

Boosting for Comparison-Based Learning

Michaël Perrot¹, Ulrike von Luxburg^{1,2}

¹Max-Planck-Institute for Intelligent Systems, Tübingen, Germany

²University of Tübingen, Department of Computer Science, Tübingen, Germany
michael.perrot@tuebingen.mpg.de, ulrike.luxburg@uni-tuebingen.de

Abstract

We consider the problem of classification in a comparison-based setting: given a set of objects, we only have access to triplet comparisons of the form *object x_i is closer to object x_j than to object x_k* . In this paper we introduce TripletBoost, a new method that can learn a classifier just from such triplet comparisons. The main idea is to aggregate the triplets information into weak classifiers, which can subsequently be boosted to a strong classifier. Our method has two main advantages: (i) it is applicable to data from any metric space, and (ii) it can deal with large scale problems using only passively obtained and noisy triplets. We derive theoretical generalization guarantees and a lower bound on the number of necessary triplets, and we empirically show that our method is both competitive with state of the art approaches and resistant to noise.

1 Introduction

In the past few years the problem of comparison-based learning has attracted growing interest in the machine learning community [Agarwal *et al.*, 2007; Jamieson and Nowak, 2011; Tamuz *et al.*, 2011; Tschopp *et al.*, 2011; Van Der Maaten and Weinberger, 2012; Heikinheimo and Ukkonen, 2013; Amid and Ukkonen, 2015; Kleindessner and Luxburg, 2015; Jain *et al.*, 2016; Haghiri *et al.*, 2017; Kazemi *et al.*, 2018]. The motivation is to relax the assumption that an explicit representation of the objects or a distance metric between pairs of examples are available. Instead one only has access to a set of ordinal distance comparisons that can take several forms depending on the problem at hand. In this paper we focus on triplet comparisons of the form *object x_i is closer to object x_j than to object x_k* , that is on relations of the form $d(x_i, x_j) < d(x_i, x_k)$ where d is an unknown metric¹.

We address the problem of classification with noisy triplets that have been obtained in a passive manner: the examples lie in an unknown metric space, not necessarily Euclidean, and

we are only given a small set of triplet comparisons — there is no way in which we could actively ask for more triplets. Furthermore we assume that the answers to the triplet comparisons can be noisy. To deal with this problem one can try to first recover an explicit representation of the examples, a task that can be solved by ordinal embedding approaches [Agarwal *et al.*, 2007; Van Der Maaten and Weinberger, 2012; Terada and von Luxburg, 2014; Jain *et al.*, 2016], and then apply standard machine learning approaches. However, such embedding methods assume that the examples lie in a Euclidean space and do not scale well with the number of examples: typically they are too slow for datasets with more than 10^3 examples. As an alternative, it would be desirable to have a classification algorithm that can work with triplets directly, without taking a detour via ordinal embedding. To the best of our knowledge, for the case of passively obtained triplets, this problem has not yet been solved in the literature.

Another interesting question in this context is that of the minimal number of triplets required to successfully learn a classifier. It is known that to exactly recover an ordinal embedding one needs of the order $\Omega(n^3)$ passively queried triplets in the worst case (essentially all of them), unless we make stronger assumptions on the underlying metric space [Jamieson and Nowak, 2011]. However, classification is a problem which seems simpler than ordinal embedding, and thus it might be possible to obtain better lower bounds.

In this paper we propose TripletBoost, a method for classification that is able to learn using only passively obtained triplets while not making any assumptions on the underlying metric space. To the best of our knowledge this is the first approach that is able to solve this problem. Our method is based on the idea that the triplets can be aggregated into simple *triplet classifiers*, which behave like decision stumps and are well-suited for boosting approaches [Schapire and Freund, 2012]. From a theoretical point of view we prove that our approach learns a classifier with low training error, and we derive generalization guarantees that ensure that its error on new examples is bounded. Furthermore we derive a new lower bound on the number of triplets that are necessary to ensure useful predictions. From an empirical point of view we demonstrate that our approach can be applied to datasets that are several order of magnitudes larger than the ones that can currently be handled by ordinal embedding methods. Furthermore we show that our method is quite resistant to noise.

¹Note that this kind of ordinal information is sometimes just used as side information [Bellet *et al.*, 2015; Kane *et al.*, 2017], but, as the references in the main text, we focus on the setting where ordinal comparisons are the sole information available.

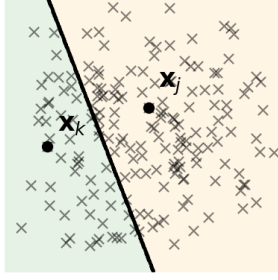


Figure 1: Example of a triplet classifier. Given two reference points x_j and x_k , the space is divided in two half-spaces: examples closer to x_j and examples closer to x_k . Triplet information is enough to reveal which half-space a point x_i is in.

2 The TripletBoost Algorithm

In this paper we are interested in multi-class classification problems. Let (\mathcal{X}, d) be an unknown and general metric space, typically not Euclidean. Let \mathcal{Y} be a finite label space. Let $S = \{(x_i, y_i)\}_{i=1}^n$ be a set of n examples drawn i.i.d. from an unknown distribution \mathcal{D}_S defined over $\mathcal{X} \times \mathcal{Y}$. Note that we use the notation x_i as a convenient way to identify an object; it does not correspond to any explicit representation that could be used by an algorithm (such as coordinates in a vector space). Let $T = \{(x_i, x_j, x_k) : (x_i, y_i), (x_j, y_j), (x_k, y_k) \in S, x_j \neq x_k\}$ be a set of m triplets. Each ordered tuple $(x_i, x_j, x_k) \in T$ encodes the following relation between the three examples:

$$d(x_i, x_j) < d(x_i, x_k). \quad (1)$$

Given the triplets in T and the label information of all points, our goal is to learn a classifier. We make two main assumptions about the data. First, we assume that the triplets are uniformly and independently sampled from the set of all possible triplets. Second, we assume that the triplets in T can be noisy (that is the inequality has been swapped, for example $(x_i, x_k, x_j) \in T$ while the true relation is $d(x_i, x_j) < d(x_i, x_k)$), but the noise is uniform and independent from one triplet to another. In the following \mathbb{I}_a denotes the indicator function returning 1 when a property a is verified and 0 otherwise, \mathcal{W}_c is an empirical distribution over $S \times \mathcal{Y}$, and $w_{c, x_i, y}$ is the weight associated to object x_i and label y .

2.1 Weak Triplet Classifiers

Rather than considering triplets as individual pieces of information we propose to aggregate them into decision stumps that we call *triplet classifiers*. The underlying idea is to select two reference examples, x_j and x_k , and to divide the space into two *half-spaces*: examples that are closer to x_j and examples that are closer to x_k . This is illustrated in Figure 1. In principle, this can be achieved with triplets only. However, a major difficulty in our setting is that our triplets are passively obtained: for most training points x_i we do not know whether they are closer to x_j or x_k . In particular, it is impossible to evaluate the classification accuracy of such a simple classifier on the whole training set. To deal with this problem we propose to use an abstention scheme where a triplet classifier

Algorithm 1 TripletBoost: Boosting with triplet classifiers.

Input: $S = \{(x_i, y_i)\}_{i=1}^n$ a set of n examples,
 $T = \{(x_i, x_j, x_k) : x_j \neq x_k\}$ a set of m triplets.

Output: $H(\cdot)$ a strong classifier.

- 1: Let \mathcal{W}_1 be the empirical uniform distribution:
 $\forall (x_i, y_i) \in S, \forall y \in \mathcal{Y}, w_{1, x_i, y} = \frac{1}{n|\mathcal{Y}|}$.
- 2: **for** $c = 1, \dots, C$ **do**
- 3: Choose a triplet classifier h_c according to \mathcal{W}_c (Equation (2)).
- 4: Compute the weight α_c of h_c according to its performance on \mathcal{W}_c (Equation (3)).
- 5: Update the weights of the examples to obtain a new distribution \mathcal{W}_{c+1} (Equation (4)).

6: **end for**

7: **return** $H(\cdot) = \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(\cdot) \neq y \wedge y \in h_c(\cdot)} \right)$

abstains if it does not know on which side of the hyperplane the considered point lies. Given a set T of triplets and two reference points x_j and x_k , we define a triplet classifier as:

$$h_{j,k}(x) = \begin{cases} o_j & \text{if } (x, x_j, x_k) \in T, \\ o_k & \text{if } (x, x_k, x_j) \in T, \\ \vartheta & \text{otherwise.} \end{cases}$$

In our multi-class setting, o_j and o_k will be sets of labels, that is $o_j, o_k \subseteq \mathcal{Y}$. In Section 2.2 we describe how we choose them in a data dependent fashion to obtain classifiers with minimal error on the training set. The prediction ϑ simply means that the triplet classifier abstains on the example. Let \mathcal{H} denote the set of all possible triplet classifiers.

The triplet classifiers are very simple and, in practice, we do not expect them to perform well at all. But we prove in Section 3.1 that, for appropriate choices of o_j and o_k , they are at least as good as random predictors. This is all we need to ensure that they can be used successfully in a boosting framework. The next section describes how this works.

2.2 TripletBoost

Boosting is based on the insight that weak classifiers (that is classifiers marginally better than random predictors) are usually easy to obtain and can be combined in a weighted linear combination to obtain a strong classifier. This weighted combination can be obtained in an iterative fashion where, at each iteration, a weak-classifier is chosen and weighted so as to minimize the error on the training examples. The weights of the points are then updated to put more focus on hard-to-classify examples [Schapire and Freund, 2012]. In this paper we use a well-known boosting algorithm called AdaBoost.MO [Schapire and Singer, 1999; Schapire and Freund, 2012]. This method can handle multi-class problems with a one-against-all approach, works with abstaining classifiers and is theoretically well founded. Algorithm 1 summarizes the main steps of our approach that we detail below.

Choosing a triplet classifier. To choose a triplet classifier we proceed in two steps. In the first step, we select two reference points x_j and x_k such that $y_j \neq y_k$. This is done by randomly sampling from an empirical distribution $\mathcal{W}_{c, \mathcal{X}}$ on the

examples. Here $\mathcal{W}_{c,\mathcal{X}}$ denotes the marginal distribution of \mathcal{W}_c with respect to the examples. This distribution is updated at each iteration to put more focus on those parts of the space that are hard to classify while promoting triplet classifiers that are able to separate different classes (see Equation 4).

In the second step, we choose o_j and o_k , the sets of labels that should be predicted for each half space of the triplet classifier. Given one of the half spaces, we propose to add a label to the set of predicted labels if the weight of examples of this class is greater than the weight of examples of different classes. Formally, with $w_{c,x_i,y}$ defined as in Algorithm 1, we construct o_j as follows:

$$o_j = \left\{ y : \sum_{\substack{(x_i,y_i) \in S, \\ (x_i,x_j,x_k) \in T}} (\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) w_{c,x_i,y} > 0 \right\}. \quad (2)$$

We construct o_k in a similar way. The underlying idea is that adding h_c to the current combination of triplet classifiers should improve the predictions on the training set as much as possible. In Section 3.1 we show that this strategy is optimal and that it ensures that the selected triplet classifier is either a weak classifier or has a weight α_c of 0.

Computing the weight of the triplet classifier. To choose the weight of the triplet classifier h_c we start by computing $W_{c,+}$ and $W_{c,-}$, the weights of correctly and incorrectly classified examples:

$$\begin{aligned} W_{c,+} &= \sum_{\substack{(x_i,y_i) \in S, \\ h_c(x_i) \neq \vartheta}} \left(\mathbb{I}_{y_i \in h_c(x_i)} w_{c,x_i,y_i} + \sum_{y \neq y_i} \mathbb{I}_{y \notin h_c(x_i)} w_{c,x_i,y} \right), \\ W_{c,-} &= \sum_{\substack{(x_i,y_i) \in S, \\ h_c(x_i) \neq \vartheta}} \left(\mathbb{I}_{y_i \notin h_c(x_i)} w_{c,x_i,y_i} + \sum_{y \neq y_i} \mathbb{I}_{y \in h_c(x_i)} w_{c,x_i,y} \right). \end{aligned} \quad (3)$$

We then set $\alpha_c = \log \left(\frac{w_{c,+} + \frac{1}{n}}{w_{c,-} + \frac{1}{n}} \right) / 2$. The term $\frac{1}{n}$ is a smoothing constant [Schapire and Singer, 2000]: in our setting with few, passively queried triplets it helps to avoid numerical problems that might arise when $W_{c,+}$ or $W_{c,-} = 0$. In Theorem 2 we show that this choice of α_c leads to a decrease in training error as the number of iterations increases.

Updating the weights of the examples. In each iteration of our algorithm, a new triplet classifier h_c is added to the weighted combination of classifiers, and we need to update the empirical distribution \mathcal{W}_c over the examples for the next iteration. The idea is (i) to reduce the weights of correctly classified examples, (ii) to keep constant the weights of the examples for which the current triplet classifier abstains, and (iii) to increase the weights of incorrectly classified examples. The weights are then normalized by a factor Z_c so that \mathcal{W}_{c+1} remains an empirical distribution over the examples. Formally, $\forall (x_i, y_i) \in S, \forall y \in \mathcal{Y}$, if $h_c(x_i) = \vartheta$ then $w_{c+1,x_i,y} = \frac{w_{c,x_i,y}}{Z_c}$ and if $h_c(x_i) \neq \vartheta$ then

$$w_{c+1,x_i,y} = \frac{w_{c,x_i,y}}{Z_c} \exp \left\{ -\alpha_c (\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) \times (\mathbb{I}_{y \in h_c(x_i)} - \mathbb{I}_{y \notin h_c(x_i)}) \right\}. \quad (4)$$

Using H for prediction. Given a new example x , TripletBoost predicts its label as

$$H(x) = \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta \wedge y \in h_c(x)} \right) \quad (5)$$

that is the label with the highest weight as predicted by the weighted combination of selected triplet classifiers. However, recall that we are in a passive setting, and thus we assume that we are given a set of triplets $T_x = \{(x, x_j, x_k) : (x_j, y_j), (x_k, y_k) \in S, x_j \neq x_k\}$ associated with the example x (but there is no way to choose them). Hence, some of the triplets in T_x correspond to triplet classifiers in H (that is $(x, x_j, x_k) \in T_x$ and the reference points for h_c were x_j and x_k) and some do not. In particular, it might happen that none of the triplets in T_x corresponds to a triplet classifier in H and, in this case, H can only randomly predict a label. In Section 3.3 we provide a lower bound on the number of triplets necessary to avoid this behaviour. The main computational bottleneck when predicting the label of a new example x is to check whether the triplets in T_x match a triplet classifier in H . A naive implementation would compare each triplet in T_x to each triplet classifier, which can be as expensive as $O(|T_x|C)$. Fortunately, by first sorting the triplets and the triplet classifiers, a far more reasonable complexity of $O(|T_x| \log(|T_x|) + C \log(C))$ can be achieved.

3 Theoretical Analysis

In this section we show that our approach is theoretically well founded. First we prove that the triplet classifiers with non-zero weights are weak learners: they are slightly better than random predictors (Theorem 1). Building upon this result we show that, as the number of iterations increases, the training error of the strong classifier learned by TripletBoost is decreased (Theorem 2). Then, to ensure that TripletBoost does not over-fit, we derive a generalization bound showing that, given a sufficient amount of training examples, the test error is bounded (Theorem 3). Finally, we derive a lower bound on the number of triplets necessary to ensure that TripletBoost does not learn a random predictor (Theorem 4).

3.1 Triplet Classifiers and Weak Learners

We start this theoretical analysis by showing that the strategy to choose the predicted labels of triplet classifiers described in Equation (2) is optimal: it ensures that their error is minimal on the training set (compared to any other labelling strategy). We also show that the triplet classifiers are never worse than random predictors and in fact, that only those triplets classifiers that are weak classifiers (strictly better than random classifiers) are affected a non-zero weight. This is summarized in the next theorem.

Theorem 1 (Triplet classifiers and weak learners). *Let \mathcal{W}_c be an empirical distribution over $S \times \mathcal{Y}$ and h_c be the corresponding triplet classifier chosen as described in Section 2.2. It holds that:*

1. *the error of h_c on \mathcal{W}_c is at most the error of a random predictor and is minimal compared to other labelling strategies,*

2. the weight α_c of the classifier is non-zero if and only if, h_c is a weak classifier, that is strictly better than a random predictor.

Proof. The proof is given in Appendix E. \square

3.2 Boosting guarantees

From a theoretical point of view the boosting framework has been extensively investigated and it has been shown that most AdaBoost-based methods decrease the training error at each iteration [Freund and Schapire, 1997]. Another question that has attracted a lot of attention is the problem of generalization. It is known that when the training error has been minimized, AdaBoost-based methods often do not over-fit and it might even be beneficial to further increase the number of weak learners. A popular explanation is the margin theory which says that as the number of iterations increases, the confidence of the algorithm in its predictions increases and thus, the test accuracy is improved [Schapire *et al.*, 1998; Breiman, 1999; Wang *et al.*, 2011; Gao and Zhou, 2013]. TripletBoost is based on AdaBoost.MO, and thus it inherits the theoretical guarantees presented above. In this section, we provide two theorems which show (i) that TripletBoost reduces the training error as the number of iterations increases, and (ii) that it generalizes well to new examples.

The following theorem shows that, as the number of iterations increases, TripletBoost decreases the training error.

Theorem 2 (Reduction of the Training Error). *Let S be a set of n examples and T be a set of m triplets (obtained as described in Section 2). Let $H(\cdot)$ be the classifier obtained after C iterations of TripletBoost (Algorithm 1) using S and T as input. It holds that:*

$$\mathbb{P}_{(x,y) \in S} [H(x) \neq y] \leq \frac{|\mathcal{Y}|}{2} \prod_{c=1}^C Z_c$$

$$\text{with } Z_c = (1 - W_{c,+} - W_{c,-}) + (W_{c,+}) \cdot \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + (W_{c,-}) \cdot \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \leq 1.$$

Proof. This result, inherited from AdaBoost.MO [Schapire and Freund, 2012], is proven in Appendix F. \square

The next theorem shows that the true error of a classifier learned by TripletBoost can be bounded by a quantity related to the confidence of the classifier on the training examples, with respect to a margin, plus a term which decreases as the number of examples increases. The confidence of the classifier on the training examples is also bounded and decreases for sufficiently small margins.

Theorem 3 (Generalization Guarantees). *Let \mathcal{D}_S be a distribution over $\mathcal{X} \times \mathcal{Y}$, let S be a set of n examples drawn i.i.d. from \mathcal{D}_S , and let T be a set of m triplets (obtained as described in Section 2). Let $H(\cdot)$ be the classifier obtained after C iterations of TripletBoost (Algorithm 1) using S and T as input. Let \mathcal{H} be a set of triplet classifiers as defined in Section 2.1. Then, given a margin parameter $\theta > \sqrt{\frac{\log |\mathcal{H}|}{16|\mathcal{Y}|^2 n}}$*

and a measure of the confidence of $H(\cdot)$ in its predictions $\theta_H(x, y)$, with probability at least $1 - \delta$, we have that

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) \neq y] \leq \mathbb{P}_{(x,y) \in S} [\theta_H(x, y) \leq \theta] + \mathcal{O} \left(\sqrt{\frac{\log \left(\frac{1}{\delta} \right)}{n}} + \log \left(\frac{|\mathcal{Y}|^2 n \theta^2}{\log |\mathcal{H}|} \right) \frac{\log |\mathcal{H}|}{n \theta^2} \right)$$

Furthermore we have that

$$\mathbb{P}_{(x,y) \in S} [\theta_H(x, y) \leq \theta] \leq \frac{|\mathcal{Y}|}{2} \prod_{c=1}^C Z_c \sqrt{\left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right)^\theta}.$$

Proof. This result, inherited from AdaBoost.MO [Schapire and Freund, 2012], is proven in Appendix G. \square

At a first glance it seems that this bound does not depend on m , the number of available triplets. However, this dependency is implicit: m impacts the probability that the training examples are well classified with a large margin θ . If the number of triplets is small, the probability that the training examples are well classified with a given margin is small. This probability increases when the number of triplets increases. We illustrate this behaviour in Appendix A.

