# Chapter 28 Imbalanced Classification for Business Analytics

**Talayeh Razzaghi** University of Central Florida, USA

Andrea Otero University of Central Florida, USA

**Petros Xanthopoulos** University of Central Florida, USA

### INTRODUCTION

In pattern recognition, classification is a crucial task for automated data driven knowledge discovery. The objective of classification is to separate a set of data into classes or sub-categories and then to identify the classes that a new observation belongs to according to a training set of data. The mathematical model trained by a classification algorithm is termed *classifier*. When the class size of given examples is not equal for all classes, the classification problem is known as *imbalanced* (Japkowicz, 2000). For instance, in a cancer diagnostic problem the main objective is to identify individuals stricken with cancer and such events are relatively rare compared to normal cases. Imbalanced classification problems are also known as skewed class distribution problems or as small/rare class learning problems (He & Garcia, 2009; Lemnaru & Potolea, 2012; Sun, Wong, & Mohamed, 2009). In binary classification, the class with fewer examples is known as the *minority class* and the other class as the *majority class*. In many applications (e.g. fraud detection, computer intrusion detection, oil spill detection, defect product detection), detection of minority class examples is more important than the majority class. Therefore, there is a need for efficient classification algorithms to address such problems. A preferred classification algorithm is the one that yields higher identification rate on rare events especially for applications where their misclassification yields to high losses. For instance in automated credit card fraud detection, a fraud event misclassification might result in high monetary losses for the credit card vendor. On the other side misclassification of non-fraudulent events will worsen the customer satisfaction experience.

DOI: 10.4018/978-1-5225-1759-7.ch028

The emerging nature of *imbalanced* classification problems has led to the development of modified algorithms and new performance metrics. Standard performance measures such as classification accuracy are not appropriate when the data is imbalanced (N. V. Chawla, 2010; He & Garcia, 2009). In this chapter, we analyze the theoretical framework of imbalanced classification, the main algorithmic approaches proposed in the literature and some of the most prominent applications in business. These business applications include customer relationship management (CRM) (Kim, Chae, & Olson, 2013), fraud detection (Wei, Li, Cao, Ou, & Chen, 2012), and risk management (Groth & Muntermann, 2011).

## BACKGROUND

Advances in science and technology accelerate the accessibility of raw data and create new opportunities for knowledge discovery. *Imbalanced* problems can be found in a wide variety of applications, including security surveillance (Wu, Wu, Jiao, Wang, & Chang, 2003), medical diagnosis (Mena & JESUS, 2009; You, Zhao, Li, & Hu, 2011), bioinformatics (Al-Shahib, Breitling, & Gilbert, 2005), geomatics (Kubat, Holte, & Matwin, 1998), telecommunications (Tang, Krasser, Judge, & Zhang, 2006), risk management (Ezawa, Singh, & Norton, 1996), manufacturing (Adam et al., 2011), quality estimation (Lee, Song, Song, & Yoon, 2005), and power management (Hu, Zhu, & Ren, 2008). Imbalanced classification has been studied in a number of studies (N. V. Chawla, 2010; Guo, Yin, Dong, Yang, & Zhou, 2008; He & Garcia, 2009; Su, Mao, Zeng, Li, & Wang, 2009; Sun et al., 2009). Previous works on the classification of *imbalanced data* (N. V. Chawla, 2010; Kubat et al., 1998; Ngai, Hu, Wong, Chen, & Sun, 2011; Su et al., 2009; Sun et al., 2009) address that many standard classification algorithms achieve poor performance. Therefore, despite the existing amounts of literature there is room for improvement and future contribution.

### MAIN FOCUS

In this part, we present (1) the appropriate performance measures for *imbalanced data*; (2) imbalanced classification techniques and (3) the most popular business analytics applications.

### Performance Measures

Classification performance measures can be obtained, directly or indirectly, from the confusion matrix. For a classification problem with k classes, the confusion matrix is a square matrix  $C \in R^k$ , with each of its entries  $c_{ij}$ , denoting the percentage of the samples that belong to the class i and classified to the class j. For the special case of binary classification (positive and negative), the confusion matrix is as follows (Table 1).

Clearly, the confusion matrix of an ideal classifier is diagonal. In this matrix, diagonal elements represent accurately classified examples and the off-diagonal elements the misclassified data for each class. A typical performance measure for classification is the so-called accuracy, which is calculated as the correctly classified samples over the total number of training samples. Since the majority class dominates the behavior of this metric, it might not be an appropriate performance indicator for imbal-

9 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: <u>www.igi-global.com/chapter/imbalanced-classification-for-business-</u> analytics/173356

## **Related Content**

## Debating About, Against, and With ChatGPT: Redesigning Academic Debate Pedagogy for the World of Generative Artificial Intelligence

John Joseph Riefand Brian J. Schrader (2024). *The Role of Generative AI in the Communication Classroom (pp. 87-105).* 

www.irma-international.org/chapter/debating-about-against-and-with-chatgpt/339064

#### Leveraging OpenAI for Enhanced Multifactor Productivity in Chinese Businesses

Mohamad Zreik (2024). *Generative AI and Multifactor Productivity in Business (pp. 55-77).* www.irma-international.org/chapter/leveraging-openai-for-enhanced-multifactor-productivity-in-chinesebusinesses/345467

## An Improved TOPSIS Method Based on a New Distance Measure and Its Application to the House Selection Problem

You En Wangand Xiao Guo Chen (2023). *International Journal of Fuzzy System Applications (pp. 1-19)*. www.irma-international.org/article/an-improved-topsis-method-based-on-a-new-distance-measure-and-its-application-tothe-house-selection-problem/328529

#### Complex Systems Modeling by Cellular Automata

Jirí Krocand Peter M.A. Sloot (2009). *Encyclopedia of Artificial Intelligence (pp. 353-360).* www.irma-international.org/chapter/complex-systems-modeling-cellular-automata/10271

### Multiagent Based Selection of Tutor-Subject-Student Paradigm in an Intelligent Tutoring System

Kiran Mishraand R. B. Mishra (2012). Insights into Advancements in Intelligent Information Technologies: Discoveries (pp. 47-71).

www.irma-international.org/chapter/multiagent-based-selection-tutor-subject/64370