

# An improved multiclass LogitBoost using adaptive-one-vs-one

Peng Sun · Mark D. Reid · Jie Zhou

Received: 25 April 2013 / Accepted: 15 January 2014 / Published online: 28 February 2014  
© The Author(s) 2014

**Abstract** LogitBoost is a popular Boosting variant that can be applied to either binary or multi-class classification. From a statistical viewpoint LogitBoost can be seen as additive tree regression by minimizing the Logistic loss. Following this setting, it is still non-trivial to devise a sound multi-class LogitBoost compared with to devise its binary counterpart. The difficulties are due to two important factors arising in multiclass Logistic loss. The first is the invariant property implied by the Logistic loss, causing the optimal classifier output being not unique, i.e. adding a constant to each component of the output vector won't change the loss value. The second is the density of the Hessian matrices that arise when computing tree node split gain and node value fittings. Oversimplification of this learning problem can lead to degraded performance. For example, the original LogitBoost algorithm is outperformed by ABC-LogitBoost thanks to the latter's more careful treatment of the above two factors. In this paper we propose new techniques to address the two main difficulties in multiclass LogitBoost setting: (1) we adopt a vector tree model (i.e. each node value is vector) where the unique classifier output is guaranteed by adding a sum-to-zero constraint, and (2) we use an adaptive block coordinate descent that exploits the dense Hessian when computing tree split gain and node values. Higher classification accuracy and faster convergence rates are observed for a range of public data sets when compared to both the original and the ABC-LogitBoost implementations. We also discuss another possibility to cope with LogitBoost's dense Hessian matrix. We derive a loss similar to the multi-class Logistic loss but which

---

Editor: Shai Shalev-Shwartz.

---

P. Sun (✉) · J. Zhou  
Tsinghua National Laboratory for Information Science and Technology(TNList),  
Department of Automation, Tsinghua University, Beijing100084, China  
e-mail: sunp08@mails.tsinghua.edu.cn

J. Zhou  
e-mail: jzhou@tsinghua.edu.cn

M. D. Reid  
Research School of Computer Science,  
The Australian National University and NICTA, Canberra, ACT, Australia  
e-mail: mark.reid@anu.edu.au

guarantees a diagonal Hessian matrix. While this makes the optimization (by Newton descent) easier we unfortunately observe degraded performance for this modification. We argue that working with the dense Hessian is likely unavoidable, therefore making techniques like those proposed in this paper necessary for efficient implementations.

**Keywords** LogitBoost · Boosting · Ensemble · Supervised learning · Convex optimization

## 1 Introduction

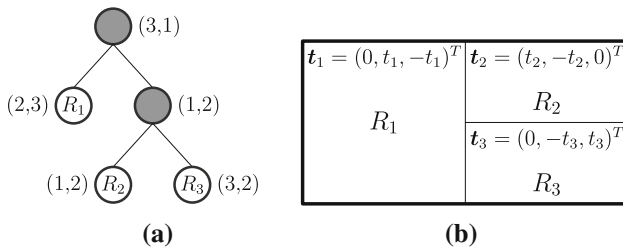
Boosting is a successful technique for training classifier for both binary and multi-class classification (Freund and Schapire 1995; Schapire and Singer 1999). In this paper, our focus is on multiclass LogitBoost (Friedman et al. 1998), one of the popular boosting variants. Originally, LogitBoost was motivated by a statistical perspective (Friedman et al. 1998), where boosting algorithm consists of three key components: the loss, the function model, and the optimization algorithm. In the case of LogitBoost, these are the multi-class Logistic loss, the use of additive tree models, and a stage-wise optimization, respectively. For  $K$ -nary classification, LogitBoost directly learns a  $K$  dimensional vector as the classifier output, each component representing the confidence of predicting the corresponding class.

There are two important factors in LogitBoost's settings. Firstly, an "invariant property" is implied by the Logistic loss, i.e., adding a constant to each component of the classifier output won't change the loss value. Therefore, the classifier output minimizing the total loss is not unique, making the optimization procedure a bit difficult. Secondly, the Logistic loss produces a dense Hessian matrix, causing troubles when deriving the tree node split gain and node value fitting. It is challenging to design a tractable optimization algorithm that fully handles both these factors. Consequently, some simplification and/or approximation is needed.

In Friedman et al. (1998) the Hessian is diagonally approximated. In this way, the minimizer becomes unique, while the optimization—essentially a quadratic problem when using one step Newton—is substantially simplified. Consequently, at each Boosting iteration the tree model updating collapses to  $K$  independent Weighted Regression Tree fittings, each tree outputting a scalar.

Unfortunately, Friedman's prescription turns out to have some drawbacks. The over simplified quadratic loss even doesn't satisfy the invariant property, and is thus a very crude approximation. A later improvement, ABC-LogitBoost, is shown to outperform LogitBoost in terms of both classification accuracy and convergence rate (Li 2008, 2010b). This is due to ABC-LogitBoost's careful handling of the above key problems of the LogitBoost setting. The invariant property is addressed by adding a sum-to-zero constraint to the output vector. To make the minimizer unique, at each iteration one variable is eliminated while the other ( $K - 1$ ) are free, so that the loss is re-written with ( $K - 1$ ) variables. At this point, the diagonal Hessian approximation is, again, adopted to permit  $K - 1$  scalar trees being independently fitted for each of the  $K - 1$  classes. The remaining class—called the base class—is recovered by taking minus the sum of the other  $K - 1$  classes. The base class—i.e., the variable to be eliminated—is selected adaptively per iteration (or every several iterations), hence the acronym ABC (Adaptive Base Class). Note the diagonal Hessian approximation in ABC-LogitBoost is taken for the ( $K - 1$ ) dimensional problem yielded by eliminating a redundant variable. Li (2008, 2010b) considers it a more refined approximation than that of original LogitBoost (Friedman et al. 1998).

In this paper, we propose two novel techniques to address the challenging aspects of the LogitBoost setting. In our approach, one vector tree is added per iteration. To make the minimizer unique, we allow a  $K$  dimensional sum-to-zero vector to be fitted for each tree



**Fig. 1** A newly added tree at some boosting iteration for a 3-class problem. **a** A class pair (shown in brackets) is selected for each tree node. For each internal node (filled), the pair is for computing split gain; For terminal nodes (unfilled), it is for node vector updating. **b** The feature space (the outer black box) is partitioned by the tree in (a) into regions  $\{R_1, R_2, R_3\}$ . On each region only two coordinates are updated based on the corresponding class pair shown in (a)

node. This permits us to explicitly formulate the computation for both node split gain and node value fitting as a  $K$  dimensional constrained quadratic optimization, which arises as a subproblem in the inner loop for split seeking when fitting a new tree. To avoid the difficulty of a dense Hessian, we propose that for each of these subproblems, only two coordinates (i.e., two classes or a class pair) are adaptively selected for updating, hence we call the modified algorithm AOSO-LogitBoost (Adaptive One vs One). Figure 1 gives an overview of our approach. In Sect 4.4 we show that first order and second order approximation of loss reduction can be a good measure for the quality of selected class pair.

Following the above formulation, ABC-LogitBoost (derived using a somewhat different framework than Li 2010b) can thus be shown to be a special case of AOSO-LogitBoost with a less flexible tree model. In Sect. 5 we compare the differences between the two approaches in detail and provide some intuition for AOSO’s improvement over ABC.

Both ABC or AOSO are carefully devised to address the difficulties of dense Hessian matrix arising in Logistic loss. In other words, the tree model learning would be easy if we were encountering diagonal Hessian matrix. Based on loss design in Primal/Dual view (Masnadi-Shirazi and Vasconcelos 2010; Reid and Williamson 2010), we show that the dense Hessian matrix essentially results from the sum-to-one constraint of class probability in Logistic loss. We thus investigate the possibility of diagonal Hessian matrix by removing such a sum-to-one constraint. The LogitBoost variant we produce does work. However, it is still inferior to AOSO LogitBoost (see Sect. 6 for detailed discussion). We therefore conclude that modifications of the original LogitBoost, such as AOSO and ABC, are necessary for efficient model learning since the dense Hessian matrix seems unavoidable.

The rest of this paper is organised as follows: In Sect. 2 we first formulate the problem setting for LogitBoost. In Sect. 3 we briefly review/discuss original/ABC-/AOSO- LogitBoost in regards of the interplay between the tree model and the optimisation procedure. In Sect. 4 we give the details of our approach. In Sect. 5 we compare our approach to (ABC-)LogitBoost. In Sect. 6 we show how to design a loss that produces diagonal Hessian matrices and assess its implications for model accuracy. In Sect. 7, experimental results in terms of classification errors and convergence rates are reported on a range of public datasets.

## 2 The problem setup

For  $K$ -class classification ( $K \geq 2$ ), consider an  $N$  example training set  $\{x_i, y_i\}_{i=1}^N$  where  $x_i$  denotes a feature value and  $y_i \in \{1, \dots, K\}$  denotes a class label. From the training

set a prediction function  $F(x) \in \mathbb{R}^K$  is learned. When clear from context, we omit the dependence on  $x$  and simply denote  $F(x)$  by  $F$  (We will do the same to other related variables.). Given a test example with known  $x$  and unknown  $y$ , we predict a class label by taking  $\hat{y} = \arg \max_k F_k, k = 1, \dots, K$ , where  $F_k$  is the  $k$ -th component of  $F$ .

The function  $F$  is obtained by minimizing a target function on training data:

$$\text{Loss} = \sum_{i=1}^N L(y_i, F_i). \tag{1}$$

Each  $L(y, F)$  term is a *loss* for a single training example  $(x, y)$ . We will make the loss concrete shortly.

To make the optimization of (1) feasible, a *function model* is needed to describe how  $F$  depends on  $x$ . For example, linear model  $F = W^T x + b$  is used in traditional Logistic regression, while Generalized Additive Model is adopted in LogitBoost where

$$F(x) = \sum_{m=1}^M f^{(m)}(x), \tag{2}$$

where  $f^{(m)}(x) \in \mathbb{R}^K$  is further expressed by tree(s), as will be seen shortly.

Each  $f^{(m)}(x)$ , or simply  $f$ , is learned by the so-called greedy stage-wise *optimization*. That is, at each Boosting iteration  $f^{(m)}$  is added only based on  $F = \sum_{j=1}^{m-1} f^{(j)}$ .

In summary, the learning procedure, called training, consists of three ingredients: the *loss*, the *function model* and the *optimization* algorithm. In what follows we discuss them in greater details.

### 2.1 The logistic loss

LogitBoost adopts the Logistic loss (Friedman et al. 1998), which is an implicit function of  $y$  and  $F$ :

$$L(y, F) = - \sum_{k=1}^K r_k \log(p_k), \tag{3}$$

where  $r_k = 1$  if  $k = y$  and 0 otherwise,  $\{p_k\}_{k=1}^K$  is connected to  $F$  via the so-called Link function:

$$p_k = \frac{\exp(F_k)}{\sum_{j=1}^K \exp(F_j)}. \tag{4}$$

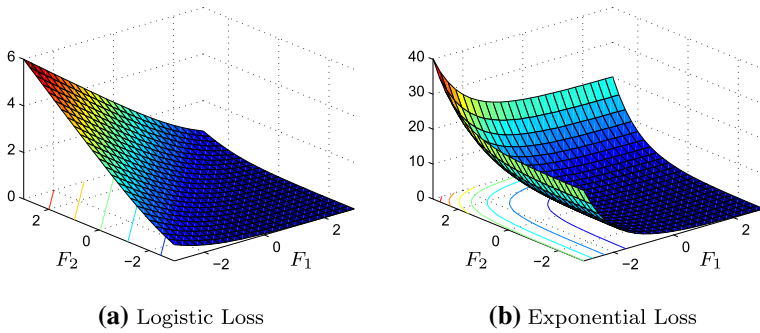
Therefore,  $p = (p_1, \dots, p_K)^T$  can be seen as probability estimate since  $p_k \geq 0$  and  $\sum_{k=1}^K p_k = 1$ .

#### 2.1.1 The invariant property

It is easy to verify that an “invariant property” holds for Logistic loss:

$$L(y, F) = L(y, F + c\mathbf{1}), \tag{5}$$

where  $c$  is an arbitrary constant,  $\mathbf{1}$  is  $K$ -dimensional all 1 vector. That is to say, the Logistic loss won’t change its value when moving  $F$  along a all-1 vector  $\mathbf{1}$ , which implies that  $F$  is “equivalent to”  $F + c\mathbf{1}$  because they both lead to the same loss (see Fig. 2b). Such an invariant property is, in turn, consistent with the multi-class prediction rule  $\arg \max_k F_k, k = 1, \dots, K$  which says adding a constant to each of  $F_k$  does not alter the prediction.



**Fig. 2** Multiclass loss  $L(y = 1, \mathbf{F})$  for  $K$ -nary problem. In this figure  $K = 2$ , thus  $\mathbf{F} \in \mathbb{R}^2$ . The Logistic loss is defined in (3); The Exponential loss is defined in (10). Note the parallel contour plots in the case of Logistic loss

2.1.2 Derivative and quadratic approximation

It is difficult to directly optimize the Logistic loss. To permit numeric method, we need the first order and second order derivatives of example wise loss (3). Simple matrix calculus shows that the gradient  $\mathbf{g} \in \mathbb{R}^K$  and Hessian  $\mathbf{H} \in \mathbb{R}^{K \times K}$  w.r.t.  $\mathbf{F}$  are:

$$\mathbf{g} = -(\mathbf{r} - \mathbf{p}), \tag{6}$$

$$\mathbf{H} = \hat{\mathbf{P}} - \mathbf{p}\mathbf{p}^T, \tag{7}$$

where the response  $\mathbf{r} = (r_1, \dots, r_K)^T$  and probability estimate  $\mathbf{p} = (p_1, \dots, p_K)^T$  are as defined before, the diagonal matrix  $\hat{\mathbf{P}} = \text{diag}(p_1, \dots, p_K)$ .

