# Improving ANNs Performance on Unbalanced Data with an AUC-Based Learning Algorithm

Cristiano L. Castro and Antônio P. Braga

Federal University of Lavras
Department of Computer Science 37200-000, Lavras, MG, Brazil
`ccastro@dcc.ufla.br`,
`apbraga@ufmg.br`

**Abstract.** This paper investigates the use of the Area Under the ROC Curve (AUC) as an alternative criteria for model selection in classification problems with unbalanced datasets. A novel algorithm, named here as AUCMLP, which incorporates AUC optimization into the Multi-layer Perceptron (MLPs) learning process is presented. The basic principle of AUCMLP is the solution of an optimization problem that aims at ranking quality as well as the separability of class distributions with respect to the threshold decision. Preliminary results achieved on real data, point out that our approach is promising, and can lead to better decision surfaces, specially under more severe unbalance conditions.

**Keywords:** unbalanced datasets, classification, Area Under the ROC Curve, parameter estimation criteria.

## 1 Introduction

Global squared error functions are often used in error-correction learning since they yield simplification of the optimization problem, specially of those algorithms which are based on gradient descent. Many of current learning algorithms for Artificial Neural Networks (ANNs) have inherited this learning principle from Backpropagation [1]. Nevertheless, the use of a global error function may fail to represent properly the true error of unbalanced classification problems. In such kind of problems, the discrimination function tend to favor the majority class since the global error function assumes uniform losses for all training samples, regardless of the prior probability of the corresponding class [2]. Model performance on each separate class is often considered as a final criteria for model assessment, but it is not usually embodied in the adaptive learning procedures.

Performance assessment and model selection for unbalanced learning problems have been often accomplished with the aid of the ROC (*Receiver Operating Characteristic*) curve [3] which represents the relationship between the true positive rate (*TPrate*) and the false positive rate (*FPrate*) of a family of classifiers resulted from different output thresholds. A more robust criteria extracted from the ROC curve is the AUC (*Area Under the ROC Curve*) which is a global metric for all thresholds regardless of class prior probabilities. Because of that,

the AUC has been applied to ranking quality estimation [4] and also to highly unbalanced learning problems [5].

Despite being a more robust metric for unbalanced classification problems, AUC maximization is not guaranteed by global error minimization learning algorithms [6]. In order to guarantee AUC maximization, learning algorithms are expected to incorporate AUC optimization into the learning procedures, approach that has been adopted by some learning algorithms [7,8,9]. It has been also shown [6,4] that the *RankBoost* algorithm under specific conditions computes a function that is equivalent to the AUC.

Although the aforementioned algorithms have been proposed to maximize ranking in specific domains, such as Information Retrieval, their application in the context of unbalanced learning have not been investigated yet in the literature. Since the inherent properties of the AUC metric motivates its use for model selection in the presence of uneven data, it is natural to suppose that AUC optimization-based algorithms could represent an alternative to the well known sampling and cost-sensitive learning approaches, which have been commonly used to increase ANNs discrimination ability [10,11]. This work aims to investigate this hypotheses with a novel algorithm for Multi-layer Perceptrons (MLPs), named here as AUCMLP, which embodies AUC optimization in the learning process.

Our goal is to adopt AUC as a general cost function in order to improve MLPs' performance for representing classification functions, particularly those induced from unbalanced datasets. The main principles of AUCMLP are the solution of the optimization problem independently of the prior distributions yielded by the AUC, as well as its relationship with the quality of classification ranking. In contrast with a global error cost function, it may yield better performances in (highly) unbalanced datasets, specially those with class overlapping.

The paper is organized as follows: Section 2 describes the foundations of our AUC-based learning approach for MLPs. Section 3 presents the methodology of the empirical study conducted to evaluate the effectiveness of our approach. Also presented are the discussions on the results obtained. At last, the final conclusions are provided in Section 4.

## 2   Area under the ROC Curve

Let $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^{N}$ be a dataset with $N$ examples belonging to two classes $C_1$ and $C_2$, where $y_i \in \{+1, -1\}$ denotes the label (target output) of each input $\mathbf{x}_i \in \mathbb{R}^n$. The dataset $D$ consists of $N_1$ examples of the minority (or positive) class $C_1$ and $N_2$ examples of the majority (or negative) class $C_2$. The union of the two sets $D_1 = \{(\mathbf{x}_p, y_p)\}_{p=1}^{N_1}$ and $D_2 = \{(\mathbf{x}_q, y_q)\}_{q=1}^{N_2}$ corresponds to the dataset $D$ ($N = N_1 + N_2$ and $D = D_1 \cup D_2$).

The Area Under the ROC Curve (AUC) of a classifier $f$ can be expressed by the probability $P(f(\mathbf{X}^+) > f(\mathbf{X}^-))$, where $f(\mathbf{X}^+)$ is the random variable corresponding to the distribution of the classifier outputs for the positive examples and $f(\mathbf{X}^-)$ the one corresponding to the negative examples. A discrete