to the minimisation of the regularized logistic loss over examples using such a wide spectrum of regularizers. Experiments with a readily available code display that regularization significantly improves rado-based learning and compares favourably with example-based learning.

1 Paper Body

What kind of data should we use to train a supervised learner ? A recent result has shown that minimising the popular logistic loss over examples with linear classifiers (in supervised learning) is equivalent to the minimisation of the exponential loss over sufficient statistics about the class known as Rademacher

observations (rados, [Nock et al., 2015]), for the same classifier. In short, we fit a classifier over data that is different from examples, and the same classifier generalizes well to new observations. It has been shown that rados offer solutions for two problems for which the state of the art involving examples can be comparatively significantly inferior: ? protection of the examples? privacy from various algebraic, geometric, statistical and computational standpoints, and learning from private data [Nock et al., 2015]; ? learning from a large number of distributed datasets without having to perform entity resolution between datasets [Patrini et al., 2016]. Quite remarkably, the training time of the algorithms involved can be smaller than it would be on examples, by orders of magnitude [Patrini et al., 2016]. Two key problems remain however: the 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

accuracy of learning from rados can compete experimentally with that of learning from examples, yet there is a gap to reduce for rados to be not just a good material to learn from in a privacy/distributed setting, but also a serious alternative to learning from examples at large, yielding new avenues to supervised learning. Second, theoretically speaking, it is now known that two widely popular losses over examples admit an equivalent loss in the rado world: the logistic loss and the square loss [Nock et al., 2015, Patrini et al., 2016]. This inevitably suggests that this property may hold for more losses, yet barely anything displays patterns of generalizability in the existing proofs. Our contributions: in this paper, we provide answers to these two questions, with three main contributions. Our first contribution is to show that this generalization indeed holds: other example losses admit equivalent losses in the rado world, meaning in particular that their minimiser classifier is the same, regardless of the dataset of examples. The technique we use exploits a two-player zero sum game representation of convex losses, that has been very useful to analyse boosting algorithms [Schapire, 2003, Telgarsky, 2012], with one key difference: payoffs are non-linear convex, eventually non-differentiable. These also resemble the entropic dual losses [Reid et al., 2015], with the difference that we do not enforce conjugacy over the simplex. The conditions of the game are slightly different for examples and rados. We provide necessary and sufficient conditions for the resulting losses over examples and rados to be equivalent. Informally, equivalence happens iff the convex functions of the games satisfy a symmetry relationship and the weights satisfy a linear system of equations. Some popular losses fit in the equivalence [Nair and Hinton, 2010, Gentile and Warmuth, 1998, Nock and Nielsen, 2008, Telgarsky, 2012, Vapnik, 1998, van Rooyen et al., 2015]. Our second contribution came unexpectedly through this equivalence. Regularizing a loss is standard in machine learning [Bach et al., 2011]. We show a sufficient condition for the equivalence under which regularizing the example loss is equivalent to regularizing the rados in the equivalent rado loss, i.e. making a Minkowski sum of the rado set with a classifier-based set. This property is independent of the regularizer, and incidentally happens to hold for all our cases of equivalence (Cf first contribution). A regularizer added to a loss over examples thus transfers to data in the rado world, in essentially the same way for all reg-

ularizers, and if one can solve the non-trivial computational and optimization problem that poses this data modification for one regularized rado loss, then, basically, "A good optimization algorithm for this regularized rado loss may fit to other regularizers as well? Our third contribution exemplifies this. We propose an iterative boosting algorithm, λ -R.A DA B OOST, that learns a classifier from rados using the exponential regularized rado loss, with regularization choice belonging to the ridge, lasso, ℓ_1 , or the recently coined SLOPE [Bogdan et al., 2015]. Since rado regularization would theoretically require to modify data at each iteration, such schemes are computationally non-trivial. We show that this modification can in fact be bypassed for the exponential rado loss, and the algorithm, λ -R.A DA B OOST, is as fast as A DA B OOST. λ -R.A DA B OOST has however a key advantage over A DA B OOST that to our knowledge is new in the boosting world: for any of these four regularizers, λ -R.A DA B OOST is a boosting algorithm λ thus, because of the equivalence between the minimization of the logistic loss over examples and the minimization of the exponential rado loss, λ -R.A DA B OOST is in fact an efficient proxy to boost the regularized logistic loss over examples using whichever of the four regularizers, and by extension, linear combination of them (e.g., elastic net regularization [Zou and Hastie, 2005]). We are not aware of any regularized logistic loss formal boosting algorithm with such a wide spectrum of regularizers. Extensive experiments validate this property: λ -R.A DA B OOST is all the better vs A DA B OOST (unregularized or regularized) as the domain gets larger, and is able to rapidly learn both accurate and sparse classifiers, making it an especially good contender for supervised learning at large on big domains. The rest of this paper is as follows. Sections 2, 3 and 4 respectively present the equivalence between example and rado losses, its extension to regularized learning and λ -R.A DA B OOST. 5 and 6 respectively present experiments, and conclude. In order not to laden the paper's body, a Supplementary Material (SM) contains the proofs and additional theoretical and experimental results.