With  $\mathbf{g}$  and  $\mathbf{H}$  given above, we can write down the Taylor expansion of (3) at  $\mathbf{F}$  up to second order:

$$q(\mathbf{f}|\mathbf{F}) = L(y, \mathbf{F}) + \mathbf{g}^T \mathbf{f} + \frac{1}{2} \mathbf{f}^T \mathbf{H} \mathbf{f}, \tag{8}$$

which locally approximates  $L(y, \mathbf{F} + \mathbf{f})$  in quadratic sense. Note that in (8)  $\mathbf{f}$  is the variable, while the “constants”  $\mathbf{g}$  and  $\mathbf{H}$  depend on  $\mathbf{F}$ . When clear from context, we also omit the dependence on  $\mathbf{F}$  and simply write  $q(\mathbf{f}|\mathbf{F})$  as  $q(\mathbf{f})$ .

It is easy to verify that the invariant property carries over to  $q(\mathbf{f})$ :

$$q(\mathbf{f}) = q(\mathbf{f} + c\mathbf{1}) \tag{9}$$

by noting  $\mathbf{1}^T \mathbf{g} = 0$  and  $\mathbf{1}^T \mathbf{H} = \mathbf{H}\mathbf{1} = 0$ .

2.1.3  $(K - 1)$  Degrees of freedom and the sum-to-zero constraint

Due to the invariant property, the minimizer  $\mathbf{f}^*$  of (8) is not unique since any  $\mathbf{f}^* = c\mathbf{1}$  would be a minimizer. To pin-down the value, we can add a constraint  $\mathbf{1}^T \mathbf{f} = 0$ , which in effect restricts  $\mathbf{f}$  to vary just in the linear subspace defined by  $\mathbf{1}^T \mathbf{f} = 0$ . Obviously, now we need only  $K - 1$  coordinates to express the vector living in the subspace, i.e., the degrees of freedom is  $K - 1$ .

In Sect. 4.2 we will discuss the rank of the Hessian matrix  $\mathbf{H}$ , which provides another perspective to why the minimizer of (8) is not unique.

Conceptually, the invariant property is primary, causing the minimizer being not unique; The sum-to-zero constraint is secondary, serving as a mathematical tool to make the minimizer unique.

### 2.1.4 More on $(K - 1)$ degrees of freedom

We don't claim that the invariant property be a necessity. In the literature, there exists other multiclass loss not satisfying this condition. In this case, it doesn't need a sum-to-zero constraint and  $f$  has  $K$  degrees of freedom. For instance, consider the  $K$ -nary exponential loss for AdaBoost.MH:

$$L(y, \mathbf{F}) = \sum_{k=1}^K \exp(-y_k^* \mathbf{F}_k), \tag{10}$$

where  $y_k^* = +1$  if  $y_k = k$  and  $-1$  otherwise. See Fig. 2(b) for its graph. In general,  $L(y, \mathbf{F}) \neq L(y, \mathbf{F} + c\mathbf{1})$  for Exponential loss.

It is out of this paper's scope to compare Exponential loss and Logistic loss in multiclass case. Instead, we focus on the Logistic loss and show how it is applied in LogitBoost. Also, in Sect. 6 we discuss a modified Logistic loss not satisfying the invariant property (but simplifying the Hessian and easing the quadratic solver) and show its degraded performance when comparing with original Logistic loss.

## 2.2 The tree model

As mentioned previously,  $\mathbf{F}(\mathbf{x}) = \sum_{m=1}^M \mathbf{f}^{(m)}(\mathbf{x})$  is additive tree model. However, the way each  $\mathbf{f}(\mathbf{x})$  (we have omitted the subscript  $m$  for simplicity and without confusion) being expressed by tree model is not unique.

In this paper, we adopt a single vector-valued tree. Further, we let it be a binary tree (i.e., only binary splits on internal node are allowed) with splits that are vertical to coordinate axis, as in Friedman et al. (1998) and Li (2010b). Formally,

$$\mathbf{f}(\mathbf{x}) = \sum_{j=1}^J \mathbf{t}_j I(\mathbf{x} \in R_j) \tag{11}$$

where  $\{R_j\}_{j=1}^J$  are a rectangular partition of the feature space,  $I(P)$  is the indicator function which is 1 when the predicate  $P$  is true and 0 otherwise, and each  $\mathbf{t}_j \in \mathbb{R}^K$  is the node value associated with  $R_j$ . See Fig. 1 for an illustration. The vector-valued tree model is also adopted in other Boosting implementation, say, Real AdaBoost.MH described in Kégl and Busa-Fekete (2009).

$\mathbf{f}(\mathbf{x})$  can be also represented by  $K$  scalar regression trees as in the original LogitBoost of Friedman et al. (1998) and the Real AdaBoost.MH implementation of Friedman et al. (1998):

$$\mathbf{f}(\mathbf{x}) = \sum_{k=1}^K \sum_{j=1}^J t_{k,j} I(\mathbf{x} \in R_{k,j}), \tag{12}$$

where  $t_{k,j}$  is a scalar and each tree can have its own partition  $\{R_{k,j}\}_{j=1}^J$ .

Finally, it is possible to express  $f v(\mathbf{x})$  with just  $K - 1$  trees by adding to  $\mathbf{f}(\mathbf{x})$  a sum-to-zero constraint, as is adopted in Li's series work (Li 2008, 2009, 2010b).

### 2.3 Stage wise optimization

Formally, using  $F_i$  and  $f_i$  as shorthand for  $F(x_i)$  and  $f(x_i)$  respectively, the stage wise optimization is:

$$f^{(m)}(x) = \arg \min_{f(x)} \sum_{i=1}^N L(y_i, F_i + f_i), \tag{13}$$

where  $f(x)$  is expressed by tree model as explained in last subsection. This procedure repeats  $M$  times with initial condition  $F = 0$ . Owing to its iterative nature, we only need to know how to solve (13) in order to implement the optimization.

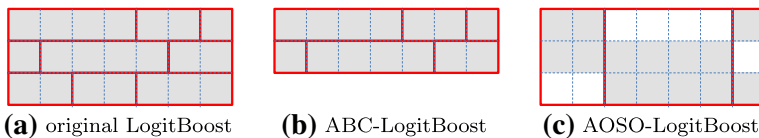
In (13) it is difficult to directly work on the Logistic loss; numeric method must be relied on. Friedman et al. (1998) suggests to use one step Newton, i.e., in (13) replace example wise Logistic loss  $L()$  with its corresponding second order Taylor expansion:

$$f^{(m)}(x) = \arg \min_{f(x)} \sum_{i=1}^N q(f_i | F_i). \tag{14}$$

### 3 Interplay between tree model and optimization

In last section we have reviewed the problem setup. In particular, we made concrete the three ingredients in LogitBoost: the loss, the function model and the optimization algorithm are the Logistic loss, the additive tree model and the stage wise optimization, respectively. To train a LogitBoost classifier, now the only thing left unexplained is how to optimize the quadratic optimization (14), which is still a bit complicated.

Basically,  $f(x)$  on a training set  $\{x_i, y_i\}_{i=1}^N$  can be seen as an  $K \times N$  matrix, as in Fig. 3, where the  $i$ -th column represents the  $K$  values of  $f(x_i)$ . However, the values in the matrix can not vary independently. Instead, the values must be subject to a tree model, which determines a type of partition on the matrix such that the elements falling into the same cell take a common, unique value. Thus the tree model has an impact on the optimization procedure. Things become even subtle when considering the invariant property of the Logistic loss. The original/ABC-/AOSO-LogitBoost makes different choices to address this tree model-optimization interplay, as will be briefly discussed in the rest of this section. The details are deferred to next section.



**Fig. 3** Viewing  $f(x)$  as a  $K \times N$  matrix on a training dataset. A  $J$ -leaf tree corresponds a  $J$ -cell partition on the matrix such that the values falling into the same cell must have a common value. In this figure,  $K = 3$ ,  $N = 7$  and  $J = 3$ . The squares in dashed lines denote the matrix elements, while the solid lines denote the partitions. Note that each cell needs not be continuous, although we intentionally plot it a continuous one for illustrative purpose. **a**  $K$  scalar trees. Each tree fits a class (a row) and has its own partition. **b**  $(K - 1)$  scalar trees. Each tree fits a class (a row) and has its own partition. The values on the remaining class (i.e. the base class) is recovered by the other classes. **c** only one vector-valued tree. Any column belongs to only one cell. At each cell, only two classes (two rows) can vary (shaded squares) and the others keep zero values (blank squares)

### 3.1 Original LogitBoost

To solve problem (14), LogitBoost (Friedman et al. 1998) takes further simplification. In the example wise quadratic loss (8), the Hessian matrix is diagonally approximated. Thus (8) collapses into  $K$  independent 1-dimensional loss, i.e., (8) is written as the sum of  $K$  terms each involving just one component of  $f$ :

$$\begin{aligned}
 q(f|\mathbf{F}) &= \sum_{k=1}^K q_k(f_k|\mathbf{F}), \\
 q_k(f_k|\mathbf{F}) &= g_k f_k + \frac{1}{2} h_{kk} f_k^2,
 \end{aligned}
 \tag{15}$$

where  $h_{kk}$  is the  $k$ -th diagonal elements in  $\mathbf{H}$  and we have without confusion omitted the constants that are irrelevant to  $f$ . Noting the  $N$ -additive form and substituting back (15), we can also collapse (14) into  $K$  1-dimensional, independent minimization problems:

$$f_k^{(m)}(\mathbf{x}) = \arg \min_{f_k(\mathbf{x})} \sum_{i=1}^N q_k(f_{i,k}|\mathbf{F}_i), \quad k = 1, \dots, K,
 \tag{16}$$

where  $f_{i,k} = f_k(\mathbf{x}_i)$  means the  $f_k$  value corresponding to the  $i$ -th training examples.

Then  $K$  scalar trees are grown, the  $k$ -th tree approximately minimizing the  $k$ -th problem. Note that each tree can have its own partition, as in Fig. 3(a).

To grow the scalar tree, Friedman et al. (1998) borrowed a traditional regression model in statistics, namely the Weighted Regression Tree. The formulation becomes: for the  $k$ -th problem, fit a Weighted Regression Tree on the training examples  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  with targets  $\{-g_{i,k}/h_{i,kk}\}_{i=1}^N$  and weights  $\{h_{i,kk}\}_{i=1}^N$ , where (in a slightly abuse of notation) we use an additional subscript  $i$  for  $g_k$  and  $h_{kk}$  to denote they correspond to the  $i$ -th training example.

### 3.2 ABC-LogitBoost

LogitBoost adopts a rather crude approximation to (8), where the Hessian matrix  $\mathbf{H} \in \mathbb{R}^{K \times K}$  is diagonally approximated. ABC-LogitBoost (Li 2010b) considers an intuitively more refined way. Recall that  $(K - 1)$  degrees of freedom suffices to express (8), i.e., (8) can be re-written with just  $(K - 1)$  variables by eliminating one redundant variable. Then the new  $(K - 1) \times (K - 1)$  Hessian matrix is, again, diagonally approximated. Li (2010b) shows that this does make a difference, as what follows.

In Sect. 2.1.3 we have explained that adding a sum to zero constraint to (8) doesn't change the problem and actually makes the minimizer unique. We can thus restrict  $f$  to vary in the subspace defined by  $\mathbf{1}^T f = 0$ . Consequently, it is convenient to use a new coordinate system  $\tilde{f} \in \mathbb{R}^{K-1}$ . In order to do this, Li (2010b) introduces a concept "base class". Suppose the base class  $b = K$ , then the new coordinates are that  $f_1 = \tilde{f}_1, \dots, f_{K-1} = \tilde{f}_{K-1}, f_K = -(\tilde{f}_1 + \dots + \tilde{f}_{K-1})$ . Writing them in compact matrix notation, we have:

$$\begin{aligned}
 f &= A \tilde{f} \\
 A &= \begin{bmatrix} \mathbf{I} \\ -\mathbf{1} \end{bmatrix},
 \end{aligned}
 \tag{17}$$



where  $\mathbf{A} \in \mathbb{R}^{K \times (K-1)}$ ,  $\mathbf{I} \in \mathbb{R}^{(K-1) \times (K-1)}$  and  $\mathbf{1}$  is  $(K - 1)$ -dim all one vector. Substituting it back to (8) we have:

$$\begin{aligned} \tilde{q}(\tilde{\mathbf{f}}) &\triangleq q(\mathbf{A}\tilde{\mathbf{f}}) \\ &= (\mathbf{A}^T \mathbf{g})^T \tilde{\mathbf{f}} + \frac{1}{2} \tilde{\mathbf{f}}^T \mathbf{A}^T \mathbf{H} \mathbf{A} \tilde{\mathbf{f}} + C \\ &= \tilde{\mathbf{g}}^T \tilde{\mathbf{f}} + \frac{1}{2} \tilde{\mathbf{f}}^T \tilde{\mathbf{H}} \tilde{\mathbf{f}} + C \end{aligned} \tag{18}$$

where  $C$  is a irrelevant constant and we have let

$$\begin{aligned} \tilde{\mathbf{g}} &= \mathbf{A}^T \mathbf{g} \\ \tilde{\mathbf{H}} &= \mathbf{A}^T \mathbf{H} \mathbf{A}. \end{aligned} \tag{19}$$

Li (2010b) then diagonally approximates the  $(K - 1) \times (K - 1)$  matrix  $\tilde{\mathbf{H}} = \mathbf{A}^T \mathbf{H} \mathbf{A}$ . After doing this, the rest optimization procedure is essentially the same with original LogitBoost: the loss (14) collapses to  $(K - 1)$  1-dimensional, independent minimization problems; the  $k$ -th tree is fitted by a Weighted Regression Tree using the new targets and weights derived from  $\tilde{\mathbf{g}}$  and  $\tilde{\mathbf{H}}$  as are defined in (19).

Li (2010b) shows that the choice of  $b$  impacts on how well the diagonal approximation is. Two intuitive methods are proposed to select  $b$ : (1) The “worst class”, i.e., the  $b$  having the biggest loss (before minimizing (14)) on that class is selected. (2) The “best class”, i.e., all possible  $b$ , up to  $K$  choices, are tried and the one leading to lowest loss (after minimizing (14)) is selected.