To prove a bound that explicitly depends on m would be of significant interest. However this is a difficult problem, as it requires to use an explicit measure of complexity for general metric spaces, which is beyond the scope of this paper.

3.3 Lower bound on the number of triplets

In this section we investigate the minimum number of triplets that are necessary to ensure that our algorithm performs well. Ideally, we would like to obtain a lower bound on the number of triplets that are necessary to achieve a given accuracy. In this paper we take a first step in this direction by deriving a bound on the number of triplets that are necessary to ensure that the learned classifier does not abstain on any unseen example. Theorems 4 and 5 show that it abstains with high probability if it is learned using too few triplets or if it combines too few triplet classifiers.

Theorem 4 (Lower bound on the probability that a strong classifier abstains). *Let $n \geq 2$ be the number of training examples, $p = \frac{2n^k}{n^3}$ with $k \in \left[0, 3 - \frac{\log(2)}{\log(n)}\right)$ be the probability that a triplet is available in the triplet set T and $C = \frac{n^\beta}{2}$ with $\beta \in \left[0, 1 + \frac{\log(n-1)}{\log(n)}\right]$ be the number of classifiers combined in the learned classifier. Let \mathcal{A} be any algorithm learning a classifier $H(\cdot) = \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{y \in h_c(\cdot)} \right)$ that combines several triplet classifiers using some weights $\alpha_c \in \mathbb{R}$. Assume that triplet classifiers that abstain on all the training examples have a weight of 0 (that is if $h_c(x_i) = \vartheta$ for all the examples $(x_i, y_i) \in S$ then $\alpha_c = 0$). Then the probability that H abstains on a test example is bounded as follows:*

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \left(1 - p + p(1 - p)^n \right)^C. \quad (6)$$

Proof. The proof is given in Appendix H. \square

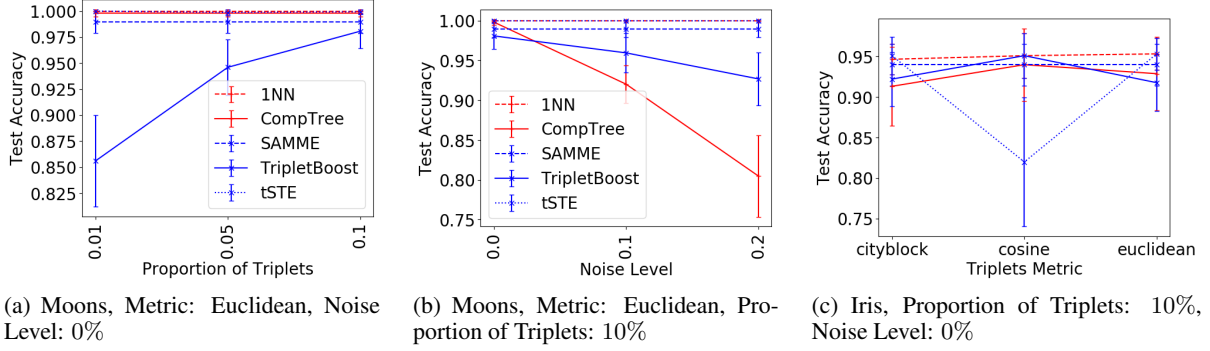


Figure 2: Moons is a small scale dataset with 500 examples and 2 dimensions. We fix the triplets metric to the Euclidean distance. In Figure 2a we consider the noise free setting and we vary the proportion of available triplets from 1 to 10% of all the triplets. In Figure 2b we fix this proportion to 10% and we vary the noise level from 0 to 20%. Iris is a small scale dataset with 150 examples and 4 dimensions. In Figure 2c we fix the proportion of available triplets to 10%, the noise level to 0% and we vary the metric considered to generate the triplets.

To understand the implications of this theorem we consider a concrete example.

Example 1. Assume that we build a linear combination of all possible triplet classifiers, that is $C = \frac{n(n-1)}{2}$. Then we have

$$\lim_{n \rightarrow +\infty} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \begin{cases} 1 & \text{if } k < \frac{3}{2}, \\ \exp(-2) & \text{if } k = \frac{3}{2}, \\ 0 & \text{if } \frac{3}{2} < k, \end{cases} \quad (7)$$

where k is the parameter that controls the probability p that a particular triplet is available in the triplets set T . The bottom line is that when $k < \frac{3}{2}$, that is when we do not have at least $\Omega(n\sqrt{n})$ random triplets, the learned classifier abstains on all the examples.

Proof. The proof is given in Appendix I. \square

Theorem 4 shows that when p and C are too small, then the strong classifier abstains with high probability. However, the theorem does not guarantee that the strong classifier does not abstain when p and C are large. The next theorem takes care of this other direction under slightly stronger assumptions on the weights learned by the algorithm.

Theorem 5 (Exact bound on the probability that a strong classifier abstains). In Theorem 4, further assume that each triplet classifier that does not abstain on at least one training example has a weight different from 0 (if for at least one example $(x_i, y_i) \in S$ we have that $h_c(x_i) \neq \vartheta$ then $\alpha_c \neq 0$). Then equality holds in Equation (6).

Proof. The proof is given in Appendix J. \square

Theorem 5 implies that equality holds in Example 1, thus when $C = \frac{n(n-1)}{2}$ we need at least $k > \frac{3}{2}$, that is at least $\Omega(n\sqrt{n})$ random triplets, to obtain a classifier that never abstains. In Appendix I we extend Example 1 and we study the limit as $n \rightarrow \infty$ for general values of C and p . We also provide a graphical illustration of the bound and a discussion on how this lower bound compares to existing results [Ailon, 2012; Jamieson and Nowak, 2011; Jain et al., 2016].

4 Experiments

We propose an empirical evaluation of TripletBoost. We consider six datasets of varying scales and four baselines.

Baselines. First, we consider an embedding approach. We use tSTE [Van Der Maaten and Weinberger, 2012] to embed the triplets in a Euclidean space and we use the 1-nearest neighbour algorithm for classification. We also would like to compare to alternative approaches able to learn directly from triplets (without embedding as a first step). However, to the best of our knowledge, TripletBoost is the only method able to do classification using only passively obtained triplets. The only option is to choose competing methods that have access to more information (providing them an unfair advantage). We settled for a method that uses actively chosen triplets to build a comparison-tree to retrieve nearest neighbours (CompTree) [Haghiry et al., 2017]. Finally, to put the results obtained in the triplet setting in perspective, we consider two methods that use the original Euclidean representations of the data, the 1-nearest neighbour algorithm (1NN) and AdaBoost.SAMME (SAMME) [Hastie et al., 2009].

Implementation details. For tSTE we used the implementation distributed on the authors' website and we set the embedding dimension to the original dimension of the data. This method was only considered for small datasets with less than 10^3 examples as it does not scale well to bigger datasets (Figure 3c). For CompTree we used our own implementation and the leaf size of the comparison tree is set to 1 as this is the only value for which this method can handle noise. For 1NN and SAMME we used sk-learn [Pedregosa et al., 2011]. The number of boosting iterations for SAMME is set to 10^3 . Finally for TripletBoost we set the number of iterations to 10^6 .

Datasets and performance measure. We consider six datasets: Iris, Moons, Gisette, Cod-rna, MNIST, and kMNIST. For each dataset we generate some triplets as in Equation (1) using three metrics: the Euclidean, Cosine, and City-block distances (details provided in Appendix C). Given a set of n examples there are $n^2(n-1)/2$ possible triplets. We consider three different regimes where 1%, 5% or 10% of them are available, and we consider three noise levels where 0%,

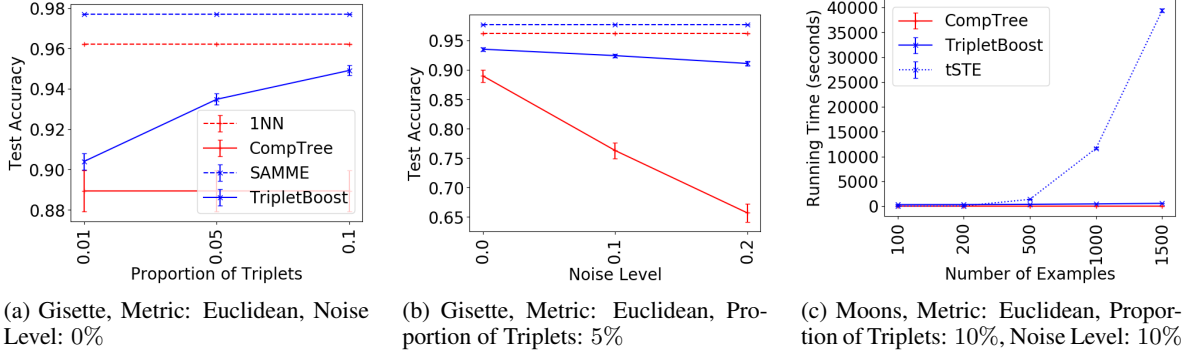


Figure 3: The Gisette dataset has 7000 examples and 5000 dimensions. In Figure 3a we consider the noise free setting and we vary the proportion of triplets available from 1 to 10% of all the triplets. In Figure 3b we fix the proportion of available triplets to 5% and we vary the noise level from 0 to 20%. Figure 3c presents the training time of the triplet-based methods with training samples of increasing sizes on the Moons dataset. The proportion of triplets and the noise level were both set to 10%. The results were obtained on a single core @3.40GHz.

10% or 20% of them are incorrect. We measure performances in terms of test accuracy (higher is better). For all the experiments we report the mean and standard deviation of 10 repetitions. Since the results are mostly consistent across the datasets we present some representative ones here and defer the others to Appendix C.

Small scale regime. We first consider the datasets with less than 10^3 training examples (Figure 2). In this setting our method does not perform well when the number of triplets is too small, but gets closer to the baselines when the number of triplets increases (Figure 2a). This behaviour can be easily explained: when only 1% of the triplets are available, the triplet classifiers abstain on all but 3 or 4 examples on average and thus their performance evaluations are not reliable. Consequently their weights cannot be chosen in a satisfactory manner. This problem vanishes when the number of triplets increases. With increasing noise levels (Figure 2b) one can notice that TripletBoost is more robust than CompTree. Indeed, CompTree generates a comparison-tree based on individual triplet queries, and the greedy decisions in the tree building procedure can easily be misleading in the presence of noise. Finally, our approach is less sensitive than tSTE to changes in the metric that generated the triplets (Figure 2c). Indeed, tSTE assumes that this metric is the Euclidean distance while our approach does not make any assumptions.

Large scale regime. On larger datasets (Figure 3), our method does not reach the accuracy of 1NN and SAMME, who exploit a significant amount of extra information. Still, it performs quite well and is competitive with CompTree, the method that uses active rather than passive queries. The ordinal embedding methods cannot compete in this regime, as they are too slow to even finish (Figure 3c). Once again TripletBoost is quite resistant to noise (Figure 3b).

5 Conclusion

In this paper we proposed TripletBoost to address the problem of comparison-based classification. It is particularly designed for situations where triplets cannot be queried actively,

and we have to live with whatever set of triplets we get. We do not make any geometric assumptions on the underlying space. From a theoretical point of view we have shown that TripletBoost is well founded and we proved guarantees on both the training error and the generalization error of the learned classifier. Furthermore we derived a new lower bound showing that to avoid learning a random predictor, at least $\Omega(n\sqrt{n})$ triplets are needed. In practice we have shown that, given a sufficient amount of triplets, our method is competitive with state of the art methods and that it is quite resistant to noise.

To the best of our knowledge, TripletBoost is the first algorithm that is able to handle large scale datasets using only passively obtained triplets. This means that the comparison-based setting could be considered for problems which were, until now, out of reach. As an illustration, consider a platform where users can watch, comment and rate movies. It is reasonable to assume that triplets of the form *movie m_i is closer to m_j than to m_k* can be automatically obtained using the ratings of the users, their comments, or their interactions. In this scenario, active learning methods are not applicable since the users might be reluctant to answer solicitations. Similarly, embedding methods are too slow to handle large numbers of movies or users. However, we can use TripletBoost to solve problems such as predicting the genres of the movies. As a proof of concept we considered the 1m movielens dataset [Harper and Konstan, 2016]. It contains 1 million ratings from 6040 users on 3706 movies. We used the users' ratings to obtain some triplets about the movies and TripletBoost to learn a classifier able to predict the genres of a new movie (more details are given in Appendix D). Given a new movie, in $\sim 83\%$ of the cases the genre predicted as the most likely one is correct and, on average, the 5 genres predicted as the most likely ones cover $\sim 92\%$ of the true genres.

Acknowledgments

Ulrike von Luxburg acknowledges funding by the DFG through the Institutional Strategy of the University of Tübingen (DFG, ZUK 63) and the Cluster of Excellence EXC 2064/1, project number 390727645.

References

- Sameer Agarwal, Josh Wills, Lawrence Cayton, Gert Lanckriet, David Kriegman, and Serge Belongie. Generalized non-metric multidimensional scaling. In *Artificial Intelligence and Statistics*, 2007.
- Nir Ailon. An active learning algorithm for ranking from pairwise preferences with an almost optimal query complexity. *Journal of Machine Learning Research*, 13(Jan), 2012.
- Ehsan Amid and Antti Ukkonen. Multiview triplet embedding: Learning attributes in multiple maps. In *International Conference on Machine Learning*, 2015.
- Aurélien Bellet, Amaury Habrard, and Marc Sebban. Metric learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 9(1), 2015.
- Leo Breiman. Prediction games and arcing algorithms. *Neural computation*, 11(7), 1999.
- Tarin Clauwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature. *arXiv preprint arXiv:1812.01718v1*, 2018.
- Dheeru Dua and Efi Karra Taniskidou. UCI machine learning repository, 2017.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1), 1997.
- Wei Gao and Zhi-Hua Zhou. On the doubt about margin explanation of boosting. *Artificial Intelligence*, 203, 2013.
- Isabelle Guyon, Steve Gunn, Asa Ben-Hur, and Gideon Dror. Result analysis of the nips 2003 feature selection challenge. In *Neural Information Processing Systems*, 2005.
- Siavash Haghir, Debarghya Ghoshdastidar, and Ulrike von Luxburg. Comparison-based nearest neighbor search. In *Artificial Intelligence and Statistics*, 2017.
- F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems*, 5(4), 2016.
- Trevor Hastie, Saharon Rosset, Ji Zhu, and Hui Zou. Multi-class adaboost. *Statistics and its Interface*, 2(3), 2009.
- Hannes Heikinheimo and Antti Ukkonen. The crowd-median algorithm. In *First AAAI Conference on Human Computation and Crowdsourcing*, 2013.
- Lalit Jain, Kevin G Jamieson, and Rob Nowak. Finite sample prediction and recovery bounds for ordinal embedding. In *Neural Information Processing Systems*, 2016.
- Kevin G Jamieson and Robert D Nowak. Low-dimensional embedding using adaptively selected ordinal data. In *Conference on Communication, Control, and Computing*, 2011.
- Daniel M Kane, Shachar Lovett, Shay Moran, and Jiapeng Zhang. Active classification with comparison queries. In *Foundations of Computer Science*, 2017.
- Ehsan Kazemi, Lin Chen, Sanjoy Dasgupta, and Amin Karbasi. Comparison based learning from weak oracles. *arXiv preprint arXiv:1802.06942v1*, 2018.
- Matthäus Kleindessner and Ulrike Luxburg. Dimensionality estimation without distances. In *Artificial Intelligence and Statistics*, 2015.
- Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 1998.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2011.
- Robert E Schapire and Yoav Freund. *Boosting: Foundations and algorithms*. MIT press, 2012.
- Robert E Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine learning*, 37(3), 1999.
- Robert E Schapire and Yoram Singer. Boostexter: A boosting-based system for text categorization. *Machine learning*, 39(2-3), 2000.
- Robert E Schapire, Yoav Freund, Peter Bartlett, and Wee Sun Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. *Annals of statistics*, 1998.
- Omer Tamuz, Ce Liu, Serge Belongie, Ohad Shamir, and Adam Tauman Kalai. Adaptively learning the crowd kernel. In *International Conference on Machine Learning*, 2011.
- Yoshikazu Terada and Ulrike von Luxburg. Local ordinal embedding. In *International Conference on Machine Learning*, 2014.
- Dominique Tschopp, Suhas Diggavi, Payam Delgosha, and Soheil Mohajer. Randomized algorithms for comparison-based search. In *Neural Information Processing Systems*, 2011.
- Andrew V Uzilov, Joshua M Keegan, and David H Mathews. Detection of non-coding rnas on the basis of predicted secondary structure formation free energy change. *BMC bioinformatics*, 7(1), 2006.
- Laurens Van Der Maaten and Kilian Weinberger. Stochastic triplet embedding. In *Machine Learning for Signal Processing*, 2012.
- Liwei Wang, Masashi Sugiyama, Zhaoxiang Jing, Cheng Yang, Zhi-Hua Zhou, and Jufu Feng. A refined margin analysis for boosting algorithms via equilibrium margin. *Journal of Machine Learning Research*, 12(Jun), 2011.

A Discussion on Theorem 3

At a first glance it seems that the bound presented in Theorem 3 does not depend on m , the number of available triplets. However, this dependency is implicit: m impacts the probability that the training examples are well classified with a large margin θ . If the number of triplets is small, the probability that the training examples are well classified with a given margin is small. This probability increases when the number of triplets increases.

To illustrate this behaviour, consider a fixed margin θ . Assume that each of the n training examples is classified by exactly one triplet classifier and that each triplet classifier abstains on all but one example. In this case, the only way to have no error on the training set is to combine the n triplet classifiers that do not abstain. For simplicity consider the case where all the triplet classifiers have uniform weight $\frac{1}{n}$ in the final classifier. Then the fixed margin will not be achieved when the number of training examples increases. The first term on the right hand side of the generalization bound will be 1, that is the bound predicts that the learned classifier might not generalize. This fact remains true for any weighting scheme. In this example, the best weighting scheme classifies with a margin at least θ , at most $\frac{1}{\theta}$ training examples.

When the number of triplets increases, the value of Z_c decreases, because the proportion of examples on which the selected triplet classifier abstains, $1 - W_{c,+} - W_{c,-}$, decreases. Similarly, the proportion of training examples that are classified with a margin at least θ increases. Hence the first term on the right hand side of the generalization bound is greatly reduced and the learned classifier generalizes well.

B Illustration and Discussion on the Lower Bound of Theorems 4 and 5

When $n = 100$ we illustrate the bound obtained in Theorems 4 and 5 in Figure 4. This figure shows that the transition between abstaining and non-abstaining classifier depends on both the proportion of triplets available and the number of classifiers considered. In particular, when the number of combined classifiers increases, one needs a smaller number of triplets. Conversely when the number of triplets available increases, one can consider combining fewer classifiers. This illustration confirms the result obtained in Example 1 that shows that when $C = \frac{n(n-1)}{2}$ the number of available triplets should at least scale as $\Omega(n\sqrt{n})$ to obtain a non-trivial classifier that does not abstain on most of the test examples. The number of triplets that are necessary to achieve good classification accuracy might be higher but is lower bounded by this value.

This lower bound does not contradict existing results [Ailon, 2012; Jamieson and Nowak, 2011; Jain *et al.*, 2016]. They were developed in the different context of triplet recovery, where the goal is not classification, but to predict the outcome of unobserved triplet questions. For example it has been shown that to exactly recover all the triplets, the number of passively available triplets should scale in $\Omega(n^3)$ [Jamieson and Nowak, 2011]. Similarly Jain *et al.* [2016] derive a finite error bound for approximate recovery of the Euclidean Gram matrix. Our bound shows that, in a classification setting, it might be possible to do better than that. To set a complete picture, one would need to derive an upper bound on the number of triplets necessary for good classification accuracy.

