Generalized Boosting

Anonymous Author(s)

Affiliation Address email

Abstract

Boosting is a widely used learning technique in machine learning for solving classification problems. In boosting, one predicts the label of an example using an ensemble of weak classifiers. While boosting has shown tremendous success on many classification problems involving tabular data, it performs poorly on complex classification tasks involving low-level features such as image classification tasks. This drawback stems from the fact that boosting builds an additive model of weak classifiers, each of which has very little predictive power. Often, the resulting additive models are not powerful enough to approximate the complex decision boundaries of real-world classification problems. In this work, we present a general framework for boosting where, similar to traditional boosting, we aim to boost the performance of a weak learner and transform it into a strong learner. However, unlike traditional boosting, our framework allows for more complex forms of aggregation of weak learners. In this work, we specifically focus on one form of aggregation - function composition. We show that many popular greedy algorithms for learning deep neural networks (DNNs) can be derived from our framework using function compositions for aggregation. Moreover, we identify the drawbacks of these greedy algorithms and propose new algorithms that fix these issues. Using thorough empirical evaluation, we show that our learning algorithms have superior performance over traditional additive boosting algorithms, as well as existing greedy learning techniques for DNNs. An important feature of our algorithms is that they come with strong theoretical guarantees.

Introduction

8

9

10

11

12

13

14

15

16

17

18

19

20

21

30

31

32

33

34

36

38

39

Boosting is a widely used learning technique in machine learning for solving classification problems. 23 Boosting aims to improve the performance of a weak learner by combining multiple weak classifiers to produce a strong classifier with good predictive performance. Since the seminal works of Freund [13], Schapire [30], a number of practical algorithms such as AdaBoost [16], gradient boosting [24], XGBoost [9], have been proposed for boosting. Over the years, boosting based methods such as XGBoost in particular, have shown tremendous success in many real-world classification problems, as well as competitive settings such as Kaggle competitions. However, this success is mostly limited to classification tasks involving structured or tabular data with hand-engineered features. On classification problems involving low-level features and complex decision boundaries, boosting tends to perform poorly [3, 28] (also see Section 5). One example where this is evident is the image classification task, where the decision boundaries are often complex and the features are low-level pixel intensities. This drawback stems from the fact that boosting builds an additive model of weak classifiers, each of which has very little predictive power. Since such additive models with any 35 reasonable number of weak classifiers are usually not powerful enough to approximate complex decision boundaries, the models' output by boosting tend to have poor performance. 37

In this work, we aim to overcome this drawback of traditional boosting by considering a generalization of boosting which allows for more complex forms of aggregation than linear combinations of weak classifiers. To achieve this goal, we work in the feature representation space and boost the

performance of weak feature transformers. Working in the representation space allows for more flexible combinations of weak feature transformers. This is unlike traditional boosting which works in the label space and builds an additive model on the predictions of the weak classifiers. The starting point for our approach is the greedy view of boosting, originally studied by Friedman et al. [18], Mason et al. [24]. Let $R_S(f)$ be the risk of a classifier f on training samples S, boosting techniques aim to approximate the minimizer of \hat{R}_S in terms of linear combinations of elements from a set of weak classifiers F. Many popular boosting algorithms including AdaBoost, XGBoost, rely on greedy techniques to find such an approximation. In our generalized framework for boosting, we take this greedy view, but differ in how we aggregate the weak learners. We approximate the minimizer of \hat{R}_S using models of the form $f_T = W\phi_T$, where $\phi_T = \sum_{t=0}^T g_t$, and $\{g_t\}_{t=0}^T$ are feature transformations learned in each iteration of the greedy algorithm, and W is the linear classifier on top of the feature transformation. Unlike additive boosting, where each g_t comes from a fixed weak feature transformer class \mathcal{G} , in our framework each g_t comes from a class \mathcal{G}_t which evolves over time t and is allowed to depend on the past iterates $\{\phi_i\}_{i=0}^{t-1}$. Some potential choices for \mathcal{G}_t that could be of interest are $\{g \circ \phi_{t-1} \text{ for } g \in \mathcal{G}\}$, $\{g \circ ([\phi_0, \dots, \phi_{t-1}]) \text{ for } g \in \mathcal{G}\}$, where $g \circ \phi(\mathbf{x}) = g(\phi(\mathbf{x}))$ denotes function of g and ϕ , and \mathcal{G} is a weak feature transformer class. Note that the former choice of \mathcal{G}_t is connected to layer-by-layer training of models with ResNet architecture [19].

As one particular instantiation of our framework, we consider weak feature transformers that are neural networks and use function compositions to combine them; that is, we use \mathcal{G}_t 's constructed using function compositions. We show that for certain choices of \mathcal{G}_t , our framework recovers the layer-by-layer training techniques developed in deep learning [6, 20]. Greedy layer-by-layer training techniques have seen a revival in recent years [5, 8, 20, 23, 27]. One reason for this revival is that greedy techniques consume less memory than end-to-end training of deep networks, and can hence accommodate much larger models in limited memory. As a primary contribution of the paper, we identify several drawbacks of existing layer-by-layer training techniques, and show that the choice of \mathcal{G}_t used by these algorithms can lead to a drop in performance. We propose alternative choices for \mathcal{G}_t which fix these issues and empirically demonstrate that the resulting algorithms have superior performance over existing layer-by-layer training techniques, and in some cases achieve performance close to that of end-to-end trained DNNs. Moreover, we show that the proposed algorithms perform much better than traditional additive boosting algorithms, on a variety of classification tasks.

As the second contribution of the paper, we provide excess risk bounds for models learned using our generalized boosting framework. Our results depend on a certain weak learning condition on feature transformer classes $\{\mathcal{G}_t\}_{t=1}^T$, which is a natural generalization of the weak learning condition that is typically imposed in traditional boosting. The resulting risk bounds are modular and depend on the generalization bounds of $\{\mathcal{G}_t\}_{t=1}^T$. An advantage of such modular bounds is that one can rely on the best-known generalization bounds for weak transformation classes $\{\mathcal{G}_t\}_{t=1}^T$ and obtain tight risk bounds for boosting. As an immediate consequence of this result, we obtain excess risk bounds for existing greedy layer-by-layer training techniques.