### 3.3 AOSO-LogitBoost

In AOSO-LogitBoost, we also adopts the loss (18) with  $K - 1$  free coordinates. However, we update only one coordinate a time. An equivalent formulation goes in the following. We add a sum-to-zero constraint  $\mathbf{1}^T \mathbf{f} = 0$  to the  $K$ -dim loss (8) and let only two coordinates of  $\mathbf{f}$ , say,  $f_r$  and  $f_s$ , vary and the other  $K - 2$  keep zero. Due to the sum-to-zero constraint, we further let  $f_r = +t$  and  $f_s = -t$  where  $t$  is a real number.

Obviously, the choice of  $(r, s)$  impacts on the goodness of approximation. However, it is unlikely to select the best class pair  $(r, s)$  for each training example. To permit the class pair selection as adaptive as possible, we adopt vector-valued tree model in AOSO-LogitBoost, i.e., any column in Fig. 3(c) can not be assigned to two cells. Then, for each node (i.e., each cell in the matrix as in Fig. 3c) we adaptively select the class pair  $(r, s)$ .

Superficially, AOSO is inferior to original/ABC- LogitBoost, because many of the  $f$  values in the matrix are untouched (zeros), as shown in Fig. 3(c). However, we should recall the “big picture”: the untouched values might hopefully receive better updating in later Boosting iterations. Consequently, AOSO is still “on average” more efficient than original/ABC-LogitBoost. We will further discuss this issue in Sect. 5.

## 4 The AOSO-LogitBoost algorithm

In this section we describe the details of AOSO-LogitBoost. Specifically, we will focus on how to build a tree in Boosting iteration. Some key ingredients of tree building improving previous techniques were firstly introduced by Li. However, we will re-derive them in our own language. The credits will be made clear when they are explained in the following.

### 4.1 Details of tree building

Solving (14) with the tree model (11) is equivalent to determining the parameters  $\{t_j, R_j\}_{j=1}^J$  at the  $m$ -th iteration. In this subsection we will show how this problem reduces to solving a collection of quadratic subproblems for which we can use standard numerical methods-Block Coordinate Descent.<sup>1</sup> Also, we will show how the gradient and Hessian can be computed incrementally.

We begin by defining some notation for the *node loss*:

$$NodeLoss(\mathbf{t}; \mathcal{I}) = \sum_{i \in \mathcal{I}} q(\mathbf{t} | \mathbf{F}_i) \tag{20}$$

where  $\mathcal{I}$  denotes the index set of the training examples on some either internal or terminal node (i.e., those falling into the corresponding region). Minimizing (20) allows us to obtain a set of nodes  $\{t_j, R_j\}_{j=1}^J$  using the following procedures:

1. To obtain the values  $t_j$  for a given  $R_j$ , we simply take the minimizer of (20):

$$t_j = \arg \min_t NodeLoss(\mathbf{t}; \mathcal{I}_j), \tag{21}$$

where  $\mathcal{I}_j$  denotes the index set for  $R_j$ .

2. To obtain the partition  $\{R_j\}_{j=1}^J$ , we recursively perform binary splitting until there are  $J$ -terminal nodes.

The key to the second procedure is the computation of the node split gain. Suppose an internal node with  $n$  training examples ( $n = N$  for the root node), we fix on some feature and re-index all the  $n$  examples according to their sorted feature values. Now we need to find the index  $n'$  with  $1 < n' < n$  that maximizes the node gain defined as loss reduction after a division between the  $n'$ -th and  $(n' + 1)$ -th examples:

$$NodeGain(n') = NodeLoss(\mathbf{t}^*; \mathcal{I}) - (NodeLoss(\mathbf{t}_L^*; \mathcal{I}_L) + NodeLoss(\mathbf{t}_R^*; \mathcal{I}_R)) \tag{22}$$

where  $\mathcal{I} = \{1, \dots, n\}$ ,  $\mathcal{I}_L = \{1, \dots, n'\}$  and  $\mathcal{I}_R = \{n' + 1, \dots, n\}$ ;  $\mathbf{t}^*$ ,  $\mathbf{t}_L^*$  and  $\mathbf{t}_R^*$  are the minimizers of (20) with index sets  $\mathcal{I}$ ,  $\mathcal{I}_L$  and  $\mathcal{I}_R$ , respectively. Generally, this search applies to all features. The division yielding the largest value of (22) is then recorded to perform the actual node split.

Note that (22) arises in the context of an  $O(N \times D)$  outer loop, where  $D$  is number of features. However, a naïve summing of the losses for (20) incurs an additional  $O(N)$  factor in complexity, which finally results in an unacceptable  $O(N^2 D)$  complexity for a single boosting iteration.

Fortunately, the gradient and Hessian can be incrementally computed. To see why, simply rewrite the (20) in standard quadratic form:

$$loss(\mathbf{t}; \mathcal{I}) = \bar{\mathbf{g}}^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \bar{\mathbf{H}} \mathbf{t}, \tag{23}$$

where  $\bar{\mathbf{g}} = -\sum_{i \in \mathcal{I}} \mathbf{g}_i$ ,  $\bar{\mathbf{H}} = \sum_{i \in \mathcal{I}} \mathbf{H}_i$ . Thanks to the additive form, both  $\bar{\mathbf{g}}$  and  $\bar{\mathbf{H}}$  can be incrementally/decrementally computed in constant time when the split searching proceeds

<sup>1</sup> We use Newton descent as there is evidence in Li (2010b) that gradient descent, i.e., in Friedmans’s MART (Friedman 2001), leads to decreased classification accuracy.

from one training example to the next. Therefore, the computation of (23) eliminates the  $O(N)$  complexity in the naïve summing of losses.<sup>2</sup>

#### 4.2 Properties of approximated node loss

To minimise (23), we make use of some properties for (23) that can be exploited when finding a solution. First, the invariant property carries over to the node loss (23):

*Property 1*  $loss(\mathbf{t}; \mathcal{I}) = loss(\mathbf{t} + c\mathbf{1}; \mathcal{I})$ .

*Proof* This is obvious by noting the additive form. □

For the Hessian  $\mathbf{H}$ , we have  $rank(\mathbf{H}) \leq rank(\mathbf{H}_i)$  by noting the additive form in (23). In Li (2010b) it is shown that  $\det \mathbf{H}_i = 0$  by brute-force determinant expansion. Here we give a stronger property:

*Property 2* Each  $\mathbf{H}_i$  is a positive semi-definite matrix such that (1)  $rank(\mathbf{H}_i) = \kappa - 1$ , where  $\kappa$  is the number of non-zero elements in  $\mathbf{p}_i$ ; (2)  $\mathbf{1}$  is the eigenvector for eigenvalue 0.

The proof can be found in Appendix 1.

The properties shown above indicate that (1)  $\mathbf{H}$  is singular so that unconstrained Newton descent is not applicable here and (2)  $rank(\mathbf{H})$  could be as high as  $K - 1$ , which prohibits the application of the standard fast quadratic solver designed for low rank Hessian. In the following we propose to address this problem via block coordinate descent, a technique that has been successfully used in training SVMs (Bottou and Lin 2007).

#### 4.3 Block coordinate descent

For the variable  $\mathbf{t}$  in (23), we only choose two coordinates, i.e., a class pair, to update while keeping the others fixed. We note that single coordinates cannot be updated independently due to the sum-to-zero constraint. Suppose we have chosen the  $r$ -th and the  $s$ -th coordinate (we explain precisely how to do so in the next section). Let  $t_r = t$  and  $t_s = -t$  be the free variables (such that  $t_r + t_s = 0$ ) and  $t_k = 0$  for  $k \neq r, k \neq s$ . Plugging these into (23) yields an unconstrained one dimensional quadratic problem with regards to the scalar variable  $t$ :

$$loss(t) = \bar{g}^T t + \frac{1}{2} \bar{h} t^2 \tag{24}$$

where the gradient and Hessian collapse to scalars:

$$\bar{g} = - \sum_{i \in \mathcal{I}} ((r_{i,r} - p_{i,r}) - (r_{i,s} - p_{i,s})) \tag{25}$$

$$\bar{h} = \sum_{i \in \mathcal{I}} (p_{i,r}(1 - p_{i,r}) + p_{i,s}(1 - p_{i,s}) + 2p_{i,r}p_{i,s}), \tag{26}$$

The derivatives (25) and (26) involving two classes were firstly given in Li (2008), where the class subscript  $r, s$  are fixed for the nodes in a tree. In this work, we allow  $r, s$  vary from node to node.

---

<sup>2</sup> In Real AdaBoost.MH, such a second order approximation is not necessary (although possible, cf. Zou et al. 2008). Due to the special form of the exponential loss and the absence of a sum-to-zero constraint, there exists analytical solution for the node loss (20) by simply setting the derivative to  $\mathbf{0}$ . Here also, the computation can be incremental/decremental. Since the loss design and AdaBoost.MH are not our main interests, we do not discuss this further.

To this extent, we are able to obtain the analytical expression for the minimizer and minimum of (24):

$$t^* = \arg \min_t \text{loss}(t) = -\frac{\bar{g}}{\bar{h}} \tag{27}$$

$$\text{loss}(t^*) = -\frac{\bar{g}^2}{2\bar{h}} \tag{28}$$

by noting the non-negativity of (26).

Based on (27), node vector (21) can be approximated by  $t_j$  with components

$$t_{j,k} = \begin{cases} +(-\bar{g}/\bar{h}) & k = r \\ -(-\bar{g}/\bar{h}) & k = s \\ 0 & \text{otherwise} \end{cases} \tag{29}$$

where  $g$  and  $h$  are computed using (25) and (26), respectively, with index set  $\mathcal{I}_j$ . Based on (28), the node gain (22) can be approximated as

$$\text{NodeGain}(n') = \frac{\bar{g}_L^2}{2\bar{h}_L} + \frac{\bar{g}_R^2}{2\bar{h}_R} - \frac{\bar{g}^2}{2\bar{h}}, \tag{30}$$

where  $\bar{g}$  (or  $\bar{g}_L, \bar{g}_R$ ) and  $\bar{h}$  (or  $\bar{h}_L, \bar{h}_R$ ) are computed by using (25) and (26) with index set  $\mathcal{I}$  (or  $\mathcal{I}_L, \mathcal{I}_R$ ). We note that (30) was firstly derived in Li (2010b) in a slightly different way.

#### 4.4 Class pair selection

Bottou and Lin (2007) provide two methods for selecting  $(r, s)$  for an SVM solver that we consider for our purposes. The first is based on a first-order approximation: let  $t_r$  and  $t_s$  be the free variables and fix the rest to 0. For a  $t$  with sufficiently small norm, let  $t_r = \epsilon$  and  $t_s = -\epsilon$  where  $\epsilon > 0$  is some small enough constant. The first order approximation of (23) is then:

$$\text{loss}(t) \approx \text{loss}(\mathbf{0}) + \bar{\mathbf{g}}^T t = \text{loss}(\mathbf{0}) - \epsilon(-\bar{g}_r - (-\bar{g}_s)), \tag{31}$$

where we have denoted the vector  $\bar{\mathbf{g}} = (\bar{g}_1, \dots, \bar{g}_r, \dots, \bar{g}_s, \dots, \bar{g}_K)^T$ . It follows that the indices  $r, s$  resulting in largest decrement to (31) are:

$$\begin{aligned} r &= \arg \max_k \{-\bar{g}_k\} \\ s &= \arg \min_k \{-\bar{g}_k\}. \end{aligned} \tag{32}$$

The second method, derived in a similar way, takes into account the second-order information:

$$\begin{aligned} r &= \arg \max_k \{-\bar{g}_k\} \\ s &= \arg \max_k \left\{ \frac{(\bar{g}_r - \bar{g}_k)^2}{\bar{h}_{rr} + \bar{h}_{kk} - 2\bar{h}_{rk}} \right\}, \end{aligned} \tag{33}$$

where we have denoted the matrix  $\bar{\mathbf{H}} = \{\bar{h}_{ij}\}$ .

Both methods are  $O(K)$  procedures that are better than the  $K \times (K - 1)/2$  naïve enumeration of all possible pairs. However, in our implementation we find that (33) achieves better results for AOSO-LogitBoost.

The selection of a class pair  $(r, s)$  here is somewhat similar to the selection of base class  $b$  in ABC Boost. Actually, Li proposed in Li (2008) that the “worst class” with the largest loss be selected as  $b$ . Clearly, the max derivative is another indicator for how “worst” it is. In

this sense, our class pair selection extends the base class idea and can be seen as a concrete implementation of the general “worst class” idea.

The pseudocode for AOSO-LogitBoost is given in Algorithm 1 and shows how all the above approximations are used together to iteratively find a model. The algorithm runs on the sample  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  for  $M$  rounds. In each round  $m$  it updates the values of  $F_i = F(\mathbf{x}_i)$  for each of the instances  $\mathbf{x}_i$  by first computing a rectangular partition of the feature space  $\{R_j^m\}_{j=1}^J$  and the corresponding node values  $\{t_j^m\}_{j=1}^J$ , and then incrementing  $F_i$  by  $t_j^m$  where  $j$  is the index of the region  $R_j^m$  such that  $\mathbf{x}_i \in R_j^m$ .

---

**Algorithm 1** AOSO-LogitBoost for a sample  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ . Parameters:  $v$  = shrinkage factor (controls learning rate);  $J$  = number of terminal nodes;  $M$  = number of rounds

---

- 1:  $F_i = \mathbf{0}, i = 1, \dots, N$
  - 2: **for**  $m = 1$  to  $M$  **do**
  - 3:  $p_{i,k} = \frac{\exp(F_{i,k})}{\sum_{j=1}^K \exp(F_{i,j})}, k = 1, \dots, K, i = 1, \dots, N.$
  - 4: Obtain  $\{R_j^m\}_{j=1}^J$  by recursive region partition. Node split gain is computed as (30), where the class pair  $(r, s)$  is selected using (33).
  - 5: Compute  $\{t_j^m\}_{j=1}^J$  by (29), where the class pair  $(r, s)$  is selected using (33).
  - 6:  $F_i = F_i + v \sum_{j=1}^J t_j^m I(\mathbf{x}_i \in R_j^m), i = 1, \dots, N.$
  - 7: **end for**
- 

### 5 Comparison to (ABC-)LogitBoost

In this section we compare the derivations of LogitBoost and ABC-LogitBoost and provide some intuition for observed behaviours in the experiments in Sect. 7.

