

Compact Multi-Class Boosted Trees

Natalia Ponomareva, Thomas Colthurst, Gilbert Hendry, Salem Haykal, Soroush Radpour
Google, Inc.

tfbt-public@google.com

Abstract—Gradient boosted decision trees are a popular machine learning technique, in part because of their ability to give good accuracy with small models. We describe two extensions to the standard tree boosting algorithm designed to increase this advantage. The first improvement extends the boosting formalism from scalar-valued trees to vector-valued trees. This allows individual trees to be used as multiclass classifiers, rather than requiring one tree per class, and drastically reduces the model size required for multiclass problems. We also show that some other popular vector-valued gradient boosted trees modifications fit into this formulation and can be easily obtained in our implementation. The second extension, layer-by-layer boosting, takes smaller steps in function space, which is empirically shown to lead to a faster convergence and to a more compact ensemble. We have added both improvements to the open-source TensorFlow Boosted trees (TFBT) package, and we demonstrate their efficacy on a variety of multiclass datasets. We expect these extensions will be of particular interest to boosted tree applications that require small models, such as embedded devices, applications requiring fast inference, or applications desiring more interpretable models.

Keywords—Multiclass gradient boosting, TensorFlow, large-scale machine learning, tree-based methods, ensemble methods

I. INTRODUCTION

There are many reasons to use and study gradient boosted decision trees (or boosted trees for short). They have outstanding accuracy, as demonstrated by their winning performance in numerous surveys of machine learning models (such as [1]) and Kaggle competitions [2]. They are easy to use, as input features do not need to be whitened or otherwise normalized. They are flexible: by supporting custom loss functions, they can attack arbitrary classification and regression problems, including ranking or regression tasks. They are backed by solid theory, which allows them to be viewed as doing gradient descent in the space of functions by taking steps in the form of decision trees [3]. And last but not least, they often produce compact models, with fewer parameters for the same accuracy when compared to say random forests. Those compact models in turn lead to faster inference speeds, less memory consumption (important for embedded devices and cellphones), and better interpretability.

This paper describes two extensions we’ve made to TensorFlow Boosted Trees (TFBT) [4] that are designed to increase the model compactness. As its name suggests, TFBT is built on top of TensorFlow; it is open-source and available on github under

`tensorflow/contrib/boosted_trees`. The first extension, described in section III, extends the usual boosting theory to apply to vector-valued outputs. Using the derived update equations, we can attack multi-class classification or multidimensional regression problems directly with trees storing vectors in the leaves, rather than the 1-vs-rest or 1-vs-1 approaches that are commonly used [5]. Significantly fewer trees are required for good performance with this approach, yielding corresponding reductions in model complexity.

The second extension, layer-by-layer boosting, described in section V, can be thought of as taking smaller steps in function space when doing the gradient descent; steps that correspond to tree layers (i.e., all the nodes of equal depth) rather than entire trees. These smaller steps let us converge faster with fewer trees, especially with custom loss functions for which a second-order Taylor expansion is inexact.

Finally, in section VII, we evaluate these two extensions. On ensemble sizes up to 100 trees, we demonstrate that vector-valued trees lead to much faster convergence and smaller ensembles on multiclass datasets, and that the combination of layer-by-layer boosting and vector-valued trees often produces significant performance improvements.

II. RELATED WORK

Many popular algorithms are inherently binary, for example, the extremely popular *AdaBoost* [6]. Historically if a multi-class dataset was used, this was addressed by decomposing the multiclass problem into a set of binary subproblems, namely either as 1-vs-rest (sometimes referred to as 1-vs-all, as well as tree-per-class) or 1-vs-1 (also referred to as “all-pairs” in some sources).

Even though binary reformulations achieve good performance on some datasets and are easy to implement, these approaches have a number of drawbacks. Firstly, the number of classifiers required is at least C (the number of classes). Secondly, as Friedman et al. point out, even if decision boundaries between classes are simple, the decision boundaries to be learnt when the problem is reformulated into binary problems can become hard, thus making these boundaries difficult to approximate [7]. Additionally, common 1-vs-rest can make the subproblems unbalanced, complicating the learning further [8], since some losses like log-loss are susceptible to class imbalances. Furthermore, when applied to boosting, theoretical guarantees state that to achieve good

performance, each of the weak learners must achieve an accuracy of at least 0.5 [9]. This is comparable to random guess in case of binary classification with balanced classes, but for a multiclass problem with balanced classes, random guess would result only in $1/C$ accuracy [9].

Several attempts to tackle multiclass problems directly, by optimizing a loss that tries to classify all of the classes correctly at the same time, were made. One of the first implementation was *LogitBoost* [7] - a generalization of *AdaBoost*. Friedman et al. showed that it produces results superior to those achieved by 1-vs-rest models on a simulated example with complicated intra-classes boundaries. *LogitBoost* takes the multinomial logistic regression (cross-entropy loss) and decomposes it in a standard way of running C independent binary logistic regressions, in which one label m is chosen as “pivot” or base class, and the other $C - 1$ labels are separately regressed against this pivot label. So for each class j , given a base class m , one can calculate a quasi-Newton update to improve the loss of class j vs class m . Now Friedman et. al note that the choice of pivot m is arbitrary, so they build a tree for j by approximating an average of steps over all choices of base classes. *LogitBoost* still works however by fitting C trees during each boosting iteration, where each tree is derived to minimize the overall loss with regard to the j -th class - so essentially an ensemble of trees is built to predict scores for each class. An alternative view would be that Friedman et al. approximates the full Hessian matrix of the expanded (up to a second derivative) loss via diagonal approximation [10]. Also note that at each boosting step, the predictions are fixed - a new tree built during an iteration i does not affect the other trees built during the same boosting iteration.

MART [11] is another popular variant of boosting with trees, which differs from *LogitBoost* in that it uses only the first order gradient for finding the splits. Both *MART* and *LogitBoost* use second-order information to derive the weights on terminal nodes [12]. For multiclass, *MART* adopts the same approach as *LogitBoost*, building C trees during each boosting iteration. An in-depth comparison between *MART* and *LogitBoost* is available in [12].

An extension to *LogitBoost* approach of multiclass handling is proposed in [13], [14]. Ping Li names it *ABC-Boost* (with two implementations *ABC-MART* and *ABC-LogitBoost*, depending on whether the second order gradient was used for finding tree splits or not). He notes that for multiclass classification, for each point, the sum of scores for all classes can be required to equal a constant (due to the fact that adding a constant to the scores of all classes does not change the overall winning class and class probabilities when common softmax is applied to the scores). This requirement on the sum of the scores does not necessarily hold true for the *LogitBoost* algorithm. If the constraint is enforced, the scores for only $C - 1$ classes need to be calculated and the derivatives can be redefined.