Related Work. Several works have proposed generalizations of traditional boosting. Cortes et al. [10] propose a boosting algorithm where the hypothesis set of weak classifiers is chosen adaptively. However, the resulting models are still additive models of weak classifiers and usually perform poorly on hard classification problems. Several recent works have attempted to learn neural networks greedily based on boosting theory. Cortes et al. [11] propose a boosting-style algorithm to learn both the structure and weights of neural networks in an adaptive way. However, the algorithms developed are restricted to feed forward neural networks and are mostly theoretical in nature. The experimental evidence in the paper is a proof-of-concept and only considers small scale binary classification tasks. Huang et al. [20], Nitanda and Suzuki [27] use ideas from classical boosting to learn neural networks in a layer-by-layer fashion. As we show later, these algorithms are specific instances of our generalized framework, and have certain drawbacks arising from the choice of \mathcal{G}_t they use.

2 Preliminaries

In this section, we set up the notation and review the necessary background on additive boosting. A consolidated list of notations can be found in Appendix A.

Notation. Let $(X,Y) \in \mathcal{X} \times \mathcal{Y}$ denote a feature-label pair following a probability distribution P. Let P^X, P^Y denote the marginal distributions of X and Y. In this work, we consider the multi-class classification problem where $\mathcal{Y} = \{0, \dots K-1\}$, and assume $\mathcal{X} \subseteq \mathbb{R}^d$. Let $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be n i.i.d samples drawn from P. Let P_n be the empirical distribution of S and P_n^X, P_n^Y be the marginal

distributions of $\{\mathbf{x}_i\}_{i=1}^n$, $\{y_i\}_{i=1}^n$. 97

In classification, our goal is to find a predictor that can well predict the label of any feature from just the 98 samples S. Let $f: \mathcal{X} \to \mathbb{R}^K$ denote a score-based classifier which assigns X to class $\operatorname{argmax}_i f_i(X)$. 99 The expected classification risk of f is defined as $\mathbb{E}_{X,Y}[\ell_{0-1}(f(X),Y)]$, where $\ell_{0-1}(f(X),Y)=0$ if $\operatorname{argmax}_i f_i(X)=Y$, and 1 otherwise. Since optimizing 0/1 risk is computationally intractable, 100 101 we consider convex surrogates of $\ell_{0-1}(f(X),Y)$, which we denote by $\ell(f(X),Y)$; typical choices for ℓ include the logistic loss and the exponential loss. The population risk of f is then defined as $R(f) = \mathbb{E}_{X,Y} \left[\ell(f(X),Y) \right]$. Since directly optimizing the population risk is impossible, we 102 103 104 approximate it with the empirical risk $\hat{R}_S(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i)$ and try to find its minimizer. 105

We consider classifiers of the form $f(X) = W\phi(X)$, where $\phi: \mathcal{X} \to \mathbb{R}^D$ is the feature transformer 106 and $W \in \mathbb{R}^{K \times D}$ is the linear classifier on top. A typical choice for ϕ is a neural network. We denote 107 the population and empirical risks of such an f as $R(W, \phi)$, $\hat{R}_S(W, \phi)$. We usually work in the space 108 of feature transforms. Let $L_2(P)$ denote the space of square integrable functions w.r.t P, and define the inner product between $\phi_1, \phi_2 \in L_2(P)$ as $\langle \phi_1, \phi_2 \rangle_P = \mathbb{E}_{X \sim P} \left[\langle \phi_1(X), \phi_2(X) \rangle \right]$. We denote with $\nabla_{\phi} R(W, \phi)$ the functional gradient of $R(W, \phi)$ w.r.t ϕ in the $L_2(P^X)$ space, which is defined as $\nabla_{\phi} R(W, \phi)(\mathbf{x}) = \mathbb{E}_{Y|\mathbf{x}} \left[W^T \nabla \ell(W\phi(\mathbf{x}), Y) \right]$, where $\nabla \ell(W\phi(\mathbf{x}), y)$ denotes the gradient of ℓ 109 110 111 112 w.r.t its first argument, evaluated at $W\phi(\mathbf{x})$. Similarly, we let $\nabla_{\phi}\hat{R}_{S}(W,\phi)$ denote the functional 113 gradient of $\widehat{R}_S(W,\phi)$ in the $L_2(P_n^X)$ space 114

$$\nabla_{\phi} \widehat{R}_{S}(W, \phi)(\mathbf{x}) = \begin{cases} W^{T} \nabla \ell(W \phi(\mathbf{x}_{i}), y_{i}), & \text{if } \mathbf{x} = \mathbf{x}_{i}, \\ 0 & \text{otherwise} \end{cases}$$

Additive Boosting. In this work, we refer to traditional boosting as additive boosting, as it constructs 115 additive models of weak classifiers. Let \mathcal{F} be a hypothesis class of weak classifiers, a typical example 116 being decision trees of bounded depth. Additive boosting aims to find an element in the linear span 117 of \mathcal{F} which minimizes the empirical risk $\hat{R}_S(f)$. As previously mentioned, there exists a duality between boosting and greedy algorithms [18, 24]. Many popular boosting algorithms use a greedy 118 119 forward stagewise approach to find a minimizer of $\hat{R}_S(f)$, and solve the following in each iteration: $\eta_t, f_t = \operatorname{argmin}_{\eta \in \mathbb{R}, f \in \mathcal{F}} \hat{R}_S\left(\sum_{i=1}^{t-1} \eta_i f_i + \eta f\right)$, 120

$$\eta_t, f_t = \operatorname{argmin}_{\eta \in \mathbb{R}, f \in \mathcal{F}} \widehat{R}_S \left(\sum_{i=1}^{t-1} \eta_i f_i + \eta f \right),$$

where η is the learning rate. Various algorithms differ in how they solve this optimization problem. In gradient boosting, one uses a linear approximation of \hat{R}_S around $\sum_{i=1}^{t-1} \eta_i f_i$ [24]. In this work, we take this greedy view of boosting to design the generalized boosting framework.