To solve (14), ABC-LogitBoost uses  $K - 1$  independent trees defined by

$$f_k(x) = \begin{cases} \sum_j t_{j,k} I(x \in R_{j,k}) & , k \neq b \\ - \sum_{l \neq b} f_l(x) & , k = b \end{cases} \tag{34}$$

where  $b$  is the *base class*. As explained in Sect. 3.2,  $b$  can be selected as either the “worst class” or the “best class”.

Li (2010b) gives how to compute the node value and the node split gain for building the  $k$ -th tree ( $k \neq b$ ). Although derived from different motivation as ours, they are actually the same with (29) and (30) in this paper, where the class pair  $(r, s)$  is replaced with  $(k, b)$ . We should not be surprised at this coincidence by noting that AOSO’s vector tree has only one freely altering coordinate on each node and thus “behaves” like a scalar tree. In this sense, AOSO and ABC is comparable. Actually, ABC can be viewed as a special form of AOSO with two differences: (1) For each tree, the class pair is fixed for every node in ABC, while it is selected adaptively in AOSO, and (2)  $K - 1$  trees are added per iteration in ABC, while only one tree is added per iteration by AOSO.

It seems unappealing to add just one tree as in AOSO, since many  $f(x)$  values are untouched (i.e., set to zeros, as illustrated in Fig. 3c); In the meanwhile, ABC would be better since it updates all the  $f(x)$  values. Considering the Boosting context, we argue, however, that AOSO should still be preferred in an “on average” sense. After adding one tree in AOSO, the  $F(x)$  values are updated and the gradient/Hessian is immediately recomputed

for each training example, which impacts on how to build the tree at next iteration. Thus the  $F(x)$  values might still receive good enough updating after several iterations, due to the adaptive class pair selection for every node at current iteration. In contrast, the  $K - 1$  trees in ABC use the same set of gradients and Hessians, which are not recomputed until adding all the  $K - 1$  trees.

Therefore, it is fair to compare ABC and AOSO in regards of number of trees, rather than number of iterations. AOSO's "on average" better performance is confirmed by the experiments in Sect. 7.2.

To evaluate whether the adaptive class pair selection is critical, we considered a variant of AOSO-LogitBoost that adopts a *fixed* class pair selection. Specifically, we still add one tree per iteration, but select a single class pair root node and let it be fixed for all other nodes, which is very similar to ABC's choice. This variant was tried but unfortunately, **degraded performance** was observed so the results are not reported here.

From the above analysis, we believe that AOSO-LogitBoost's more flexible model obtained from the adaptive split selection (as well as its immediate model updating after adding one tree per iteration) is what contributes to its improvement over ABC.

## 6 Sum-to-one probability and dense Hessian matrix

As in previous discussion, both ABC or AOSO improve on the original LogitBoost method by dealing with dense Hessian matrix due to the Logistic loss as given in (3). An immediate question is that whether we can derive an effective alternative surrogate loss that has a diagonal Hessian matrix "by design"—i.e., can we define a modified Logistic loss that guarantees a diagonal Hessian matrix? In this section, we show how the original multi-class Logistic loss can be derived from a maximum entropy argument via convex duality, in a manner similar to derivations of boosting updates by [Shen and Li \(2010\)](#), [Shen and Hao \(2011\)](#), [Lafferty \(1999\)](#), and [Kivinen and Warmuth \(1999\)](#) and results connecting entropy and loss by [Masnadi-Shirazi and Vasconcelos \(2010\)](#) and [Reid and Williamson \(2010\)](#).

In contrast to earlier work, our analysis focuses on the role of the constraint in defining the loss and the effect it has on the form of its gradient and Hessian. In particular, we'll see that the original Logistic loss's dense Hessian matrix essentially results from sum-to-one constraint on class probabilities. Moreover, we are able to obtain a diagonal Hessian matrix by dropping this constraint when deriving an alternative to the Logistic loss from the same maximum entropy formulation. By doing so we show that the optimization (i.e., via Newton descent) becomes straightforward, however lower classification accuracy and slower convergence rate are observed for the new loss. Therefore, we argue the techniques used by ABC/AOSO seem necessary for dealing with the dense Hessian matrix for the original, more effective, Logistic loss.

### 6.1 Logistic loss in primal/dual view

We now examine the duality between entropy and loss in a manner similar to that of the more general treatment of [Masnadi-Shirazi and Vasconcelos \(2010\)](#); [Reid and Williamson \(2010\)](#). By starting with a "trivial" maximum entropy problem, we show how consideration of its dual problem recovers a "composite representation" [Reid and Williamson \(2010\)](#) of a loss function in terms of a loss defined on the simplex and corresponding link function. In case of Shannon entropy, we show that such a construction results in the original Logistic loss

function. Appendix 2 describes the matrix calculus definitions and conventions we adopt, mainly from [Magnus and Neudecker \(2007\)](#).

For each feature vector  $\mathbf{x}$ , let  $\boldsymbol{\eta}(\mathbf{x}) = Pr(\mathbf{y}|\mathbf{x}) \in \Delta^K$  be the true conditional probability for the  $K$  classes where  $\Delta^K = \{\mathbf{p} \in [0, 1]^K : \sum_{k=1}^K p_k = 1\}$  denotes the  $K$ -simplex. Let  $H : \Delta^K \rightarrow \mathbb{R}$  be an *entropy function*—a concave function such that  $H(\mathbf{e}^k) = 0$  for each vertex  $\mathbf{e}^k$  of  $\Delta^K$ .

Given some set of probabilities  $S \subset \Delta^K$ —typically defined by a set of constraints—the MaxEnt criterion ([Jaynes 1957](#)) advises choosing the  $\mathbf{p} \in S$  such that  $H(\mathbf{p})$  is maximized. We consider a “trivial” instance of the MaxEnt problem in which the set  $S$  contains only the single point  $\boldsymbol{\eta}$ , i.e.,  $S = \{\boldsymbol{\eta}\}$ , yielding the following problem:

$$\begin{aligned} \max \quad & H(\mathbf{p}) \\ \text{s.t.} \quad & \mathbf{p} - \boldsymbol{\eta} = 0 \\ & \mathbf{1}^T \mathbf{p} - \mathbf{1} = 0, \end{aligned} \tag{35}$$

where the argument  $\mathbf{x}$  to  $\boldsymbol{\eta}()$ ,  $\mathbf{p}()$  has been omitted. Despite its apparent triviality, the Lagrangian of this problem has an interesting form:

$$\mathcal{L}(\mathbf{p}, \mathbf{F}, \lambda) = H(\mathbf{p}) + \mathbf{F}^T (\mathbf{p} - \boldsymbol{\eta}) + \lambda(\mathbf{1}^T \mathbf{p} - \mathbf{1}), \tag{36}$$

where  $\mathbf{F} \in \mathbb{R}^K$ ,  $\lambda \in \mathbb{R}$  are the so-called dual variables. The dual function is obtained by maximizing the Lagrangian over the domain of the primal variable  $\mathbf{p}$ . By setting to zero the derivative of (36) w.r.t.  $\mathbf{p}$ , we can implicitly express  $\mathbf{p}$  as a function of  $\mathbf{F}$  and  $\lambda$  via

$$\nabla H(\mathbf{p}) + \lambda \mathbf{1} + \mathbf{F} = 0. \tag{37}$$

This gives us the dual function  $g$  which depends only on  $\mathbf{F}$ ,  $\lambda$ :

$$g(\mathbf{F}, \lambda) = H(\mathbf{p}) + \mathbf{F}^T (\mathbf{p} - \boldsymbol{\eta}) + \lambda(\mathbf{1}^T \mathbf{p} - \mathbf{1}), \tag{38}$$

where  $\mathbf{p} = \mathbf{p}(\mathbf{F}, \lambda)$  as per (37), and the unconstrained convex dual problem

$$\min \quad g(\mathbf{F}, \lambda).$$

Now further eliminating the dual variable  $\lambda$  in (38) by taking derivative of (38) w.r.t.  $\lambda$  as

$$\begin{aligned} D_\lambda g(\mathbf{F}, \lambda) &= (D_\lambda H(\mathbf{p}))(D_\lambda \mathbf{p}) + \mathbf{F}^T (D_\lambda \mathbf{p}) + \lambda \mathbf{1}^T (D_\lambda \mathbf{p}) + (\mathbf{1}^T \mathbf{p} - \mathbf{1}) \\ &= (D_\lambda H(\mathbf{p}) + \lambda \mathbf{1}^T + \mathbf{F}^T)(D_\lambda \mathbf{p}) + \mathbf{1}^T \mathbf{p} - \mathbf{1} \\ &= \mathbf{1}^T \mathbf{p} - \mathbf{1}. \end{aligned}$$

and setting to zero, we obtain the “partial” dual function:

$$\ell(\mathbf{F}) = H(\mathbf{p}) + \mathbf{F}^T (\mathbf{p} - \boldsymbol{\eta}), \tag{39}$$

where  $\mathbf{p} = \mathbf{p}(\mathbf{F}) \in \mathbb{R}^K$ ,  $\lambda = \lambda(\mathbf{F}) \in \mathbb{R}$  are given by the  $(K + 1) \times (K + 1)$  equation array

$$\begin{cases} \nabla H(\mathbf{p}) + \lambda \mathbf{1} + \mathbf{F} &= 0 \\ \mathbf{1}^T \mathbf{p} - \mathbf{1} &= 0. \end{cases} \tag{40}$$

In Machine Learning, (39) and (40) are precisely the “loss function” and the “link function”, respectively (e.g., the Logistic loss (3) and Logistic link (4)). Training boosting classifier is essentially a numerical optimization procedure, where the gradient and Hessian of the loss (39) are needed. We give them as follows:

**Theorem 1** *The gradient of (39) is  $\nabla\ell(\mathbf{F}) = \mathbf{p} - \boldsymbol{\eta}$ .*

See Appendix 2 for its proof. Note that the gradient is precisely the left-hand-side of the constraint in primal (35), which is a common result in primal-dual theory in convex optimization (Bertsekas 1982).

**Theorem 2** *The Hessian of (39) is  $\nabla^2\ell(\mathbf{F}) = -A^{-1} + \frac{1}{\mathbf{1}^T A^{-1} \mathbf{1}} A^{-1} \mathbf{1} \mathbf{1}^T A^{-1}$ , where the shorthand  $A = \nabla^2\mathbf{H}(\mathbf{p})$ .*

See Appendix 2 for its proof.

In both of the above two theorems, the dependence on  $\mathbf{F}$  is implicit where  $\mathbf{p}$  is given by the Link function (40).

To this extent, we can show some concrete results in case of Logistic settings. The starting point is that we adopt Shannon entropy in (35) such that  $\mathbf{H}(\mathbf{p}) = -\sum_{k=1}^K p_k \log p_k$ . Thus its gradient and Hessian are  $\nabla\mathbf{H}(\mathbf{p}) = -(1 + \log p_1, \dots, 1 + \log p_K)^T$ ,  $\nabla^2\mathbf{H}(\mathbf{p}) = -\text{diag}(1/p_1, \dots, 1/p_K)$ , respectively. With these two expressions, it follows that the link function by solving equations (40) is precisely the Logistic link in (4). Furthermore, according to Theorem 1 and 2 the gradient and Hessian for the loss (39) are  $\nabla\ell = \mathbf{p} - \boldsymbol{\eta}$  and  $\nabla^2\ell = \text{diag}(p_1, \dots, p_K) - \mathbf{p}\mathbf{p}^T$ , which are respectively no more than (6) and (7) by noting that  $\boldsymbol{\eta}$  reduces to 1-of- $K$  response  $\mathbf{r}$  when taking value on the vertex of the probability simplex. These are shown in the left half of Table 1.

### 6.2 Diagonal Hessian matrix by dropping the sum-to-one constraint

In last subsection we derived the LogitBoost via a duality argument. An immediate observation is that the sum-to-one constraint  $\mathbf{1}^T \mathbf{p} - 1 = 0$  in (35) seems redundant, since it is already guaranteed by the constraint  $\mathbf{p} - \boldsymbol{\eta} = 0$  where  $\boldsymbol{\eta}$  itself is sum-to-one. This means that the primal problem can be rewritten as:

$$\max \mathbf{H}(\mathbf{p}) \tag{41}$$

$$s.t. \quad \mathbf{p} - \boldsymbol{\eta} = 0 \tag{42}$$

Deriving its dual form, we obtain the corresponding loss

$$\ell(\mathbf{F}) = \mathbf{H}(\mathbf{p}) + \mathbf{F}^T (\mathbf{p} - \boldsymbol{\eta}), \tag{43}$$

and link

$$\nabla\mathbf{H}(\mathbf{p}) + \mathbf{F} = 0 \tag{44}$$

Still, with  $\mathbf{H}(\mathbf{p}) = -\sum_{k=1}^K p_k \log p_k$  we can obtain the link, the gradient and Hessian for the loss, as shown in the right half of Table 1. One attractive feature of this alternative derivation is the diagonal Hessian matrix that is yielded. When calculating the node value or node gain, we can obtain the precise Newton Step and Newton Decrement since the inversion of Hessian is much easier to compute. Consequently, no effort (e.g., base class selection in ABC or class pair selection in AOSO) is needed to deal with the difficulty of dense Hessian and all classes can be updated for each terminal node.