2

Games and equivalent example/rado losses

To avoid notational load, we briefly present our learning setting to point the key quantity in our . . formulation of the general two players game. Let $[m] = \{1, 2, \dots, m\}$ and $\mathbb{R}_m = \{\mathbb{R}^1, \mathbb{R}^m\}$, for $m \geq 0$. The classical (batch) supervised learner is example-based: it is given a set of examples $S = \{(x_i, y_i), i \in [m]\}$ where $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^1$, $i \in [m]$. It returns a classifier $h : \mathbb{R}^d \rightarrow \mathbb{R}$ from 2

a predefined set H . Let $z_i(h) = y_i h(x_i)$ and abbreviate $z(h)$ by z for short. The learner fits h to the minimization of a loss. Table 1, column 'e', presents some losses that can be used: we remark that h appears only through z , so let us consider in this section that the learner rather fits vector $z \in \mathbb{R}^m$. We can now define our two players game setting. Let $\phi_e : \mathbb{R} \rightarrow \mathbb{R}$ and $\phi_r : \mathbb{R} \rightarrow \mathbb{R}$ two convex and m lower-semicontinuous generators. We define functions $L_e : \mathbb{R}^m \rightarrow \mathbb{R}$ and $L_r : \mathbb{R}^2 \rightarrow \mathbb{R}^m \rightarrow \mathbb{R} : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$. $L_e(p, z) = \sum_{i \in [m]} p_i z_i + \phi_e(\sum_{i \in [m]} p_i)$, $(1) i \in [m]$

$$L_r(q, z) =$$

$$\begin{aligned}
& i^?[m] \\
& X \\
& qI \\
& I^?[m] \\
& X \\
& z_i + ?r \\
& X \\
& ?r(qI) , \\
& (2) \\
& I^?[m] \\
& i^?I
\end{aligned}$$

where $?e, ?r \in [0, 1]$ do not depend on z . For the notation to be meaningful, the coordinates in q are assumed (wlog) to be in bijection with $2[m]$. The dependence of both problems in their respective generators is implicit and shall be clear from context. The adversary's goal is to fit

$$\begin{aligned}
& p^?(z) = . \\
& . \\
& (3) \\
& \arg \min_m L_r(q, z) , \\
& (4) \\
& p^?R \\
& q^?(z) = m \\
& \arg \min_m L_e(p, z) , q^?H2 \\
& m \\
& \text{with } H2 = \{q^?R2 : 1 \leq q \leq m\}, \text{ so as to attain } . \\
& L_e(z) = L_e(p^?(z), z) , . L_r(z) = L_r(q^?(z), z) , \\
& (5) (6)
\end{aligned}$$

and let $?Le(z)$ and $?Lr(z)$ denote their subdifferentials. We view the learner's task as the problem of maximising the corresponding problems in eq. (5) (with examples; this is already sketched above) or (6) (with what we shall call Rademacher observations, or rados), or equivalently minimising negative the corresponding function, and then resort to a loss function. The question of when these two problems are equivalent from the learner's standpoint motivates the following definition. Definition 1 Two generators $?e, ?r$ are said proportionate iff $?m \in [0, 1]$, there exists $(?e, ?r)$ such that $L_e(z)$

$$\begin{aligned}
& = L_r(z) + b , ?z \in [0, 1] \\
& (b \text{ does not depend on } z) ?m \in [0, 1] , \text{ let} \\
& . 0 \leq 2m \leq 1 \text{ } G_m = \\
& G_m \leq 1 \\
& 1 \leq 2m \leq 1 \text{ } G_m \leq 1 \\
& (7) \\
& m \\
& (? \in \{0, 1\}^m) \\
& (8) \\
& .
\end{aligned}$$

if $m \leq 1$, and $G1 = [0, 1]$ otherwise (notation $z \in \mathbb{R}^d$ indicates a vector in \mathbb{R}^d). Theorem 2 ϕ_e, ϕ_r are proportionate iff the optima $p^*(z)$ and $q^*(z)$ to eqs (3) and (4) satisfy: $p^*(z) \in \text{Le}(z)$, $Gm q^*(z) \in \text{Lr}(z)$.