Additionally, instead of averaging over a choice of all “pivot” classes during each iteration, *ABC-boost* adaptively chooses the best base class (by first considering all base classes and then choosing the one that maximally reduces the training loss). Experiments show improvements over the baseline algorithms, however a large number of boosting iterations (for example, around 4,000 on the covtype dataset) is still required to achieve good performance [14]. Additionally experiments demonstrate that in some cases, *ABC-MART* requires even more training iterations than the original *MART* algorithm [14]. Nevertheless, it was shown that *ABC-MART* and *ABC-LogitBoost* improve over their respective original algorithms [14], as well as that *ABC-LogitBoost* outperforms *ABC-MART* on most datasets. An alternative view of *ABC-LogitBoost* is presented in [10]: Peng et al. show that the difference between *LogitBoost* and *ABC-LogitBoost* is essentially in the Hessian matrix approximation. In *ABC-LogitBoost*, the Hessian is approximated by first choosing a dimension (“base” or “pivot” class) and then again approximating the remaining $(C - 1) \times (C - 1)$ matrix using diagonal approximation, where in *LogitBoost* the full $C \times C$ Hessian is approximated via diagonal. The results suggest that the better approximation in *ABC-LogitBoost* results in improved performance [10]. Note that *ABC-Boost* and its modifications still build $C - 1$ trees at every iteration. *AOSO-LogitBoost* is an adaptation of *LogitBoost* to multiclass problems where only one vector-leaf tree is built during an iteration. It essentially approximates the Hessian in a block-diagonal fashion, at each iteration selecting only scores for 2 classes to be updated, leading to a one-vs-one classifier [10]. *AOSO-LogitBoost* is similar to *ABC-LogitBoost* in that it also works only on a $(C - 1) \times (C - 1)$ matrix, where the final dimension, the “pivot”, is fixed. The pair of classes is chosen based on the magnitude of the derivatives (i.e. choose the class you do the worst on). Experiments show that *AOSO-LogitBoost* outperforms *ABC-LogitBoost* and, not surprisingly, requires a smaller number of trees to reach convergence.

SAMME [9] is another attempt at tackling multiclass with boosting: it is a modification of *AdaBoost* exponential loss (as opposed to multinomial logistic loss in *LogitBoost*). *SAMME* works on vector encoded labels, building only one tree per iteration. The class encoding scheme for the labels is as follows: instead of one-hot vector encoding, the authors use an encoding where 1 is put in position of the real class, and $-1/(C-1)$ is set for all other classes positions. With this label encoding scheme, a new multi-class exponential loss is introduced. Optimizing this loss with a constraint that scores should sum to 0, *SAMME* derives a closed-form solution by using a Lagrange multiplier and using only first-order derivatives (vs Newton update in *LogitBoost*) with respect to the prediction so far and a Lagrange multiplier. *SAMME* also shows that in this framework, weak classifiers are required to do better than C -random guessing.

GD-MCBoost is another attempt at extending AdaBoost to a multiclass version [15], that also uses gradient descent step but a different encoding scheme for the labels. *GD-MCBoost* has solid theory behind its encoding scheme and was shown to produce larger-margin models than *SAMME*.

GAMBLE [16] also extends *AdaBoost* to multiclass classification. Authors use the same multiclass exponential loss and label encoding scheme as in *SAMME* and also produces one vector-leaf tree per boosting iteration. The difference between *GAMBLE* and *SAMME* lies in derivations for loss optimization - *GAMBLE* uses quasi-Newton step (with second order derivative). The original paper [16] provides evidence that *GAMBLE* performs better on a number of datasets than both 1-vs-rest and *SAMME*, which might be due to the fact that second order derivatives are used.

Many more variants of multiclass vector-form boosting exists. However, most successful variations differ only in the following choices:

- Choice of loss (cross-entropy/multinomial logistic, multiclass exponential etc.)
- Choice of loss optimization (to use or not to use second order information,
- If second order gradient is used, the choice of approximation of the Hessian.
- The choice of enforcing the constraint on scores (via a penalty term added to the Loss, or not enforced, or enforced via softmax over the scores etc.)
- And finally, choice of label encoding scheme (one hot, 1 vs -1/(C-1), etc.)

Many papers show that models that handle a multiclass loss directly and build trees with vector-leaves result in better convergence rates and better performing models. However, the adoption of these methods is unfortunately lacking. Many popular libraries like XGBoost [2], Scikit-learn [17], R GBM [18], LightGBM [19] support multiclass only as 1-vs-rest, whereas Spark MLlib [20] (at the time of writing this paper) does not support multiclass at all. This might be due to the fact that each modification that tackles multiclass problems uses different losses, weight update scheme and re-derives gradients and Hessians (if used), complicating implementations.

III. MULTICLASS HANDLING

In the conventional formulation that is used by many libraries like XGBoost [2], gradient boosted trees store only scalar values in their leaves. In order to handle vector regression or multiclass classification problems, multiple scalar-leaved trees must be used. In this section, we show how the usual derivation of the update formula for gradient boosted trees effortlessly generalizes to handling vector values which can handle vector regression and multiclass problems directly. For clarity, we follow the derivation given in [2].

Assuming m is the number of features, n is the number of instances and C is the number of classes, our model maps the input $x_i \in \mathcal{R}^m, i \in 1..n$ to the output $\hat{y}_i \in \mathcal{R}^C$ using the sum of K trees:

$$\hat{y}_i = \varphi(x_i) = \sum_{k=1}^K f_k(x_i)$$

Assuming T is the number of leaves in a tree $f(x)$, we represent each tree $f(x)$ as the combination of a structure function $q(x) : \mathcal{R}^m \rightarrow 1..T$, mapping an instance to the tree leaf where it ends up, and a set of leaf weights $\{w_j \in \mathcal{R}^C \mid j \in 1..T\}$ so that $f(x) = w_{q(x)}$. We seek to minimize the regularized objective function

$$L(\varphi) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

which is the sum of the non-regularized loss l and the ensemble regularizer Ω . In what follows, we will specifically consider an Ω regularizer of the form

$$\Omega(f) = \alpha T + \frac{1}{2} \lambda \sum_{j=1}^T \|w_j\|^2$$

that penalizes both the number of the tree leaves and the L2 norms of its leaf weight vectors.

The minimization is done iteratively, and the training happens in an additive manner: at the start of each boosting iteration K we have $K - 1$ fixed trees built so far and we are looking add in the new tree f_K as to minimize

$$L^{(K)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(K-1)} + f_K(x_i)) + \Omega(f_K) \quad (1)$$

For a given function $l(x)$, the vector form Taylor expansion (up to second order derivative) can be written as:

$$l(x + \Delta x) = l(x) + (\Delta x)^\top \nabla g(x) + \frac{1}{2} (\Delta x)^\top H(x) \Delta x \quad (2)$$

where $\nabla g(x)$ - is the vector of gradients, $H_g(x)$ is the matrix of second order derivatives. If x is C -dimensional, then $\nabla g(x)$ will be of size C and $H(x)$ will be of size $C \times C$.