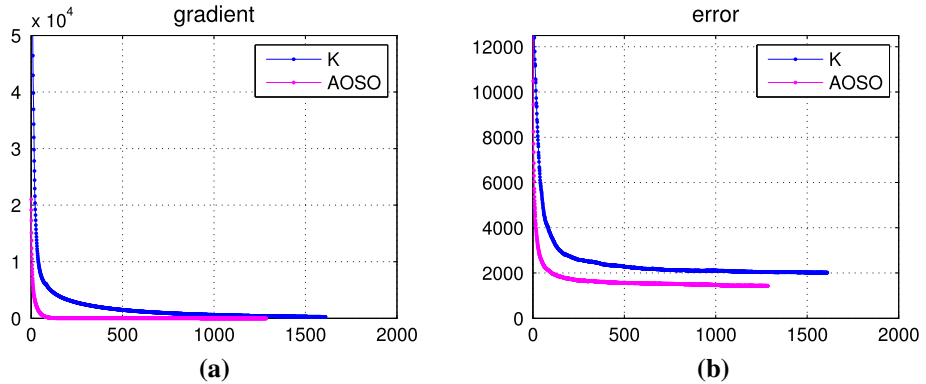
### 6.3 Degraded performance

We refer to the modified LogitBoost without sum-to-one constraint as “K-LogitBoost”, where all the  $K$  classes are updated at each terminal node. We compare K-LogitBoost with AOSO-LogitBoost as in Fig. 4. Noting that only 2 classes are updated at each terminal node in



**Table 1** The Link, gradient/Hessian for the loss with or without the sum-to-one constraint

	with sum-to-one	without sum-to-one
Link $p = p(F)$	$p_k = \frac{e^{F_k}}{\sum_j^K e^{F_j}}, k = 1, \dots, K$	$p_k = e^{F_k-1}, k = 1, \dots, K$
Gradient $\nabla \ell(F)$	$p - \eta$	$p - \eta$
Hessian $\nabla^2 \ell(F)$	$\text{diag}(p_1, \dots, p_K) - p p^T$	$\text{diag}(p_1, \dots, p_K)$



**Fig. 4** Comparison between K LogitBoost and AOSO LogitBoost on dataset M-Basic. **a** How the L2 norm of gradient  $\|p - \eta\|_2^2$  decreases when iteration proceeds. **b** Test error v.s. iteration

AOSO, we scale the horizontal axis (i.e., the number of iterations) by a factor  $2/K$  for AOSO to make it a fair comparison. As can be seen in the results, a slower convergence than AOSO is incurred by K LogitBoost with regard to both test error and L2 norm of gradient, although K LogitBoost is still competitive with AOSO LogitBoost. More comparative results are provided in Sect. 7.3.

We provide an intuitive explanation for this phenomenon. At convergence, the gradient  $p - \eta$  vanishes, and consequently,  $p$  satisfies the sum-to-one constraint since  $\eta$  is a probability. For the original Logistic loss setting, the apparently redundant sum-to-one constraint for  $p$  in (35) enforces that the primal variable  $p$  approaches  $\eta$  on the plane  $\mathbf{1}^T p - 1 = 0$  containing the simplex  $\Delta^K$ . In contrast, for K LogitBoost  $p$  may reside outside that plane during the optimization procedure. The latter optimization is intuitively slower since the L2 norm  $\|p - \eta\|_2^2$  never increases when  $p - \eta$  is projected onto a plane.

### 7 Experiments

In this section we compare AOSO-LogitBoost with ABC-LogitBoost, which was shown to outperform original LogitBoost in Li’s experiments (Li 2010b, 2009a). We test AOSO on all the datasets used in Li (2010b, 2009a), as listed in Table 2. In the top section are UCI datasets and in the bottom are Mnist datasets with many variations (see Li (2010a) for detailed descriptions).<sup>3</sup>

<sup>3</sup> Code and data are available at <http://ivg.au.tsinghua.edu.cn/index.php?n=People.PengSun>.

**Table 2** Datasets used in our experiments

Datasets	K	#features	#training	#test
Poker525k	10	25	5,25,010	5,00,000
Poker275k	10	25	2,75,010	5,00,000
Poker150k	10	25	1,50,010	5,00,000
Poker100k	10	25	1,00,010	5,00,000
Poker25kT1	10	25	25,010	5,00,000
Poker25kT2	10	25	25,010	5,00,000
Coverttype290k	7	54	2,90,506	2,90,506
Coverttype145k	7	54	1,45,253	2,90,506
Letter	26	16	16,000	4,000
Letter15k	26	16	15,000	5,000
Letter2k	26	16	2,000	18,000
Letter4K	26	16	4,000	16,000
Pendigits	10	16	7,494	3,498
Zipcode	10	256	7,291	2,007
(a.k.a. USPS)				
Isolet	26	617	6,238	1,559
Optdigits	10	64	3,823	1,797
Mnist10k	10	784	10,000	60,000
M-Basic	10	784	12,000	50,000
M-Image	10	784	12,000	50,000
M-Rand	10	784	12,000	50,000
M-Noise1	10	784	10,000	2,000
M-Noise2	10	784	10,000	2,000
M-Noise3	10	784	10,000	2,000
M-Noise4	10	784	10,000	2,000
M-Noise5	10	784	10,000	2,000
M-Noise6	10	784	10,000	2,000

To exhaust the learning ability of (ABC-)LogitBoost, Li let the boosting stop when either the training loss is small (implemented as  $\leq 10^{-16}$ ) or a maximum number of iterations,  $M$ , is reached. Test errors at last iteration are simply reported since no obvious over-fitting is observed. By default,  $M = 10,000$ , while for those large datasets (Coverttype290k, Poker525k, Pokder275k, Poker150k, Poker100k)  $M = 5,000$ . We adopt the same criteria, except that our maximum iterations  $M_{AOSO} = (K - 1) \times M_{ABC}$ , where  $K$  is the number of classes. Note that only one tree is added at each iteration in AOSO, while  $K - 1$  are added in ABC. Thus, this correction compares the same maximum number of trees for both AOSO and ABC.

The most important tuning parameters in LogitBoost are the number of terminal nodes  $J$ , and the shrinkage factor  $v$ . In (Li 2010b, 2009a), Li reported results of (ABC-)LogitBoost for a number of  $J$ - $v$  combinations. We report the corresponding results for AOSO-LogitBoost for the same combinations. In the following, we intend to show that **for nearly all  $J$ - $v$  combinations, AOSO-LogitBoost has lower classification error and faster convergence rates than ABC-LogitBoost.**

**Table 3** Summary of test classification errors

Datasets	#tests	ABC	AOSO	R	$pv$	ABC*	AOSO*	R	$pv$
Poker525k	500000	1736	<b>1537</b>	0.1146	0.0002	–	–	–	–
Poker275k	500000	2727	<b>2624</b>	0.0378	0.0790	–	–	–	–
Poker150k	500000	5104	<b>3951</b>	0.2259	0.0000	–	–	–	–
Poker100k	500000	13707	<b>7558</b>	0.4486	0.0000	–	–	–	–
Poker25kT1	500000	37345	<b>31399</b>	0.1592	0.0000	37345	<b>31399</b>	0.1592	0.0000
Poker25kT2	500000	36731	<b>31645</b>	0.1385	0.0000	36731	<b>31645</b>	0.1385	0.0000
Covertyp290k	290506	9727	<b>9586</b>	0.0145	0.1511	–	–	–	–
Covertyp145k	290506	13986	<b>13712</b>	0.0196	0.0458	–	–	–	–
Letter	4000	<b>89</b>	92	−0.0337	0.5892	89	<b>88</b>	0.0112	0.4697
Letter15k	5000	<b>109</b>	116	−0.0642	0.6815	–	–	–	–
Letter4k	16000	1055	<b>991</b>	0.0607	0.0718	1034	<b>961</b>	0.0706	0.0457
Letter2k	18000	2034	<b>1862</b>	0.0846	0.0018	1991	<b>1851</b>	0.0703	0.0084
Pendigits	3498	100	<b>83</b>	0.1700	0.1014	90	<b>81</b>	0.1000	0.2430
Zipcode	2007	<b>96</b>	99	−0.0313	0.5872	<b>92</b>	94	−0.0217	0.5597
Isolet	1559	65	<b>55</b>	0.1538	0.1759	55	<b>50</b>	0.0909	0.3039
Optdigits	1797	55	<b>38</b>	0.3091	0.0370	38	<b>34</b>	0.1053	0.3170
Mnist10k	60000	2102	<b>1948</b>	0.0733	0.0069	2050	<b>1885</b>	0.0805	0.0037
M-Basic	50000	1602	<b>1434</b>	0.1049	0.0010	–	–	–	–
M-Rotate	50000	5959	<b>5729</b>	0.0386	0.0118	–	–	–	–
M-Image	50000	4268	<b>4167</b>	0.0237	0.1252	4214	<b>4002</b>	0.0503	0.0073
M-Rand	50000	4725	<b>4588</b>	0.0290	0.0680	–	–	–	–
M-Noise1	2000	234	<b>228</b>	0.0256	0.3833	–	–	–	–
M-Noise2	2000	237	<b>233</b>	0.0169	0.4221	–	–	–	–
M-Noise3	2000	238	<b>233</b>	0.0210	0.4031	–	–	–	–
M-Noise4	2000	238	<b>233</b>	0.0210	0.4031	–	–	–	–
M-Noise5	2000	227	<b>214</b>	0.0573	0.2558	–	–	–	–
M-Noise6	2000	201	<b>191</b>	0.0498	0.2974	–	–	–	–

Lower one is in bold. Middle panel  $J = 20, v = 0.1$  except for `Poker25kT1` and `Poker25kT2` on which  $J, v$  are chosen by validation (See the text in 7.1); Right panel the overall best. Dash means unavailable in Li (2010b, 2009a). Relative improvements ( $R$ ) and  $p$  values ( $pv$ ) are given

## 7.1 Classification errors

### 7.1.1 Summary

In Table 3 we summarize the results for all datasets. Li (2010b) reports that ABC-LogitBoost is insensitive to  $J, v$  on all datasets except for `Poker25kT1` and `Poker25kT2`. Therefore, Li reports classification errors for ABC simply with  $J = 20$  and  $v = 0.1$ , except that on `Poker25kT1` and `Poker25kT2` errors are reported by using each other’s test set as a validation set. Based on the same criteria we report AOSO in the middle panel of Table 3 where the test errors as well as the improvement relative to ABC are given. In the right panel

**Table 4** Summary of test error rates for (ABC- and AOSO-)LogitBoost and deep learning algorithms on variants of Mnist

	M-Basic (%)	M-Rotate (%)	M-Image (%)	M-Rand (%)
SVM-RBF	3.05	11.11	22.61	14.58
SVM-POLY	3.69	15.42	24.01	16.62
NNET	4.69	18.11	27.41	20.04
DBN-3	3.11	<b>10.30</b>	16.31	<b>6.73</b>
SAA-3	3.46	<b>10.30</b>	23.00	11.28
DBN-1	3.94	14.69	16.15	9.80
ABC	3.20	11.92	8.54	9.45
AOSO	<i>2.87</i>	<i>11.46</i>	<b>8.33</b>	<i>9.18</i>

*SVM-RBF* SVM with RBF kernel, *SVM-POLY* SVM with polynomial kernel, *DBN-1(-3)* Deep Belief Network with 1(3) hidden layers, *NNET* Neural Network with 1 hidden layer, *SAA-3* Stacked Auto Associator with 3 hidden layers. See (Larochelle et al 2007) for details. For each dataset, the lowest error among all algorithms is in bold, while the lower one between ABC and AOSO is in italic. Bold and italic result means both

**Table 5** Test classification errors on Mnist10k

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	2,630 <b>2,515</b>	2,600 <b>2,414</b>	2,535 <b>2,414</b>	2,522 <b>2,392</b>
$J = 6$	2,263 <b>2,133</b>	2,252 <b>2,146</b>	2,226 <b>2,146</b>	2,223 <b>2,134</b>
$J = 8$	2,159 <b>2,055</b>	2,138 <b>2,046</b>	2,120 <b>2,046</b>	2,143 <b>2,055</b>
$J = 10$	2,122 <b>2,010</b>	2,118 <b>1,980</b>	2,091 <b>1,980</b>	2,097 <b>2,014</b>
$J = 12$	2,084 <b>1,968</b>	2,090 <b>1,965</b>	2,090 <b>1,965</b>	2,095 <b>1,995</b>
$J = 14$	2,083 <b>1,945</b>	2,094 <b>1,938</b>	2,063 <b>1,938</b>	2,050 <b>1,935</b>
$J = 16$	2,111 <b>1,941</b>	2,114 <b>1,928</b>	2,097 <b>1,928</b>	2,082 <b>1,966</b>
$J = 18$	2,088 <b>1,925</b>	2,087 <b>1,916</b>	2,088 <b>1,916</b>	2,097 <b>1,920</b>
$J = 20$	2,128 <b>1,930</b>	2,112 <b>1,917</b>	2,095 <b>1,917</b>	2,102 <b>1,948</b>
$J = 24$	2,174 <b>1,901</b>	2,147 <b>1,920</b>	2,129 <b>1,920</b>	2,138 <b>1,903</b>
$J = 30$	2,235 <b>1,887</b>	2,237 <b>1,885</b>	2,221 <b>1,885</b>	2,177 <b>1,917</b>
$J = 40$	2,310 <b>1,923</b>	2,284 <b>1,890</b>	2,257 <b>1,890</b>	2,260 <b>1,912</b>
$J = 50$	2,353 <b>1,958</b>	2,359 <b>1,910</b>	2,332 <b>1,910</b>	2,341 <b>1,934</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

of Table 3 we provide the comparison for the best results achieved over all  $J-v$  combinations when the corresponding results for ABC are available in Li (2010b) or Li (2009a).

We also tested the statistical significance between AOSO and ABC. We assume the classification error rate is subject to some Binomial distribution. Let  $z$  denote the number of errors and  $n$  the number of tests, then the estimate of error rate  $\hat{p} = z/n$  and its variance is  $\hat{p}(1 - \hat{p})/n$ . Subsequently, we approximate the Binomial distribution by a Gaussian distribution and perform a hypothesis test. The  $p$  values are reported in Table 3.