C Details on the Experiments

The characteristics of the different datasets are given in Table 1. To generate the triplets we used three different metrics:

- Euclidean distance: $d(x, y) = \|x - y\|_2$,
- Cityblock distance: $d(x, y) = \|x - y\|_1$,
- Cosine distance: $d(x, y) = 1 - \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$.

The set \mathcal{T} of all triplets is then defined as follows:

$$\mathcal{T} = \{(x_i, x_j, x_k) : d(x_i, x_j) < d(x_i, x_k)\}.$$

In all the experiments we considered subsets T of \mathcal{T} by selecting uniformly at random without replacement 1%, 5%, or 10% of the triplets. Similarly we added some noise by randomly swapping 0%, 10%, or 20% of the triplets, that is $(x_i, x_k, x_j) \in T$ while $(x_i, x_j, x_k) \in \mathcal{T}$.

The results that were omitted in the main paper are given in Figures 5 to 16.

D Details on the Movielens experiment

As a proof of concept we considered the 1m movielens dataset [Harper and Konstan, 2016]. This dataset contains 1 million ratings from 6040 users on 3706 movies and each movie has one or several genres (there is 18 genres in total). To demonstrate the interest of our approach we proposed (i) to use the users' ratings to obtain some triplets of the form *movie m_i is closer to movie m_j than to movie m_k* , and (ii) to use TripletBoost to learn a classifier predicting the genres of the movies.

To generate the triplets we propose to consider that movie m_i is closer to m_j than to m_k if, on average, users that rated all three movies rated m_i and m_j more similarly than m_i and m_k . The underlying intuition is that users like and dislike genres of movies — for example a user that dislikes horror movies but likes comedy movies will probably give low ratings to The Ring

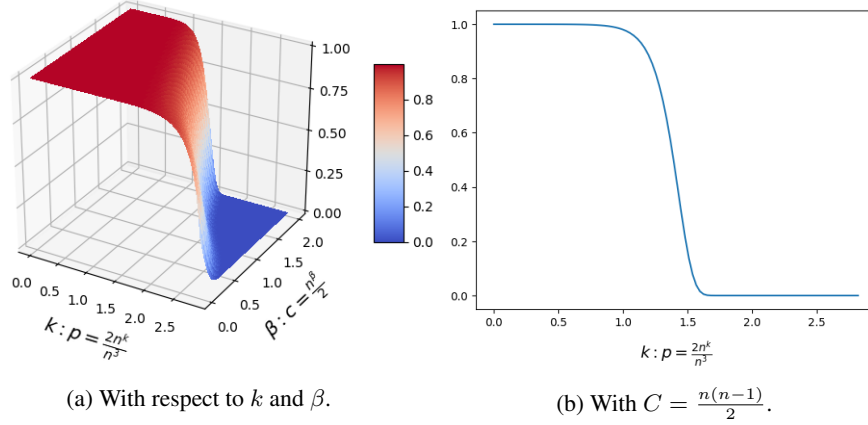


Figure 4: Illustration of the bound in Theorems 4 and 5 when $n = 100$.

Dataset	Dimension	Train / Test	Classes	Results
Iris	4	105 / 45	3	Figures 5 and 6
Moons	2	350 / 150	2	Figures 7 and 8
Gisette	5000	6000 / 1000	2	Figures 9 and 10
CodRna	8	59535 / 271617	2	Figures 11 and 12
MNIST	784	60000 / 10000	10	Figures 13 and 14
kMNIST	784	60000 / 10000	10	Figures 15 and 16

Table 1: Summary of the different datasets.

(2002) and Scream (1996) and a higher rating to The Big Lebowski (1998). Formally let $r_{u,i}$ be the rating of user u on movie m_i and $r_{u,i,j} = |r_{u,i} - r_{u,j}|$ then the triplet set T is

$$T = \left\{ (m_i, m_j, m_k) : \left(\sum_{u \in U_{i,j,k}} \frac{\mathbb{I}_{r_{u,i,j} < r_{u,i,k}} - \mathbb{I}_{r_{u,i,j} > r_{u,i,k}}}{|U_{i,j,k}|} \right) > 0 \right\}$$

where $U_{i,j,k}$ is the set of users that rated all three movies. Each user has only rated a small number of movies and might give a high, respectively low, rating to a movie with a genre that he usually rates lower, respectively higher. Thus we only have access to a noisy subset of all the possible triplets.

We used a random sample of 2595 movies, and their corresponding triplets, to learn a multi-label classifier using TripletBoost (with 10^6 boosting iterations). Since we are in a multi-label setting we would like to predict how relevant each genre is for a new movie rather than a single genre. To obtain such a quantity we can simply ignore the $\arg \max$ in Equation (5) in the main paper to obtain a classifier $H(m, y)$ that predicts the weight of a genre y for a movie m :

$$H(m, y) = \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(m) \neq y \wedge y \in h_c(m)}.$$

In the test phase we used the 1111 remaining movies to measure the performance of the learned classifier. First we considered the precision of the genre predicted with the highest weight and obtained a value of 0.83168. It means that in $\sim 83\%$ of the cases the genre predicted with the highest weight is correct. We also considered the recall of the 5 genres predicted with the highest weights and obtained a value of 0.92943. It implies that, on average, the 5 genres predicted with the highest weights cover $\sim 92\%$ of the genres of the considered movie.

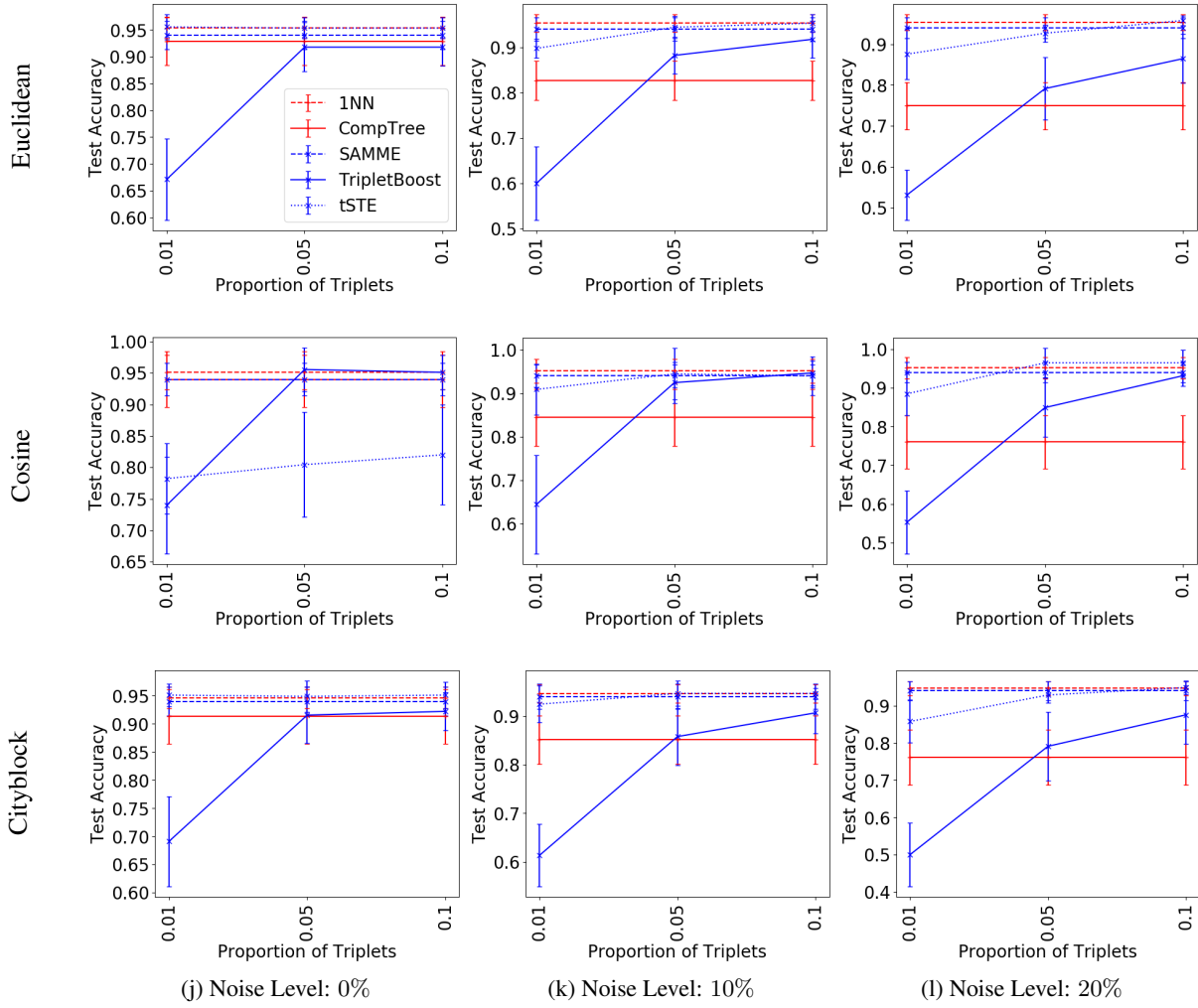


Figure 5: Results on the Iris dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different noise level from 0 to 20%. In each plot we vary the proportion of triplets available from 1 to 10%.

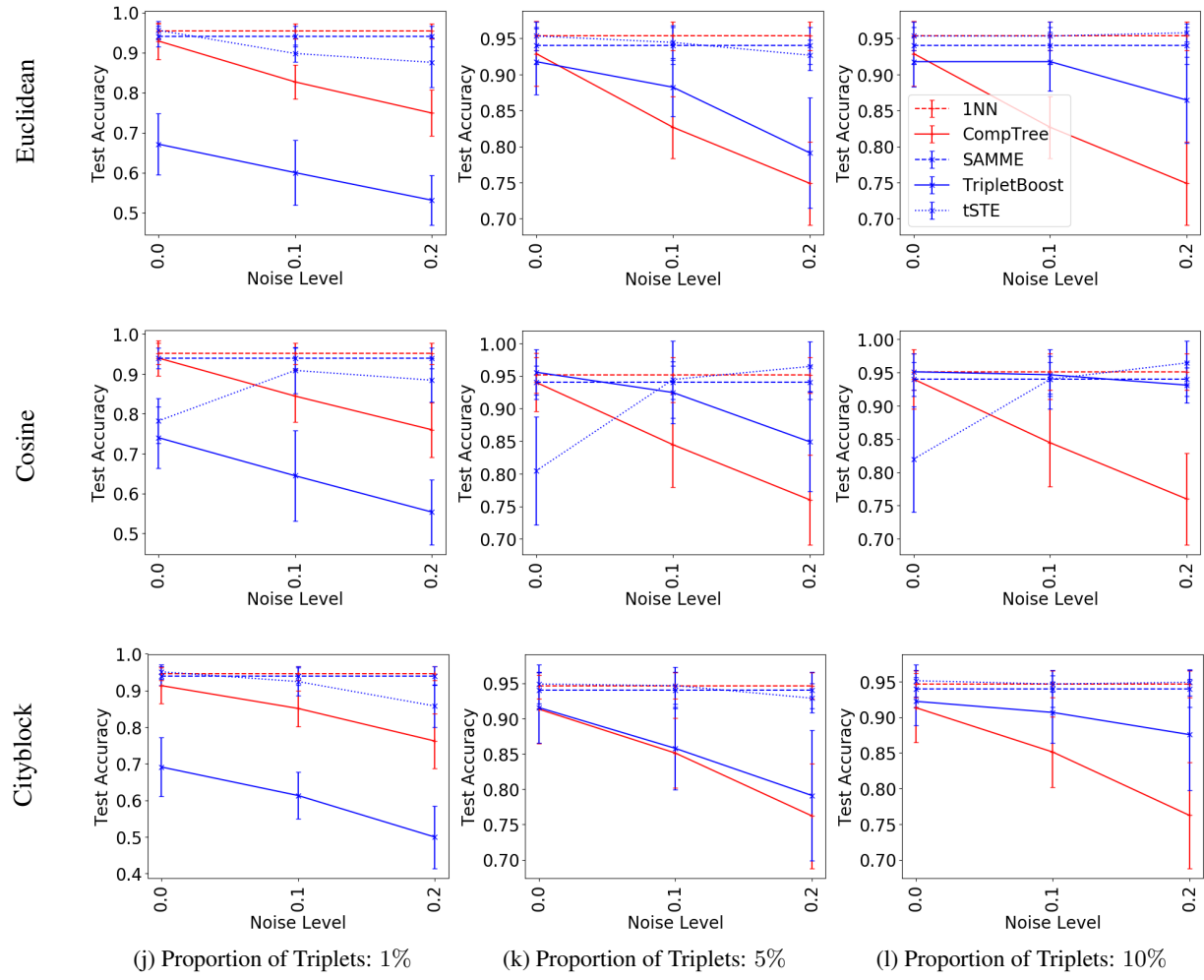


Figure 6: Results on the Iris dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different proportion of triplets available from 1 to 10%. In each plot we vary the noise level from 0 to 20%.

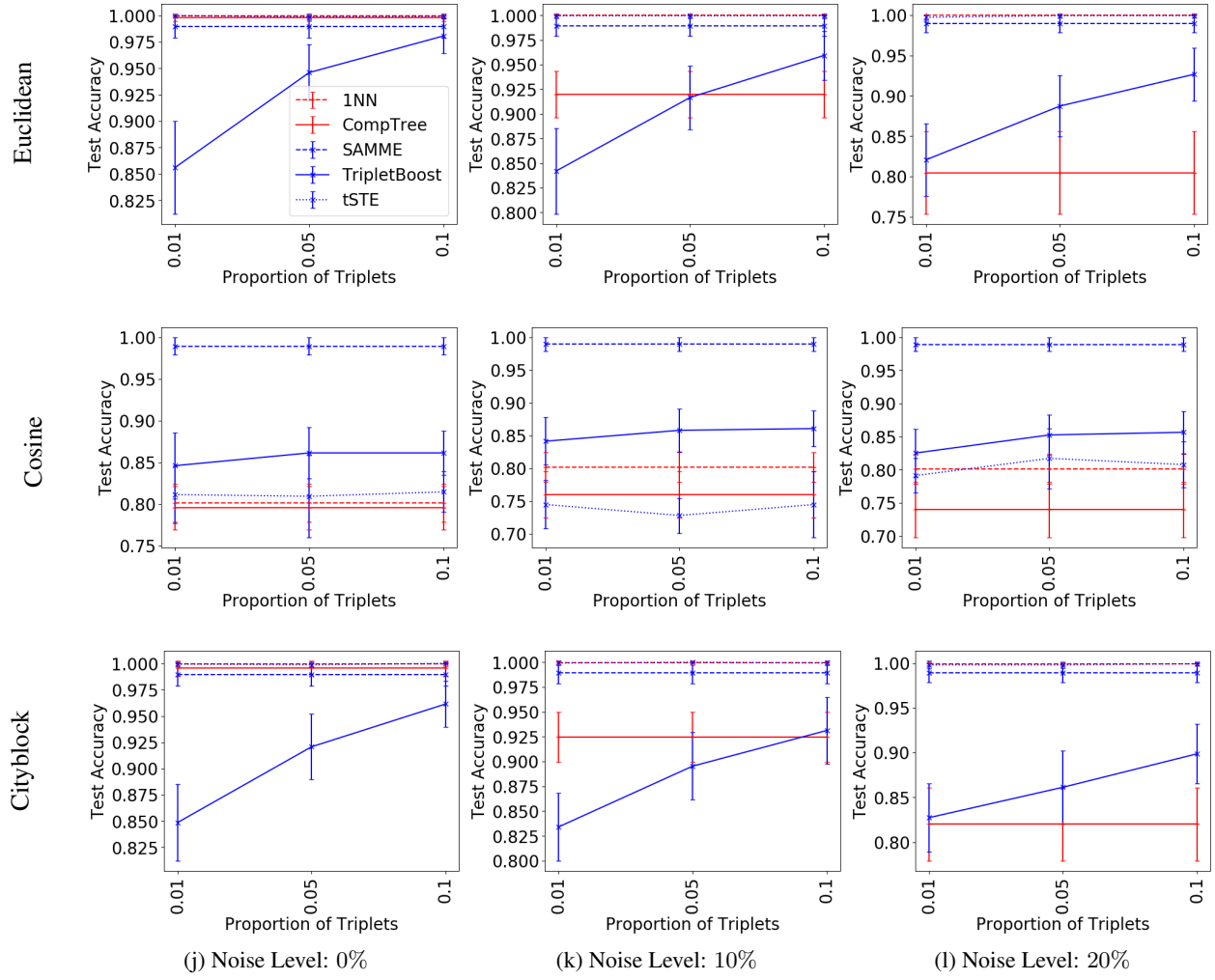


Figure 7: Results on the Moons dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different noise level from 0 to 20%. In each plot we vary the proportion of triplets available from 1 to 10%.

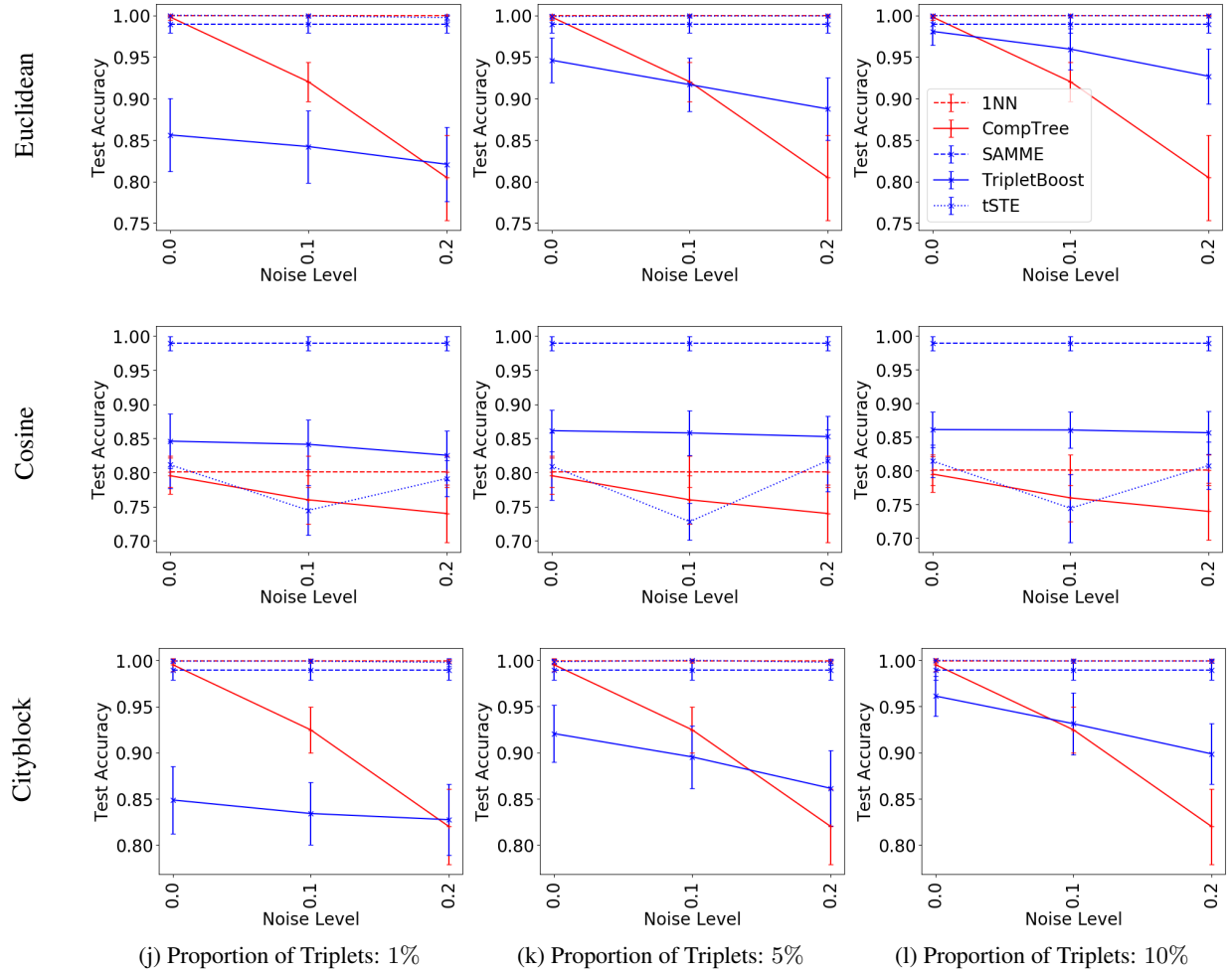


Figure 8: Results on the Moons dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different proportion of triplets available from 1 to 10%. In each plot we vary the noise level from 0 to 20%.