(9) (10)

If ϕ_e, ϕ_r are differentiable and strictly convex, they are proportionate iff $p^*(z) = Gm q^*(z)$. We can alleviate the fact that convexity is strict, which results in a set-valued identity for ϕ_e, ϕ_r to be proportionate. This gives a necessary and sufficient condition for two generators to be proportionate. It does not say how to construct one from the other, if possible. We now show that it is indeed possible and prune the search space: if ϕ_e is proportionate to some ϕ_r , then it has to be a "symmetrized" version of ϕ_r , according to the following definition.

Definition 3 Let ϕ_r s.t. $\text{dom} \phi_r \subset (0, 1)$. $\phi_s(r)(z) = \phi_r(z) + \phi_r(1 - z)$ is the symmetrisation of ϕ_r . Lemma 4 If ϕ_e and ϕ_r are proportionate, then $\phi_e(z) = (\phi_r / \phi_e) \in \phi_s(r)(z) + (b / \phi_e)$ (b is in (7)). 3

I II III IV

P

$\phi_e(z, \phi_e) \log(1 + \exp(z \cdot e)) \leq 2(1 + z \cdot e) \leq \max\{0, z \cdot e\} \leq \phi_e(z)$

$\phi_r(z, \phi_r) \leq \exp(z \cdot r) \leq \text{EI}[\phi_r] + \phi_r \leq \text{VI}[\phi_r]$

$\max\{0, \max_{I \in \mathcal{I}} \{z \cdot I\} \leq \text{EI}[z \cdot I]\}$

$\phi_r(z) \leq z \log z + z(1/2) \leq z^2 \leq [0, 1](z) \leq [1m, 1](z) \leq 2$

2

$\phi_e \phi_e / 4 \phi_e \phi_e$

ϕ_e and ϕ_r $\phi_e = \phi_r \phi_e = \phi_r \phi_e, \phi_r \phi_e, \phi_r$

Table 1: Examples of equivalent example and rado losses. Names of the rado-losses $\phi_r(z, \phi_r)$ are respectively the Exponential (I), Mean-variance (II), ReLU P (III) and Unhinged (IV) rado loss. We use shorthands $z \cdot e = (1/\phi_e) \cdot z$ and $z \cdot r = (1/\phi_r) \cdot z$. Parameter ϕ_e appears in eq. (14). Column ϕ_e and ϕ_r gives the constraints for the equivalence to hold. EI and VI are the expectation and variance over uniform sampling in sets $I \in \mathcal{I}$ (see text for details). To summarize, ϕ_e and ϕ_r are proportionate iff (i) they meet the structural property that ϕ_e is (proportional to) the symmetrized version of ϕ_r (according to Definition 3), and (ii) the optimal solutions $p^*(z)$ and $q^*(z)$ to problems (1) and (2) satisfy the conditions of Theorem 2. Depending on the direction, we have two cases to craft proportionate generators. First, if we have ϕ_r , then necessarily $\phi_e = \phi_s(r)$ so we merely have to check Theorem 2. Second, if we have ϕ_e , then it matches Definition 3. In this case, we have to find $\phi_r = f + g$ where $g(z) = \phi_g(1 - z)$ and $\phi_e(z) = f(z) + f(1 - z)$. We now come back to $\text{Le}(z)$, $\text{Lr}(z)$ (Definition 1), and make the connection with example and rado losses. In the next definition, an e-loss $\phi_e(z)$ is a function defined over the coordinates of z , and a r-loss $\phi_r(z)$ is a function defined over the subsets of sums of coordinates. Functions can depend on other parameters as well. Definition 5 Suppose e-loss $\phi_e(z)$ and r-loss $\phi_r(z)$ are such that there exist (i) $f_e : \mathbb{R} \rightarrow \mathbb{R}$

and $f_r(z) : \mathbb{R} \rightarrow \mathbb{R}$ both strictly increasing and such that $\forall z \in \mathbb{R}_m, \ell_e(z) = \ell_r(z) =$

$$f_e(e(z)), f_r(r(z)), \quad (11) \quad (12)$$

where $\ell_e(z)$ and $\ell_r(z)$ are defined via two proportionate generators e and r (Definition 1). Then the couple (e, r) is called a couple of equivalent example-rado losses. Following is the main Theorem of this Section, which summarizes all the cases of equivalence between example and rado losses, and shows that the theory developed on example / rado losses with proportionate generators encompasses the specific proofs and cases already known [Nock et al., 2015, Patrini et al., 2016]. Table 1 also displays generator r . Theorem 6 In each row of Table 1, $e(z, e)$ and $r(z, r)$ are equivalent for e and r as indicated. The proof (SM, Subsection 2.3) details for each case the proportionate generators e and r .

3

Learning with (rado) regularized losses.