**Table 6** Test classification errors on M-Image

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	5,539 <b>5,445</b>	5,480 <b>5,404</b>	5,408 <b>5,387</b>	5,430 <b>5,310</b>
$J = 6$	5,076 <b>4,909</b>	4,925 <b>4,851</b>	4,950 <b>4,801</b>	4,919 <b>4,822</b>
$J = 8$	4,756 <b>4,583</b>	4,748 <b>4,587</b>	4,678 <b>4,599</b>	4,670 <b>4,594</b>
$J = 10$	4,597 <b>4,436</b>	4,572 <b>4,445</b>	4,524 <b>4,448</b>	4,537 <b>4,467</b>
$J = 12$	4,432 <b>4,332</b>	4,455 <b>4,350</b>	4,416 <b>4,331</b>	4,389 <b>4,332</b>
$J = 14$	4,378 <b>4,264</b>	4,338 <b>4,277</b>	4,356 <b>4,240</b>	4,299 <b>4,245</b>
$J = 16$	4,317 <b>4,163</b>	4,307 <b>4,184</b>	4,279 <b>4,202</b>	4,313 <b>4,258</b>
$J = 18$	4,301 <b>4,119</b>	4,255 <b>4,176</b>	4,230 <b>4,167</b>	4,287 <b>4,188</b>
$J = 20$	4,251 <b>4,084</b>	4,231 <b>4,101</b>	4,214 <b>4,093</b>	4,268 <b>4,167</b>
$J = 24$	4,242 <b>4,049</b>	4,298 <b>4,077</b>	4,226 <b>4,067</b>	4,250 <b>4,120</b>
$J = 30$	4,351 <b>4,022</b>	4,307 <b>4,025</b>	4,311 <b>4,026</b>	4,286 <b>4,099</b>
$J = 40$	4,434 <b>4,002</b>	4,426 <b>4,033</b>	4,439 <b>4,047</b>	4,388 <b>4,090</b>
$J = 50$	4,502 <b>4,067</b>	4,534 <b>4,077</b>	4,487 <b>4,064</b>	4,479 <b>4,121</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

**Table 7** Test classification errors on Letter4k

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	2,630 <b>2,515</b>	2,600 <b>2,414</b>	2,535 <b>2,414</b>	2,522 <b>2,392</b>
$J = 6$	2,263 <b>2,133</b>	2,252 <b>2,146</b>	2,226 <b>2,146</b>	2,223 <b>2,134</b>
$J = 8$	2,159 <b>2,055</b>	2,138 <b>2,046</b>	2,120 <b>2,046</b>	2,143 <b>2,055</b>
$J = 10$	2,122 <b>2,010</b>	2,118 <b>1,980</b>	2,091 <b>1,980</b>	2,097 <b>2,014</b>
$J = 12$	2,084 <b>1,968</b>	2,090 <b>1,965</b>	2,090 <b>1,965</b>	2,095 <b>1,995</b>
$J = 14$	2,083 <b>1,945</b>	2,094 <b>1,938</b>	2,063 <b>1,938</b>	2,050 <b>1,935</b>
$J = 16$	2,111 <b>1,941</b>	2,114 <b>1,928</b>	2,097 <b>1,928</b>	2,082 <b>1,966</b>
$J = 18$	2,088 <b>1,925</b>	2,087 <b>1,916</b>	2,088 <b>1,916</b>	2,097 <b>1,920</b>
$J = 20$	2,128 <b>1,930</b>	2,112 <b>1,917</b>	2,095 <b>1,917</b>	2,102 <b>1,948</b>
$J = 24$	2,174 <b>1,901</b>	2,147 <b>1,920</b>	2,129 <b>1,920</b>	2,138 <b>1,903</b>
$J = 30$	2,235 <b>1,887</b>	2,237 <b>1,885</b>	2,221 <b>1,885</b>	2,177 <b>1,917</b>
$J = 40$	2,310 <b>1,923</b>	2,284 <b>1,890</b>	2,257 <b>1,890</b>	2,260 <b>1,912</b>
$J = 50$	2,353 <b>1,958</b>	2,359 <b>1,910</b>	2,332 <b>1,910</b>	2,341 <b>1,934</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

### 7.1.2 Comparisons with SVM and deep learning

For some problems, we note LogitBoost (both ABC and AOSO) outperforms other state-of-the-art classifier such as SVM or Deep Learning.

On dataset `Poker`, Li (2009a) reports that linear SVM works poorly (the test error rate is about 40%), while ABC-LogitBoost performs far better (i.e. <10% on `Poker25kT1` and `Poker25kT2`). AOSO-LogitBoost proposed in this paper has even lower test error than ABC-LogitBoost, see Table 3.

**Table 8** Test classification errors on Letter2k

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	2,347 <b>2,178</b>	2,299 <b>2,162</b>	2,256 <b>2,179</b>	2,231 <b>2,189</b>
$J = 6$	2,136 <b>2,049</b>	2,120 <b>2,058</b>	2,072 <b>2,021</b>	2,077 <b>2,005</b>
$J = 8$	2,080 <b>1,975</b>	2,049 <b>1,977</b>	2,035 <b>1,985</b>	2,037 <b>1,997</b>
$J = 10$	2,044 <b>1,954</b>	2,003 <b>1,935</b>	2,021 <b>1,926</b>	2,002 <b>1,936</b>
$J = 12$	2,024 <b>1,924</b>	1,992 <b>1,905</b>	2,018 <b>1,916</b>	2,018 <b>1,906</b>
$J = 14$	2,022 <b>1,912</b>	2,004 <b>1,874</b>	2,006 <b>1,888</b>	2,030 <b>1,881</b>
$J = 16$	2,024 <b>1,889</b>	2,004 <b>1,899</b>	2,005 <b>1,888</b>	1,999 <b>1,899</b>
$J = 18$	2,044 <b>1,873</b>	2,021 <b>1,882</b>	1,991 <b>1,862</b>	2,034 <b>1,894</b>
$J = 20$	2,049 <b>1,870</b>	2,021 <b>1,875</b>	2,024 <b>1,884</b>	2,034 <b>1,862</b>
$J = 24$	2,060 <b>1,854</b>	2,037 <b>1,865</b>	2,021 <b>1,860</b>	2,047 <b>1,860</b>
$J = 30$	2,078 <b>1,851</b>	2,057 <b>1,865</b>	2,041 <b>1,860</b>	2,045 <b>1,867</b>
$J = 40$	2,121 <b>1,882</b>	2,079 <b>1,870</b>	2,090 <b>1,879</b>	2,110 <b>1,861</b>
$J = 50$	2,174 <b>1,915</b>	2,155 <b>1,915</b>	2,133 <b>1,851</b>	2,150 <b>1,868</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

**Table 9** Test classification errors on Poker25kT1

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	90,323 1,02,905 <b>89,561</b>	<b>67,417</b> 71,450 69,674	<b>49,403</b> 51,226 56,702	<b>42,126</b> 42,140 47,165
$J = 6$	38,017 43,156 <b>34,460</b>	36,839 39,164 <b>32,716</b>	35,467 37,954 <b>31,714</b>	34,879 37,546 <b>31,399</b>
$J = 8$	39,220 46,076 <b>33,188</b>	37,112 40,162 <b>32,566</b>	36,407 38,422 <b>32,031</b>	35,777 37,345 <b>31,562</b>
$J = 10$	39,661 44,830 <b>34,203</b>	38,547 40,754 <b>33,251</b>	36,990 40,486 <b>32,873</b>	36,647 38,141 <b>32,316</b>
$J = 12$	41,362 48,412 <b>35,448</b>	39,221 44,886 <b>34,489</b>	37,723 42,100 <b>33,641</b>	37,345 39,798 <b>33,241</b>
$J = 14$	42,764 52,479 <b>36,495</b>	40,993 48,093 <b>37,122</b>	40,155 44,688 <b>35,465</b>	37,780 43,048 <b>34,840</b>
$J = 16$	44,386 53,363 <b>39,087</b>	43,360 51,308 <b>38,111</b>	41,952 47,831 <b>37,766</b>	40,050 46,968 <b>37,275</b>
$J = 18$	46,463 57,147 <b>42,047</b>	45,607 55,468 <b>40,906</b>	45,838 50,292 <b>41,227</b>	43,040 47,986 <b>41,859</b>
$J = 20$	49,577 62,345 <b>43,394</b>	47,901 57,677 <b>44,614</b>	45,725 53,696 <b>43,065</b>	<b>44,295</b> 49,864 46,479

In each  $J-v$  entry, the first is for ABC-MART, the second for ABC-LogitBoost and the third for AOSO-LogitBoost. Lower one is in bold

On those variants of Mnist specially designed for testing various Deep Learning algorithms, Li reported ABC’s results by simply setting  $J = 20, v = 0.1$  and concluded that ABC are comparable to or better than Deep Learning. We provide AOSO’s results in Table 4.

### 7.1.3 Detailed results

We provide one-on-one comparison between ABC and AOSO over a number of  $J-v$  combinations, as follows.

Mnist10k, M-Image, Letter4k and Letter2k: For these four datasets, classification errors are reported in Li (2010b) with every combination of  $J \in \{4, 6, 8, 10, 12, 14, 16,$

**Table 10** Test classification errors on Poker25kT2

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	<b>89,608</b> 1,02,014 90,120	<b>67,071</b> 70,886 69,921	<b>48,855</b> 50,783 57,028	<b>41,688</b> 41,551 47,447
$J = 6$	37,567 42,699 <b>34,602</b>	36,345 38,592 <b>32,982</b>	34,920 37,397 <b>31,840</b>	34,326 36,914 <b>31,645</b>
$J = 8$	38,703 45,737 <b>33,456</b>	36,586 39,648 <b>32,883</b>	35,836 37,935 <b>32,288</b>	35,129 36,731 <b>31,790</b>
$J = 10$	39,078 44,517 <b>34,429</b>	38,025 40,286 <b>33,570</b>	36,455 40,044 <b>33,125</b>	36,076 37,504 <b>32,547</b>
$J = 12$	40,834 47,948 <b>35,725</b>	38,657 44,602 <b>34,773</b>	37,203 41,582 <b>33,930</b>	36,781 39,378 <b>33,417</b>
$J = 14$	42,348 52,063 <b>36,666</b>	40,363 47,642 <b>37,371</b>	39,613 44,296 <b>35,711</b>	37,243 42,720 <b>35,021</b>
$J = 16$	44,067 52,937 <b>39,238</b>	42,973 50,842 <b>38,208</b>	41,485 47,578 <b>37,851</b>	39,446 46,635 <b>37,454</b>
$J = 18$	46,050 56,803 <b>42,233</b>	45,133 55,166 <b>40,978</b>	45,308 49,956 <b>41,252</b>	42,485 47,707 <b>42,047</b>
$J = 20$	49,046 61,980 <b>43,506</b>	47,430 57,383 <b>44,754</b>	45,390 53,364 <b>43,317</b>	<b>43,888</b> 49,506 46,666

In each  $J-v$  entry, the first is for ABC-MART, the second for LogitBoost and the third for AOSO-LogitBoost. Lower one is in bold

**Table 11** Test classification errors on Letter

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	<b>125</b> 133	<b>121</b> 134	122 <b>121</b>	119 <b>114</b>
$J = 6$	112 <b>103</b>	107 <b>104</b>	<b>101</b> 106	<b>102</b> 109
$J = 8$	104 <b>98</b>	102 <b>100</b>	<b>93</b> 96	<b>95</b> 99
$J = 10$	101 <b>94</b>	100 <b>93</b>	<b>96</b> 98	<b>93</b> 88
$J = 12$	<b>96</b> 88	100 <b>91</b>	95 <b>92</b>	95 <b>92</b>
$J = 14$	<b>96</b> 88	98 <b>89</b>	<b>94</b> 88	89 <b>88</b>
$J = 16$	97 <b>88</b>	94 <b>93</b>	93 <b>90</b>	<b>95</b> 96
$J = 18$	95 <b>91</b>	<b>92</b> 95	96 <b>91</b>	93 <b>91</b>
$J = 20$	95 <b>90</b>	97 <b>89</b>	93 <b>87</b>	89 <b>92</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

**Table 12** Test classification errors on Pendigits

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	92 <b>86</b>	93 <b>85</b>	90 <b>84</b>	92 <b>85</b>
$J = 6$	98 <b>86</b>	97 <b>85</b>	96 <b>84</b>	93 <b>85</b>
$J = 8$	97 <b>84</b>	94 <b>84</b>	95 <b>84</b>	93 <b>86</b>
$J = 10$	100 <b>84</b>	98 <b>83</b>	97 <b>84</b>	97 <b>83</b>
$J = 12$	98 <b>84</b>	98 <b>82</b>	98 <b>83</b>	98 <b>83</b>
$J = 14$	100 <b>84</b>	101 <b>81</b>	99 <b>84</b>	98 <b>83</b>
$J = 16$	100 <b>83</b>	97 <b>84</b>	98 <b>82</b>	96 <b>82</b>
$J = 18$	102 <b>88</b>	97 <b>83</b>	99 <b>83</b>	97 <b>81</b>
$J = 20$	106 <b>83</b>	102 <b>82</b>	100 <b>83</b>	100 <b>83</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

18, 20, 24, 30, 40, 50} and  $v \in \{0.04, 0.06, 0.08, 0.1\}$ . The comparison with AOSO-LogitBoost is listed in Tables 5, 6, 7 and 8.