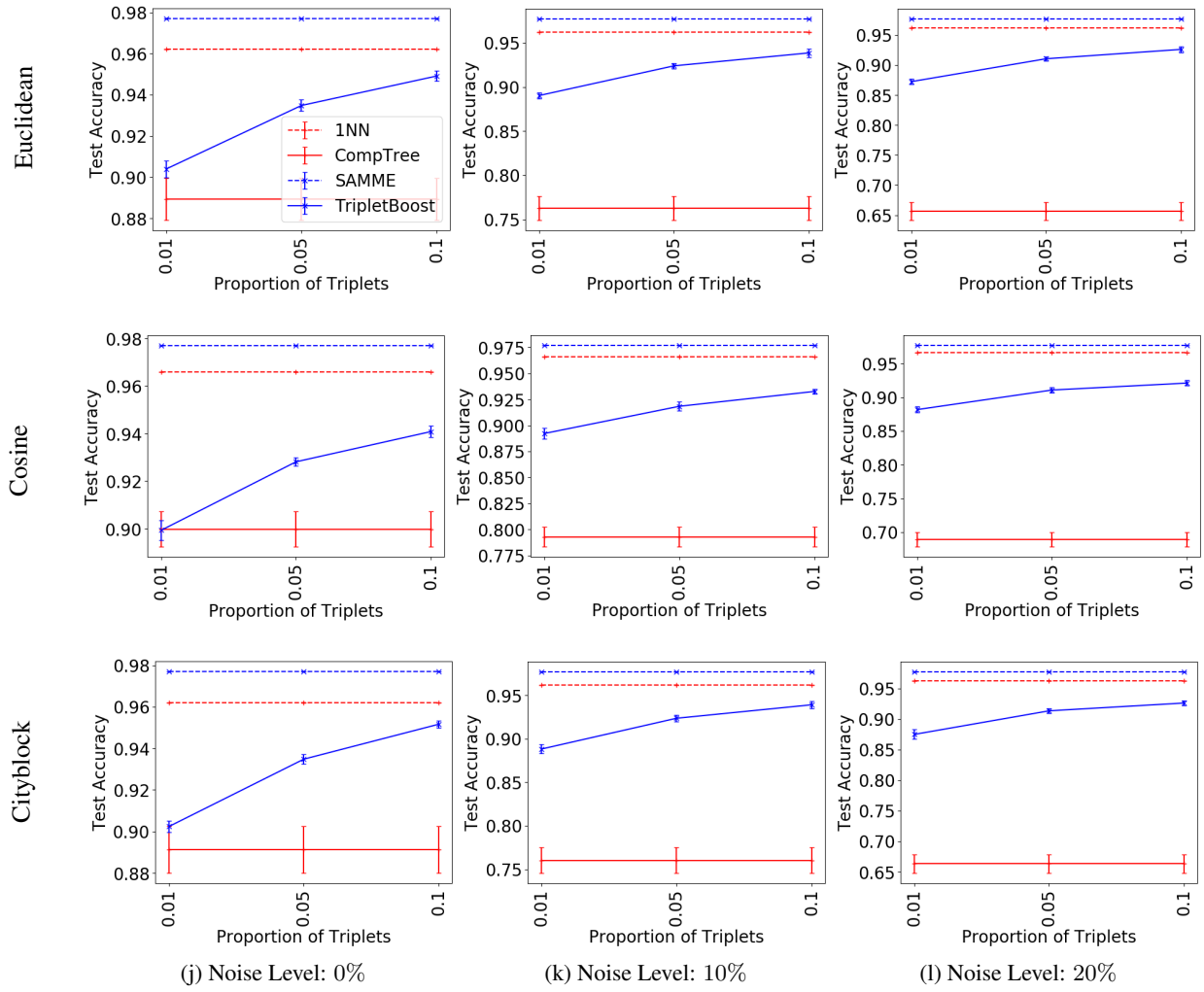


Figure 9: Results on the Gisette dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different noise level from 0 to 20%. In each plot we vary the proportion of triplets available from 1 to 10%.

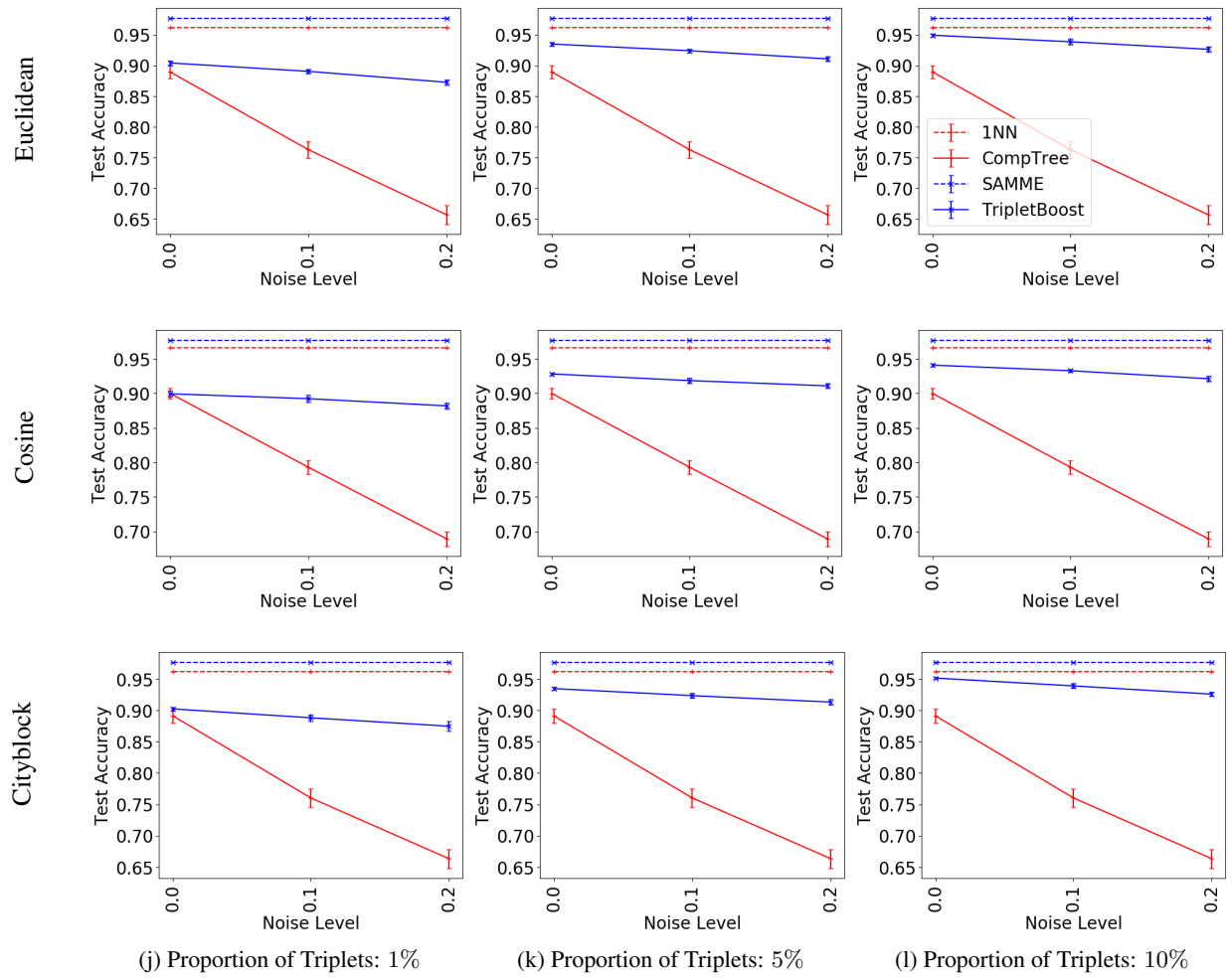


Figure 10: Results on the Gisette dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different proportion of triplets available from 1 to 10%. In each plot we vary the noise level from 0 to 20%.

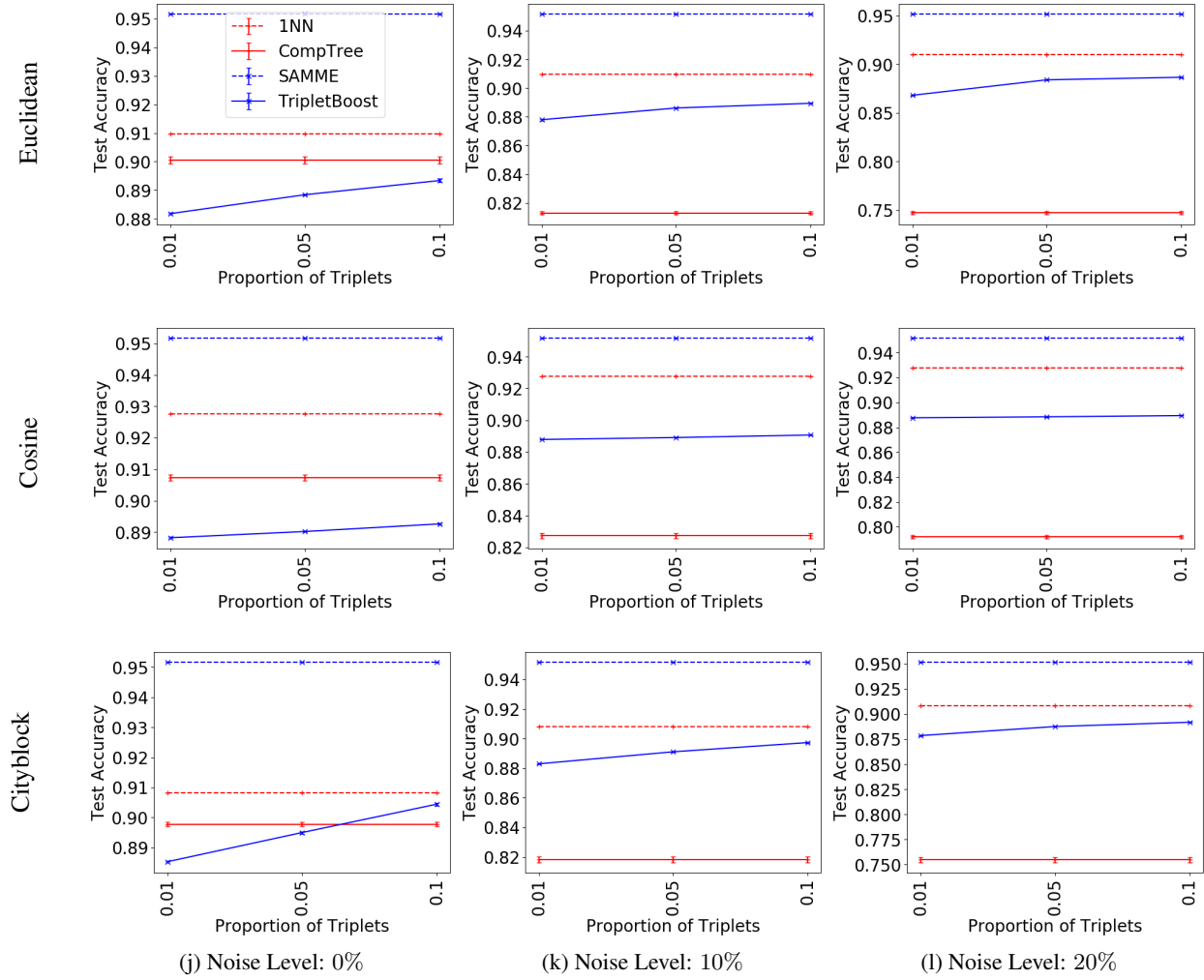


Figure 11: Results on the Cod-rna dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different noise level from 0 to 20%. In each plot we vary the proportion of triplets available from 1 to 10%.

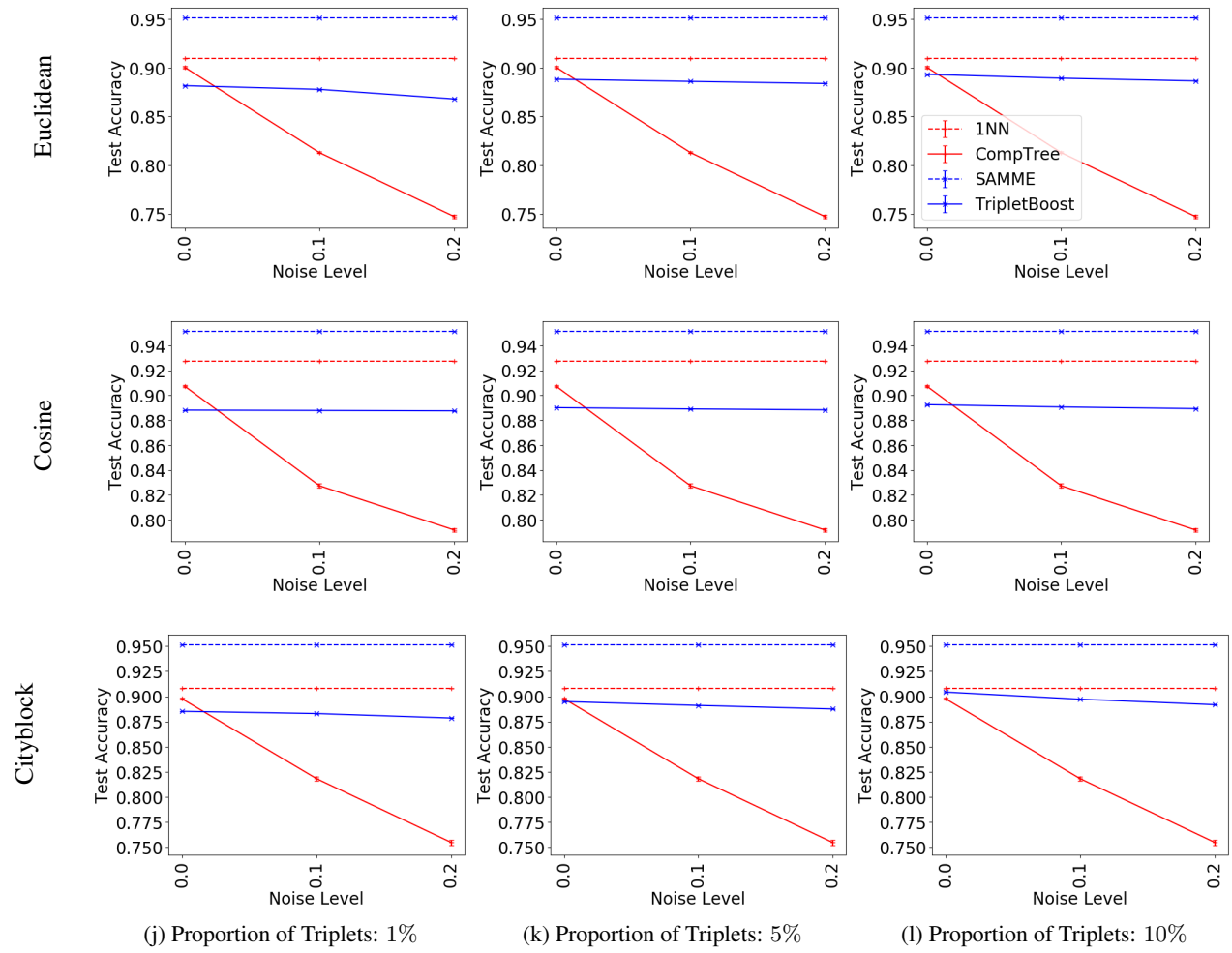


Figure 12: Results on the Cod-rna dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different proportion of triplets available from 1 to 10%. In each plot we vary the noise level from 0 to 20%.

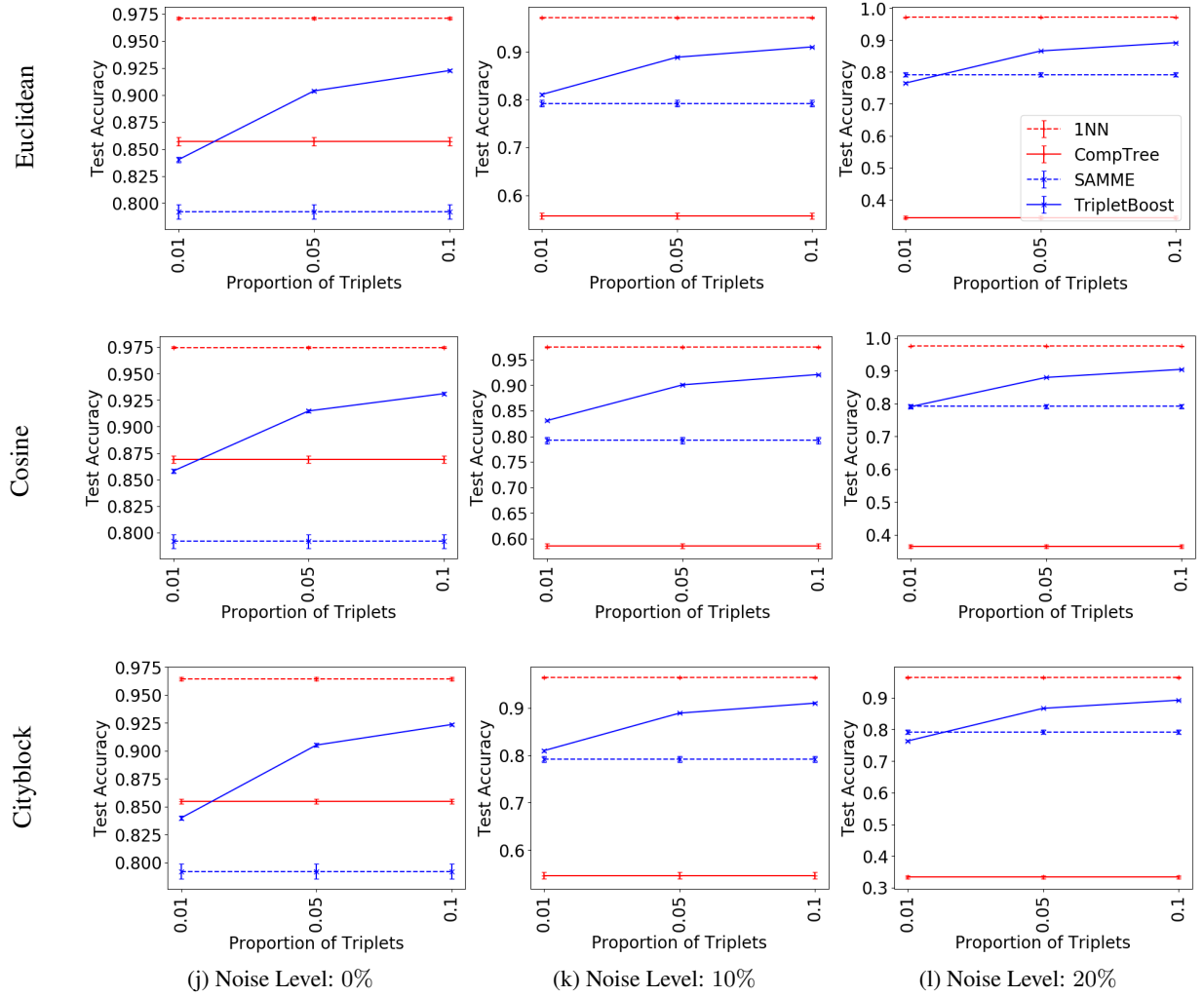


Figure 13: Results on the MNIST dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different noise level from 0 to 20%. In each plot we vary the proportion of triplets available from 1 to 10%.