Additive Representation Boosting. In this work, we perform boosting in the representation space, contrasting with traditional boosting which works in the output space. Let \mathcal{G} be a hypothesis class of weak feature transformers, whose examples include the set of one layer neural networks of bounded width and a set of vector-valued polynomials of bounded degree. More generally, \mathcal{G} can be any set of non-linear transformations. In additive representation boosting, we aim to find a strong feature transform ϕ in the linear span of \mathcal{G} , and a linear predictor $W \in \mathcal{W} \subseteq \mathbb{R}^{K \times D}$ that minimizes $\hat{R}_S(W,\phi)$. To this end, we consider greedy algorithms that solve the following problem each iteration:

$$W_{t}, g_{t} = \operatorname{argmin}_{W \in \mathcal{W}, q \in \mathcal{G}} \hat{R}_{S} \left(W, \phi_{t-1} + \eta_{t} g \right), \tag{1}$$

where $\phi_t = \phi_0 + \sum_{i=1}^t \eta_i g_i$ with ϕ_0 being the initial feature transformation, and $\{\eta_i\}_{i=1}^{\infty}$ is a predefined learning rate schedule. 132 133

Generalized Boosting 3

121

122 123

124

125

126

127

128

129

130 131

134

The starting point for our generalized boosting framework is the additive representation boosting 135 described in Section 2. Typically, linear combinations of weak feature transformations are not 136 powerful enough to model complex decision boundaries. Consequently, the minimizer of $R_S(W,\phi)$ 137 over the linear span of \mathcal{G} tends to have a high risk. A simple workaround for this issue would be 138 to perform additive boosting with a complex hypothesis class \mathcal{G} . For example, if the weak feature 139 transformers are one layer neural networks, then one could increase the complexity of \mathcal{G} by using deeper networks. However, such an alternative has several drawbacks both from an optimization and generalization perspective and defeats the purpose of boosting, which aims to convert weak learners into strong learners. From an optimization perspective, moving to complex \mathcal{G} makes each greedy step harder to optimize. For example, compared to deep neural networks, shallow networks are easier to optimize, require fewer resources, and are easier to analyze or interpret [5]. From a generalization perspective, since the generalization bounds of boosting depend on the complexity of \mathcal{G} , larger hypothesis classes can lead to overfitting and poor performance on unseen data.

In this work, we are interested in other approaches for increasing the complexity of models produced 148 by boosting, while ensuring the boosting/greedy steps are easy to implement. One way to achieve 149 this is by considering more complex combinations of weak feature transformers than the linear 150 combinations considered in additive representation boosting. Formally, let \mathcal{G}_t denote the hypothesis 151 class of feature transformations used in the t^{th} iteration of boosting. In additive boosting, $\mathcal{G}_t = \mathcal{G}$ for all t. In our generalized boosting framework, we increase the complexity of \mathcal{G}_t by letting it depend on the past iterates $\{\phi_i\}_{i=0}^{t-1}$. Here are some potential choices for \mathcal{G}_t , other than the ones stated in the introduction: $\{g \circ (\sum_{i=0}^{t-1} \alpha_i \phi_i), \text{ for } g \in \mathcal{G}, \alpha_i \in \mathbb{R}\}, \{g \circ \phi_{t-1} \circ \phi_{t-2} \cdots \circ \phi_0, \text{ for } g \in \mathcal{G}\}$. Depending on the problem domain, one could consider several other ways of constructing \mathcal{G}_t using the past 152 153 154 155 156 iterates. Note that even with these complex choices of \mathcal{G}_t , the greedy steps are easy to implement and 157 only need a weak learner which can identify an element in \mathcal{G} that best fits the data. As a result, this 158 remains in the spirit of boosting and at the same time ensures the models learned are complex enough 159 for real world problems. 160

We now present our algorithm for generalized boosting (see Algorithm 1). Similar to additive representation boosting, our algorithm proceeds in a greedy fashion. In the t^{th} iteration of the algorithm, we aim to solve the following optimization problem:

$$W_t, g_t = \underset{W \in \mathcal{W}, g \in \mathcal{G}_t}{\operatorname{argmin}} \widehat{R}_S \left(W, \phi_{t-1} + \eta_t g \right). \tag{2}$$

We provide two approaches for solving this problem. One is the *exact greedy approach*, which directly solves the optimization problem (Algorithm 2). For problems where direct optimization of Equation (2) is difficult¹, we provide an approximate technique which performs functional gradient descent on the objective. In this approach, which we call *gradient greedy approach*, we approximate the objective with the linear approximation of \hat{R}_S around ϕ_{t-1} (Algorithm 3):

$$\widehat{R}_{S}\left(W,\phi_{t-1}+\eta_{t}g\right)\approx\widehat{R}_{S}\left(W,\phi_{t-1}\right)+\eta_{t}\left\langle\nabla_{\phi}\widehat{R}_{S}(W,\phi_{t-1}),g\right\rangle_{P_{S}^{X}}.$$

To optimize the linear approximation, we first fix W to W_{t-1} and find a minimizing $g_t \in \mathcal{G}_t$. Intuitively, this step can be seen as finding a g which best aligns with the negative functional gradient of empirical risk at the current iterate. For appropriate choice of learning rate η , moving along g_t results in reduction of \hat{R}_S . Next, we fix g_t and find a linear predictor W which minimizes the empirical risk $\hat{R}_S(W, \phi_t)$. This alternating optimization of g and W makes the algorithm easy to implement in practice. Moreover, this algorithm is more stable than joint optimization of g and W. We note that such gradient greedy approaches have been developed for traditional boosting [24].