By taking f_K as a formal vector-valued symbol and applying formula 2 to formula 1, we get

$$L^{(K)} \approx \sum_{i=1}^n (l(y_i, \hat{y}_i^{(K-1)}) + f_K(x_i)^\top g_i + \frac{1}{2} f_K(x_i)^\top H_i f_K(x_i)) + \Omega(f_K) \quad (3)$$

where g_i is a C size vector of gradients (with respect to the predicted score for the first class, second class etc.), H_i - Hessian matrix (with respect to pairs of predicted scores, of size $C \times C$).

Now fix a structure $q(x_i)$ on f_K , and consider a particular weight vector w_j of leaf j . We can rewrite the objective so the summation happens on leaf index j . Dropping the loss on predictions so far, which is constant, we get

$$L^{(K)} = \sum_{j=1}^T (w_j^\top \tilde{g}_j + \frac{1}{2} (w_j^\top \tilde{H}_j w_j + \lambda w_j^\top w_j)) + \gamma T$$

$$\begin{aligned} Leaf_set_j &= \{i \mid q(x_i) = j\} \\ \tilde{g}_j &= \sum_{i \in Leaf_set_j} g_i \\ \tilde{H}_j &= \sum_{i \in Leaf_set_j} H_i \end{aligned} \quad (4)$$

where \tilde{g}_j is the vector representing the sum of gradients of instances that fall into the leaf j , \tilde{H}_j is the matrix representing the sum of Hessians of instances that fall into that leaf j . This objective is a sum of independent objectives per leaf. For each given leaf, we have

$$\bar{L}^{(K)}(w_j) = w_j^\top \tilde{g}_j + \frac{1}{2} w_j^\top (\tilde{H}_j + \lambda I) w_j \quad (5)$$

where I is the identity matrix.

This approximation is a quadratic function of the vector w_j and has a global minimum at

$$w_j = -(\lambda I + \tilde{H}_j)^{-1} \tilde{g}_j \quad (6)$$

if the matrix $\lambda I + \tilde{H}_j$ is symmetric positive definite. This is easy to guarantee: by Schwarz's theorem, the individual Hessian matrices H_i will be symmetric if the second-order derivatives are continuous, and a convex loss function will make the H_i 's and their sum \tilde{H}_j positive semi-definite. Finally, if $\lambda > 0$, adding λI to the sum will make the final matrix positive definite.

Under those conditions, adding the leaf j will decrease the loss by

$$\text{Gain}_j = \frac{1}{2} \tilde{g}_j^\top (\lambda I + \tilde{H}_j)^{-1} \tilde{g}_j \quad (7)$$

To evaluate the quality of the split, the contributions of both the new left and right leaves L and R must be compared against the contribution of the removed parent P , along with any penalty the regularizer might impose on the increased tree complexity:

$$\text{Gain}_{P \rightarrow L, R} = \text{Gain}_L + \text{Gain}_R - \text{Gain}_P - \gamma \quad (8)$$

Similarly to XGBoost, we build our trees greedily based on this gain, and always pick the split with the highest gain. We also offer the usual option of only picking a split if its gain is greater than zero ("pre-pruning"), or allowing the gain to be negative (which can happen because of the regularization) and then post-pruning afterwards. From an implementation point of view, it is worth noting that $\tilde{g}_P = \tilde{g}_L + \tilde{g}_R$ and $\tilde{H}_P = \tilde{H}_L + \tilde{H}_R$.

A. Matrix inversion

Both accumulation of Hessian matrices and matrix inversion in Formulas (7 and 6) are potentially expensive. We implement two strategies in TFBT:

1) *Full Hessian*: Note that if we want to calculate $x = A^{-1}v$, or alternatively $Ax = v$, where A is some matrix, instead of explicitly calculating A^{-1} and multiplying by v , we can treat it as linear least squares system [21]. In general, the usual methods to solve such a system are SVD decomposition, the QR decomposition and normal equations. SVD decomposition is accurate but slow, normal equations tend to be the least accurate but the fastest, and the QR decomposition is in between [22]. We go with QR decomposition with column pivoting of A using the *Eigen* library.

2) *Diagonal Hessian*: One potential simplification would be to assume that matrix \tilde{H}_j is diagonal, which makes both Hessian accumulation for each split and matrix inverse $O(C)$ (vs. $O(C^3)$ for full Hessian inverse and $O(C^2)$ for Hessians accumulation). In our experiments we show that such simplification actually does not result in decreased performance. On the contrary, it serves as additional regularization and results in faster convergence and better results.

IV. CHOICE OF LOSSES AND EXTENSIONS

For classification problems with $C > 2$ classes, our default model stores length C vectors in the leaves, turns those vectors into a vector of C probabilities using the softmax function, and then evaluates those probabilities using the cross-entropy error function. This model does have the drawback of using C parameters to describe the $C - 1$ dimensional space of C probabilities summing to 1, but it has two out-weighing benefits: it maintains the symmetry between the class weights, and it allows the regularization function Ω to maintain its usual form. In particular, the zero vector of weights, which minimizes the L2 penalty, maps to the maximally uninformative vector of equal probabilities, which is desirable.

This default model and loss is the one used for the experiments in section VII. However, as previously mentioned, our formulation and implementation is easily extensible and can be used to obtain many of the multiclass approaches discussed in Section II, in the following ways:

- A compact (because we build 1 tree per iteration for all of the classes) version of *LogitBoost*-like model can be obtained when multinomial logistic/cross entropy loss is chosen, and labels are encoded as 1-hot, with a diagonal Hessian mode.
- *SAMME*-like model can be approximated by not using the Hessian at all, and modifying the loss to add Lagrange to enforce the constraints of scores summing to 0.
- *GAMBLE*-like model can be obtained by using multiclass exponential loss, vector-form label encoding

scheme of 1 in position of the real class, and $-1/(C-1)$, and no regularization. The label encoding can be done in a TensorFlow input function.

- *AOSO LogitBoost* can be simulated by using cross-entropy loss and adding new block-coordinate strategy for Hessian approximation.
- *GD-MCBoost* can be obtained by using the label encoding scheme from [15], no second order gradient in loss optimization and exponential multiclass loss.

One thing to keep in mind is that for a large number of classes, the full Hessian approximation might become prohibitively slow. Since our experiments show no significant difference between full and diagonal Hessian, we recommend that diagonal Hessian is used as default for a larger number of classes.