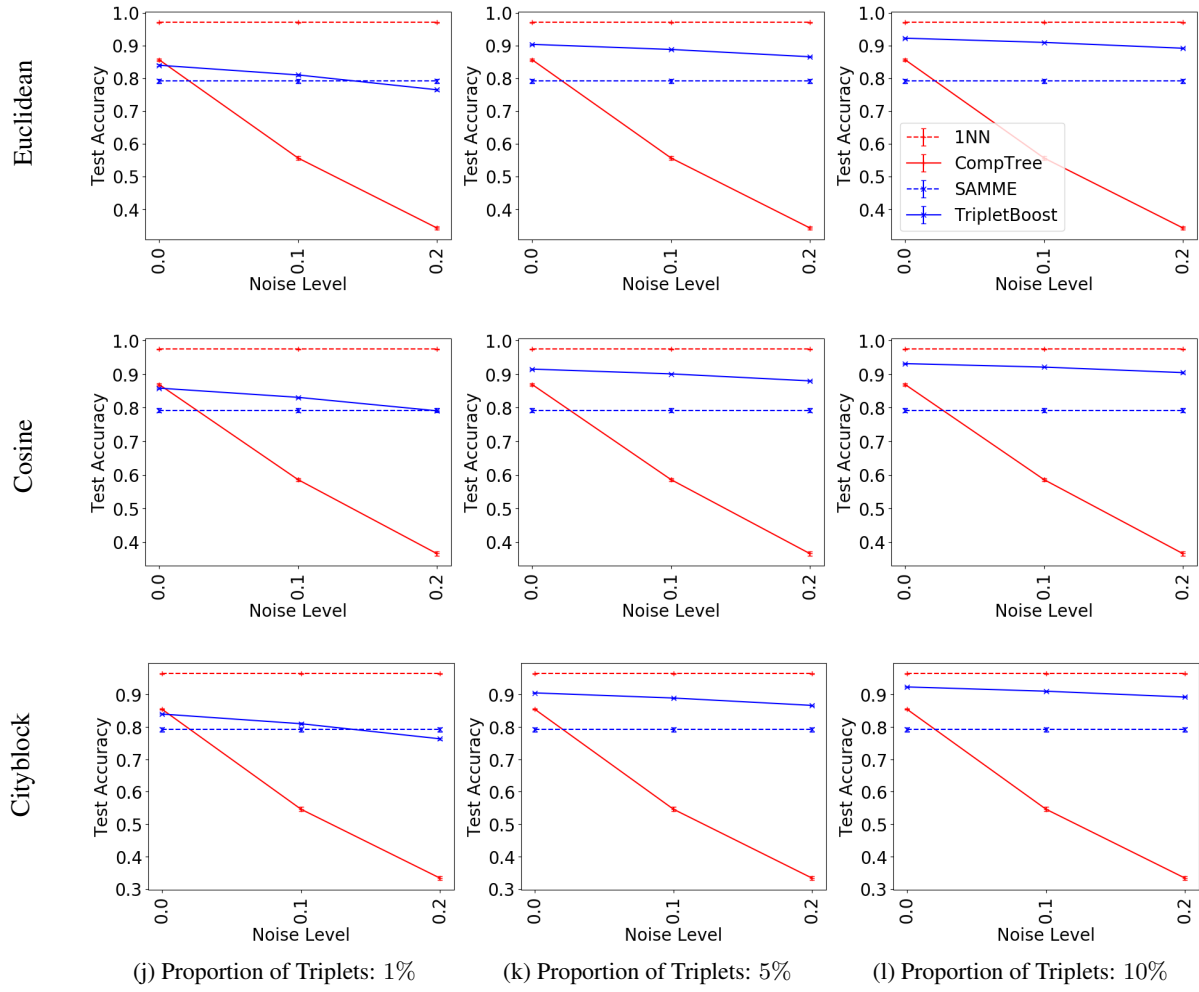


Figure 14: Results on the MNIST dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different proportion of triplets available from 1 to 10%. In each plot we vary the noise level from 0 to 20%.

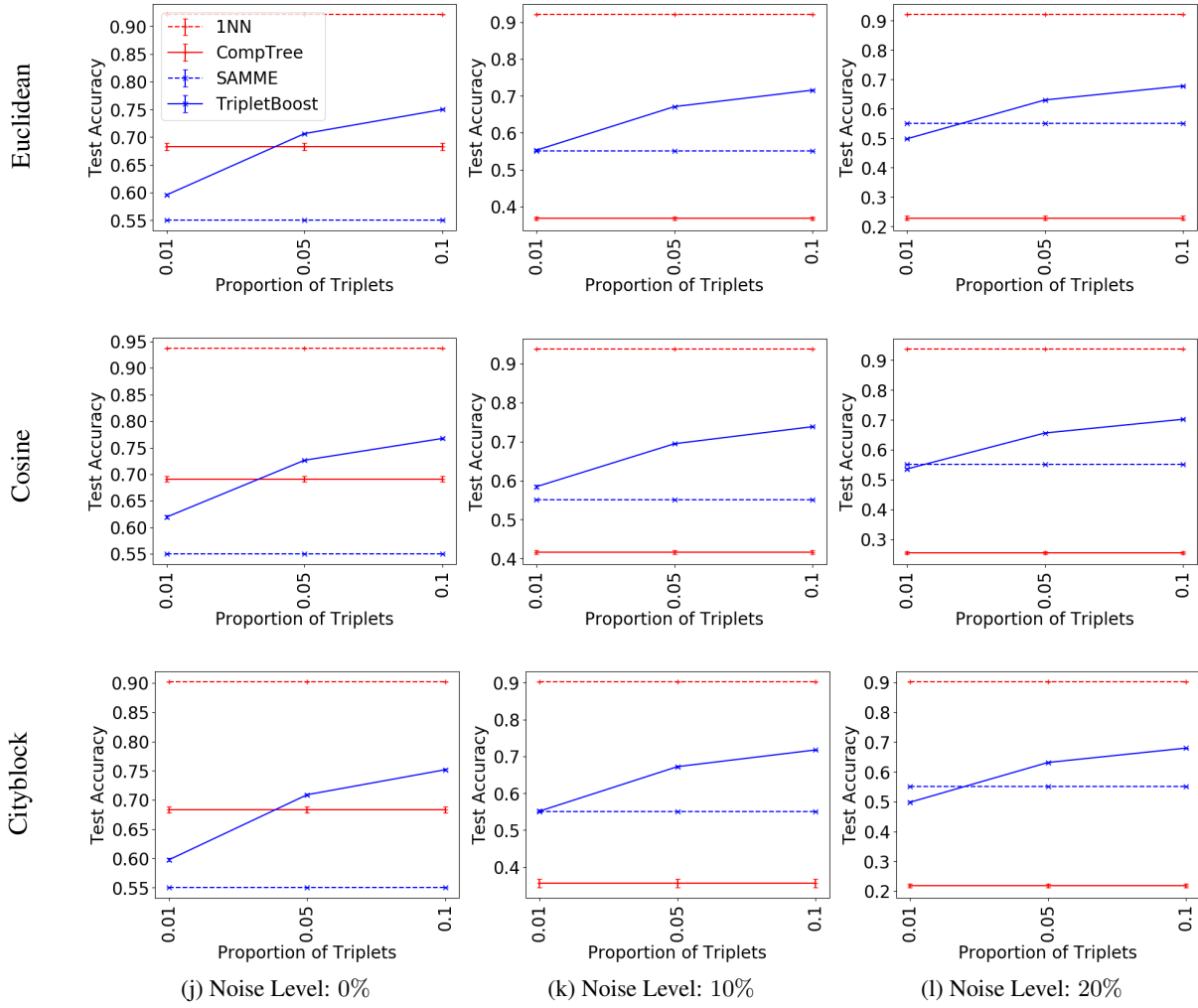


Figure 15: Results on the kMNIST dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different noise level from 0 to 20%. In each plot we vary the proportion of triplets available from 1 to 10%.

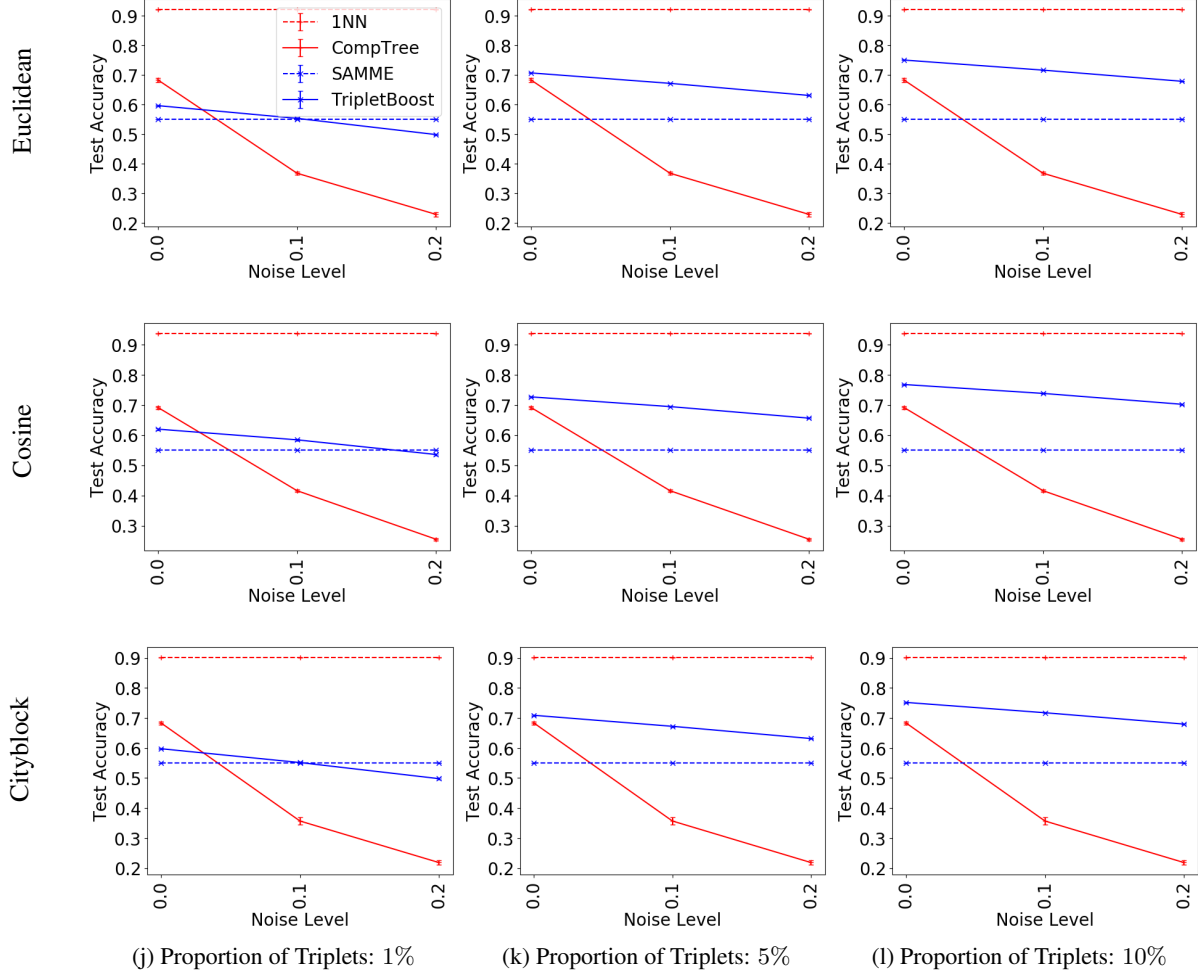


Figure 16: Results on the kMNIST dataset. In each row we consider a different metric to generate the triplets. In each column we consider a different proportion of triplets available from 1 to 10%. In each plot we vary the noise level from 0 to 20%.

E Proof of Theorem 1

Theorem 1 (Triplet classifiers and weak learners). *Let \mathcal{W}_c be an empirical distribution over $S \times \mathcal{Y}$ and h_c be the corresponding triplet classifier chosen as described in Section 2.2. It holds that:*

1. *the error of h_c on \mathcal{W}_c is at most the error of a random predictor and is minimal compared to other labelling strategies,*
2. *the weight α_c of the classifier is non-zero if and only if, h_c is a weak classifier, that is is strictly better than a random predictor.*

Proof. To prove the first claim of the theorem we first study the error of h_c on \mathcal{W}_c and then we show that any labelling strategy different from our would increase the proportion of incorrectly classified examples. To prove the second claim we simply use the definition of the weights α_c .

First claim First of all notice that the probability that a triplet classifier h_c makes an error on \mathcal{W}_c is

$$\mathbb{P}_{(x_i, y) \sim \mathcal{W}_c} [(\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) \neq (\mathbb{I}_{y \in h_c(x_i)} - \mathbb{I}_{y \notin h_c(x_i)})]$$

that is, by the law of total probabilities:

$$\begin{aligned} & \mathbb{P}_{\substack{(x_i, y) \sim \mathcal{W}_c \\ h_c(x_i) \neq \vartheta}} [(\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) \neq (\mathbb{I}_{y \in h_c(x_i)} - \mathbb{I}_{y \notin h_c(x_i)})] \mathbb{P}_{(x_i, y) \sim \mathcal{W}_c} [h_c(x_i) \neq \vartheta] \\ & + \mathbb{P}_{\substack{(x_i, y) \sim \mathcal{W}_c \\ h_c(x_i) = \vartheta}} [(\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) \neq (\mathbb{I}_{y \in h_c(x_i)} - \mathbb{I}_{y \notin h_c(x_i)})] \mathbb{P}_{(x_i, y) \sim \mathcal{W}_c} [h_c(x_i) = \vartheta]. \end{aligned}$$

Notice that when h_c abstains the best possible behaviour is to randomly predict the labels and thus

$$\mathbb{P}_{\substack{(x_i, y) \sim \mathcal{W}_c \\ h_c(x_i) = \emptyset}} \left[(\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) \neq (\mathbb{I}_{y \in h_c(x_i)} - \mathbb{I}_{y \notin h_c(x_i)}) \right] = \frac{1}{2}.$$

Hence to show that h_c has at most the error of a random predictor it is sufficient to show that

$$\mathbb{P}_{\substack{(x_i, y) \sim \mathcal{W}_c \\ h_c(x_i) \neq \emptyset}} \left[(\mathbb{I}_{y=y_i} - \mathbb{I}_{y \neq y_i}) \neq (\mathbb{I}_{y \in h_c(x_i)} - \mathbb{I}_{y \notin h_c(x_i)}) \right] = \frac{W_{c,-}}{W_{c,+} + W_{c,-}} \leq \frac{1}{2}$$

where $W_{c,+}$ and $W_{c,-}$ are respectively the proportions of correctly and incorrectly classified examples.

On the one hand assume that h_c uses x_j and x_k as reference points. Let $S_{c,j}$ be the set of all the training examples for which h_c does not abstain and which are closer to x_j than to x_k . If $y \in o_j$ then we have that:

$$\sum_{z_i \in S_{c,j}} \mathbb{I}_{y=y_i} w_{c,x_i,y} > \sum_{z_i \in S_{c,j}} \mathbb{I}_{y \neq y_i} w_{c,x_i,y}. \quad (8)$$

Similarly, if $y \notin o_j$ then it implies that:

$$\sum_{z_i \in S_{c,j}} \mathbb{I}_{y \neq y_i} w_{c,x_i,y} \geq \sum_{z_i \in S_{c,j}} \mathbb{I}_{y=y_i} w_{c,x_i,y}. \quad (9)$$

The same result holds for o_k .

On the other hand recall that

$$\begin{aligned} W_{c,+} &= \sum_{\substack{(x_i, y_i) \in S \\ h_c(x_i) \neq \emptyset}} \left(\mathbb{I}_{y_i \in h_c(x_i)} w_{c,x_i,y_i} + \sum_{y \neq y_i} \mathbb{I}_{y \notin h_c(x_i)} w_{c,x_i,y} \right), \\ W_{c,-} &= \sum_{\substack{(x_i, y_i) \in S \\ h_c(x_i) \neq \emptyset}} \left(\mathbb{I}_{y_i \notin h_c(x_i)} w_{c,x_i,y_i} + \sum_{y \neq y_i} \mathbb{I}_{y \in h_c(x_i)} w_{c,x_i,y} \right). \end{aligned} \quad (10)$$

Combining Equations (8), (9), and (10) we obtain

$$\begin{aligned} & W_{c,-} \leq W_{c,+} \\ \Leftrightarrow & 2W_{c,-} \leq W_{c,+} + W_{c,-} \\ \Leftrightarrow & \frac{W_{c,-}}{W_{c,+} + W_{c,-}} \leq \frac{1}{2} \end{aligned}$$

which proves that h_c has at most the error of a random predictor.

To see that our labelling strategy is optimal just notice that any other labelling strategy would either predict some labels in o_j or o_k that do not conform with Equation (8) or not predict some labels from o_j or o_k that do not conform with Equation (9). As a consequence $W_{c,-}$ would be increased and thus the error of h_c would also be increased. This concludes the proof of the first claim.

Second claim Recall that α_c is computed as

$$\alpha_c = \frac{1}{2} \log \left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right).$$

Also notice that h_c is a weak classifier, that is better than a random predictor, when $W_{c,-} < W_{c,+}$. We have that:

$$\begin{aligned} & W_{c,-} < W_{c,+} \\ \Leftrightarrow & W_{c,-} + \frac{1}{n} < W_{c,+} + \frac{1}{n} \\ \Leftrightarrow & 1 < \frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \\ \Leftrightarrow & 0 < \frac{1}{2} \log \left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right) \end{aligned}$$

and it proves the second claim. □

F Proof of Theorem 2

Theorem 2 (Reduction of the Training Error). *Let S be a set of n examples and T be a set of m triplets (obtained as described in Section 2). Let $H(\cdot)$ be the classifier obtained after C iterations of TripletBoost (Algorithm 1) using S and T as input. It holds that:*

$$\mathbb{P}_{(x,y) \in S} [H(x) \neq y] \leq \frac{|\mathcal{Y}|}{2} \prod_{c=1}^C Z_c$$

$$\text{with } Z_c = (1 - W_{c,+} - W_{c,-}) + (W_{c,+}) \cdot \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + (W_{c,-}) \cdot \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \leq 1.$$

Proof. This result is inherited from AdaBoost.MO and can be obtained directly by applying Theorem 10.4 in Chapter 10 of the book of Schapire and Freund [2012] using a so-called loss-based decoding and noticing that in our case $\bar{K} = |\mathcal{Y}|$ and $\rho = 2$. For the sake of completeness we detail the complete proof below.