3.1 Compositional Boosting

169

170

171

172

173

174

175

176

188

189

190

As one particular instantiation of our framework, we consider \mathcal{G}_t 's constructed by composing elements from a weak feature transformer class \mathcal{G} with the past iterates $\{\phi_i\}_{i=0}^{t-1}$ and study the resulting boosting 177 178 algorithms. We refer to such boosting algorithms as compositional boosting algorithms since the 179 strong feature transformer is constructed from weak feature transformer via function composition. 180 When $\mathcal{G}_t = \{g \circ \phi_{t-1} \text{ for } g \in \mathcal{G}\}$, the models in our framework have the ResNet architecture and can 181 be defined recurrently as $\phi_t = \phi_{t-1} + \eta_t g_t \circ \phi_{t-1}$. Moreover, Algorithm 1 with this choice of \mathcal{G}_t and 182 Algorithm 2 as update routine give us the greedy layer-wise supervised training technique proposed 183 by Bengio et al. [6] and recently revisited by Belilovsky et al. [5]. In another recent work, Huang 184 et al. [20] propose a boosting-based algorithm for learning ResNets. We now show that their approach 185 is equivalent to the greedy technique of Bengio et al. [6], and thus can be seen as an instance of our 186 general framework. We note that such a connection is not known previously. 187

Proposition 3.1. Suppose the classification loss ℓ is the exponential loss. Then the greedy technique of Huang et al. [20] for learning ResNets is equivalent to the greedy layer-wise supervised training technique of Bengio et al. [6].

In another recent work, Nitanda and Suzuki [27] propose a gradient boosting technique to greedily learn a ResNet. This algorithm is closely related to the gradient greedy approach described in Algorithm 3, with $\mathcal{G}_t = \{g \circ \phi_{t-1} \text{ for } g \in \mathcal{G}\}$.

¹Such scenarios arise if the feature transformations are non-differentiable functions such as decision trees.

Algorithm 1 Generalized Boosting

```
1: Input: Training data S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n, iterations T, initial linear predictor W_0, initial feature transformer \phi_0, learning rates \{\eta_i\}_{i=1}^T, Update-routine: UPDATE

2: t \leftarrow 1

3: while t \leq T do

4: Construct feature transformer class \mathcal{G}_t based on past iterates \{(W_i, \phi_i)\}_{i=0}^{t-1}

5: W_t, \phi_t \leftarrow \text{UPDATE}(S, W_{t-1}, \phi_{t-1}, \eta_t, \mathcal{G}_t)

6: t \leftarrow t + 1

7: end while

8: Return: W_T, \phi_T
```

Algorithm 2 Exact Greedy Update

Algorithm 3 Gradient Greedy Update

```
1: Input: Training data S, previous iterate (W, \phi), learning rate \eta, feature transformer class \mathcal{G}
2: W^+, g^+ \leftarrow \underset{\widetilde{W} \in \mathcal{W}, \widetilde{g} \in \mathcal{G}}{\operatorname{argmin}} \hat{R}_S(\widetilde{W}, \phi + \eta \widetilde{g})
4: Q^+ \leftarrow Q^+ + \eta g^+
4: Return: Q^+ \leftarrow Q^+ + \eta g^+
5: Q^+ \leftarrow Q^+ + \eta g^+
6: Q^+ \leftarrow Q^+ + \eta g^+
7: Return: Q^+ \leftarrow Q^+ + \eta g^+
7: Return: Q^+ \leftarrow Q^+ + \eta g^+
```

We now highlight certain drawbacks of the existing greedy layer-wise training techniques, which arise from the particular choice of \mathcal{G}_t used by these algorithms. Since $\{g \circ \phi_{t-1} \text{ for } g \in \mathcal{G}\}$ is constructed solely based on the past iterate ϕ_{t-1} , any mistake in ϕ_{t-1} is propagated to all the future iterates. As a result, these algorithms can not recover from their past mistakes. As an example, consider the following scenario where two points $\mathbf{x}_1, \mathbf{x}_2$ belonging to two different classes are placed close to each other in the feature space, after 1^{st} iteration of greedy; that is $\phi_1(\mathbf{x}_1) \approx \phi_1(\mathbf{x}_2)$. In such a scenario, the future iterates $\{\phi_t\}_{t=2}^\infty$ generated by existing greedy algorithms will always place $\mathbf{x}_1, \mathbf{x}_2$ close to each other in the representation space. As a result, the algorithm will always misclassify at least one of $\mathbf{x}_1, \mathbf{x}_2$. Another issue with existing greedy techniques is that they do not guarantee that the complexity of \mathcal{G}_t increases with time t. In such scenarios, Algorithm 1 doesn't make much progress in each iteration and can result in poor models. As an example, consider the setting where \mathcal{G} is the set of all linear transformations. Suppose ϕ_0 is the identity transform and ϕ_1 is such that its range lies in a low dimensional subspace. Then, it is evident that $\mathcal{G}_1 \supseteq \mathcal{G}_t$ for all $t \geqslant 2$.