We now detail further the learning setting. In the preceeding Section, we have defined $z_i(h) = y_i h(x_i)$, which we plug in the losses of Table 1 to obtain the corresponding example and rado losses. Losses simplify conveniently when H consists of linear classifiers, $h(x) = \sum_i x_i$ for some $\{x_i\} \in \mathbb{R}^d$. In this case, the example loss can be described using edge vectors $S_e = \{y_i x_i, i = 1, 2, \dots, m\}$ since $z_i = \sum_j (y_j x_j)$, and the rado loss can be described using rademacher observations [Nock et al., 2015], since $z_i = \sum_j x_j$ for $y_j = y_i$ iff $j \in I$ (and y_j otherwise) and $\sum_j x_j = (1/2) \sum_j (x_j + y_j x_j)$. Let us define $S_r = \{x_j, j \in [m]\}$ the set of all rademacher observations. We rewrite any couple of equivalent example and rado losses as $e(S_e, ?)$ and $r(S_r, ?)$ respectively, omitting parameters e and r , assumed to be fixed beforehand for the equivalence to hold (see Table 1). Let us regularize the example loss, so that the learner's goal is to minimize $e(S_e, ?, ?) + \lambda$

$$= e(S_e, ?) + \lambda, \quad (13)$$

Alternatively, e is permissible [Kearns and Mansour, 1999]. To prevent notational overload, we blend notions of (pointwise) loss and (samplewise) risk, as just losses.

4

Algorithm 1 R.A DA B OOST. Input set of rados $S_r = \{x_1, x_2, \dots, x_n\}$; $T \in \mathbb{N}$; parameters $\eta \in (0, 1)$, $\beta \in \mathbb{R}_+$; Step 1 : let $\eta_0 = 0$, $w_0 = (1/n) \mathbf{1}$; Step 2 : for $t = 1, 2, \dots, T$ Step 2.1 : call the weak learner: $(\ell(t), r_t) \leftarrow \text{WL}(S_r, w_t, \eta, \beta, \ell(t-1))$; Step 2.2 : compute update parameters $\eta(t)$ and $\beta(t)$ (here, $\eta(t) = \max_j \sum_k \beta(t) \ell(t, x_k)$): $\eta(t) = (1/(2\sum_k \ell(t, x_k))) \log((1 + r_t)/(1 - r_t))$ and $\beta(t) = \eta(t) / (1 - \eta(t))$;

(16)

Step 2.3 : update and normalize weights: for $j = 1, 2, \dots, n$, $w_{t,j} \leftarrow w_{t-1,j} \exp(\beta(t) \ell(t, x_j)) / Z_t$;

(17)

Return \hat{T} ;

with ϕ a regularizer [Bach et al., 2011]. The following shows that when f_e in eq. (11) is linear, there is a rado-loss equivalent to this regularized loss, regardless of ϕ . Theorem 7 Suppose H contains linear classifiers. Let $(\phi(\mathbf{S}_e, \mathbf{r}), \mathbf{r}(\mathbf{S}_e \mathbf{r}, \mathbf{r}))$ be any couple of equivalent example-rado losses such that f_e in eq. (11) is linear: $f_e(\mathbf{z})$

$$= a_e \phi(\mathbf{z}) + b_e ,$$

(14)

for some $a_e \in \mathbb{R}$, $b_e \in \mathbb{R}$. Then for any regularizer $\phi(\cdot)$ (assuming wlog $\phi(0) = 0$), the regularized example loss $\phi(\mathbf{S}_e, \mathbf{r}, \mathbf{r})$ is equivalent to rado loss $\mathbf{r}(\mathbf{S}_e, \mathbf{r}, \mathbf{r})$ computed over regularized rados: $\mathbf{r}(\mathbf{S}_e, \mathbf{r}, \mathbf{r})$

$$. \phi = \mathbf{S}_e \mathbf{r} \phi \{ \phi(\mathbf{r}) \} ,$$

(15)

. Here, ϕ is Minkowski sum and $\phi(\mathbf{r}) = a_e \phi(\mathbf{r}) / \|\mathbf{r}\|_2^2$ if $\mathbf{r} \neq 0$ (and 0 otherwise).

Theorem 7 applies to all rado losses (I-IV) in Table 1. The effect of regularization on rados is intuitive from the margin standpoint: assume that a "good" classifier ϕ is one that ensures lowerbounded inner products $\phi(\mathbf{z}) \geq \gamma$ for some margin γ threshold γ . Then any good classifier on a regularized rado ϕ shall actually meet, over examples, $\phi(\mathbf{y}_i - \mathbf{x}_i) \geq \gamma + a_e \phi(\mathbf{r})$. This inequality ties an "accuracy" of ϕ (edges, left hand-side) and its sparsity (right-hand side). Clearly, Theorem 7 has an unfamiliar shape since regularisation modifies data in the rado world: a different ϕ , or a different \mathbf{r} , yields a different $\mathbf{S}_e, \mathbf{r}, \mathbf{r}$, and therefore it may seem very tricky to minimize such a regularized \mathbf{r} loss. Even more, iterative algorithms like boosting algorithms look at first glance a poor choice, since any update on \mathbf{r} implies an update on the rados as well. What we show in the following Section is essentially the opposite for the exponential rado loss, and a generalization of the RADO BOOST algorithm of Nock et al. [2015], which does not modify rados, is a formal boosting algorithm for a broad set of regularizers. Also, remarkably, only the high-level code of the weak learner depends on the regularizer; that of the strong learner is not affected.