**Table 13** Test classification errors on Zipcode

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	<b>111</b> 104	108 <b>105</b>	114 <b>104</b>	107 <b>101</b>
$J = 6$	101 <b>99</b>	102 <b>100</b>	<b>98</b> 106	<b>99</b> 104
$J = 8$	<b>99 99</b>	<b>95</b> 98	<b>96</b> 102	98 <b>95</b>
$J = 10$	<b>97</b> 98	<b>94</b> 99	97 <b>96</b>	<b>94</b> 102
$J = 12$	<b>98</b> 99	98 <b>97</b>	<b>99</b> 100	<b>93</b> 100
$J = 14$	100 <b>96</b>	99 <b>98</b>	<b>97</b> 100	<b>92</b> 94
$J = 16$	98 <b>97</b>	<b>95</b> 97	99 <b>96</b>	<b>98</b> 103
$J = 18$	<b>96 96</b>	<b>98</b> 100	<b>101</b> 104	<b>8</b> 100
$J = 20$	<b>97</b> 98	<b>96</b> 97	100 <b>95</b>	<b>96</b> 99

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

**Table 14** Test classification errors on Isolet

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	- 56	55 <b>54</b>	- 56	<b>55</b> 57
$J = 6$	- 55	59 <b>56</b>	- 54	58 <b>54</b>
$J = 8$	- 54	57 <b>54</b>	- 53	60 <b>55</b>
$J = 10$	- 54	61 <b>54</b>	- 55	62 <b>56</b>
$J = 12$	- 52	63 <b>54</b>	- 54	64 <b>50</b>
$J = 14$	- 48	65 <b>51</b>	- 54	60 <b>55</b>
$J = 16$	- 55	64 <b>57</b>	- 57	62 <b>58</b>
$J = 18$	- 55	67 <b>57</b>	- 53	62 <b>57</b>
$J = 20$	- 51	63 <b>57</b>	- 56	65 <b>55</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold. Dash means unavailable

**Table 15** Test classification errors on Optdigits

	$v = 0.04$	$v = 0.06$	$v = 0.08$	$v = 0.1$
$J = 4$	<b>41</b> 43	42 <b>41</b>	<b>40 40</b>	41 <b>39</b>
$J = 6$	43 <b>40</b>	45 <b>39</b>	44 <b>39</b>	<b>38</b> 39
$J = 8$	44 <b>38</b>	44 <b>37</b>	45 <b>41</b>	45 <b>41</b>
$J = 10$	50 <b>37</b>	50 <b>38</b>	46 <b>36</b>	42 <b>38</b>
$J = 12$	50 <b>37</b>	48 <b>38</b>	47 <b>35</b>	46 <b>41</b>
$J = 14$	48 <b>38</b>	46 <b>41</b>	51 <b>40</b>	48 <b>40</b>
$J = 16$	54 <b>39</b>	51 <b>37</b>	49 <b>38</b>	46 <b>38</b>
$J = 18$	54 <b>41</b>	55 <b>38</b>	53 <b>40</b>	51 <b>41</b>
$J = 20$	61 <b>37</b>	56 <b>34</b>	55 <b>36</b>	55 <b>38</b>

In each  $J-v$  entry, the first is for ABC-LogitBoost and the second for AOSO-LogitBoost. Lower one is in bold

Poker25kT1 and Poker25kT2: These are the only two datasets on which ABC-MART<sup>4</sup> outperforms ABC-LogitBoost in experiments by Li (2010b). Thus we cite the results for both ABC-MART and ABC-LogitBoost in Tables 9 and 10, with

<sup>4</sup> To put it simple, one gets MART (Multiple Additive Regression Trees Friedman 2001) by replacing the Hessian (23) with identity matrix. ABC-MART is its adaptive base class version (Li 2009).



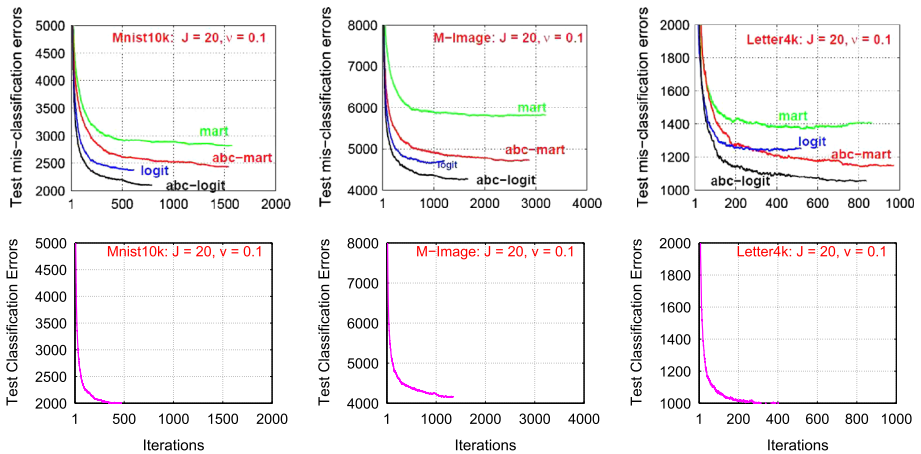
**Table 16** #trees added when convergence on selected datasets.  $R$  stands for the ratio AOSO/ABC

	Mnist10k	M-Rand	M-Image	Letter15k	Letter4k	Letter2k
ABC	7,092	15,255	14,958	45,000	20,900	13,275
$R$	0.7689	0.7763	0.8101	0.5512	0.5587	0.5424

**Table 17** #trees added when convergence on Mnist10k for a number of  $J-v$  combinations

	$v = 0.04$	$v = 0.06$	$v = 0.1$
$J = 4$	90,000 1.0	90,000 1.0	90,000 1.0
$J = 6$	90,000 0.7740	63,531 0.7249	38,223 0.7175
$J = 8$	55,989 0.7962	38,223 0.7788	22,482 0.7915
$J = 10$	39,780 0.8103	27,135 0.7973	16,227 0.8000
$J = 12$	31,653 0.8109	20,997 0.8074	12,501 0.8269
$J = 14$	26,694 0.7854	17,397 0.8047	10,449 0.8160
$J = 16$	22,671 0.7832	11,704 1.0290	8,910 0.8063
$J = 18$	19,602 0.7805	13,104 0.7888	7,803 0.7933
$J = 20$	17,910 0.7706	11,970 0.7683	7,092 0.7689
$J = 24$	14,895 0.7514	9,999 0.7567	6,012 0.7596
$J = 30$	12,168 0.7333	8,028 0.7272	4,761 0.7524
$J = 40$	9,846 0.6750	6,498 0.6853	3,870 0.6917
$J = 50$	8,505 0.6420	5,571 0.6448	3,348 0.6589

For each  $J-v$  entry, the first number is for ABC, the second for the ratio AOSO/ABC



**Fig. 5** Errors versus iterations on selected datasets and parameters. *Top row* ABC [copied from Li (2010b)]; *Bottom row* AOSO (*horizontal axis scaled to compensate the  $K - 1$  factor*)

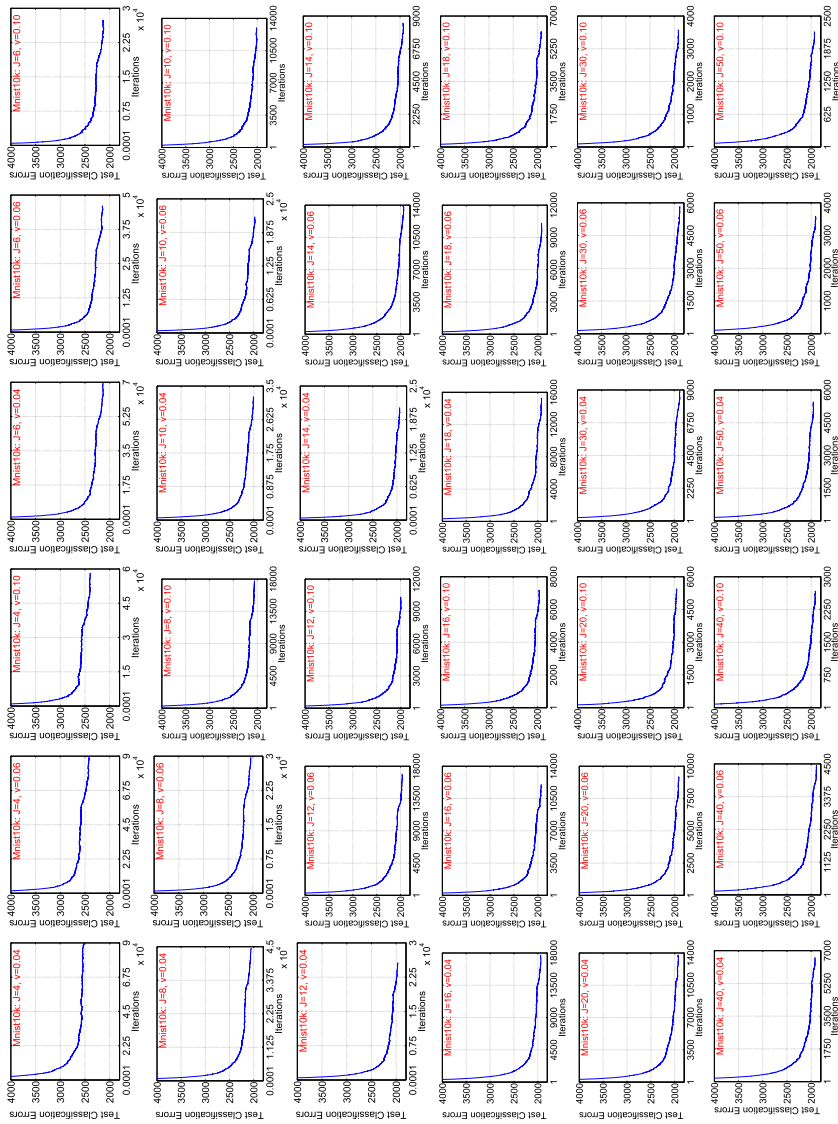
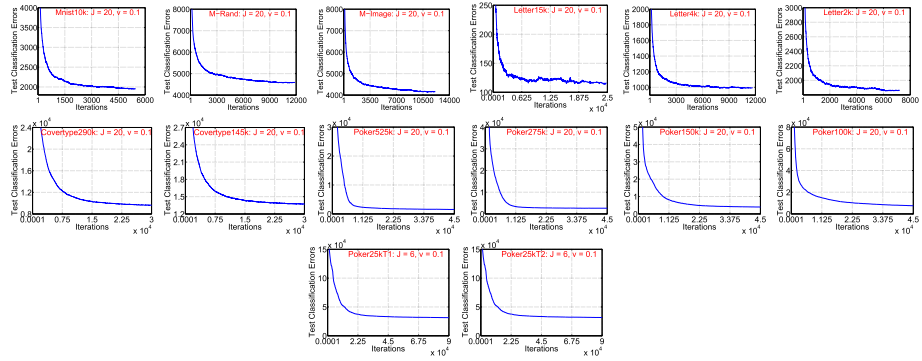


Fig. 6 Errors versus iterations on  $\text{Minst10k}_k, J \in \{4, 6, 8, 10, 12, 14, 16, 18, 24, 30, 40, 50\}$



**Fig. 7** Errors versus iterations on selected datasets and parameters

$J \in \{4, 6, 8, 10, 12, 14, 16, 18, 20\}$  and  $v \in \{0.04, 0.06, 0.08, 0.1\}$ . The comparison with AOSO-LogitBoost is also listed. Unlike on previous datasets, AOSO-LogitBoost is a bit sensitive to parameters, which is also observed for ABC-MART and ABC-LogitBoost by Li (2010b).

**Letter, Pendigits, Zipcode, Isolet and Optdigits:** For these five datasets, classification errors are reported by Li (2010b) with every combination of  $J \in \{4, 6, 8, 10, 12, 14, 16, 18, 20\}$  and  $v \in \{0.04, 0.06, 0.08, 0.1\}$  (except that  $v \in \{0.06, 0.1\}$  for Isolet). The comparison with AOSO-LogitBoost is listed in Tables 11, 12, 13, 14 and 15.

### 7.2 Convergence rate

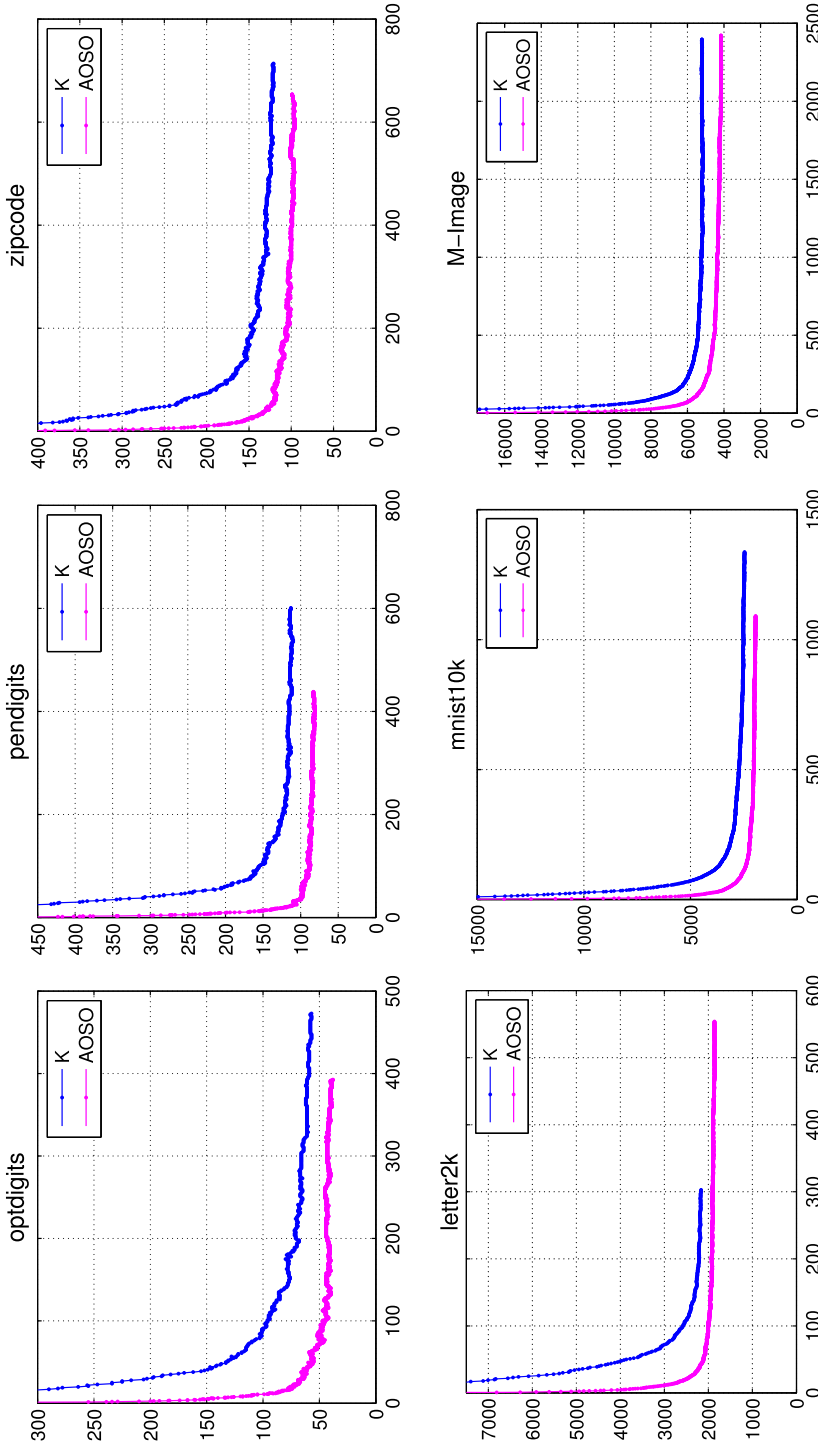
Recall that we stop the boosting procedure if either the maximum number of iterations is reached or the training loss is small (i.e. the loss  $(1) \leq 10^{-16}$ ). The fewer trees added when boosting stops, the faster the convergence and the lower the time cost for either training or testing. We compare AOSO with ABC in terms of the number of trees added when boosting stops for the results of ABC available in Li (2010b, 2009a). Note that simply comparing number of boosting iterations is unfair to AOSO, since at each iteration only one tree is added in AOSO and  $K - 1$  in ABC.