To fix these issues, we propose two new compositional boosting algorithms obtained with a more careful choice of \mathcal{G}_t . In our first algorithm, which we call DenseCompBoost, we choose \mathcal{G}_t as follows

$$\mathcal{G}_t = \left\{ g \circ (\operatorname{Id} + \sum_{i=0}^{t-1} \alpha_i \phi_i), \text{ for } g \in \mathcal{G}, \alpha_i \in \mathbb{R} \right\},$$
(3)

where $\mathrm{Id}(\cdot)$ is the identify function. Such a choice of \mathcal{G}_t helps us recover from the past mistakes. For example, if ϕ_1 is a constant function, then the algorithm can still learn a good feature transformer by relying on the input x and the initial feature transform ϕ_0 . Moreover, our choice of \mathcal{G}_t ensures its complexity grows with t and satisfies: $\mathcal{G}_{t-1} \subseteq \mathcal{G}_t$, for all t. We call our algorithm DenseCompBoost, since the resulting model for this choice of \mathcal{G}_t resembles a DenseNet [21], where each layer is allowed to be connected to all the previous layers. That being said, the models output by DenseCompBoost differ from DenseNet in how they aggregate the previous layers. DenseNet concatenates the features from previous layers, whereas DenseCompBoost adds the features. Our second algorithm, which we call CmplxCompBoost, tries to increase the complexity of \mathcal{G}_t in each iteration as follows

$$\mathcal{G}_t = \left\{ g \circ \phi_{t-1}, \text{ for } g \in \widetilde{\mathcal{G}}_t \right\},\tag{4}$$

where $\widetilde{\mathcal{G}}_t$ is a weak feature transformer class and satisfies $\widetilde{\mathcal{G}}_{t-1} \subset \widetilde{\mathcal{G}}_t$ for all t. In the case of one layer neural networks, such $\widetilde{\mathcal{G}}_t$'s can be constructed by increasing the layer width with t. We note that the $\widetilde{\mathcal{G}}_t$ in this algorithm is independent of the past iterates. By increasing the complexity of $\widetilde{\mathcal{G}}_t$ with t, we expect the complexity of \mathcal{G}_t to increase and Algorithm 1 to make more progress in each iteration. While not immediately evident, we note that this technique can also fix the mistakes made by past iterates. For example, suppose ϕ_1 is such that it places two points $\mathbf{x}_1, \mathbf{x}_2$ from different classes, close to each other in the feature space. Then having a more complex $\widetilde{\mathcal{G}}_2$ can help recover from this mistake, as one can potentially find a $g \in \widetilde{\mathcal{G}}_2$ which can separate these two points. In Section 5, we present empirical evidence showing that our new boosting algorithms have superior performance over

existing additive and compositional boosting algorithms. Further empirical evidence corroborating the issues we identified with existing layer-wise training techniques can be found in Appendix J.1.

4 Excess Risk Bounds

In this section, we provide excess risk bounds for the models' output by the generalized boosting framework. Our results depend on a *weak learning condition* on the hypothesis class \mathcal{G}_t used in the t^{th} iteration of Algorithm 1. This condition is a way to quantify the relative strength of \mathcal{G}_t and roughly says that there always exists an element in \mathcal{G}_t which has an acute angle with the negative functional gradient at the current iterate. Such a condition ensures progress in each iteration of boosting.

Definition 4.1. Let $\beta \in (0,1]$, $\epsilon \ge 0$ be constants. \mathcal{G}_{t+1} is said to satisfy the (β,ϵ) -weak learning condition for a dataset S, if there exists a $g \in \mathcal{G}_{t+1}$ such that

$$\left\langle g, -\nabla_{\phi} \widehat{R}_{S}(W_{t}, \phi_{t}) \right\rangle_{P_{n}^{X}} \geqslant \beta B(\mathcal{G}_{t+1}) \|\nabla_{\phi} \widehat{R}_{S}(W_{t}, \phi_{t})\|_{P_{n}^{X}} - \epsilon,$$

where $B(\mathcal{G}_{t+1}) = \sup_{g \in \mathcal{G}_{t+1}} \|g\|_{P_n^X}$, and P_n is the empirical distribution of S.

In traditional boosting, such conditions are typically referred to as the edge of a weak learner and play a crucial role in the convergence analysis. For example, Freund and Schapire [14] assume that for any set of weights over the training set S, there exists a classifier in the hypothesis class of weak classifiers which has better than random accuracy on the weighted samples. The following proposition shows that their condition is closely related to Definition 4.1.

Proposition 4.1. For binary classification, the weak learning condition of Freund and Schapire [14] satisfies the empirical weak learning condition in Definition 4.1, albeit in the label space.

For binary classification problems, it is well known that the weak learning condition of [14] is the weakest condition under which boosting is possible [15, 29]. This, together with the above proposition, suggests that our weak learning condition in Definition 4.1 cannot be weakened for binary classification problems.

To begin with, we derive excess risk bounds for the gradient greedy approach. Our analysis crucially relies on the observation that it can be viewed as performing inexact gradient descent on the population risk R. Several recent works have analyzed inexact gradient descent on convex objectives [2, 12, 31, 32]. However, the condition on the inexact gradient imposed by these works is different from ours and in many cases is stronger than our condition. For example, the condition of Balakrishnan et al. [2] translates to $\|g + \nabla_{\phi} R(W, \phi)\|_{P^X} \leqslant \epsilon$ in our setting, which is stronger than our weak learning condition. So the core of our analysis focuses on understanding inexact gradient descent with descent steps satisfying the weak learning condition in Definition 4.1. In our analysis, we consider a sample-splitting variant of the algorithm, where in each iteration we use a fresh batch of samples. This is mainly done to simplify the analysis by avoiding complex statistical dependencies between the iterates of the algorithm. Let $\tilde{n} = \lfloor \frac{n}{T} \rfloor$, we split the training dataset S into T subsets $\{S_t\}_{t=1}^T$ of size \tilde{n} , where $S_t = \{(\mathbf{x}_{t,i}, y_{t,i})\}_{i=1}^{\tilde{n}}$. We work with the subset S_t in the t^{th} iteration of Algorithm 1. We are now ready to state our main result on the excess risk bounds of the iterates of Algorithm 3. Our results depend on the Rademacher complexity terms related of the hypothesis sets W, \mathcal{G}_t

$$\mathcal{R}\left(\mathcal{W}, \mathcal{G}_{t}\right) = \mathbb{E}\left[\sup_{\substack{W \in \mathcal{W}, \\ g \in \mathcal{G}_{t}}} \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \sum_{k=1}^{K} \rho_{ik} [Wg(\mathbf{x}_{t,i})]_{k}\right], \ \mathcal{R}\left(\mathcal{G}_{t}\right) = \mathbb{E}\left[\sup_{g \in \mathcal{G}_{t}} \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \sum_{j=1}^{D} \rho_{ij} [g(\mathbf{x}_{t,i})]_{j}\right],$$

where $[\mathbf{u}]_k$ denotes the k^{th} entry of a vector \mathbf{u} , and the expectation is taken w.r.t the randomness from S_t and the Rademacher random variables ρ_{ij} 's.

