

Chapter 30

Learning From Imbalanced Data

Lincy Mathews

M. S. Ramaiah Institute of Technology, India

Seetha Hari

Vellore Institute of Technology, India

ABSTRACT

A very challenging issue in real-world data is that in many domains like medicine, finance, marketing, web, telecommunication, management, etc. the distribution of data among classes is inherently imbalanced. A widely accepted researched issue is that the traditional classifier algorithms assume a balanced distribution among the classes. Data imbalance is evident when the number of instances representing the class of concern is much lesser than other classes. Hence, the classifiers tend to bias towards the well-represented class. This leads to a higher misclassification rate among the lesser represented class. Hence, there is a need of efficient learners to classify imbalanced data. This chapter aims to address the need, challenges, existing methods, and evaluation metrics identified when learning from imbalanced data sets. Future research challenges and directions are highlighted.

INTRODUCTION

Pattern Identification on various domains have become one of the most researched fields. Accuracy of all traditional and standard classifiers is highly proportional to the completeness or quality of the training data. Completeness is bound by various parameters such as noise, highly representative samples of the real world population, availability of training data, dimensionality etc.

Another very pressing and domineering issue identified in real world data sets is that the data is well-dominated by typical occurring examples but with only a few rare or unusual occurrences. This distribution among classes make the real world data inherently imbalanced in many domains like medicine, finance, marketing, web, fault detection, anomaly detection etc.

This chapter aims to highlight the existence of imbalance in all real world data and the need to focus on the inherent characteristics present in imbalanced data that can degrade the performance of classifiers. It provides an overview of the existing effective methods and solutions implemented towards the significant problems of imbalanced data for improvement in the performance of standard classifiers. Efficient metrics for evaluating the performance of imbalanced learning models followed by future directions for research is been highlighted.

DOI: 10.4018/978-1-5225-7598-6.ch030

BACKGROUND

The field of data mining has identified learning from data that suffer from imbalance distribution as one of the top problems of today (Yang & Wu, 2006). However, a widely accepted issue is that the traditional learning algorithms assume a balanced distribution among the classes. It does not address nor recognize the presence of imbalance in the data. Data imbalance is evident when the number of instances representing the class of concern is much lesser than other classes.

To cite an example, the 1999 KDD Cup data set (UCI machine repository) is considered. The information collected by a simulated LAN environment consists of normal traffic with a relatively small number of intrusion attempts. The original data set consists of 23 classes, of which one of the classes belonged to normal traffic. When the data set was grouped down to a total of 2 classes, 'normal' and 'attack', the KDD data had 972,781 minority 'attack' class examples and 3,925,650 majority 'normal' class examples, which is approximately 80.14% majority examples. The training data thus will have only very few samples from the 'attack' class, due to which the classifier will be biased to the normal cases. This under representation of the interested class is evident in many applications such as intrusion detection, pollution detection, fault monitoring, biomedical, bioinformatics and remote sensing.

The under-represented class and well-represented class are known as the positive class (denoted by +1) and negative class (denoted by -1) respectively. As the class of interest indicates a positive case and is rarer by nature, it represents the minority data. The research community addresses the other well-defined classes as majority class. The ratio between the instances of majority versus minority is termed as imbalance ratio.

The skew distribution present in the training data, leads to the bias by most classifiers. Studies have however shown that the base classifiers perform well when presented with balanced data than with imbalanced data (Weiss & Provost, 2001). This justifies the need for learning models that can address the challenges posed by imbalanced data.

CHARACTERISTICS OF IMBALANCED DATA

The imbalance ratio between the majority and minority instances need not necessarily affect the performance of classifiers if the degree of imbalance is moderate (Chen & Wasikowski, 2008). The inherent characteristics within minority data however; can cause degrade in performance by the learning models. Two basic categorization of minority instances exist; Safe and unsafe minority instances. Safe minority instances are instances where the misclassification is minimal by the base learners. These instances exist much away from the borderline of majority instances. Unsafe minority instances are so called because the misclassifications occur highly with these kinds of minority instances.

The causes of unsafe instances in imbalanced data sets are contributed by four significant occurrences. They are the presence of small disjuncts, lack of density in the sample space, noisy data and data shift. Addressing these issues alone can sometimes bring a positive effect on the accuracy of the classifier without having to address the imbalance factor (Japkowicz N, 2003).

Small disjuncts exist when there is a small cluster of similar instances amidst cluster of majority or minority instances. However, in case of imbalanced datasets, the presence of sub concepts (disjuncts) in the majority class will be rare as they are represented well. The occurrence of small disjuncts is frequent

10 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/learning-from-imbalanced-data/214631

Related Content

Mobility and Multimodal User Interfaces

C. Pavlovski and S. Mitchell (2007). *Encyclopedia of Mobile Computing and Commerce* (pp. 644-650).

www.irma-international.org/chapter/mobility-multimodal-user-interfaces/17150

Location-Based Multimedia Content Delivery System for Monitoring Purposes

A. Sotiriou (2007). *Encyclopedia of Mobile Computing and Commerce* (pp. 381-386).

www.irma-international.org/chapter/location-based-multimedia-content-delivery/17105

Deep Reinforcement Learning for Mobile Video Offloading in Heterogeneous Cellular Networks

Nan Zhao, Chao Tian, Menglin Fan, Minghu Wu, Xiao He and Pengfei Fan (2018). *International Journal of Mobile Computing and Multimedia Communications* (pp. 34-57).

www.irma-international.org/article/deep-reinforcement-learning-for-mobile-video-offloading-in-heterogeneous-cellular-networks/214042

Mobile Virtual Communities of Commuters

Jalal Kawash, Christo El Morr, Hamza Taha and Wissam Charaf (2009). *Mobile Computing: Concepts, Methodologies, Tools, and Applications* (pp. 1771-1779).

www.irma-international.org/chapter/mobile-virtual-communities-commuters/26624

Mobile File-Sharing over P2P Networks

L. Yan (2007). *Encyclopedia of Mobile Computing and Commerce* (pp. 492-496).

www.irma-international.org/chapter/mobile-file-sharing-over-p2p/17123