First of all, let $F(x, y) = \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(x) \neq y} (\mathbb{I}_{y \in h_c(x)} - \mathbb{I}_{y \notin h_c(x)})$ and notice that the strong classifier $H(\cdot)$ can be equivalently rewritten as:

$$\begin{aligned} H(\cdot) &= \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(\cdot) \neq y \wedge y \in h_c(\cdot)} \right) && (x > y \Leftrightarrow x - (B - x) > y - (B - y).) \\ \Leftrightarrow H(\cdot) &= \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(\cdot) \neq \varnothing} (\mathbb{I}_{y \in h_c(\cdot)} - \mathbb{I}_{y \notin h_c(\cdot)}) \right) \\ &&& (F(\cdot, y) = \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(\cdot) \neq \varnothing} (\mathbb{I}_{y \in h_c(\cdot)} - \mathbb{I}_{y \notin h_c(\cdot)})) \\ \Leftrightarrow H(\cdot) &= \arg \max_{y \in \mathcal{Y}} F(\cdot, y) \\ \Leftrightarrow H(\cdot) &= \arg \min_{y \in \mathcal{Y}} -F(\cdot, y) \\ \Leftrightarrow H(\cdot) &= \arg \min_{y \in \mathcal{Y}} \exp(-F(\cdot, y)) && (\arg \min_x \exp(-f(x)) = \arg \min_x -\exp(f(x)).) \\ \Leftrightarrow H(\cdot) &= \arg \min_{y \in \mathcal{Y}} \exp(-F(\cdot, y)) - \exp(F(\cdot, y)) \\ &&& (\text{Adding a constant does not change the value of arg min.}) \\ \Leftrightarrow H(\cdot) &= \arg \min_{y \in \mathcal{Y}} \exp(-F(\cdot, y)) - \exp(F(\cdot, y)) + \sum_{y' \in \mathcal{Y}} \exp(F(\cdot, y')) \\ \Leftrightarrow H(\cdot) &= \arg \min_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(\cdot, y')). \end{aligned}$$

Given an example $(x, y) \in S$, the classifier $H(\cdot)$ makes an error if and only if there exists a label $\ell \neq y$ such that

$$\sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) F(x, y')) \leq \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y')).$$

It implies that, when $H(\cdot)$ makes an error

$$\begin{aligned} &\sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y')) \\ &\geq \frac{1}{2} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) F(x, y')) + \frac{1}{2} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y')) \\ &\quad (\text{Dropping all the terms that do not depend on } y \text{ and } \ell \text{ since } \forall x \in \mathbb{R}, \exp x > 0.) \\ &\geq \frac{1}{2} \sum_{y' \in \{y, \ell\}} [\exp(-(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) F(x, y')) + \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y'))] \\ &\quad (\forall y' \in \{y, \ell\}, (\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) = -(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'})). \end{aligned}$$

$$\begin{aligned}
&\geq \frac{1}{2} \sum_{y' \in \{y, \ell\}} \left[\frac{1}{\exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'})F(x, y'))} + \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'})F(x, y')) \right] \\
&\quad (\forall x > 0, x + \frac{1}{x} \geq 2.) \\
&\geq 2.
\end{aligned}$$

Assuming that the classifier $H(\cdot)$ makes M errors, we have that:

$$2M \leq \sum_{(x, y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'})F(x, y')). \quad (11)$$

Now, we need to show that the right hand side of Inequality (11) is bounded by $n|\mathcal{Y}| \prod_{c=1}^C Z_c$ where n is the number of training examples and Z_c is the normalization factor used in TripletBoost. Recall that, given an example $(x, y) \in S$ and a label $y' \in \mathcal{Y}$, the weight $w_{c+1, x, y'}$, obtained after c iterations of our algorithm, is

$$\begin{aligned}
w_{c+1, x, y'} &= \begin{cases} \frac{w_{c, x, y'}}{Z_c} & \text{if } h_c(x) = \vartheta, \\ \frac{w_{c, x, y'}}{Z_c} \exp(-\alpha_c (\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)})) & \text{if } h_c(x) \neq \vartheta. \end{cases} \\
\Leftrightarrow w_{c+1, x, y'} &= \frac{w_{c, x, y'}}{Z_c} \exp(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta} (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)})) \quad (12)
\end{aligned}$$

By recursively using Equation (12), we obtain

$$w_{C+1, x, y'} = \frac{w_{1, x, y'}}{\prod_{c=1}^C Z_c} \prod_{c=1}^C \exp(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta} (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)})) \quad (13)$$

($\exp(x) \exp(y) = \exp(x+y)$.)

$$\begin{aligned}
\Leftrightarrow w_{C+1, x, y'} &= \frac{w_{1, x, y'}}{\prod_{c=1}^C Z_c} \exp\left(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta} (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)})\right) \quad (14) \\
&\quad (F(x, y) = \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta} (\mathbb{I}_{y \in h_c(x)} - \mathbb{I}_{y \notin h_c(x)}).)
\end{aligned}$$

$$\begin{aligned}
\Leftrightarrow w_{C+1, x, y'} &= \frac{w_{1, x, y'}}{\prod_{c=1}^C Z_c} \exp(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) F(x, y')) \quad (15) \\
&\quad (\sum_{(x, y) \in S} \sum_{y' \in \mathcal{Y}} w_{C+1, x, y'} = 1 \text{ and } \forall (x, y) \in S, \forall y' \in \mathcal{Y}, w_{1, x, y'} = \frac{1}{n|\mathcal{Y}|}).
\end{aligned}$$

$$\begin{aligned}
\Leftrightarrow \prod_{c=1}^C Z_c &= \frac{1}{n|\mathcal{Y}|} \sum_{(x, y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) F(x, y')). \quad (16)
\end{aligned}$$

Combining this result with Inequality (11) we obtain:

$$2M \leq n|\mathcal{Y}| \prod_{c=1}^C Z_c.$$

Noticing that $\mathbb{P}_{(x, y) \in S} [H(x) \neq y] = \frac{M}{n}$ gives the first part of the theorem:

$$\mathbb{P}_{(x, y) \in S} [H(x) \neq y] \leq \frac{|\mathcal{Y}|}{2} \prod_{c=1}^C Z_c.$$

In the last part of the proof we show that $Z_c = (1 - W_{c,+} - W_{c,-}) + (W_{c,+}) \cdot \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + (W_{c,-}) \cdot \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \leq 1$. Using the result obtained in Equation (12) we have that

$$\begin{aligned}
w_{c+1, x, y'} &= \frac{w_{c, x, y'}}{Z_c} \exp(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta} (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)})) \\
&\quad (\sum_{(x, y) \in S} \sum_{y' \in \mathcal{Y}} w_{c+1, x, y'} = 1.) \\
\Leftrightarrow Z_c &= \sum_{(x, y) \in S} \sum_{y' \in \mathcal{Y}} w_{c, x, y'} \exp(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \alpha_c \mathbb{I}_{h_c(x) \neq \vartheta} (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)}))
\end{aligned}$$

$$(\exp(0) = 1.)$$

$$\Leftrightarrow Z_c = \sum_{\substack{(x,y) \in S \\ h_c(x) = \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} + \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} \exp \left(- (\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \alpha_c (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)}) \right) \quad (17)$$

Notice that $W_{c,+}$ and $W_{c,-}$, the weights of correctly and incorrectly classified examples, can be rewritten as

$$W_{c,+} = \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} (\mathbb{I}_{y'=y \wedge y' \in h_c(x)} + \mathbb{I}_{y' \neq y \wedge y' \notin h_c(x)}) w_{c,x,y'},$$

$$W_{c,-} = \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} (\mathbb{I}_{y'=y \wedge y' \notin h_c(x)} + \mathbb{I}_{y' \neq y \wedge y' \in h_c(x)}) w_{c,x,y'}.$$

Replacing $\alpha_c = \frac{1}{2} \log \left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right)$ in Equation (17) gives

$$\begin{aligned} & \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} \exp \left(- (\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \alpha_c (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)}) \right) \\ &= \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} \exp \left(- (\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \frac{1}{2} \log \left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right) (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)}) \right) \\ & \quad (\exp(a \log(b)) = b^a) \\ &= \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} \left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right)^{(-(\mathbb{I}_{y'=y} - \mathbb{I}_{y' \neq y}) \frac{1}{2} (\mathbb{I}_{y' \in h_c(x)} - \mathbb{I}_{y' \notin h_c(x)}))} \\ & \quad (\exp(a \log(b)) = b^a) \\ &= \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} (\mathbb{I}_{y'=y \wedge y' \in h_c(x)} + \mathbb{I}_{y' \neq y \wedge y' \notin h_c(x)}) \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} \\ & \quad + \sum_{\substack{(x,y) \in S \\ h_c(x) \neq \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} (\mathbb{I}_{y'=y \wedge y' \notin h_c(x)} + \mathbb{I}_{y' \neq y \wedge y' \in h_c(x)}) \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \\ &= (W_{c,+}) \cdot \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + (W_{c,-}) \cdot \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \end{aligned}$$

Similarly, using the fact that $\sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} = 1$, we obtain

$$\sum_{\substack{(x,y) \in S \\ h_c(x) = \emptyset}} \sum_{y' \in \mathcal{Y}} w_{c,x,y'} = 1 - W_{c,+} - W_{c,-}.$$

Combining these two results with Equation (17) gives the value of Z_c :

$$Z_c = (1 - W_{c,+} - W_{c,-}) + (W_{c,+}) \cdot \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + (W_{c,-}) \cdot \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}}.$$

To show that $Z_c \leq 1$ note that

$$\Leftrightarrow \left(\sqrt{W_{c,+} + \frac{1}{n}} - \sqrt{W_{c,-} + \frac{1}{n}} \right)^2 \geq 0$$

$$W_{c,+} + \frac{1}{n} + W_{c,-} + \frac{1}{n} - 2\sqrt{W_{c,+} + \frac{1}{n}}\sqrt{W_{c,-} + \frac{1}{n}} \geq 0$$

$$\begin{aligned}
&\Leftrightarrow W_{c,+} + W_{c,-} + \frac{2}{n} - \left(W_{c,+} + \frac{1}{n}\right) \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} - \left(W_{c,-} + \frac{1}{n}\right) \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \geq 0 \\
&\Leftrightarrow W_{c,+} + W_{c,-} + \frac{2}{n} - W_{c,+} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} - \frac{1}{n} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} - W_{c,-} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} - \frac{1}{n} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \geq 0 \\
&\Leftrightarrow -W_{c,+} - W_{c,-} - \frac{2}{n} + W_{c,+} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + \frac{1}{n} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + W_{c,-} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} + \frac{1}{n} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \leq 0 \\
&\Leftrightarrow 1 - W_{c,+} - W_{c,-} - \frac{2}{n} + W_{c,+} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + \frac{1}{n} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + W_{c,-} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} + \frac{1}{n} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \leq 1 \\
&\Leftrightarrow Z_c - \frac{2}{n} + \frac{1}{n} \sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + \frac{1}{n} \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \leq 1 \\
&\Leftrightarrow Z_c - \frac{2}{n} + \frac{1}{n} \left(\sqrt{\frac{W_{c,-} + \frac{1}{n}}{W_{c,+} + \frac{1}{n}}} + \sqrt{\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}}} \right) \leq 1 \\
&\hspace{15em} (\forall x > 0, x + \frac{1}{x} \geq 2.) \\
&\Rightarrow Z_c \leq 1.
\end{aligned}$$

It concludes the proof of the theorem. \square

G Proof of Theorem 3

Theorem 3 (Generalization Guarantees). Let \mathcal{D}_S be a distribution over $\mathcal{X} \times \mathcal{Y}$, let S be a set of n examples drawn i.i.d. from \mathcal{D}_S , and let T be a set of m triplets (obtained as described in Section 2). Let $H(\cdot)$ be the classifier obtained after C iterations of TripletBoost (Algorithm 1) using S and T as input. Let \mathcal{H} be a set of triplet classifiers as defined in Section 2.1. Then, given a margin parameter $\theta > \sqrt{\frac{\log |\mathcal{H}|}{16|\mathcal{Y}|^2 n}}$ and a measure of the confidence of $H(\cdot)$ in its predictions $\theta_H(x, y)$, with probability at least $1 - \delta$, we have that

$$\begin{aligned}
&\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) \neq y] \leq \mathbb{P}_{(x,y) \in S} [\theta_H(x, y) \leq \theta] \\
&\quad + \mathcal{O} \left(\sqrt{\frac{\log \left(\frac{1}{\delta} \right)}{n}} + \log \left(\frac{|\mathcal{Y}|^2 n \theta^2}{\log |\mathcal{H}|} \right) \frac{\log |\mathcal{H}|}{n \theta^2} \right)
\end{aligned}$$

Furthermore we have that

$$\mathbb{P}_{(x,y) \in S} [\theta_H(x, y) \leq \theta] \leq \frac{|\mathcal{Y}|}{2} \prod_{c=1}^C Z_c \sqrt{\left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right)^\theta}.$$

Proof. This result is inherited from AdaBoost.MO and can be obtained directly by following the steps of Exercise 10.3 in Chapter 10 of the book of Schapire and Freund [2012]. Following our notations $\theta_{f,\eta}(x, y)$ is equal to $\mathcal{M}_{f,\eta}(x, y)$ and the free parameter n in their proof has been chosen as $\left\lceil \frac{16}{\theta^2} \log \left(\frac{4|\mathcal{Y}|^2 n \theta^2}{\log |\mathcal{H}|} \right) \right\rceil$. The second part of the theorem can be obtained in a similar way by following the steps of Exercise 10.4 in Chapter 10 of the book of Schapire and Freund [2012]. For the sake of completeness we write the complete proof below.

First, we define several quantities used throughout. Recall that our approach outputs a strong classifier $H(\cdot)$ defined as

$$H(x) = \arg \max_{y \in \mathcal{Y}} \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(x) \neq y \wedge y \in h_c(x)} \quad (18)$$

As in the proof of Theorem 2, let $F(x, y) = \sum_{c=1}^C \alpha_c \mathbb{I}_{h_c(x) \neq y} (\mathbb{I}_{y \in h_c(x)} - \mathbb{I}_{y \notin h_c(x)})$ and let $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ be the normalized version of F :

$$f(x, y) = \sum_{c=1}^C a_c \mathbb{I}_{h_c(x) \neq y} (\mathbb{I}_{y \in h_c(x)} - \mathbb{I}_{y \notin h_c(x)}). \quad (19)$$

where $a_c = \frac{\alpha_c}{\sum_{c=1}^C \alpha_c} \in [0, 1]$ since we have that $\alpha_c \geq 0$ because of our labelling strategy for the triplet classifiers. From the proof of Theorem 2, recall that we have $H(x) = \arg \max_{y \in \mathcal{Y}} F(x, y) = \arg \max_{y \in \mathcal{Y}} f(x, y)$.

Notice that f is part of the following convex set:

$$\text{co}(\mathcal{H}) = \left\{ f : x, y \mapsto \sum_{c=1}^C a_c \mathbb{I}_{h_c(x) \neq \varnothing} (\mathbb{I}_{y \in h_c(x)} - \mathbb{I}_{y \notin h_c(x)}) \mid a_1, \dots, a_C \geq 0; \sum_{c=1}^C a_c = 1; h_1, \dots, h_C \in \mathcal{H} \right\}. \quad (20)$$

Let us define the following set of unweighted averages of size T :

$$\text{un}_T(\mathcal{H}) = \left\{ f : x, y \mapsto \frac{1}{T} \sum_{t=1}^T \mathbb{I}_{h_t(x) \neq \varnothing} (\mathbb{I}_{y \in h_t(x)} - \mathbb{I}_{y \notin h_t(x)}) \mid h_1, \dots, h_T \in \mathcal{H} \right\}. \quad (21)$$

For $f \in \text{co}(\mathcal{H})$, $\eta > 0$, and $(x, y) \in \mathcal{X} \times \mathcal{Y}$, let

$$\nu_{f, \eta}(x, y) = -\frac{1}{\eta} \log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \right). \quad (22)$$

We use this quantity to define the following confidence margin for $(x, y) \in \mathcal{X} \times \mathcal{Y}$:

$$\theta_H(x, y) = \theta_{f, \eta}(x, y) = \frac{1}{2} \left(\nu_{f, \eta}(x, y) - \max_{y' \neq y} \nu_{f, \eta}(x, y') \right). \quad (23)$$

Bound on the true risk of $H(\cdot)$. Note that $\theta_{f, \eta}(x, y) \in [-1, 1]$ since $f(x, y) \in [-1, 1]$. Similarly, given an example $(x, y) \in \mathcal{X} \times \mathcal{Y}$, if $H(\cdot)$ makes an error then it implies that there exists a label $\ell \neq y$ such that

$$\begin{aligned} & F(x, y) \leq F(x, \ell) \\ & f(x, y) \leq f(x, \ell) \\ & \Leftrightarrow \eta f(x, y) \leq \eta f(x, \ell) \\ & \Leftrightarrow \exp(-\eta f(x, y)) - \exp(\eta f(x, y)) \geq \exp(-\eta f(x, \ell)) - \exp(\eta f(x, \ell)) \\ & \Leftrightarrow \exp(-\eta f(x, y)) - \exp(\eta f(x, y)) + \sum_{y' \in \mathcal{Y}} \exp(\eta f(x, y')) \geq \exp(-\eta f(x, \ell)) - \exp(\eta f(x, \ell)) + \sum_{y' \in \mathcal{Y}} \exp(\eta f(x, y')) \\ & \Leftrightarrow \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \geq \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') \right) \\ & \Leftrightarrow \log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \right) \geq \log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') \right) \right) \\ & \Leftrightarrow \nu_{f, \eta}(x, y) \leq \nu_{f, \eta}(x, \ell) \\ & \Leftrightarrow \nu_{f, \eta}(x, y) - \nu_{f, \eta}(x, \ell) \leq 0 \\ & \Leftrightarrow \theta_{f, \eta}(x, y) \leq 0. \end{aligned} \quad (\eta > 0.)$$

It implies that

$$\mathbb{P}_{(x, y) \sim \mathcal{D}_S} [H(x) \neq y] = \mathbb{P}_{(x, y) \sim \mathcal{D}_S} [\theta_{f, \eta}(x, y) \leq 0]$$

and hence, in the following, we prove an upper bound on $\mathbb{P}_{(x, y) \sim \mathcal{D}_S} [\theta_{f, \eta}(x, y) \leq 0]$.

To this end, let us first define $\tilde{f}(x, y) = \frac{1}{T} \sum_{t=1}^T \mathbb{I}_{\tilde{h}_t(x) \neq \varnothing} (\mathbb{I}_{y \in \tilde{h}_t(x)} - \mathbb{I}_{y \notin \tilde{h}_t(x)}) \in \text{un}_T$ with $\tilde{h}_1, \dots, \tilde{h}_T$ drawn from \mathcal{H} independently at random following the probability distribution defined by the weights a_1, \dots, a_C . Notice that $\tilde{f}(x, y) \in [-1, 1]$ and that given $x \in \mathcal{X}$ and $y' \in \mathcal{Y}$ fixed, $\mathbb{E}_{\tilde{f}} [\tilde{f}(x, y')] = f(x, y')$. Hence, using Hoeffding's inequality, it holds that

$$\mathbb{P}_{\tilde{f}} \left[\left| \tilde{f}(x, y') - f(x, y') \right| \geq \frac{\theta}{4} \right] \leq 2 \exp \left(\frac{-T\theta^2}{32} \right).$$

Then, by the union bound, it holds that

$$\mathbb{P}_{\tilde{f}} \left[\exists y' \in \mathcal{Y} : \left| \tilde{f}(x, y') - f(x, y') \right| \geq \frac{\theta}{4} \right] \leq 2 |\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right). \quad (24)$$

Consider the following technical fact from Schapire and Freund [2012]. Define the grid

$$\varepsilon_\theta = \left\{ \frac{4 \log |\mathcal{Y}|}{i\theta} : i = 1, \dots, \left\lceil \frac{8 \log |\mathcal{Y}|}{\theta^2} \right\rceil \right\}.$$

For any $\eta > 0$, let $\hat{\eta}$ be the closest value in ε_θ to η . Then for all $f \in \text{co}(\mathcal{H})$ and for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$,

$$|\nu_{f,\eta}(x, y) - \nu_{f,\hat{\eta}}(x, y)| \leq \frac{\theta}{4}. \quad (25)$$

We will now show that, with high probability, $\theta_{f,\eta}(x, y)$ is close to $\theta_{\tilde{f},\hat{\eta}}(x, y)$ for any distribution $\mathcal{D}_\mathcal{T}$ over $\mathcal{X} \times \mathcal{Y}$.