Theorem 4.1 (Gradient Greedy). Suppose the classification loss ℓ is L-Lipschitz and M-smooth w.r.t the first argument. Let the hypothesis set of linear predictors W be s.t. any $W \in W$ satisfies $\lambda_{min}(WW^T) \geqslant \sigma_{min}^2 > 0$ and $\lambda_{max}(WW^T) \leqslant \sigma_{max}^2$. Moreover, suppose for all t, \mathcal{G}_t satisfies the (β, ϵ_t) -weak learning condition of Definition 4.1 for any dataset S_t . Finally, suppose any $g \in \mathcal{G}_t$ is bounded with $\sup_X \|g(X)\|_2 \leqslant B$. Let the learning rates $\{\eta_t\}_{t=1}^{\infty}$ be chosen as $\eta_t = ct^{-s}$, for some $s \in \left(\frac{\beta+1}{\beta+2},1\right)$ and positive constant c. If Algorithm 1 is run for T iterations with Algorithm 3 as update routine, then (W_T, ϕ_T) , the T^{th} iterate output by the algorithm, satisfies the following risk bound for any W^* , ϕ^* and $\alpha \in (0, \beta(1-s))$, with probability at least $1-\delta$ over datasets of size n

$$R(W_T, \phi_T) \leqslant R(W^*, \phi^*) + O\left(\frac{1}{T^{\alpha}} + T^{2-s} \sqrt{\frac{\log \frac{T}{\delta}}{\tilde{n}}}\right) + 2\sum_{t=1}^{T} \eta_t \left(L\mathcal{R}\left(\mathcal{W}, \mathcal{G}_t\right) + L\mathcal{R}\left(\mathcal{G}_t\right) + \epsilon_t\right).$$

The $T^{-\alpha}$ term above corresponds to the *optimization error*, the $\eta_t \epsilon_t$ term corresponds to the *approximation error*, and the rest correspond to the *generalization error*. As T increases, $T^{-\alpha}$ goes down, and as \tilde{n} increases, the generalization error goes down. If there is no approximation error, that is if $\epsilon_t = 0$ for all t, then the excess risk goes down to 0 as $\tilde{n}, T \to \infty$ at an appropriate rate. Further discussion on this result can be found in Appendix D. We now extend the analysis of Theorem 4.1 to the exact greedy approach.

Corollary 4.1 (Exact Greedy). Consider the setting of Theorem 4.1. Suppose Algorithm 1 is run with Algorithm 2 as update routine. Then (W_T, ϕ_T) , the T^{th} iterate output by the algorithm, satisfies the same risk bounds as gradient greedy algorithm in Theorem 4.1.

In the rest of the section, we instantiate Theorem 4.1 for specific choices of \mathcal{G}_t . We first consider the additive representation boosting algorithm.

Corollary 4.2. Consider the setting of Theorem 4.1 and consider the additive representation boosting algorithm, where $\mathcal{G}_t = \mathcal{G}$ for all t. Suppose \mathcal{G} is the set of one layer neural networks with sigmoid activation functions: $\mathcal{G} = \{\sigma(C\mathbf{x}), \text{ for } C \in \mathbb{R}^{D \times d}, \|C_{i,*}\|_1 \leq \Lambda, \forall i\}$. Moreover, suppose the feature domain \mathcal{X} is a subset of $[0,1]^d$. Then the T^{th} iterate output by Algorithm 1, with Algorithm 2 or 3 as update routine, satisfies the following risk bound for any (W^*, ϕ^*) , with probability at least $1 - \delta$

$$R(W_T, \phi_T) \leqslant R(W^*, \phi^*) + O\left(\frac{1}{T^{\alpha}}\right) + 2\sum_{t=1}^T \eta_t \epsilon_t + O\left(\frac{KD\Lambda T^{1-s} \log D}{\sqrt{\tilde{n}}} + T^{2-s} \sqrt{\frac{\log \frac{T}{\delta}}{\tilde{n}}}\right).$$

Next, we consider the layer-by-layer fitting technique of Bengio et al. [6].

Corollary 4.3. Consider the setting of Corollary 4.2 and consider the layer-by-layer training technique of Bengio et al. [6], where $\mathcal{G}_t = \{g \circ \phi_{t-1} \text{ for } g \in \mathcal{G}\}$. Suppose \mathcal{G} is the set of one layer neural networks with sigmoid activation functions: $\mathcal{G} = \{\sigma(C\mathbf{x}), \text{ for } C \in \mathbb{R}^{D \times D}, \|C_{i,*}\|_1 \leq \Lambda, \forall i\}$. Then the T^{th} iterate output by Algorithm 1, with Algorithm 2 or 3 as update routine, satisfies the following risk bound for any (W^*, ϕ^*) with probability at least $1 - \delta$

$$R(W_T, \phi_T) \leqslant R(W^*, \phi^*) + O\left(\frac{1}{T^{\alpha}}\right) + 2\sum_{t=1}^T \eta_t \epsilon_t + O\left(\frac{KD\Lambda T^{2-2s} \log D}{\sqrt{\tilde{n}}} + T^{2-s} \sqrt{\frac{\log \frac{T}{\delta}}{\tilde{n}}}\right).$$

Note that the generalization and optimization errors for both additive feature boosting and layer-by-layer fitting have similar dependence on T, \tilde{n} . However, the latter tends to have a smaller approximation error (ϵ_t) as it is able to build complex \mathcal{G}_t 's over time. So one would expect layer-by-layer fitting to output models with a better population risk, which our empirical results in fact verify.