4

Boosting with (rado) regularized losses

ϕ -RADO BOOST presents our approach to learning with rados regularized with regularizer ϕ to . Pt exp minimise loss $\mathbf{r}(\mathbf{S}_e, \mathbf{r}, \mathbf{r})$ in eq. (45). Classifier ϕ_t is defined as $\phi_t = \mathbf{t}_0 + \phi(\mathbf{t}_0) - \mathbf{1}(\mathbf{t}_0)$, where $\mathbf{1}_k$ is the k th canonical basis vector. The expected edge \mathbf{r}_t used to compute ϕ_t in eq. (16) is based on the following basis assignation: $\mathbf{r}(\mathbf{t})$

ϕ

1

\mathbf{X}

$\phi(\mathbf{t})$

$\mathbf{j}=1$

$\mathbf{w}_t \mathbf{j} \phi(\mathbf{t}) (\mathbf{r}[\mathbf{j}, 1])$.

(19)

The computation of r_t is eventually tweaked by the weak learner, as displayed in Algorithm 2. We investigate four choices for γ . For each of them, we prove the boosting ability of γ -R.A DA B OOST (γ is symmetric positive definite, S_d is the symmetric group of order d , γ is the 5

Algorithm 2 γ -WL, for $\gamma \in \{k.k1, k.k2, k.k, k.k\}$. Input set of n samples $S = \{x_1, x_2, \dots, x_n\}$; weights $w \in \mathbb{R}^n$; parameters $\gamma \in (0, 1)$, $\eta \in \mathbb{R}^+$; classifier $f \in \mathbb{R}^d$; Step 1: pick weak feature $j \in [d]$; Optional γ use preference order: $\gamma_0 \gamma \rightarrow \gamma \rightarrow \gamma \rightarrow \gamma_0 \rightarrow \gamma \rightarrow \gamma_0$. // $\gamma = \gamma \gamma (\gamma + \gamma \gamma 1) \gamma \gamma (\gamma)$, r_j is given in (19) and γ_j is given in (16) Step 2: if $\gamma = k.k2$ then

r_j if $r_j \gamma \gamma [\gamma, \gamma] r_j \gamma$; (18) $\text{sign}(r_j \gamma) \gamma \gamma$ otherwise else $r_j \gamma r_j \gamma$; Return (γ, r_j) ;

vector whose coordinates are the absolute values of the coordinates of γ): $\gamma \cdot k.k1 = \gamma \gamma 1$ Lasso $\gamma \gamma \cdot 2 k.k \gamma = \gamma \gamma \gamma$ Ridge $\gamma \gamma = k.k = \max \gamma \gamma \gamma \gamma k k \gamma \cdot k.k \gamma = \max M \gamma S_d (M \gamma \gamma) \gamma \gamma$ SLOPE

(20)

[Bach et al., 2011, Bogdan et al., 2015, Duchi and Singer, 2009, Su and Candès, 2015]. The γ coordinates of γ in SLOPE are $\gamma_k = \gamma_1 (1 \gamma kq/(2d))$ where $\gamma_1 (\cdot)$ is the quantile of the standard normal distribution and $q \in (0, 1)$; thus, the largest coordinates (in absolute value) of γ are more penalized. We now establish the boosting ability of γ -R.A DA B OOST. We give no direction for Step 1 in γ -WL, which is consistent with the definition of a weak learner in the boosting theory: all we require from the weak learner is γ . γ no smaller than some weak learning threshold $\gamma_{WL} \gamma 0$. Definition 8 Fix any constant $\gamma_{WL} \in (0, 1)$. γ -WL is said to be a γ_{WL} -Weak Learner iff the feature $\gamma(t)$ it picks at iteration t satisfies $\gamma \gamma(t) \gamma \gamma \gamma_{WL}$, for any $t = 1, 2, \dots, T$. We also provide an optional step for the weak learner in γ -WL, which we exploit in the experimentations, which gives a total preference order on features to optimise further γ -R.A DA B OOST. Theorem 9 (boosting with ridge). Take $\gamma(\cdot) = k.k2$. Fix any $0 \gamma a \gamma 1/5$, and suppose that γ and the number of iterations T of γ -R.A DA B OOST are chosen so that γ

