

# Chapter 4

## Mitigating Data Imbalance Issues in Medical Image Analysis

**Debapriya Banik**

*Jadavpur University, India*

**Debotosh Bhattacharjee**

*Jadavpur University, India*

### ABSTRACT

*Medical images mostly suffer from data imbalance problems, which make the disease classification task very difficult. The imbalanced distribution of the data in medical datasets happens when a proportion of a specific type of disease in a dataset appears in a small section of the entire dataset. So analyzing medical datasets with imbalanced data is a significant challenge for the machine learning and deep learning community. A standard classification learning algorithm might be biased towards the majority class and ignore the importance of the minority class (class of interest), which generally leads to the wrong diagnosis of the patients. So, the data imbalance problem in the medical image dataset is of utmost importance for the early prediction of disease, specifically cancer. This chapter attempts to explore different problems concerning data imbalance in medical diagnosis. The authors have discussed different rebalancing strategies that offer guidelines for choosing appropriate optimal procedures to train the samples by a classifier for an efficient medical diagnosis.*

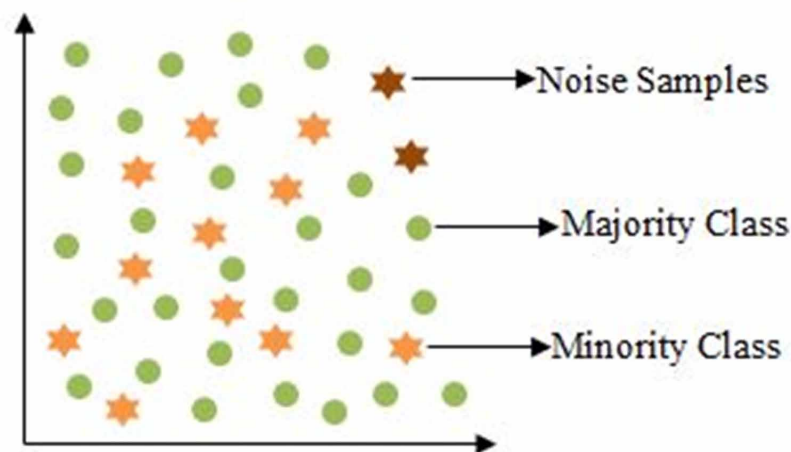
### INTRODUCTION

The data imbalance problem is prevalent in medical image analysis. The training of machine learning (ML) algorithm from an imbalanced medical data set is an inherently challenging task (Mena & Gonzalez, 2006). A classifier in ML's objective is to learn and predict the unseen output class of an unknown instance with good generalization capability. The mining of knowledge in a machine learning paradigm is accomplished by a set of  $\mathcal{D}$  input instances such as  $\eta_1, \eta_2, \eta_3, \dots, \eta_{\mathcal{D}}$  described by  $k$  features

DOI: 10.4018/978-1-7998-7371-6.ch004

$\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_k \in F$  whose intended output class labels  $\mathcal{O}_j \in C = \{c_1, c_2, \dots, c_m\}$ . A mapping function  $F^k \rightarrow C$ , implies the learning algorithm which is known as a classifier (Galar, Fernandez, Barrenechea, Bustince, & Herrera, 2011). This is a general idea for how a supervised learning algorithm performs its task. The imbalanced distribution of the data in medical image datasets happens when a specific disease type in a dataset appears in a small section of the entire dataset (C. Zhang, 2019). Hence, analyzing medical data posed severe challenges in the classification of a disease. A standard ML classifier will be skewed against the majority class and underestimate the importance of the minority class because the minority class has a lesser number of instances compared to the majority class. However, the minority class is generally referred to as the class of interest (Napierala & Stefanowski, 2016) in medical image analysis. So, the minority class is of utmost importance for the early prediction of disease. This problem influences all supervised classification algorithms. A well-balanced medical image dataset is very crucial for designing a reliable and standard prediction model. Typically, real-world medical data, specifically cancer data, usually suffer from data imbalance, leading to the degradation of ML algorithms' generalization. These eventually degrade the efficiency and accuracy of the computer-aided early prediction of cancer. The biasness of the medical data in healthcare domain due to individual diversity can cause missclassification which may affect early diagnosis of cancer and disease risk prediction (Zhao, Wong, & Tsui, 2018). However, the imbalanced class problem is generally ignored in Conventional Learning (CL) algorithms. Those algorithms give the same priority to both classes: the majority class and the minority class. However, when the majority class and the minority class are highly imbalanced, it is very challenging to build a good classifier using CL algorithms (Krawczyk, 2016). It is a significant concern in most medical datasets where patients at high-risk tend to be in the minority class, and so the cost in miss-classification of the minority classes is higher than that of the majority class. In Figure 1 a graphical representation of the distribution of majority class and the minority class is shown. The noisy data is a small part of the minority class, which significantly impacts the performance of the classifier (López, Fernández, García, Palade, & Herrera, 2013).