Additionally, we would like to point out that our formulation should work with multi-label problems as well, as long as an appropriate multi-label loss is provided. However for multi-label problems it might be beneficial to explore other than diagonal approximations of the Hessian, that can account for the connections between the labels and the sparsity of the Hessian, for example block diagonal Hessians. It is however not suitable for extreme-multilabel problems, where only few labels out of millions apply to an instance. This is due to the fact that a full dense vector of scores will be stored in leaves.

Finally, we would like to highlight that since TensorFlow does automatic differentiation, switching between losses and creating new customizable losses for multiclass setting should be very easy in our proposed framework. Any twice-differentiable loss should be easily pluggable into TFBT.

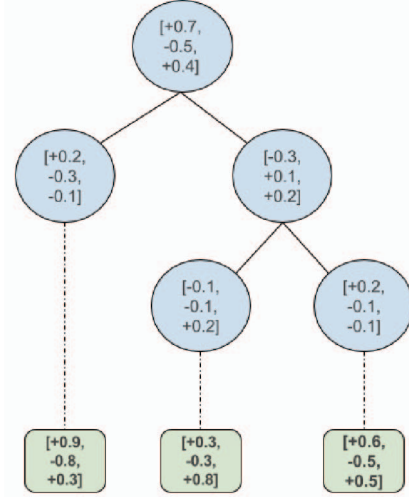
V. LAYER-BY-LAYER BOOSTING

One of the novel features of TFBT's tree building procedure is layer-by-layer boosting. In TFBT's layer-by-layer boosting, we allow internal nodes to contribute weight updates while the current tree is getting built. A leaf node's final contribution is therefore the aggregate contributions from its ancestors all the way up to the root node. This is illustrated in the following diagram: We rewrite the objective defined in Formula 1 to be a function of both K which still tracks the index of the current tree we are building and Z which tracks the layer of the tree we are currently building, we now have:

$$L^{(K)} = \sum_i l(y_i, \hat{y}_i^{K-1}) + f_K^{Z-1}(x_i) + f_K^Z(x_i) + \Omega(f_K^Z)$$

where $f_K^Z(x_i)$ is the prediction from the last Z layer of tree K that is currently being built. Notice that one boosting iteration now results in building one layer instead of a whole tree. Intuitively, the leaves grown in $f_K^Z(x_i)$ are learning a residual adjustment over the previous layer. This is useful as deeper trees define more fine grained partitions of the

Figure 1. Layered Prediction.



example space and we end up with leaves having few examples where it's desirable to learn smaller adjustments which leverage parent nodes as priors to minimize overfitting. Furthermore, recalculating gradients at every layer results in better approximating the functional space gradient which in turn typically enables TFBT to build fewer trees due to faster convergence. This is especially true for more complex loss functions, which can be user-defined in TFBT, where the second order Taylor expansion still results in relatively sizable approximation errors.

Additionally, we would like to mention that our implementation allows to choose how many instances to use to recalculate gradients per each layer. On deeper trees, each node would get fewer and fewer examples in a conventional scheme, and splits quality will deteriorate. TFBT allows users to specify the number of instances required for each layer, so the deeper layers can be built on more instances than shallower ones.

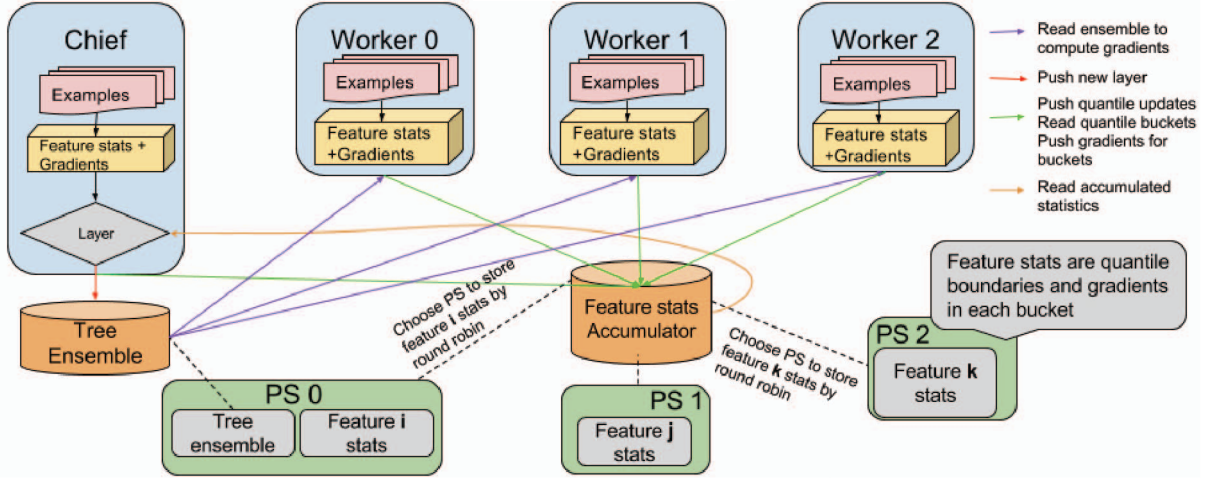
VI. TFBT SYSTEM DESIGN

Below we briefly describe our computation model and changes we had to make to support multi-class learning. For more in-depth review of TFBT please refer to [4].

Our design is similar to XGBoost [2] and TencentBoost [23] in that we build distributed quantile sketches of feature values and use them to build histograms, to be used later to find the best split. In TencentBoost [23] and PSMART [24] the full training data is partitioned and loaded in workers' memory, which can be a problem for larger datasets. To address this we instead work on mini-batches, updating quantiles in an online fashion without loading all the data into the memory. As far as we know, this approach is not implemented anywhere else.

Each worker loads a mini-batch of data, builds a local quantile sketch, pushes it to a Parameter Server (PS) and