5 Experiments

287

289

290

298

300

301 302

303

304

305

318

319

320

321

322

In this section, we present experiments comparing the performance of various boosting techniques on both simulated and benchmark datasets.

Baselines. We compare our proposed boosting techniques with XGBoost, AdaBoost, additive 306 representation boosting (discussed in Corollary 4.2) and greedy layer-by-layer training technique of 307 Bengio et al. [6] (Corollary 4.3). XGBoost uses decision stumps as weak classifiers. For AdaBoost, 308 we use 1 hidden layer neural networks as weak classifiers. We use two kinds of neural networks, 309 based on the dataset. For tabular datasets, we use fully connected networks and for image datasets, we 310 use convolutional networks (CNN) with the convolution block made up of Convolution, BatchNorm, 311 ReLU layers arranged sequentially. For additive representation boosting (Additive Feature Boost 312 from now on) and layer-by-layer fitting (StdCompBoost from now on), the weak feature transformer 313 class \mathcal{G} consists of one layer neural network transformations. Similar to AdaBoost, we use two kinds 314 of transformations: a) fully connected transformations of the form $g(\mathbf{x}) = \text{ReLU}(C\mathbf{x} + \mathbf{d})$, and b) 315 convolutional transformations with Convolution, BatchNorm, ReLU blocks arranged sequentially. 316 Finally, we also compare against end-to-end training of neural networks. 317

Proposed Techniques. For DenseCompBoost, we consider two choices for \mathcal{G} : one based on fully connected blocks and the other based on convolution blocks. For CmplxCompBoost, we again consider two choices for the weak transformer class $\tilde{\mathcal{G}}_t$ in Equation (4): a) ReLU($C\mathbf{x} + \mathbf{d}$) with $C \in \mathbb{R}^{D_t \times D_{t-1}}$, where $D_t = D_{t-1} + \Delta$ for some positive constant Δ , and b) convolution blocks with number of output channels equal to the number of input channels plus a constant Δ . This choice of feature transformers ensures the complexity of $\tilde{\mathcal{G}}_t$ increases with t. We use exact greedy updates (Algorithm 2) for both of our proposed methods and set learning rate η_t to 1.

Table 1: Test accuracy of various boosting techniques on synthetic datasets. Numbers in bold indicate the best performance among various greedy techniques. The row corresponding to *half width* shows the performance of DenseCompBoost using layers with a width equal to half of the best width.

Technique	Simulation 1	Simulation 2	Simulation 3	
XGBoost (Trees)	84.40	97.59	50.10	
AdaBoost (1 NN)	67.90	93.73	72.64	
Additive Feature Boost	88.49	93.91	73.13	
StdCompBoost	91.53	96.95	82.49	
DenseCompBoost	93.55	98.35	95.70	
DenseCompBoost (half width)	93.54	97.99	94.37	
CmplxCompBoost	91.97	97.22	82.52	
End-to-End	93.88	98.35	99.09	

Table 2: Test accuracy of various boosting techniques on benchmark datasets. We use convolution blocks for the first 5 datasets and fully connected blocks for the other datasets.

Technique	SVHN	FashionMNIST	CIFAR10	Convex	MNIST-rot- back-image	MNIST	Letter	CovType	Connect4
XGBoost (Trees)	77.72	90.34	58.34	82.29	53.89	97.96	96.16	97.46	86.63
AdaBoost (1 NN)	82.88	88	72.78	86.17	50.02	98.27	92.08	90.95	86.39
Additive Feature Boost	83.36	89.95	74.33	89.30	54.31	98.27	90.86	93.12	86.58
StdCompBoost	90.81	92.77	81.93	98.19	73.17	98.37	96.43	95.61	86.33
DenseCompBoost	91.03	93.17	82.31	98.6	73.1	98.34	96.96	96.28	86.85
CmplxCompBoost	91.25	93.18	82.43	98.52	74.32	98.34	96.66	95.92	86.49
End-to-End	94.82	93.49	86.88	98.81	82.69	98.95	97.67	96.86	87.37

Results: The baseline and proposed methods are tested on 3 simulated datasets involving binary classification tasks with polynomial decision boundaries, and 9 benchmark datasets consisting of both tabular and image data. Details of the datasets and hyperparameters can be found in Appendix J.

Table 1 presents the results on simulated datasets. Both CmplxCompBoost and StdCompBoost largely outperform the additive boosting methods, with CmplxCompBoost being slightly better due to the increasing complexity in \tilde{G}_t . Notably, DenseCompBoost performs significantly better than the rest and is able to bridge the gap between StdCompBoost and End-to-End. We attribute its success to its ability to recover from earlier mistakes: while StdCompBoost or CmplxCompBoost necessarily accumulate errors at each layer, DenseCompBoost is further connected to earlier layers, allowing it to undo its past mistakes. To verify this, we corrupt the weights of the first layers of models learned using StdCompBoost and DenseCompBoost, during the training process. We observed that DenseCompBoost is barely affected by a poor first layer (0.7% drop in accuracy), whereas StdCompBoost suffers a significant performance drop (about 10%) (see Table 3 in Appendix J.1).

Results on benchmark datasets are presented in Table 2. It can be seen that additive boosting techniques have poor performance on image classification tasks. Among compositional boosting methods, StdCompBoost performs the worst. While DenseCompBoost performs comparably to CmplxCompBoost on image datasets, it is better on tabular data. We believe a hybrid of DenseCompBoost and CmplxCompBoost algorithms can achieve better performance than either of the algorithms.