Results are shown in Table 16 and Table 17. Except for when  $J-v$  is too small, or particularly difficult datasets where both ABC and AOSO reach maximum iterations, we found that trees needed in AOSO are typically only 50–80% of those in ABC.

Figure 5 shows plots for test classification error vs. iterations in both ABC and AOSO and show that AOSO’s test error decreases faster. More plots for AOSO are given in Figs. 6 and 7. We only choose the datasets and parameters on which the plots of ABC are provided in Li (2010b).

### 7.3 Comparison between K-LogitBoost and AOSO-LogitBoost

In Fig. 8 we provide comparison between K-LogitBoost and AOSO-LogitBoost on several datasets. The parameters for both algorithms are  $J = 20, v = 0.1$ . As can be seen, K-LogitBoost converges slower than AOSO-LogitBoost, confirming our arguments in Sect. 6.3.



**Fig. 8** Comparison between K-LogitBoost and AOSO-LogitBoost on datasets *Optdigits*, *Pendigits*, *Zipcode*, *Letter2k*, *Mnist10k*, *M-Image* with parameters  $J = 20$  and  $\nu = 0.1$ : how the testing error decreases with iterations

## 8 Conclusions

We present an improved LogitBoost, namely AOSO-LogitBoost, for multi-class classification. Compared with ABC-LogitBoost, our experiments suggest that our adaptive class pair selection technique results in lower classification error and faster convergence rates.

**Acknowledgments** We appreciate Ping Li’s inspiring discussion and generous encouragement. We also thank the MLJ anonymous reviewers and the action editor for their valuable comments. This work was supported by National Natural Science Foundation of China (61020106004, 61225008) and an Australian Research Council Discovery Early Career Research Award (DE130101605). NICTA is funded by the Australian Government through the Department of Communications and the Australian Research Council through the ICT Centre of Excellence Program.

### Appendix 1: Proof for property 2

First, we give without proof the lemma below that is well known in Matrix Analysis:

**Lemma 1** *Let  $\mathbf{x}, \mathbf{y}$  be  $K$  dimensional vector,  $\mathbf{A} = \mathbf{x}\mathbf{y}^T$  be their out-product. Then the  $K$  eigen values for  $\mathbf{A}$  are  $\{\mathbf{y}^T\mathbf{x}, \underbrace{0, \dots, 0}_{K-1}\}$ .*

For positive semi definiteness

Let  $\hat{\mathbf{Q}} = \text{diag}(\sqrt{p_{i,1}}, \dots, \sqrt{p_{i,K}}), \mathbf{q} = (\sqrt{p_{i,1}}, \dots, \sqrt{p_{i,K}})^T$ . For any  $\mathbf{x} \in \mathbb{R}^K$ , let  $\mathbf{y} = \hat{\mathbf{Q}}\mathbf{x}$ . Then

$$\begin{aligned} \mathbf{x}^T \mathbf{H}_i \mathbf{x} &= \underbrace{\mathbf{x}^T \hat{\mathbf{P}} \mathbf{x}}_{\mathbf{y}^T \mathbf{y}} - \underbrace{\mathbf{x}^T \mathbf{p} \mathbf{p}^T \mathbf{x}}_{\mathbf{y}^T \mathbf{q} \mathbf{q}^T \mathbf{y}} \\ &= \mathbf{y}^T (\mathbf{I} - \mathbf{q} \mathbf{q}^T) \mathbf{y}. \end{aligned} \tag{45}$$

From Lemma 1, we have that the eigen values for  $\mathbf{q} \mathbf{q}^T$  are  $\{1, \underbrace{0, \dots, 0}_{K-1}\}$  by noting  $\mathbf{q}^T \mathbf{q} = 1$ .

Then the eigen values for its first order matrix polynomial  $\mathbf{I} - \mathbf{q} \mathbf{q}^T$  is  $\{0, \underbrace{1, \dots, 1}_{K-1}\}$ , i.e. all the eigen values are non-negative. So it follows that the quadratic  $\mathbf{y}^T (\mathbf{I} - \mathbf{q} \mathbf{q}^T) \mathbf{y} \geq 0$ .

For the rank

First, assume the number of non-zero elements in  $\mathbf{p}_i$  equals the number of classes:  $\kappa = K$ . We write  $\mathbf{H}_i$  as

$$\mathbf{H}_i = \hat{\mathbf{P}} (\mathbf{I} - \mathbf{1} \mathbf{p}_i^T) \tag{46}$$

where  $\hat{\mathbf{P}} = \text{diag}(p_{i,1}, \dots, p_{i,K})$  is full rank. So we have

$$\text{rank}(\mathbf{H}_i) = \text{rank}(\mathbf{I} - \mathbf{1} \mathbf{p}_i^T) \tag{47}$$

Noting Lemma 1, the eigenvalues for  $A = \mathbf{1}p_i^T$  are  $\{1, \underbrace{0, \dots, 0}_{K-1}\}$ . It follows that the eigenvalues for its matrix (first order) polynomial  $I - A$  is  $\{0, \underbrace{1, \dots, 1}_{K-1}\}$ , which proves that

$$\text{rank}(\hat{P} - p_i p_i^T) = K - 1.$$

Then we assume  $\kappa < K$ . We collect the non-zero elements in  $p_{i,k}$  to the upper left corner in matrix  $H_i$  and drop the corresponding zero elements by permutation. Repeat the procedure for  $\kappa = K$  and we obtain that  $\text{rank}(H_i) = \kappa - 1$ . Note that since  $p_i$  is a probability,  $\kappa \geq 1$ .

For eigenvector

It's straightforward by noting that  $p_i$  is a probability and thereby satisfies the sum-to-one constraint.

### Appendix 2: Proof for Theorem 1 and 2

Conventions for matrix calculus

Before the proof, we introduce our conventions for matrix calculus, which are borrowed from Magnus and Neudecker (2007).

Let  $\phi(x) : \mathbb{R}^K \rightarrow \mathbb{R}$  be a scalar function defined for  $K$ -dim vector  $x$ , its gradient is a  $K$ -dim vector  $\nabla_x \phi(x)$ . In particular, we use the symbol  $D$  to represent the derivative, which is no more than the transpose of gradient  $D_x \phi(x) = (\nabla_x \phi(x))^T$  and thus a  $K$ -dim row vector. Without confusion from context, we sometimes abbreviate them as  $\nabla \phi$  and  $D\phi$ .

For the vector function  $\xi(x) : \mathbb{R}^K \rightarrow \mathbb{R}^K$ , the matrix  $D\xi \in \mathbb{R}^{K \times K}$  is precisely the Jacobian matrix. Therefore, the Hessian matrix of scalar function  $\phi()$  can be expressed by the Jacobian of its gradient  $\nabla^2 \phi = D(\nabla \phi)$ .

For Theorem 1

From first equation of (40) we have

$$DH(p) + F^T = -\lambda \mathbf{1}. \tag{48}$$

By applying derivative w.r.t  $F$  on both sides of the second equation of (40) we have

$$\mathbf{1}^T (Dp) = 0. \tag{49}$$

By noting both (48) and (49), it follows that

$$D\ell(F) = (DH(p))(Dp) + (p - \eta)^T (DF) + F^T (Dp) \tag{50}$$

$$D\ell(F) = -\lambda \mathbf{1}^T (Dp) + (p - \eta)^T \tag{51}$$

$$D\ell(F) = (p - \eta)^T, \tag{52}$$

which completes the proof.

For Theorem 2

The  $(K + 1) \times (K + 1)$  equations (40) defines an implicit function  $\mathbf{p} = \mathbf{p}(\mathbf{F}) \in \mathbb{R}^K$ ,  $\lambda = \lambda(\mathbf{F}) \in \mathbb{R}$ . We apply to it the derivative w.r.t.  $\mathbf{F}$  on both sides and have

$$\begin{cases} (\mathbf{D}\mathbf{F}) + (\nabla^2\mathbf{H}(\mathbf{p}))(\mathbf{D}\mathbf{p}) + \mathbf{1}(\mathbf{D}\lambda) = 0 \\ \mathbf{1}^T(\mathbf{D}\mathbf{p}) = 0 \end{cases} \tag{53}$$

Noting  $\mathbf{D}\mathbf{F} = \mathbf{I}$  as well as the formula for inverse of block matrix, we obtain from (53) the following equation:

$$\begin{aligned} \begin{bmatrix} \mathbf{D}\mathbf{p} \\ \mathbf{D}\lambda \end{bmatrix} &= - \begin{bmatrix} \nabla^2\mathbf{H}(\mathbf{p}) & \mathbf{1} \\ \mathbf{1}^T & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{I} \\ 0 \end{bmatrix} \\ &= - \begin{bmatrix} A^{-1} + A^{-1}\mathbf{1}(-\mathbf{1}^T A^{-1}\mathbf{1})\mathbf{1}^T A^{-1} & \times \\ \times & \times \end{bmatrix} \begin{bmatrix} \mathbf{I} \\ 0 \end{bmatrix} \end{aligned} \tag{54}$$

where the abbreviation  $A = \nabla^2\mathbf{H}(\mathbf{p})$  and the symbol “ $\times$ ” means formula we don’t care. It thus follows from above formula

$$\mathbf{D}\mathbf{p} = -(A^{-1} + A^{-1}\mathbf{1}(-\mathbf{1}^T A^{-1}\mathbf{1})\mathbf{1}^T A^{-1}). \tag{55}$$

In the meanwhile, by applying derivative to  $\nabla\ell(\mathbf{F})$  in Theorem 1 we precisely have  $\nabla^2\ell(\mathbf{F}) = (\mathbf{D}\mathbf{p})^T$ . This completes the proof.

### References

Bertsekas, D. P. (1982). *Constrained optimization and Lagrange multiplier methods*. Boston: Academic Press.

Bottou, L., & Lin, C. J. (2007). Support vector machine solvers. In L. Bottou, O. Chapelle, D. DeCoste, & J. Weston (Eds.), *Large scale Kernel machines* (pp. 301–320). Cambridge: MIT Press. <http://leon.bottou.org/papers/bottou-lin-2006>

Freund, Y., & Schapire, R. (1995). A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational learning theory* (pp. 23–37). New York: Springer.

Friedman, J. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5), 1189–1232.

Friedman, J., Hastie, T., & Tibshirani, R. (1998). Additive logistic regression: A statistical view of boosting. *Annals of Statistics*, 28(2), 337–407.

Jaynes, E. (1957). Information theory and statistical mechanics. *The Physical Review*, 106(4), 620–630.

Kégl, B., & Busa-Fekete, R. (2009). Boosting products of base classifiers. In *Proceedings of the 26th Annual International Conference on Machine Learning* (pp. 497–504). New York: ACM.

Kivinen, J., & Warmuth, M. K. (1999). Boosting as entropy projection. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory* (pp. 134–144). New York: ACM.

Lafferty, J. (1999). Additive models, boosting, and inference for generalized divergences. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory* (pp. 125–133).

Larochelle, H., Erhan, D., Courville, A., Bergstra, J., & Bengio, Y. (2007). An empirical evaluation of deep architectures on problems with many factors of variation. In *Proceedings of the 24th International Conference on Machine learning* (pp. 473–480). New York: ACM.

Li, P. (2008). Adaptive base class boost for multi-class classification. *Arxiv preprint arXiv:08111250*.

Li, P. (2009a). Abc-boost: Adaptive base class boost for multi-class classification. In *Proceedings of the 26th Annual International Conference on Machine Learning* (pp. 625–632). New York: ACM.

Li, P. (2009b). Abc-logitboost for multi-class classification. *Arxiv preprint arXiv:09084144*.

Li, P. (2010a). An empirical evaluation of four algorithms for multi-class classification: Mart, abc-mart, robust logitboost, and abc-logitboost. *Arxiv preprint arXiv:10011020*.

Li, P. (2010b). Robust logitboost and adaptive base class (abc) logitboost. In *Conference on Uncertainty in Artificial Intelligence*.

- Magnus, J. R., & Neudecker, H. (2007). *Matrix differential calculus with applications in statistics and econometrics* (3rd ed.). New York: Wiley.
- Masnadi-Shirazi, H., & Vasconcelos, N. (2010). Risk minimization, probability elicitation, and cost-sensitive svms. In *Proceedings of the International Conference on Machine Learning* (pp. 204–213).
- Reid, M. D., & Williamson, R. C. (2010). Composite binary losses. *The Journal of Machine Learning Research*, *11*, 2387–2422.
- Schapire, R., & Singer, Y. (1999). Improved boosting algorithms using confidence-rated predictions. *Machine learning*, *37*(3), 297–336.
- Shen, C., & Hao, Z. (2011). A direct formulation for totally corrective multi-class boosting. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 2585–2592).
- Shen, C., & Li, H. (2010). On the dual formulation of boosting algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *32*(12), 2216–2231.
- Zou, H., Zhu, J., & Hastie, T. (2008). New multiclass boosting algorithms based on multiclass fisher-consistent losses. *The Annals of Applied Statistics*, *2*(4), 1290–1306.