*Figure 1. Pictorial representation of a class imbalanced dataset*



22 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/mitigating-data-imbalance-issues-in-medical-image-analysis/280911](http://www.igi-global.com/chapter/mitigating-data-imbalance-issues-in-medical-image-analysis/280911)

## Related Content

---

### Internet of Things in the 5G Ecosystem and Beyond 5G Networks

Swati S. Roy, Shatarupa Dashand Bharat Jyoti Ranjan Sahu (2023). *Handbook of Research on Data Science and Cybersecurity Innovations in Industry 4.0 Technologies* (pp. 476-504).

[www.irma-international.org/chapter/internet-of-things-in-the-5g-ecosystem-and-beyond-5g-networks/331026](http://www.irma-international.org/chapter/internet-of-things-in-the-5g-ecosystem-and-beyond-5g-networks/331026)

### A Multi-Stage Framework for Classification of Unconstrained Image Data from Mobile Phones

Shashank Mujumdar, Dror Porat, Nithya Rajamaniand L.V. Subramaniam (2014). *International Journal of Multimedia Data Engineering and Management* (pp. 22-35).

[www.irma-international.org/article/a-multi-stage-framework-for-classification-of-unconstrained-image-data-from-mobile-phones/120124](http://www.irma-international.org/article/a-multi-stage-framework-for-classification-of-unconstrained-image-data-from-mobile-phones/120124)

### Applying Bibliometrics to Examine Research Output and Highlight Collaboration

Nandita S. Mani, Michelle A. Cawley, Adam Doddand Barrie E. Hayes (2022). *Handbook of Research on Academic Libraries as Partners in Data Science Ecosystems* (pp. 75-101).

[www.irma-international.org/chapter/applying-bibliometrics-to-examine-research-output-and-highlight-collaboration/302748](http://www.irma-international.org/chapter/applying-bibliometrics-to-examine-research-output-and-highlight-collaboration/302748)

### Diagnostic Device for Sustainable Medical Care Using Hyperspectral Imaging

Vinuja G., Saravanan V., Maharajan K., Jayasudha V., Ramya R.and Jothi Arunachalam S. (2024). *Emerging Advancements in AI and Big Data Technologies in Business and Society* (pp. 129-144).

[www.irma-international.org/chapter/diagnostic-device--for-sustainable-medical-care-using-hyperspectral-imaging/351262](http://www.irma-international.org/chapter/diagnostic-device--for-sustainable-medical-care-using-hyperspectral-imaging/351262)

### Brand Privacy Policy and Brand Trust Reference in Online Business Communication

Tugba Orten Tugruland Tugberk Kara (2023). *Enhancing Business Communications and Collaboration Through Data Science Applications* (pp. 157-177).

[www.irma-international.org/chapter/brand-privacy-policy-and-brand-trust-reference-in-online-business-communication/320755](http://www.irma-international.org/chapter/brand-privacy-policy-and-brand-trust-reference-in-online-business-communication/320755)