6 Conclusion

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346 347

348

349

350

351

We proposed a generalized framework for boosting, which allows for more complex forms of aggregation of weak learners than traditional boosting. Our generalized framework allows to derive learning algorithms that (a) have performance close to that of end-to-end trained DNNs, and (b) come with strong theoretical guarantees. Additive boosting algorithms do not satisfy property (a), while DNNs do not satisfy property (b). In particular, additive boosting algorithms, even with small neural networks as their weak classifiers, do not not have the strong performance of end-to-end trained DNNs. Improving their performance requires the hypothesis space to increase in complexity while not increasing sample complexity of each boosting step too greatly, which can be achieved by our generalized boosting framework. One particular instantiation of our framework is aggregation using function compositions. A number of existing greedy techniques for learning neural networks fall into our framework, and our analysis allowed us to delineate some of their key flaws, then consequently, propose new techniques which improve upon them. We believe our work opens up a new line of inquiry for greedy learning of highly flexible models with rigorous theoretical guarantees, by leveraging the theory of boosting and generalized greedy algorithms in function spaces. We moreover believe our work has the potential to bridge the gap in performance between existing greedy layer-by-layer training techniques and end-to-end training of deep networks.

360 Broader Impact

In the long term, our work can potentially lead to XGBoost like packages for efficiently learning complex machine learning models and can have similar societal consequences as XGBoost. In the short term, this work does not present any societal consequences.

364 References

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing* systems, pages 6155–6166, 2019.
- Sivaraman Balakrishnan, Martin J Wainwright, Bin Yu, et al. Statistical guarantees for the em
 algorithm: From population to sample-based analysis. *The Annals of Statistics*, 45(1):77–120,
 2017.
- [3] Pierre Baldi, Peter Sadowski, and Daniel Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature communications*, 5:4308, 2014.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [5] Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Greedy layerwise learning can scale to imagenet. *arXiv preprint arXiv:1812.11446*, 2018.
- Francisco [6] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise training of deep networks. In *Advances in neural information processing systems*, pages 153–160, 2007.
- [7] Catherine L Blake and Christopher J Merz. Uci repository of machine learning databases, 1998, 1998.
- [8] Chang Chen, Zhiwei Xiong, Xinmei Tian, and Feng Wu. Deep boosting for image denoising.
 In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018.
- [9] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [10] Corinna Cortes, Mehryar Mohri, and Umar Syed. Deep boosting. 2014.
- Corinna Cortes, Xavier Gonzalvo, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. Adanet:
 Adaptive structural learning of artificial neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 874–883. JMLR. org, 2017.
- Olivier Devolder, François Glineur, and Yurii Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146(1-2):37–75, 2014.
- 133 Yoav Freund. Boosting a weak learning algorithm by majority. *Information and computation*, 121(2):256–285, 1995.
- Yoav Freund and Robert E Schapire. A desicion-theoretic generalization of on-line learning and an application to boosting. In *European conference on computational learning theory*, pages 23–37. Springer, 1995.
- Yoav Freund and Robert E Schapire. Game theory, on-line prediction and boosting. In
 Proceedings of the ninth annual conference on Computational learning theory, pages 325–332,
 1996.
- [16] Yoav Freund, Robert E Schapire, et al. Experiments with a new boosting algorithm. In *icml*, volume 96, pages 148–156. Citeseer, 1996.
- ⁴⁰³ [17] Peter W Frey and David J Slate. Letter recognition using holland-style adaptive classifiers. ⁴⁰⁴ *Machine learning*, 6(2):161–182, 1991.
- for [18] Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
 pages 770–778, 2016.

- Furong Huang, Jordan Ash, John Langford, and Robert Schapire. Learning deep resnet blocks sequentially using boosting theory. *arXiv preprint arXiv:1706.04964*, 2017.
- 413 [21] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected 414 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern* 415 *recognition*, pages 4700–4708, 2017.
- [22] Hugo Larochelle, Dumitru Erhan, Aaron Courville, James Bergstra, and Yoshua Bengio. An
 empirical evaluation of deep architectures on problems with many factors of variation. In
 Proceedings of the 24th international conference on Machine learning, pages 473–480, 2007.
- [23] Sindy Löwe, Peter O'Connor, and Bastiaan Veeling. Putting an end to end-to-end: Gradient-isolated learning of representations. In *Advances in Neural Information Processing Systems*, pages 3033–3045, 2019.
- Llew Mason, Jonathan Baxter, Peter L Bartlett, and Marcus R Frean. Boosting algorithms as gradient descent. In *Advances in neural information processing systems*, pages 512–518, 2000.
- 424 [25] Andreas Maurer. A vector-contraction inequality for rademacher complexities. In *International Conference on Algorithmic Learning Theory*, pages 3–17. Springer, 2016.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.
 Reading digits in natural images with unsupervised feature learning. 2011.
- 428 [27] Atsushi Nitanda and Taiji Suzuki. Functional gradient boosting based on residual network perception. *arXiv preprint arXiv:1802.09031*, 2018.
- [28] Natalia Ponomareva, Thomas Colthurst, Gilbert Hendry, Salem Haykal, and Soroush Radpour.
 Compact multi-class boosted trees. In 2017 IEEE International Conference on Big Data (Big Data), pages 47–56. IEEE, 2017.
- 433 [29] Gunnar Rätsch and Manfred K Warmuth. Efficient margin maximizing with boosting. *Journal*434 of Machine Learning Research, 6(Dec):2131–2152, 2005.
- [30] Robert E Schapire. The strength of weak learnability. *Machine learning*, 5(2):197–227, 1990.
- 436 [31] Mark Schmidt, Nicolas L Roux, and Francis R Bach. Convergence rates of inexact proximal-437 gradient methods for convex optimization. In *Advances in neural information processing* 438 *systems*, pages 1458–1466, 2011.
- [32] Vladimir Nikolaevich Temlyakov. Greedy expansions in convex optimization. *Proceedings of the Steklov Institute of Mathematics*, 284(1):244–262, 2014.
- [33] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
 benchmarking machine learning algorithms, 2017.