Let $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ then we have that

$$\left| \nu_{f,\eta}(x, y) - \nu_{\tilde{f},\hat{\eta}}(x, y) \right| \leq \left| \nu_{f,\eta}(x, y) - \nu_{\tilde{f},\eta}(x, y) \right| + \left| \nu_{\tilde{f},\eta}(x, y) - \nu_{\tilde{f},\hat{\eta}}(x, y) \right| \quad (26)$$

(Using Inequality (25) since $\tilde{f} \in \text{co}(\mathcal{H})$.)

$$\leq \left| \nu_{f,\eta}(x, y) - \nu_{\tilde{f},\eta}(x, y) \right| + \frac{\theta}{4} \quad (27)$$

$$\leq \left| -\frac{1}{\eta} \log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \right) \right| \quad (28)$$

$$+ \frac{1}{\eta} \log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) \tilde{f}(x, y') \right) \right) \Big| + \frac{\theta}{4} \quad (29)$$

$$\leq \left| \frac{1}{\eta} \log \left(\frac{\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) \tilde{f}(x, y') \right)}{\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right)} \right) \right| + \frac{\theta}{4} \quad (30)$$

(By induction, one can show that $\frac{\sum_i a_i}{\sum_i b_i} \leq \max_i \frac{a_i}{b_i}$ for positive a_i and b_i .)

$$\leq \left| \frac{1}{\eta} \log \left(\max_{y' \in \mathcal{Y}} \frac{\exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) \tilde{f}(x, y') \right)}{\exp \left(-\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right)} \right) \right| + \frac{\theta}{4} \quad (31)$$

$$\leq \left| \frac{1}{\eta} \log \left(\max_{y' \in \mathcal{Y}} \exp \left(\eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') - \eta (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) \tilde{f}(x, y') \right) \right) \right| + \frac{\theta}{4} \quad (32)$$

$$\leq \left| \max_{y' \in \mathcal{Y}} (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') - (\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) \tilde{f}(x, y') \right| + \frac{\theta}{4} \quad (33)$$

$$\leq \max_{y' \in \mathcal{Y}} \left| f(x, y') - \tilde{f}(x, y') \right| + \frac{\theta}{4}. \quad (34)$$

And then it implies that:

$$\begin{aligned} \left| \theta_{f,\eta}(x, y) - \theta_{\tilde{f},\hat{\eta}}(x, y) \right| &\leq \left| \frac{1}{2} \left(\nu_{f,\eta}(x, y) - \max_{y' \neq y} \nu_{f,\eta}(x, y') \right) - \frac{1}{2} \left(\nu_{\tilde{f},\hat{\eta}}(x, y) + \max_{y' \neq y} \nu_{\tilde{f},\hat{\eta}}(x, y') \right) \right| \\ &\leq \frac{1}{2} \left| \nu_{f,\eta}(x, y) - \nu_{\tilde{f},\hat{\eta}}(x, y) \right| + \frac{1}{2} \left| \max_{y' \neq y} \nu_{\tilde{f},\hat{\eta}}(x, y') - \max_{y' \neq y} \nu_{f,\eta}(x, y') \right| \end{aligned}$$

(Noticing that if $|g(x) - g'(x)| \leq B$ for all x then $|\max_x g(x) - \max_x g'(x)| \leq B$ and applying Inequality (34) twice.)

$$\leq \max_{y' \in \mathcal{Y}} \left| f(x, y') - \tilde{f}(x, y') \right| + \frac{\theta}{4}$$

Then, using Inequality (24), for any $(x, y) \in \mathcal{X} \times \mathcal{Y}$, we have

$$\mathbb{P}_{\tilde{f}} \left[\left| \theta_{f,\eta}(x, y) - \theta_{\tilde{f},\hat{\eta}}(x, y) \right| \geq \frac{\theta}{2} \right] \leq 2 |\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right).$$

Then, given any distribution $\mathcal{D}_{\mathcal{T}}$ over $\mathcal{X} \times \mathcal{Y}$, we have that:

$$\mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{T}}} \mathbb{P}_{\tilde{f}} \left[\left| \theta_{f,\eta}(x,y) - \theta_{\tilde{f},\hat{\eta}}(x,y) \right| \geq \frac{\theta}{2} \right] \leq \mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{T}}} 2|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right)$$

which implies

$$\mathbb{P}_{\mathcal{D}_{\mathcal{T}}, \tilde{f}} \left[\left| \theta_{f,\eta}(x,y) - \theta_{\tilde{f},\hat{\eta}}(x,y) \right| \geq \frac{\theta}{2} \right] \leq 2|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right). \quad (35)$$

This last inequality shows that $\theta_{f,\eta}(x,y)$ is close to $\theta_{\tilde{f},\hat{\eta}}(x,y)$ for any distribution $\mathcal{D}_{\mathcal{T}}$ over $\mathcal{X} \times \mathcal{Y}$.

Let θ and T be fixed. Given \tilde{f} and $\hat{\eta}$, let B_i be a Bernoulli random variable that is 1 if $\theta_{\tilde{f},\hat{\eta}}(x_i, y_i) \leq \frac{\theta}{2}$ with $(x_i, y_i) \in S$. We have that

$$\begin{aligned} \mathbb{P}_{(x,y) \in S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] &= \frac{1}{n} \sum_{i=1}^n B_i, \\ \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] &= \mathbb{E}_{S \sim \mathcal{D}_S^n} \left[\mathbb{P}_{(x,y) \in S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] \right]. \end{aligned}$$

Then, using Hoeffding's inequality, it holds that:

$$\mathbb{P}_{S \sim \mathcal{D}_S^n} \left[\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] \geq \mathbb{P}_{(x,y) \in S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] + \epsilon \right] \leq \exp(-2\epsilon^2 n).$$

Using the union bound for both \tilde{f} and $\hat{\eta}$, it holds that $\forall \tilde{f} \in \text{un}_T$ and $\forall \hat{\eta} \in \varepsilon_\theta$, with probability at least $1 - \delta$ over the choice of the random training set S :

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] \leq \mathbb{P}_{(x,y) \in S} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] + \sqrt{\frac{\log \left(\frac{|\varepsilon_\theta| |\mathcal{H}|^T}{\delta} \right)}{2n}}. \quad (36)$$

Now, notice that for two events E and E' it holds that $\mathbb{P}[E] \leq \mathbb{P}[E'] + \mathbb{P}[E, \neg E']$ and thus:

$$\begin{aligned} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [\theta_{f,\eta}(x,y) \leq 0] &\leq \mathbb{P}_{(x,y) \sim \mathcal{D}_{S,\tilde{f}}} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] + \mathbb{P}_{(x,y) \sim \mathcal{D}_{S,\tilde{f}}} \left[\theta_{f,\eta}(x,y) \leq 0, \theta_{\tilde{f},\hat{\eta}}(x,y) > \frac{\theta}{2} \right] \\ &\quad \text{(Applying Inequality (35) to the rightmost term.)} \\ &\leq \mathbb{P}_{(x,y) \sim \mathcal{D}_{S,\tilde{f}}} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] + 2|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right). \end{aligned}$$

Similarly it holds that

$$\mathbb{P}_{(x,y) \in S, \tilde{f}} \left[\theta_{\tilde{f},\hat{\eta}}(x,y) \leq \frac{\theta}{2} \right] \leq \mathbb{P}_{(x,y) \in S} [\theta_{f,\eta}(x,y) \leq \theta] + 2|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right).$$

Combining these two results with Inequality (36) yields with probability at least $1 - \delta$:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [\theta_{f,\eta}(x,y) \leq 0] \leq \mathbb{P}_{(x,y) \in S} [\theta_{f,\eta}(x,y) \leq \theta] + \sqrt{\frac{\log \left(\frac{|\varepsilon_\theta| |\mathcal{H}|^T}{\delta} \right)}{2n}} + 4|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right). \quad (37)$$

Choosing $T = \left\lceil \frac{16}{\theta^2} \log \left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log |\mathcal{H}|} \right) \right\rceil$ with $\theta > \sqrt{\frac{\log |\mathcal{H}|}{16|\mathcal{Y}|^2 n}}$ and noticing that $|\varepsilon_\theta| \leq \frac{8 \log |\mathcal{Y}|}{\theta^2} + 1$ the error term becomes:

$$\begin{aligned} &\sqrt{\frac{\log \left(\frac{|\varepsilon_\theta| |\mathcal{H}|^T}{\delta} \right)}{2n}} + 4|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right) \\ &= \sqrt{\frac{\log \left(\frac{|\varepsilon_\theta|}{\delta} \right) + T \log |\mathcal{H}|}{2n}} + 4|\mathcal{Y}| \exp \left(\frac{-T\theta^2}{32} \right) \end{aligned}$$

$$\begin{aligned}
&= \sqrt{\frac{\log\left(\frac{|\varepsilon_\theta|}{\delta}\right) + \left(\frac{16}{\theta^2} \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right) + 1\right) \log|\mathcal{H}|}{2n}} + 4|\mathcal{Y}| \exp\left(\frac{-\frac{16}{\theta^2} \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right) \theta^2}{32}\right) \\
&= \sqrt{\frac{\log\left(\frac{|\varepsilon_\theta|}{\delta}\right) + \left(\frac{16}{\theta^2} \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right) + 1\right) \log|\mathcal{H}|}{2n}} + 4|\mathcal{Y}| \exp\left(-\frac{1}{2} \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right)\right) \\
&= \sqrt{\frac{\log\left(\frac{|\varepsilon_\theta|}{\delta}\right) + \left(\frac{16}{\theta^2} \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right) + 1\right) \log|\mathcal{H}|}{2n}} + 4|\mathcal{Y}| \sqrt{\frac{\log|\mathcal{H}|}{16|\mathcal{Y}|^2 n \theta^2}} \\
&= \sqrt{\frac{\log\left(\frac{|\varepsilon_\theta|}{\delta}\right) + \left(\frac{16}{\theta^2} \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right) + 1\right) \log|\mathcal{H}|}{2n}} + \sqrt{\frac{\log|\mathcal{H}|}{n \theta^2}} \\
&= \sqrt{\frac{\log\left(\frac{8 \log|\mathcal{Y}|}{\theta^2} + 1\right)}{2n} + \frac{\log\left(\frac{1}{\delta}\right)}{2n} + 8 \log\left(\frac{16|\mathcal{Y}|^2 n \theta^2}{\log|\mathcal{H}|}\right) \frac{\log|\mathcal{H}|}{n \theta^2} + \frac{\log|\mathcal{H}|}{2n} + \sqrt{\frac{\log|\mathcal{H}|}{n \theta^2}}}.
\end{aligned}$$

This concludes the first part of the proof.

Bound on the empirical margin of $H(\cdot)$. This part of the proof is very similar to the proof of Theorem 2 and the main difference is that θ has to be taken into account. To this end, let us assume that the strong classifier violates the margin, that is $\theta_{f,\eta}(x, y) \leq \theta$. Given a label $\ell \neq y$, we define the following quantities:

$$\begin{aligned}
z(y') &= \eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') - \eta\theta, \\
z_\ell(y') &= \eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') + \eta\theta.
\end{aligned}$$

Notice that:

$$\begin{aligned}
&\theta_{f,\eta}(x, y) \leq \theta \\
\Leftrightarrow &\frac{1}{2} \left(\nu_{f,\eta}(x, y) - \max_{\ell \neq y} \nu_{f,\eta}(x, \ell) \right) \leq \theta \\
\Leftrightarrow &\nu_{f,\eta}(x, y) - \theta \leq \max_{\ell \neq y} \nu_{f,\eta}(x, \ell) + \theta \\
\Leftrightarrow &-\log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \right) - \eta\theta \leq \\
&\max_{\ell \neq y} -\log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') \right) \right) + \eta\theta \\
\Leftrightarrow &\log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \right) + \eta\theta \geq \\
&\min_{\ell \neq y} \log \left(\frac{1}{|\mathcal{Y}|} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') \right) \right) - \eta\theta \\
\Leftrightarrow &\sum_{y' \in \mathcal{Y}} \exp(\eta\theta) \exp \left(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') \right) \geq \\
&\min_{\ell \neq y} \sum_{y' \in \mathcal{Y}} \exp(-\eta\theta) \exp \left(-\eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') \right) \\
\Leftrightarrow &\sum_{y' \in \mathcal{Y}} \exp \left(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') + \eta\theta \right) \geq \\
&\min_{\ell \neq y} \sum_{y' \in \mathcal{Y}} \exp \left(-\eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') - \eta\theta \right)
\end{aligned}$$

\Leftrightarrow

$$\sum_{y' \in \mathcal{Y}} \exp(-z(y')) \geq \min_{\ell \neq y} \sum_{y' \in \mathcal{Y}} \exp(-z_\ell(y'))$$

Then, when the strong classifier violates the margin, that is $\theta_{f,\eta}(x, y) \leq \theta$, we have that:

$$\begin{aligned} \sum_{y' \in \mathcal{Y}} \exp(-z(y')) &\geq \frac{1}{2} \sum_{y' \in \mathcal{Y}} \exp(-z(y')) + \frac{1}{2} \min_{\ell \neq y} \sum_{y' \in \mathcal{Y}} \exp(-z_\ell(y')) \\ &\geq \min_{\ell \neq y} \left[\frac{1}{2} \sum_{y' \in \mathcal{Y}} \exp(-z(y')) + \frac{1}{2} \sum_{y' \in \mathcal{Y}} \exp(-z_\ell(y')) \right] \\ &\geq \frac{1}{2} \min_{\ell \neq y} \sum_{y' \in \mathcal{Y}} \left[\exp(-z(y')) + \exp(-z_\ell(y')) \right] \\ &\geq \frac{1}{2} \min_{\ell \neq y} \sum_{y' \in \mathcal{Y}} \left[\exp(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') + \eta\theta) + \exp(-\eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') - \eta\theta) \right] \\ &\quad (\text{Dropping all the terms that do not depend on } y \text{ and } \ell \text{ since } \forall x \in \mathbb{R}, \exp x > 0.) \\ &\geq \frac{1}{2} \min_{\ell \neq y} \sum_{y' \in \{y, \ell\}} \left[\exp(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') + \eta\theta) + \exp(-\eta(\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) f(x, y') - \eta\theta) \right] \\ &\quad (\forall y' \in \{y, \ell\}, (\mathbb{I}_{\ell=y'} - \mathbb{I}_{\ell \neq y'}) = -(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'})) \\ &\geq \frac{1}{2} \min_{\ell \neq y} \sum_{y' \in \{y, \ell\}} \left[\exp(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') + \eta\theta) + \exp(\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') - \eta\theta) \right] \\ &\geq \frac{1}{2} \min_{\ell \neq y} \sum_{y' \in \{y, \ell\}} \left[\exp(-z(y')) + \exp(z(y')) \right] \\ &\quad (\forall x > 0, x + \frac{1}{x} \geq 2.) \\ &\geq 2 \end{aligned}$$

Now assume that the strong classifier violates the margin M times, that is $\theta_{f,\eta}(x, y) \leq \theta$ for M examples, then

$$\begin{aligned} 2M &\leq \sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-z(y')) \\ \Leftrightarrow 2M &\leq \sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y') + \eta\theta) \\ \Leftrightarrow 2M &\leq \sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-\eta(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) f(x, y')) \exp(\eta\theta) \\ &\quad (\text{Let } \eta = \sum_c \alpha_c \text{ and recall that } (\sum_c \alpha_c) f(x, y) = F(x, y).) \\ \Leftrightarrow 2M &\leq \sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y')) \exp\left(\sum_c \alpha_c \theta\right) \\ \Leftrightarrow 2M &\leq \exp\left(\sum_c \alpha_c \theta\right) \sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y')) \end{aligned}$$

Recall that $\sum_{(x,y) \in S} \sum_{y' \in \mathcal{Y}} \exp(-(\mathbb{I}_{y=y'} - \mathbb{I}_{y \neq y'}) F(x, y')) = n |\mathcal{Y}| \prod_{c=1}^C Z_c$ as shown in Equation (16) in the proof of Theorem 2.

$$\begin{aligned} 2M &\leq \exp\left(\sum_c \alpha_c \theta\right) n |\mathcal{Y}| \prod_{c=1}^C Z_c \\ \Leftrightarrow 2M &\leq \exp\left(\sum_c \alpha_c \theta\right) n |\mathcal{Y}| \prod_{c=1}^C Z_c \\ \Leftrightarrow 2M &\leq n |\mathcal{Y}| \prod_{c=1}^C Z_c \exp(\alpha_c \theta) \end{aligned}$$

$$\Leftrightarrow \frac{M}{n} \leq \frac{|\mathcal{Y}|}{2} \prod_{c=1}^C Z_c \exp(\alpha_c \theta)$$

Replacing α_c by its value, that is $\alpha_c = \frac{1}{2} \log \left(\frac{W_{c,+} + \frac{1}{n}}{W_{c,-} + \frac{1}{n}} \right)$ and noticing that $\mathbb{P}_{(x,y) \in S} [\theta_{f,\eta}(x,y) \leq \theta] = \frac{M}{n}$ concludes the proof. \square

H Proof of Theorem 4

Theorem 4 (Lower bound on the probability that a strong classifier abstains). *Let $n \geq 2$ be the number of training examples, $p = \frac{2n^k}{n^3}$ with $k \in \left[0, 3 - \frac{\log(2)}{\log(n)}\right)$ be the probability that a triplet is available in the triplet set T and $C = \frac{n^\beta}{2}$ with $\beta \in \left[0, 1 + \frac{\log(n-1)}{\log(n)}\right]$ be the number of classifiers combined in the learned classifier. Let \mathcal{A} be any algorithm learning a classifier $H(\cdot) = \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{y \in h_c(\cdot)} \right)$ that combines several triplet classifiers using some weights $\alpha_c \in \mathbb{R}$. Assume that triplet classifiers that abstain on all the training examples have a weight of 0 (that is if $h_c(x_i) = \vartheta$ for all the examples $(x_i, y_i) \in S$ then $\alpha_c = 0$). Then the probability that H abstains on a test example is bounded as follows:*

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \left(1 - p + p(1-p)^n\right)^C. \quad (38)$$

Proof. Recall that $H(\cdot) = \arg \max_{y \in \mathcal{Y}} \left(\sum_{c=1}^C \alpha_c \mathbb{I}_{y \in h_c(\cdot)} \right)$. Each component $\alpha_c \mathbb{I}_{y \in h_c(\cdot)}$ of H abstains for a given example x if $h_c(x) = \vartheta$ or $\alpha_c = 0$. Overall H abstains if all its components also abstain and we bound this probability here.

Recall that $S = \{(x_i, y_i)\}_{i=1}^n$ is a set of n examples drawn i.i.d. from an unknown distribution \mathcal{D}_S over the space $\mathcal{X} \times \mathcal{Y}$ where $|\mathcal{Y}| < \infty$ and recall that C is the number of classifiers selected in H . Finally recall that each triplet is obtained with probability p independently of the other triplets. It implies that each triplet classifier $h_c(\cdot)$ abstains with probability $q = 1 - p$ independently of the other triplet classifiers.

First we start by defining several random variables. Let $Y_{1,1}, \dots, Y_{n,C}$ be nC independent random variables encoding the event that a classifier abstains on one example:

$$Y_{i,c} = \begin{cases} 0 & \text{if } h_c \text{ abstains for example } x_i, \\ 1 & \text{otherwise.} \end{cases} \quad (39)$$

From the definition of the classifiers each of these random variables follow a Bernoulli distribution $B(1, p)$.

Let $Y_{x,1}, \dots, Y_{x,C}$ be C independent random variables encoding the event that a classifier abstains on a new example x :

$$Y_{x,c} = \begin{cases} 0 & \text{if } h_c \text{ abstains for example } x, \\ 1 & \text{otherwise.} \end{cases} \quad (40)$$

From the definition of the classifiers each of these random variables follow a Bernoulli distribution $B(1, p)$.