Figure 2. TFBT architecture



Algorithm 1 Chief and Workers' work

```

1: procedure CALCULATESTATISTICS(PS, MODEL, STAMP, BATCH_DATA, LOSS_FN)
2:    $predictions \leftarrow \text{model.predict}(BATCH\_DATA)$ 
3:    $quantile\_stats \leftarrow \text{calculate\_quantile\_stats}(BATCH\_DATA)$ 
4:    $\text{push\_stats}(PS, quantile\_stats, \text{stamp})$  ▷ PS updates quantiles
5:    $current\_boundaries \leftarrow \text{fetch\_latest\_boundaries}(PS, \text{stamp})$ 
6:    $gradients, Hessians \leftarrow \text{calculate\_derivatives}(predictions, LOSS\_FN)$ 
7:    $gradients, Hessians \leftarrow \text{aggregate}(current\_boundaries, gradients, Hessians)$ 
8:    $\text{push\_stats}(PS, gradients, Hessians, \text{size}(BATCH\_DATA), \text{stamp})$ 
9: procedure DOWORK(PS, LOSS_FN, IS_CHIEF) ▷ Runs on workers and 1 chief
10:  while true do
11:     $BATCH\_DATA \leftarrow \text{read\_data\_batch}()$ 
12:     $model \leftarrow \text{fetch\_latest\_model}(PS)$ 
13:     $stamp \leftarrow \text{model.stamp\_token}$ 
14:     $\text{CalculateStatistics}(PS, model, stamp, BATCH\_DATA, LOSS\_FN)$ 
15:    if  $is\_chief \ \& \ \text{get\_num\_examples}(PS, stamp) \geq N\_PER\_LAYER$  then
16:       $next\_stamp \leftarrow stamp + 1$ 
17:       $stats \leftarrow \text{flush}(PS, stamp, next\_stamp)$  ▷ Update stamp, returns stats
18:       $\text{build\_layer}(PS, model, next\_stamp, stats)$  ▷ PS updates the ensemble: best splits for the nodes in the layer are chosen according to Formula 8, splitting is performed

```

fetches the bucket boundaries that were built at the previous iteration. Workers then compute per bucket gradients and Hessians and push them back to the PS. One of the workers, designated as Chief, checks at each iteration if the PS have accumulated enough statistics for the current layer and if so, starts building the new layer by finding best splits for each of the nodes in the layer. Code that finds the best splits for each feature is executed on the PS that have accumulated the gradient statistics for the feature. The Chief receives the best split for every leaf from the PS and grows a new layer on the tree.

Once the Chief adds a new layer, the workers copy of the tree ensemble become stale. To avoid stale updates, we introduce an abstraction called StampedResource - a TensorFlow resource with an int64 stamp. The tree ensemble, as well as gradients and quantile accumulators are all stamped resources with such a token. When the worker fetches the model, it gets the stamp token which is then used for all the reads and writes to stamped resources until the end of the iteration. This guarantees that all the updates are

consistent and ensures that the Chief doesn't need to wait for Workers for synchronization, which is important when using preemptible VMs (Figure 4). The Chief checkpoints resources to disk and workers don't hold any state, so if they are restarted, they can load a new mini-batch and resume. In order to support multi-class with Diagonal Hessian and Full Hessian strategies, our statistics accumulators support accumulating tensors as well as scalars for each bucket. During TensorFlow graph construction we pick either a tensor accumulator for multiclass or a scalar accumulator for binary classification.

VII. EXPERIMENTS

In this paper, we make the following claims that we are trying to confirm with experiments.

- Our general vector-form multiclass handling is a better way than the conventional 1-vs-rest strategy implemented in most libraries: with a vector-form multiclass, a considerably smaller number of trees is required to reach the best performance.

- Layer-by-layer boosting when used with vector form multiclass allows for faster convergences and results in smaller ensembles.

Although not the focus of this paper, we also check (Table III) that our implementations are on par with two multiclass vector methods, *LogitBoost* and *ABC-LogitBoost*. We use larger ensembles to be able to fairly compare against publicly available results in [13].

A. Datasets

We perform the experiments on medium and large size multiclass datasets, ranging from thousands to a million instances. The details about the datasets and the preprocessing we have done can be found in Table I.

B. Experiments setup

We build trees of depth 4 and explore ensembles of 10, 25, 50 and 100 trees. As a baseline we use XGBoost [2] conventional 1-vs-rest multiclass handling. For a default XGBoost model we use the following values of hyperparameters: `max_depth:4`, `learning_rate:0.3`, `objective:multi:softprob`, `lambda:1`, `scale_pos_weight:False`. For tuned XGBoost, we tune `min_child_weight`, `learning_rate` and `lambda` using scikitlearn `RandomizedSearchCV` [17] with 20 iterations with 5-fold CV over the training data using a predefined grid of values. For the poker dataset, 3-fold CV was used due to the fact that there were only 3 instances of class 9. If during hyper parameter tuning, values from the edges of the grid are chosen, the grid is expanded and search is repeated, depth and objective remaining fixed. The grids are as follows:

```
'min_child_weight':
  [0,0.05,0.1,0.5,1,2,4,8,10,12,14,16,18,20,22],
'learning_rate':
  [0.01,0.03,0.05,0.075,0.1,0.3,0.5,1,2,3],
'reg_lambda': [0.01, 0.5, 0.1, 1,2,4]
```

Note that XGBoost's `num_round` parameter (`n_estimators` in Scikit-learn) denotes the number of boosting iterations, not the number of trees. During each boosting iteration XGBoost will build C trees. To perform the comparison over the predefined number of trees, we adjust the number of rounds according to following formula:

$$\text{num_round} = \text{int}(\text{math.ceil}(1.0 * \text{num_trees} / C))$$

In some cases it results in a slightly larger number of XGBoost trees than reported. For example, for *Letter-26* 104 trees are built (for 100 trees column). For the same reasons, it is impossible to build 10 trees for *Letter-26* dataset. We don't adjust the number of trees for TFBT implementation (so 10, 25, 50 and 100 trees are built).

For TFBT we use the same objective as XGBoost (namely Max-Ent/Cross entropy/Multinomial Logistic Loss). For TFBT experiments, we use the same hyperparameter values as in default XGBoost: all parameters apart from *lambda* (L2 regularization) translate directly into TFBT settings. We set

TFBT's L2 to, since L2 in TFBT is on per instance basis. We use full batch size (equal to the training data size), and accumulate train batch size number of instances for layer-by-layer boosting.

We compare several variations of vector multiclass handling, namely full Hessian, diagonal Hessian and diagonal Hessian with layer-by-layer boosting. We also include the results for conventional tree-per-class (1-vs-rest) implementation, which is also available in TFBT. Note that we don't tune the hyper parameters of our TFBT methods at all, since our goal is to show that vector form multiclass results in smaller ensembles no matter of how much tuning is done. It should be noted that except in tree-per-class mode, TFBT will use more parameters per tree than XGBoost, since it stores length C vectors rather than scalars in the leaves. But we feel that for fixed depth trees, the number of trees is the more relevant comparison for purposes of inference speed and understandability.

Finally, to compare with *LogitBoost* and *ABC-LogitBoost*, we replicate Mnist10k (an Mnist dataset with train and test sets reversed) experiment from [13] for depth 2, 3 and 4 and learning rates 0.04, 0.06, 0.08, 0.1. For reasons of space and time, we did not compare against all other true multiclass approaches mentioned in section II, but as previously noted, many of them are variations of our formulation.

C. Metrics

We report the accuracy and the cross entropy loss on each test set. To check for significance between the TFBT and XGBoost results, we use an unpaired t-test with a p-value threshold of 0.05. (We use an unpaired t-test because it was easier to collect the data for, but since it is less powerful than the paired t-test, if it shows significance, the paired t-test would have as well.)