$$\begin{aligned} & \gamma \\ & (2a \min \max \gamma_{2jk}) / (T \gamma \gamma) \gamma, j \\ & k \\ & (21) \end{aligned}$$

where $\gamma \gamma \gamma 0$ is the largest eigenvalue of γ . Then there exists some $\gamma \gamma 0$ (depending on a , and given to γ -WL) such that for any fixed $0 \gamma \gamma_{WL} \gamma \gamma$, if γ -WL is a γ_{WL} -Weak Learner, then γ -R.A DA B OOST returns at the end of the T boosting iterations a classifier γ_T which meets: ‘ $\exp(S_r, \gamma_T, k.k2) \gamma \gamma \exp(a \gamma_{2WL} T / 2) \cdot r$

(22)

Furthermore, if we fix $a = 1/7$, then we can fix $\gamma = 0.98$, and if $a = 1/10$, then we can fix $\gamma = 0.999$. Two remarks are in order. First, the cases $a = 1/7, 1/10$ show that γ -WL can still obtain large edges in eq. (19), so even a ‘strong’ weak learner might fit in for γ -WL, without clamping edges. Second, the right-hand side of ineq. (21) may be very large if we consider that $\max_j \gamma_{2jk}$ may be proportional to m^2 . So the constraint on γ is in fact loose.

Theorem 10 (boosting with lasso or ‘?’). Take $\eta(\cdot) = \{k.k1, k.k?\}$. Suppose η -WL is a η -WL-Weak Learner for some η -WL ≤ 0 . Suppose $\eta_0 \leq 3/11$ s. t. η satisfies:

$$\eta = \eta\text{-WL} \min_j \max_k -\eta_{jk} - \eta.k \quad (23)$$

Then η -R.A DA B OOST returns at the end of the T boosting iterations a classifier η T which meets: $\exp(Sr, \eta T, \eta) \leq \exp(\eta T \eta\text{-WL} / 2) \cdot r$

(24)

where $T\eta = \eta\text{-WL} T$ if $\eta = k.k1$, and $T\eta = (T \eta T\eta) + \eta\text{-WL} \eta T\eta$ if $\eta = k.k?$; $T\eta$ is the number of iterations where the feature computing the η norm was updated. We finally investigate the SLOPE choice. The Theorem is proven for $\eta = 1$ in η -R.A DA B OOST, for two reasons: it matches the original definition [Bogdan et al., 2015] and furthermore it unveils an interesting connection between boosting and SLOPE properties.

Theorem 11 (boosting with SLOPE). Take $\eta(\cdot) = k.k?$. Let $a = \min\{3\eta\text{-WL} / 11, \eta(1 - \eta/(2d)) / \min_j \max_k -\eta_{jk} - \eta.k\}$. Suppose $\eta \log -\eta T k - \eta - \eta T (k+1) - \eta.k$, and fix $\eta = 1$. Suppose (i) η -WL is a η -WL-Weak Learner for some η -WL ≤ 0 , and (ii) the q -value is chosen to meet:

$3\eta\text{-WL} k q \leq 2 \eta \max_j \max_k -\eta_{jk} - \eta.k$ Then classifier η T returned by η -R.A DA B OOST at the end of the T boosting iterations satisfies: $\exp(Sr, \eta T, k.k?) \leq \exp(\eta a \eta\text{-WL} T / 2) \cdot r$

(25)

Constraint (ii) on q is interesting in the light of the properties of SLOPE [Su and Candès, 2015]. Modulo some assumptions, SLOPE yields a control the false discovery rate (FDR) i.e., negligible coefficients in the “true” linear model that are found significant in the learned η . Constraint (ii) links the “small” achievable FDR (upperbounded by q) to the “boostability” of the data: the fact that each feature k can be chosen by the weak learner for a “large” η -WL, or has $\max_j -\eta_{jk} - \eta.k$ large, precisely flags potential significant features, thus reducing the risk of sparsity errors, and allowing small q , which is constraint (ii). Using the second order approximation of normal quantiles [Su and Candès, 2015], a sufficient condition for (ii) is that, for some $K \leq 0$, $p \eta\text{-WL} \min_j \max_k -\eta_{jk} - \eta.k \leq K \eta \log d + \log q$; (26)