Let V_1, \dots, V_C be C independent random variables encoding the probability that a triplet classifier abstains on all the training examples:

$$V_c = \sum_{i=1}^n Y_{i,c}. \quad (41)$$

As a sum of n independent Bernoulli trials, these random variables follow a Binomial distribution $B(n, p)$.

Using the law of total probabilities we have that:

$$\begin{aligned} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] &= \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[H(x) = \vartheta \mid \sum_{c=1}^C V_c Y_{x,c} = 0 \right] \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C V_c Y_{x,c} = 0 \right] \\ &+ \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[H(x) = \vartheta \mid \sum_{c=1}^C V_c Y_{x,c} \neq 0 \right] \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C V_c Y_{x,c} \neq 0 \right]. \end{aligned}$$

Note that from the assumption that $\forall (x_i, y_i) \in S, h_c(x_i) = \vartheta$ then $\alpha_c = 0$ and the definition of the random variables $Y_{x,\cdot}$, we have that:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[H(x) = \vartheta \mid \sum_{c=1}^C V_c Y_{x,c} = 0 \right] = 1.$$

Similarly, we have that:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[H(x) = \vartheta \left| \sum_{c=1}^C V_c Y_{x,c} \neq 0 \right. \right] \geq 0. \quad (42)$$

Hence we have that:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C V_c Y_{x,c} = 0 \right]. \quad (43)$$

Let U_1, \dots, U_C be C independent random variables such that:

$$U_c = \begin{cases} 1 & \text{if } V_c Y_{x,c} = 0, \\ 0 & \text{otherwise.} \end{cases} \quad (44)$$

These random variables follow a Bernoulli distribution $B(1, (1-p) + p(1-p)^n)$. To see that, note that V_c and $Y_{x,c}$ are independent and thus we have that:

$$\begin{aligned} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [V_c Y_{x,c} = 0] &= 1 - \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [V_c Y_{x,c} \neq 0] \\ &= 1 - \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [\{V_c \neq 0\} \cap \{Y_{x,c} \neq 0\}] \\ &= 1 - \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [V_c \neq 0] \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [Y_{x,c} \neq 0] \\ &\quad (Y_{x,c} \text{ follows a Bernoulli distribution } B(1, p).) \\ &= 1 - \left(1 - \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [V_c = 0] \right) p \\ &\quad (V_c \text{ follows a Binomial distribution } B(n, p).) \\ &= 1 - \left(1 - \binom{n}{0} p^0 (1-p)^{n-0} \right) p \\ &= 1 - (1 - (1-p)^n) p \\ &= (1-p) + p(1-p)^n. \end{aligned}$$

The r.h.s. of Inequality (43) can then be written as:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C V_c Y_{x,c} = 0 \right] = \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C U_c = c \right] \quad (45)$$

which, by definition of the random variables U , corresponds to the probability of obtaining c successes among c Bernoulli trials. In other words the random variable $U = \sum_{c=1}^C U_j$ follows a Binomial distribution $B(c, (1-p) + p(1-p)^n)$. Following this we have that:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C V_c Y_{x,c} = 0 \right] = \binom{c}{c} ((1-p) + p(1-p)^n)^C (1 - (1-p) + p(1-p)^n)^{C-C} \quad (46)$$

$$= ((1-p) + p(1-p)^n)^C. \quad (47)$$

Hence plugging this result in Inequality (43) we have that:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq ((1-p) + p(1-p)^n)^C$$

which concludes the proof. \square

I Proof of Example 1

Note that here we prove a more general version of Example 1 than the one presented in the main paper. The latter follows directly by choosing $\beta = 1 + \frac{\log(n-1)}{\log(n)}$.

Example 1. In Theorem 4, we choose different values for $p = 2n^{k-3}$ the probability of that a triplet is available and $C = \frac{n^\beta}{2}$ the number of combined classifiers. Then we take the limit as $n \rightarrow \infty$ to obtain the following results.

If $0 \leq \beta < 1$ then:

$$\lim_{n \rightarrow +\infty} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \begin{cases} 1 & \text{if } 0 \leq k < 3 - \beta, \\ \exp(-1) & \text{if } k = 3 - \beta, \\ 0 & \text{if } 3 - \beta < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $\beta = 1$ then:

$$\lim_{n \rightarrow +\infty} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \begin{cases} 1 & \text{if } 0 \leq k < 2, \\ \exp(\exp(-2) - 1) & \text{if } k = 2, \\ 0 & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $1 < \beta \leq 1 + \frac{\log(n-1)}{\log(n)}$ then:

$$\lim_{n \rightarrow +\infty} \mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] \geq \begin{cases} 1 & \text{if } 0 \leq k < \frac{5-\beta}{2}, \\ \exp(-2) & \text{if } k = \frac{5-\beta}{2}, \\ 0 & \text{if } \frac{5-\beta}{2} < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \quad (48)$$

Proof. Replacing p and C by their values in Equation (38) of Theorem 4 and applying Lemma 1 gives the example. \square

The next Lemma is a technical result used in the proof of Theorem 4.

Lemma 1 (Limit of the right hand side of Theorem 4). Given $n \geq 2$, $k \in \left[0, 3 - \frac{\log(2)}{\log(n)}\right)$ and $\beta \in \left[0, 1 + \frac{\log(n-1)}{\log(n)}\right]$, define:

$$f(n, k, \beta) \doteq \left(1 - 2n^{k-3} + 2n^{k-3} (1 - 2n^{k-3})^n\right)^{\frac{n^\beta}{2}} \quad (49)$$

then the following limits hold:

If $0 \leq \beta < 1$ then:

$$\lim_{n \rightarrow +\infty} f(n, k, \beta) = \begin{cases} 1 & \text{if } 0 \leq k < 3 - \beta, \\ \exp(-1) & \text{if } k = 3 - \beta, \\ 0 & \text{if } 3 - \beta < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $\beta = 1$ then:

$$\lim_{n \rightarrow +\infty} f(n, k, \beta) = \begin{cases} 1 & \text{if } 0 \leq k < 2, \\ \exp(\exp(-2) - 1) & \text{if } k = 2, \\ 0 & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $1 < \beta \leq 1 + \frac{\log(n-1)}{\log(n)}$ then:

$$\lim_{n \rightarrow +\infty} f(n, k, \beta) = \begin{cases} 1 & \text{if } 0 \leq k < \frac{5-\beta}{2}, \\ \exp(-2) & \text{if } k = \frac{5-\beta}{2}, \\ 0 & \text{if } \frac{5-\beta}{2} < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

Proof. We are looking for the limit when $n \rightarrow \infty$ of the function defined in Equation (49).

$$\begin{aligned} & \lim_{n \rightarrow +\infty} \left(1 - 2n^{k-3} + 2n^{k-3} (1 - 2n^{k-3})^n\right)^{\frac{n^\beta}{2}} \\ &= \lim_{n \rightarrow +\infty} \exp\left(\frac{n^\beta}{2} \log\left(1 - 2n^{k-3} + 2n^{k-3} (1 - 2n^{k-3})^n\right)\right) \\ & \quad \text{(Maclaurin series of log since } |2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n| < 1.) \\ &= \lim_{n \rightarrow +\infty} \exp\left(\frac{n^\beta}{2} \left(-\sum_{l=1}^{\infty} \frac{(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n)^l}{l}\right)\right) \end{aligned}$$

$$\begin{aligned}
&= \lim_{n \rightarrow +\infty} \exp \left(- \sum_{l=1}^{\infty} \frac{n^{\beta}}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n \right)^l \right) \\
&= \lim_{n \rightarrow +\infty} \exp \left(\underbrace{- \frac{n^{\beta}}{2} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n \right)}_{\text{Main term, with } l=1} \right. \\
&\quad \left. \underbrace{- \sum_{l=2}^{\infty} \frac{n^{\beta}}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n \right)^l}_{\text{Remaining term, with } l > 1} \right). \tag{50}
\end{aligned}$$

In the following we study the limits of the two underlined terms in Equation (50).

Main term, with $l = 1$. To obtain the limit of this term we proceed in two steps. First we decompose it in a product of two sub terms and we study their individual limits. Then we consider the indeterminate forms that arise after considering the products of the limits and we obtain the correct limits by upper and lower bounding the term and showing that the limits of the two bounds are identical.

We are interested in:

$$\lim_{n \rightarrow +\infty} -\frac{n^{\beta}}{2} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n \right) = \lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n \right).$$

- On the one hand, since $\beta \in \left[0, 1 + \frac{\log(n-1)}{\log(n)}\right]$ we have:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} = \begin{cases} 0 & \text{if } 0 \leq k < 3 - \beta, \\ -1 & \text{if } k = 3 - \beta, \\ -\infty & \text{if } 3 - \beta < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \tag{51}$$

- On the other hand, we have:

$$\lim_{n \rightarrow +\infty} \left(1 - (1 - 2n^{k-3})^n \right) = \begin{cases} 0 & \text{if } 0 \leq k < 2, \\ 1 - \exp(-2) & \text{if } k = 2, \\ 1 & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \tag{52}$$

To see that note that:

$$\lim_{n \rightarrow +\infty} (1 - 2n^{k-3})^n = \exp \left(\lim_{n \rightarrow +\infty} n \log (1 - 2n^{k-3}) \right)$$

which is an indeterminate form that can be solved using l'Hôpital's rule. Given the derivatives:

$$\begin{aligned}
\frac{\partial \log (1 - 2n^{k-3})}{\partial n} &= \frac{-2(k-3)n^{k-4}}{1 - 2n^{k-3}}, \\
\frac{\partial \frac{1}{n}}{\partial n} &= \frac{-1}{n^2},
\end{aligned}$$

it follows that:

$$\begin{aligned}
\lim_{n \rightarrow +\infty} n \log (1 - 2n^{k-3}) &= \lim_{n \rightarrow +\infty} -\frac{2(3-k)n^{k-2}}{1 - 2n^{k-3}} \\
&= \begin{cases} 0 & \text{if } 0 \leq k < 2, \\ -2 & \text{if } k = 2, \\ -\infty & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}
\end{aligned}$$

It implies:

$$\lim_{n \rightarrow +\infty} (1 - 2n^{k-3})^n = \begin{cases} 1 & \text{if } 0 \leq k < 2, \\ \exp(-2) & \text{if } k = 2, \\ 0 & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

Combining Limits (51) and (52) we obtain the following limits:

If $0 \leq \beta < 1$ then:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 0 \leq k < 2, \\ 0 & \text{if } k = 2, \\ 0 & \text{if } 2 < k < 3 - \beta, \\ -1 & \text{if } k = 3 - \beta, \\ -\infty & \text{if } 3 - \beta < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $\beta = 1$ then:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 0 \leq k < 2, \\ \exp(-2) - 1 & \text{if } k = 2, \\ -\infty & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $1 < \beta \leq 2$ then:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 0 \leq k < 3 - \beta, \\ 0 & \text{if } k = 3 - \beta, \\ \text{Indeterminate} & \text{if } 3 - \beta < k < 2, \\ -\infty & \text{if } k = 2, \\ -\infty & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \quad (53)$$

The only indeterminate form arises when $1 < \beta \leq 2$ and $3 - \beta < k < 2$. To solve this indeterminate form we propose to upper and lower bound $-n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right)$ and to show that the limits coincides. Notice that, given our assumptions on k we have that $-2n^{k-3} \in (-1, 0]$ and hence, since $n \geq 2$ we can apply Bernoulli's inequalities² to obtain:

$$\begin{aligned} (1 - 2n^{k-2}) &\leq (1 - 2n^{k-3})^n \leq \left(1 - \frac{2n^{k-2}}{1 + 2n^{k-2} - 2n^{k-3}}\right) \\ \Leftrightarrow 1 - (1 - 2n^{k-2}) &\geq 1 - (1 - 2n^{k-3})^n \geq 1 - \left(1 - \frac{2n^{k-2}}{1 + 2n^{k-2} - 2n^{k-3}}\right) \\ \Leftrightarrow 2n^{k-2} &\geq 1 - (1 - 2n^{k-3})^n \geq \frac{2n^{k-2}}{1 + 2n^{k-2} - 2n^{k-3}} \\ \Leftrightarrow -2n^{2k-5+\beta} &\leq -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) \leq -\frac{2n^{2k-5+\beta}}{1 + 2n^{k-2} - 2n^{k-3}} \end{aligned}$$

Furthermore, recalling that we are only interested in the case where $1 < \beta \leq 2$ and $3 - \beta < k < 2$, we have the following limits:

$$\begin{aligned} \lim_{n \rightarrow +\infty} -2n^{2k-5+\beta} &= \begin{cases} 0 & \text{if } 3 - \beta < k < \frac{5-\beta}{2}, \\ -2 & \text{if } k = \frac{5-\beta}{2}, \\ -\infty & \text{if } \frac{5-\beta}{2} < k < 2. \end{cases} \\ \lim_{n \rightarrow +\infty} -\frac{2n^{2k-5+\beta}}{1 + 2n^{k-2} - 2n^{k-3}} &= \begin{cases} 0 & \text{if } 3 - \beta < k < \frac{5-\beta}{2}, \\ -2 & \text{if } k = \frac{5-\beta}{2}, \\ -\infty & \text{if } \frac{5-\beta}{2} < k < 2. \end{cases} \end{aligned}$$

Hence we obtain:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 3 - \beta < k < \frac{5-\beta}{2}, \\ -2 & \text{if } k = \frac{5-\beta}{2}, \\ -\infty & \text{if } \frac{5-\beta}{2} < k < 2. \end{cases}$$

Combining this result with Equation (53) gives the limit of the main term:

If $0 \leq \beta < 1$ then:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 0 \leq k < 3 - \beta, \\ -1 & \text{if } k = 3 - \beta, \\ -\infty & \text{if } 3 - \beta < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

² $1 + \frac{nx}{1-nx+x} \geq (1+x)^n \geq 1 + nx$ for $x \in (-1, 0]$ and $n \in \mathbb{N}$. Note that these inequalities might hold for more general values of x and n but we restricted ourselves to the case of interest for the proof.

If $\beta = 1$ then:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 0 \leq k < 2, \\ \exp(-2) - 1 & \text{if } k = 2, \\ -\infty & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

If $1 < \beta \leq 2$ then:

$$\lim_{n \rightarrow +\infty} -n^{k-3+\beta} \left(1 - (1 - 2n^{k-3})^n\right) = \begin{cases} 0 & \text{if } 0 \leq k < \frac{5-\beta}{2}, \\ -2 & \text{if } k = \frac{5-\beta}{2}, \\ -\infty & \text{if } \frac{5-\beta}{2} < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \quad (54)$$

Remaining term, with $l > 1$. We can rewrite the remaining term as follows:

$$\begin{aligned} -\sum_{l=2}^{\infty} \frac{n^{\beta}}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right)^l &= -\sum_{l=2}^{\infty} \frac{n^{\beta}}{2l} 2n^{kl-3l} \left(1 - (1 - 2n^{k-3})^n\right)^l \\ &= -\sum_{l=2}^{\infty} \frac{(1 - (1 - 2n^{k-3})^n)^l}{l} n^{(k-3)l+\beta} \end{aligned}$$

Notice that $\forall k \in \left[0, 3 - \frac{\log(2)}{\log(n)}\right)$, $\forall l \geq 2, \forall n \geq 2$ we have that $(1 - (1 - 2n^{k-3})^n)^l \in [0, 1]$. Following this we have:

$$\begin{aligned} l > \frac{\beta}{3-k} &\Rightarrow \lim_{n \rightarrow +\infty} -n^{(k-3)l+\beta} = 0 \Rightarrow \lim_{n \rightarrow +\infty} -\frac{(1 - (1 - 2n^{k-3})^n)^l}{l} n^{(k-3)l+\beta} = 0 \\ l \leq \frac{\beta}{3-k} &\Rightarrow \lim_{n \rightarrow +\infty} -n^{(k-3)l+\beta} \leq 0 \Rightarrow \lim_{n \rightarrow +\infty} -\frac{(1 - (1 - 2n^{k-3})^n)^l}{l} n^{(k-3)l+\beta} \leq 0. \end{aligned}$$

From this, noticing that $k < \frac{6-\beta}{2}$ implies $\frac{\beta}{3-k} < 2$, we obtain:

$$0 \leq k < \frac{6-\beta}{2} \Rightarrow \lim_{n \rightarrow +\infty} -\sum_{l=2}^{\infty} \frac{n^2}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right)^l = 0 \quad (55)$$

$$\frac{6-\beta}{2} \leq k < 3 - \frac{\log(2)}{\log(n)} \Rightarrow \lim_{n \rightarrow +\infty} -\sum_{l=2}^{\infty} \frac{n^2}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right)^l \leq 0 \quad (56)$$

Limit of Equation (50). We now combine the limits obtained in Equations (54), (55) and (56) to obtain the following limits:

If $0 \leq \beta < 1$ then $\frac{6-\beta}{2} > 3 - \beta$ and thus the remaining term does not change the limit of the main term which implies:

$$\begin{aligned} \lim_{n \rightarrow +\infty} -\frac{n^{\beta}}{2} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right) - \sum_{l=2}^{\infty} \frac{n^{\beta}}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right)^l &= \\ &= \begin{cases} 0 & \text{if } 0 \leq k < 3 - \beta, \\ -1 & \text{if } k = 3 - \beta, \\ -\infty & \text{if } 3 - \beta < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \end{aligned}$$

If $\beta = 1$ then $\frac{6-\beta}{2} > 2$ and thus the remaining term does not change the limit of the main term which implies:

$$\begin{aligned} \lim_{n \rightarrow +\infty} -\frac{n^{\beta}}{2} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right) - \sum_{l=2}^{\infty} \frac{n^{\beta}}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right)^l &= \\ &= \begin{cases} 0 & \text{if } 0 \leq k < 2, \\ \exp(-2) - 1 & \text{if } k = 2, \\ -\infty & \text{if } 2 < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases} \end{aligned}$$

If $1 < \beta \leq 1 + \frac{\log(n-1)}{\log(n)}$ then $\frac{6-\beta}{2} > \frac{5-\beta}{2}$ and thus the remaining term does not change the limit of the main term which implies:

$$\lim_{n \rightarrow +\infty} -\frac{n^{\beta}}{2} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right) - \sum_{l=2}^{\infty} \frac{n^{\beta}}{2l} \left(2n^{k-3} - 2n^{k-3} (1 - 2n^{k-3})^n\right)^l =$$

$$\begin{cases} 0 & \text{if } 0 \leq k < \frac{5-\beta}{2}, \\ -2 & \text{if } k = \frac{5-\beta}{2}, \\ -\infty & \text{if } \frac{5-\beta}{2} < k < 3 - \frac{\log(2)}{\log(n)}. \end{cases}$$

Plugging this result back into Equation (50) gives the lemma. \square

J Proof of Theorem 5

Theorem 5 (Exact bound on the probability that a strong classifier abstains). *In Theorem 4, further assume that each triplet classifier that does not abstain on at least one training example has a weight different from 0 (if for at least one example $(x_i, y_i) \in S$ we have that $h_c(x_i) \neq \vartheta$ then $\alpha_c \neq 0$). Then equality holds in Equation (38).*

Proof. The theorem follows directly by noticing that, in the proof of Theorem 4, if we further assume that each triplet classifier that does not abstain on at least one training example has a weight different from 0 (if for at least one example $(x_i, y_i) \in S$ we have that $h_c(x_i) \neq \vartheta$ then $\alpha_c \neq 0$) then we have equality in Equation (42):

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[H(x) = \vartheta \mid \sum_{c=1}^C V_c Y_{x,c} \neq 0 \right] = 0.$$

And thus we also have that:

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_S} [H(x) = \vartheta] = \mathbb{P}_{(x,y) \sim \mathcal{D}_S} \left[\sum_{c=1}^C V_c Y_{x,c} = 0 \right].$$

\square

Note that, in the very unlikely event that, in one of the iterations of TripletBoost, the selected triplet classifier has an error of exactly $1/2$, Theorem 5 might not hold for our method.