D. Experiments results

Table II summarizes our results. Several observations are apparent:

- When results are significant, none of the tree-per-class (1-vs-rest) implementations of multiclass handling are able to beat vector multiclass TFBT variants neither in terms of Accuracy nor in terms of Cross-Entropy loss. Sensit dataset is the only dataset where results are non-significant and where tuned XGBoost is able to achieve accuracy and loss that is not significantly different based on the test size and collected standard deviations. One thing to note is that for Sensit, tuned XGBoost ends up using a high learning rate (1 for 10, 25 and 50 trees, and 0.5 for 100 trees), resulting in "faster" learning which might lead to overfitting. All TFBT methods still use a learning rate of 0.3.
- Diagonal Hessian strategy seems to be often significantly better than full Hessian strategy. Since Diagonal version is much faster to run, we recommend to use

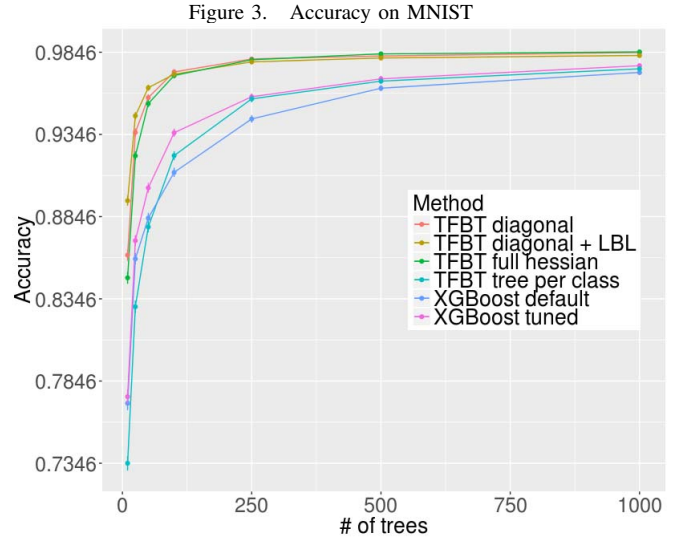
Table I
MULTICLASS DATASETS DESCRIPTIONS

| Dataset | # train | # test | # features | # classes | Comments |
|-----------------|-----------|---------|------------|-----------|---|
| Mnist [25] | 60,000 | 10,000 | 784 | 10 | Hand-written digits recognition task. Conventional train/test split. All features are dense numeric. |
| Sensit [26] | 78,823 | 19,705 | 100 | 3 | SensIT Vehicle (combined) dataset, preprocessed data obtained from LibSVM [27] |
| Covertypes [28] | 464,810 | 116,202 | 54 | 7 | Obtained from [29]. Predicting forest cover type. We split all available data into 80% train and 20% test. |
| Letter-26 [30] | 16,000 | 4,000 | 16 | 26 | Obtained from [29]. Conventional test/train split. Predicting a English letter from image features. |
| Poker [31] | 1,000,000 | 25,010 | 10 | 10 | Predicting type of "poker hand" based on information about 5 cards. Preprocessed data obtained from LibSVM [27] |
| CIFAR-10 [32] | 50,000 | 10,000 | 3072 | 10 | Classify RGB 32x32 images into a number of classes. Conventional test/train split. No convolutional features were used. |

it by default, instead of the full Hessian method. We also hypothesize that approximating the Hessian matrix via a diagonal matrix serves as a sort of additional regularizer, resulting in better performance in most of the cases.

- Layer-by-layer boosting (LBL) results in best performance for ensembles of all fixed sizes apart from a single case for Sensit 100 trees. For this particular case, it seems that 100 trees is already enough to achieve best performance even with tree-per-class method, and it suggests that since LBL speeds up convergence, it also may result in overfitting on problems that are easier.
- If extremely small ensembles are required (10-50 trees), for example due to inference time or memory size of the device requirements, Diagonal+LBL boosting should be the first choice. Even if larger ensembles can be tolerated, vector form multiclass strategies should be used instead of 1-vs-rest. Figures 3 and 4 present the accuracy convergence on larger number of trees, which shows that even for larger ensembles, vector form multiclass methods dominate. However the line between different vector-form strategies blurs as more trees are added.
- The improvements of vector form over tree-per-class, not surprisingly, are more pronounced on datasets with larger number of classes (like Letter-26, Poker, CIFAR-10 and MNIST). It is assumed that approximately C times more trees will be required for tree-per-class to achieve the same performance, and the experiments seem to confirm this folk wisdom.
- An interesting observation is that TFBT tree-per-class implementation produces different results from XGBoost's version. The difference between implementations lies in the fact that XGBoost computes the predictions for all the instances first and then uses them to construct all C trees in the same boosting iteration, whereas our implementation recalculates the

predictions and uses them to calculate the loss and gradients after each subsequent tree. This seems to result in better loss but sometimes leads to worse accuracy (for example, Mnist 10, 25 and 50 trees). Since an improvement in cross-entropy loss does not necessarily translates into a direct improvement in accuracy, we don't find these experiment results alarming. It also demonstrates that recalculating the predictions after each new trees is beneficial to a faster loss convergence, as expected.



VIII. CONCLUSION

We show that the conventional boosting formalism used by most popular open-sourced gradient boosting libraries can be easily extended from scalar-valued trees to vector-valued trees. We demonstrate, not surprisingly, that vector-valued trees lead to much faster convergence and smaller ensembles. We also show that other recent vector-valued gradient