j

but $\min_j \max_k -\eta_{jk} - \eta.k$ is proportional to m , so ineq. (26), and thus (ii), may hold even for small samples and q -values. An additional Theorem deferred to SM for space considerations shows that for any applicable choice of regularization (eq. 20), the regularized log-loss of η T over examples enjoys with high probability a monotonically decreasing upperbound with T as: $\log(\text{Se}, \eta, \eta) \leq e \log 2 \eta \eta T + \eta(m)$, with $\eta(m) \leq 0$ when $m \geq 0$ (and η does not depend on T), and $\eta \leq 0$ does not depend on T . Hence, η -R.A DA B OOST is an efficient proxy to boost the regularized log-loss over examples, using whichever of the ridge, lasso, ‘?’ or SLOPE regularization establishing the first boosting algorithm for this choice η , or linear combinations of the choices, e.g. for elastic nets. If we were to compare Theorems 9 – 11 (eqs (22, 24, 25)), then the

convergence looks best for ridge $\lambda = 2 \times 10^{-4}$) while it looks slightly worse for $\lambda = 10^{-4}$ and SLOPE (the unsigned (the unsigned exponent is $O(10^{-4})$ WL 3 λ exponent is now $O(10^{-4})$ WL)), the lasso being in between.

5

Experiments

We have implemented λ -WL4 using the order suggested to retrieve the topmost feature in the order. Hence, the weak learner returns the feature maximising $-\mathbf{r}^T \mathbf{w} - \lambda \|\mathbf{w}\|_1$. The rationale for this comes from Q 2 the proofs of Theorems 9 & 11, showing that $\frac{1}{2} \exp(\mathbf{r}^T \mathbf{w} / 2 - \lambda \|\mathbf{w}\|_1)$ is an upperbound on the exponential regularized rado-loss. We do not clamp the weak learner for $\mathbf{r}(\cdot) = \mathbf{k} \cdot \mathbf{k}^2$, so the weak learner is restricted to Step 1 in λ -WL5. The objective of these experiments is to evaluate λ -R.A DA B OOST as a contender for supervised learning per se. We compared λ -R.A DA B OOST to A DA B OOST/ λ regularized-A DA B OOST [Schapire and Singer, 1999, Xi et al., 2009]. All algorithms are run for a total of $T = 1000$ iterations, and at the end of the iterations, the classifier in the sequence that minimizes the empirical loss is kept. Notice therefore that rado-based classifiers are evaluated on the training set which computes the 3

If several features match this criterion, T is the total number of iterations for all these features. Code available at: <http://users.cecs.anu.edu.au/~rnock/> 5 the values for λ that we test, in $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 2, 3, 4, 5\}$, are small with respect to the upperbound in ineq. (21) given the number of boosting steps ($T = 1000$), and would yield on most domains a maximal $\lambda \leq 1.4$

7

rados. To obtain very sparse solutions for regularized-A DA B OOST, we pick its λ (λ in [Xi et al., 2009]) in $\{10^{-4}, 1, 10^4\}$. The complete results aggregate experiments on twenty (20) domains, all but one coming from the UCI [Bache and Lichman, 2013] (plus the Kaggle competition domain ‘Give me some credit’), with up to $d = 500+$ features and $m = 100\,000+$ examples. Two tables, in the SM (Tables 1 and 2 in Section 3) report respectively the test errors and sparsity of classifiers, whose summary is given here in Table 2. The experimental setup is a ten-folds stratified cross validation for all algorithms and each domain. A DA B OOST/regularized-A DA B OOST is trained using the complete training fold. When the domain size $m \leq 40000$, the number of rados n used for λ -R.A DA B OOST is a random subset of rados of size equal to that of the training fold. When the domain size exceeds 40000, a random set of $n = 10000$ rados is computed from the training fold. Thus, (i) there is no optimisation of the examples chosen to compute rados, (ii) we always keep a very small number of rados compared to the maximum available, and (iii) when the domain size gets large, we keep a comparatively tiny number of rados. Hence, the performances of λ -R.A DA B OOST do not stem from any optimization in the choice or size of the rado sample. Ada