Table II
MULTICLASS DATASETS RESULTS

| Data | Method | 10 trees | | 25 trees | | 50 trees | | 100 trees | |
|------------|------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | | Accuracy | Cross-Ent | Accuracy | Cross-Ent | Accuracy | Cross-Ent | Accuracy | Cross-Ent |
| Mnist | XGBoost default | 0.7711 | 1.5444 | 0.8589 | 1.0151 | 0.8837 | 0.7497 | 0.9115 | 0.4452 |
| | XGBoost tuned | 0.7751 | 0.8479 | 0.8700 | 0.4515 | 0.9021 | 0.3218 | 0.9356 | 0.2084 |
| | Tree-per-class | 0.7346 | 1.2998 | 0.8298 | 0.7830 | 0.8783 | 0.4957 | 0.9217 | 0.2827 |
| | Full Hessian | 0.8474 | 0.5465 | 0.9216 | 0.2613 | 0.9533 | 0.1555 | 0.9704 | 0.0972 |
| | Diagonal Hessian | 0.8612 | 0.4741 | 0.9358 | 0.2201 | 0.9570 | 0.1398 | 0.9726 | 0.0917 |
| | Diagonal + LBL | 0.8943 | 0.3513 | 0.9459 | 0.1775 | 0.9630 | 0.1226 | 0.9712 | 0.0899 |
| Sensit | XGBoost default | 0.7904 | 0.6582 | 0.8091 | 0.5292 | 0.8222 | 0.4691 | 0.8364 | 0.4309 |
| | XGBoost tuned | 0.8088 | 0.4921 | 0.8261 | 0.4456 | 0.8359 | 0.4230 | 0.8430 | 0.4145 |
| | Tree-per-class | 0.7836 | 0.5977 | 0.8144 | 0.4869 | 0.8331 | 0.4399 | 0.8409 | 0.4158 |
| | Full Hessian | 0.7998 | 0.5376 | 0.8245 | 0.4571 | 0.8386 | 0.4225 | 0.8454 | 0.4060 |
| | Diagonal Hessian | 0.8041 | 0.4995 | 0.8307 | 0.4407 | 0.8327 | 0.4431 | 0.8343 | 0.4409 |
| | Diagonal + LBL | 0.8110 | 0.4833 | 0.8285 | 0.4442 | 0.8391 | 0.4235 | 0.8407 | 0.4197 |
| Coverttype | XGBoost default | 0.5991 | 1.3425 | 0.6173 | 1.1343 | 0.6164 | 0.9669 | 0.6205 | 0.8884 |
| | XGBoost tuned | 0.6097 | 1.0181 | 0.6158 | 0.9916 | 0.6072 | 0.9182 | 0.6092 | 0.8982 |
| | Tree-per-class | 0.7071 | 1.1480 | 0.7264 | 0.8013 | 0.7362 | 0.6482 | 0.7577 | 0.5668 |
| | Full Hessian | 0.7113 | 0.6812 | 0.7536 | 0.5647 | 0.7776 | 0.5085 | 0.8085 | 0.4549 |
| | Diagonal Hessian | 0.7354 | 0.6364 | 0.7705 | 0.5466 | 0.7810 | 0.5361 | 0.7632 | 0.5624 |
| | Diagonal + LBL | 0.7399 | 0.6081 | 0.7696 | 0.5337 | 0.8047 | 0.4573 | 0.8348 | 0.3927 |
| Letter-26 | XGBoost default | N/A | N/A | 0.6850 | 1.7549 | 0.7282 | 1.4718 | 0.7708 | 1.1522 |
| | XGBoost tuned | N/A | N/A | 0.6855 | 1.4266 | 0.6967 | 1.1718 | 0.7813 | 0.8844 |
| | Tree-per-class | 0.2843 | 2.5872 | 0.6015 | 1.5482 | 0.6988 | 1.1599 | 0.7593 | 0.8561 |
| | Full Hessian | 0.7623 | 0.9297 | 0.8665 | 0.5191 | 0.9190 | 0.3118 | 0.9465 | 0.1879 |
| | Diagonal Hessian | 0.7595 | 0.9263 | 0.8705 | 0.4913 | 0.9223 | 0.2926 | 0.9510 | 0.1800 |
| | Diagonal + LBL | 0.8060 | 0.7339 | 0.8973 | 0.3758 | 0.9375 | 0.2165 | 0.9560 | 0.1409 |
| Poker | XGBoost default | 0.5376 | 1.7858 | 0.5471 | 1.5421 | 0.5524 | 1.2046 | 0.5913 | 1.0002 |
| | XGBoost tuned | 0.5376 | 1.1414 | 0.5568 | 0.9871 | 0.5634 | 0.9492 | 0.5880 | 0.9083 |
| | Tree-per-class | 0.5086 | 1.7671 | 0.5212 | 1.2253 | 0.5713 | 1.0170 | 0.5966 | 0.9170 |
| | Full Hessian | 0.5633 | 0.9492 | 0.6201 | 0.8567 | 0.6497 | 0.8089 | 0.6872 | 0.7426 |
| | Diagonal Hessian | 0.5625 | 0.9424 | 0.6301 | 0.8368 | 0.6606 | 0.7871 | 0.7185 | 0.7012 |
| | Diagonal + LBL | 0.5830 | 0.9299 | 0.6351 | 0.8249 | 0.6750 | 0.7634 | 0.7329 | 0.6651 |
| CIFAR-10 | XGBoost default | 0.3012 | 2.1725 | 0.3528 | 2.0195 | 0.3769 | 1.9221 | 0.4050 | 1.7724 |
| | XGBoost tuned | 0.3003 | 2.0021 | 0.3546 | 1.9215 | 0.3829 | 1.7361 | 0.4194 | 1.6361 |
| | Tree-per-class | 0.2885 | 2.1113 | 0.3285 | 1.9709 | 0.3622 | 1.8335 | 0.4053 | 1.6967 |
| | Full Hessian | 0.3450 | 1.8574 | 0.4085 | 1.6830 | 0.4507 | 1.5573 | 0.4856 | 1.4601 |
| | Diagonal Hessian | 0.3568 | 1.8418 | 0.4068 | 1.6746 | 0.4497 | 1.5674 | 0.4798 | 1.4778 |
| | Diagonal + LBL | 0.3631 | 1.7933 | 0.4148 | 1.6458 | 0.4548 | 1.5434 | 0.4769 | 1.4903 |

For each dataset and column, in **bold** are the best values of the metrics for cross entropy-loss and accuracy for this number of trees. In *italics*, we highlight the best value of cross-entropy and accuracy when Diagonal+LBL is not considered. If there is no significant difference between several best values, we highlight all the values the difference between which is insignificant.

Table III
COMPARISON OF ACCURACY ON MNIST10K

| LR | D | LogitBoost | ABC-Logit | Diagonal | Diag+LBL |
|------|---|------------|-----------|----------|----------|
| 0.04 | 2 | 0.9511 | 0.9562 | 0.9562 | 0.9585 |
| | 3 | 0.9567 | 0.9640 | 0.9643 | 0.9617 |
| | 4 | 0.9596 | 0.9648 | 0.9661 | 0.9646 |
| 0.06 | 2 | 0.9505 | 0.9567 | 0.9768 | 0.9556 |
| | 3 | 0.9564 | 0.9644 | 0.9640 | 0.9614 |
| | 4 | 0.9594 | 0.9648 | 0.9657 | 0.9637 |
| 0.08 | 2 | 0.9503 | 0.9578 | 0.9582 | 0.9561 |
| | 3 | 0.9569 | 0.9647 | 0.9642 | 0.9615 |
| | 4 | 0.9601 | 0.9651 | 0.9660 | 0.9617 |
| 0.1 | 2 | 0.9497 | 0.9580 | 0.9588 | 0.9561 |
| | 3 | 0.9567 | 0.9643 | 0.9640 | 0.9614 |
| | 4 | 0.9601 | 0.9653 | 0.9660 | 0.9635 |

D stands for depth, *LR* for learning rate, *ABC-Logit* is ABC-LogitBoost, *Diagonal* and *Diag+LBL* are TFBT implementations. Results for LogitBoost and ABC-LogitBoost are from [13] Table 2.

boosted trees formulations fit into our general framework and can be easily implemented in TFBT by changing the loss, Hessian handling strategy and label encoding scheme.

We would like to encourage researchers who work on multiclass boosting to reuse our TensorFlow based TFBT library, which is open-sourced and convenient to use due to automatic differentiation capabilities. Finally, we argue that vector-valued leaves should be the default strategy for handling multiclass problems.

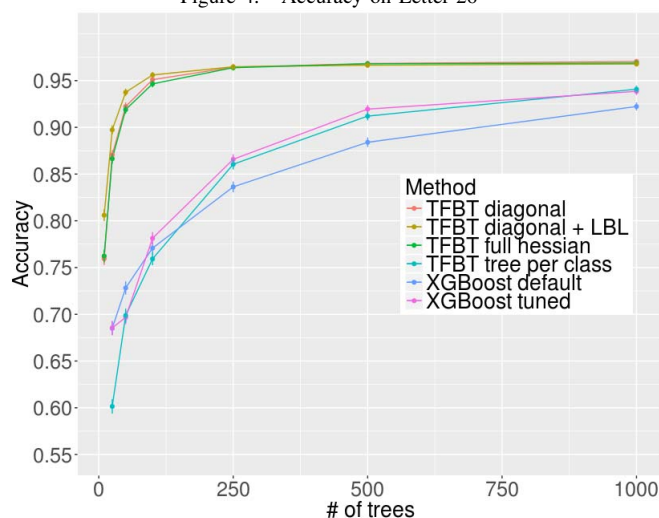
ACKNOWLEDGMENT

The authors would like to thank Boris Dadachev, Afshin Rostamizadeh, Corinna Cortes, Tal Shaked, D. Sculley, Alexander Grushetsky and Petr Mitrichev for their invaluable comments on this paper.

REFERENCES

- [1] R. Caruana and A. Niculescu-Mizil, "An empirical comparison of supervised learning algorithms using different performance metrics," in *In Proc. 23 rd Intl. Conf. Machine learning (ICML06)*, 2005, pp. 161–168.
- [2] T. Chen *et al.*, "XGBoost: A scalable tree boosting system," *CoRR*, 2016.

Figure 4. Accuracy on Letter-26



- [3] L. Mason, J. Baxter, P. L. Bartlett, and M. R. Frean, "Boosting algorithms as gradient descent," in *Advances in neural information processing systems*, 2000, pp. 512–518.
- [4] N. Ponomareva, S. Radpour, G. Hendry, S. Haykal, T. Colthurst, P. Mitrichev, and A. Grushetsky, "Tf boosted trees: A scalable tensorflow based framework for gradient boosting," in *ECML PKDD*, 2017.
- [5] E. L. Allwein, R. E. Schapire, and Y. Singer, "Reducing multiclass to binary: A unifying approach for margin classifiers," *J. Mach. Learn. Res.*, 2001.
- [6] Y. Freund and R. E. Schapire, "Experiments with a new boosting algorithm," 1996.
- [7] J. Friedman, T. Hastie, and R. Tibshirani, "Additive logistic regression: a statistical view of boosting," *Annals of Statistics*, vol. 28, p. 2000, 1998.
- [8] A. Saffari, M. Godec, T. Pock, C. Leistner, and H. Bischof, "Online multi-class lboost."
- [9] J. Zhu, H. Zou, S. Rosset, and T. Hastie, "Multi-class adaboost," 2009.
- [10] P. Sun and other, "Aoso-logitboost: Adaptive one-vs-one logitboost for multi-class problem," in *ICML-12*. ACM, 2012.
- [11] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Annals of Statistics*, vol. 29, pp. 1189–1232, 2000.
- [12] P. Sun, T. Zhang, and J. Zhou, "A convergence rate analysis for logitboost, mart and their variant," in *Proceedings of the 31st International Conference on International Conference on Machine Learning*, ser. ICML'14, 2014.
- [13] P. Li, "Robust logitboost and adaptive base class (ABC) logitboost," *CoRR*, vol. abs/1203.3491, 2012. [Online]. Available: <http://arxiv.org/abs/1203.3491>
- [14] —, "Adaptive base class boost for multi-class classification," *CoRR*, vol. abs/0811.1250, 2008. [Online]. Available: <http://arxiv.org/abs/0811.1250>
- [15] M. J. Saberian and N. Vasconcelos, "Multiclass boosting: Theory and algorithms," in *Advances in Neural Information Processing Systems* 24, 2011.
- [16] J. Huang, S. Ertekin, Y. Song, H. Zha, and C. L. Giles, *Efficient Multiclass Boosting Classification with Active Learning*, pp. 297–308. [Online]. Available: <http://epubs.siam.org/doi/abs/10.1137/1.9781611972771.27>
- [17] F. Pedregosa *et al.*, "Scikit-learn: Machine learning in Python," *JMLR*, vol. 12, 2011.
- [18] G. Ridgeway, "Generalized boosted models: A guide to the gbm package," 2005.
- [19] Microsoft, "Microsoft/dmtk," <https://github.com/microsoft/dmtk>, 2013.
- [20] X. Meng, J. Bradley *et al.*, "MLlib: Machine learning in Apache Spark," 2016.
- [21] G. H. Golub and C. F. Van Loan, *Matrix Computations (3rd Ed.)*. Baltimore, MD, USA: Johns Hopkins University Press, 1996.
- [22] D. Q. Lee, "Numerically efficient methods for solving least squares problems."
- [23] J. Jiang *et al.*, "Tencentboost: A gradient boosting tree system with parameter server."
- [24] J. Zhou *et al.*, "PSMART: Parameter server based multiple additive regression trees system," ser. WWW '17 Companion.
- [25] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [26] M. F. Duarte and Y. H. Hu, "Vehicle classification in distributed sensor networks," *J. Parallel Distrib. Comput.*, 2004.
- [27] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, 2011, datasets at <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multi-class.html>.
- [28] J. A. Blackard and D. J. Dean, "Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables," *Computers and Electronics in Agriculture*, 1999.
- [29] M. Lichman, "UCI machine learning repository," 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [30] P. W. Frey and D. J. Slate, "Letter recognition using holland-style adaptive classifiers," *Machine Learning*, 1991.
- [31] R. Cattral, F. Oppacher, and D. Deugo, "Evolutionary data mining with automatic rule generalization," in *Recent Advances in Computers, Computing and Communications*, 2002.
- [32] A. Krizhevsky, "Learning multiple layers of features from tiny images," Tech. Rep., 2009.