$\lambda = 10^{-4}$

k.k2Id 10 3

k.k1 10 3 11

k.k? 8 2 9 7

k.k? 9 1 7 4 8

Experiments support several key observations. First, regularization consistently reduces the 9 ? test error of ?-R.A DA B OOST, by more than k.k2Id 10 17 15% on Magic, and 20% on Kaggle. In Table k.k1 10 17 7 2, ?-R.A DA B OOST unregularized ("?") is virk.k? 11 18 9 9 tually always beaten by its SLOPE regularized k.k? 10 19 10 10 11 version. Second, ?-R.A DA B OOST is able to obtain both very sparse and accurate classiTable 2: Number of domains for which algorithm in fiers (Magic, Hardware, Marketing, Kaggle). row beats algorithm in column (Ada = best result of A D Third, ?-R.A DA B OOST competes or beats A B OOST , ? = ?-R.A DA B OOST not regularized, see text). A DA B OOST on all domains, and is all the better as the domain gets bigger. Even qualitatively as seen in Table 2, the best result obtained by A DA B OOST (regularized or not) does not manage to beat any of the regularized versions of ?-R.A DA B OOST on the majority of the domains. Fourth, it is important to have several choices of regularizers at hand. On domain Stat-log, the difference in test error between the worst and the best regularization of ?-R.A DA B OOST exceeds 15%. Fifth, as already remarked [Nock et al., 2015], significantly subsampling rados (e.g. Marketing, Kaggle) still yields very accurate classifiers. Sixth, regularization in ?-R.A DA B OOST successfully reduces sparsity to learn more accurate classifiers on several domains (Spectf, Transfusion, Hill-noise, Winered, Magic, Marketing), achieving efficient adaptive sparsity control. Last, the comparatively extremely poor results of A DA B OOST on the biggest domains seems to come from another advantage of rados that the theory developed so far does not take into account: on domains for which some features are significantly correlated with the class and for which we have a large number of examples, the concentration of the expected feature value in rados seems to provide leveraging coefficients that tend to have much larger (absolute) value than in A DA B OOST, making the convergence of ?-R.A DA B OOST significantly faster than A DA B OOST. For example, we have checked that it takes much more than the $T = 1000$ iterations for A DA B OOST to start converging to the results of regularized ?-R.A DA B OOST on Hardware or Kaggle. Ada

6

Conclusion

We have shown that the recent equivalences between two example and rado losses can be unified and generalized via a principled representation of a loss function in a two-player zero-sum game. Furthermore, we have shown that this equivalence extends to regularized losses, where the regularization in the rado loss is performed over the rados themselves with Minkowski sums. Our theory and experiments on ?-R.A DA B OOST with prominent regularizers (including ridge, lasso, '? , SLOPE) indicate that when such a simple regularized form of the rado loss is available, it may help to devise accurate and efficient workarounds to boost a regularized loss over examples via the rado loss, even when the regularizer is significantly more involved like e.g. for group norms [Bach et al., 2011].

Acknowledgments Thanks are due to Stephen Hardy and Giorgio Patrini for

stimulating discussions around this material. 8

2 References

F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4:1?106, 2011.

K. Bache and M. Lichman. UCI machine learning repository, 2013.

M Bogdan, E. van den Berg, C. Sabatti, W. Su, and E.-J. Candès. SLOPE : adaptive variable selection via convex optimization. *Annals of Applied Statistics*, 2015. Also arXiv:1310.1969v2.

J.-C. Duchi and Y. Singer. Efficient learning using forward-backward splitting. In *NIPS*22*, pages 495?503, 2009.

C. Gentile and M. Warmuth. Linear hinge loss and average margin. In *NIPS*11*, pages 225?231, 1998.

M.J. Kearns and Y. Mansour. On the boosting ability of top-down decision tree learning algorithms. *J. Comp. Syst. Sc.*, 58:109?128, 1999.

V. Nair and G. Hinton. Rectified linear units improve restricted boltzmann machines. In *27th ICML*, pages 807?814, 2010.

R. Nock and F. Nielsen. On the efficient minimization of classification-calibrated surrogates. In *NIPS*21*, pages 1201?1208, 2008.

R. Nock, G. Patrini, and A Friedman. Rademacher observations, private data, and boosting. In *32nd ICML*, pages 948?956, 2015.

G. Patrini, R. Nock, S. Hardy, and T. Caetano. Fast learning from distributed datasets without entity matching. In *26 th IJCAI*, 2016.

M.-D. Reid, R.-M. Frongillo, R.-C. Williamson, and N.-A. Mehta. Generalized mixability via entropic duality. In *28th COLT*, pages 1501?1522, 2015.

R.-E. Schapire. The boosting approach to machine learning: An overview. In D.-D. Denison, M.-H. Hansen, C.-C. Holmes, B. Mallick, and B. Yu, editors, *Nonlinear Estimation and Classification*, volume 171 of *Lecture Notes in Statistics*, pages 149?171. Springer Verlag, 2003.

R. E. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. *MLJ*, 37:297?336, 1999.

W. Su and E.-J. Candès. SLOPE is adaptive to unknown sparsity and asymptotically minimax. *CoRR*, abs/1503.08393, 2015.

M. Telgarsky. A primal-dual convergence analysis of boosting. *JMLR*, 13:561?606, 2012.

B. van Rooyen, A. Menon, and R.-C. Williamson. Learning with symmetric label noise: The importance of being unhinged. In *NIPS*28*, 2015.

V. Vapnik. *Statistical Learning Theory*. John Wiley, 1998.

Y.-T. Xi, Z.-J. Xiang, P.-J. Ramadge, and R.-E. Schapire. Speed and sparsity of regularized boosting. In *12th AISTATS*, pages 615?622, 2009.

H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society B*, 67:301?321, 2005.