# UNIVERSITI TEKNOLOGI MARA

# AN ENHANCED SYNTHETIC OVERSAMPLING FRAMEWORK WITH SELF-SUPERVISED CONTRASTIVE LEARNING FOR MULTI-CLASS IMAGE IMBALANCE

## GAO XIAOLING

Thesis submitted in partial fulfillment
of the requirements for the degree of
**Doctor of Philosophy**
**(Computer Science)**

**Faculty of Computer and Mathematical Sciences**

**October 2025**

# ABSTRACT

Class imbalance significantly affects the performance of machine learning and deep learning classifiers, especially in image recognition tasks where certain classes are underrepresented. Traditional synthetic oversampling methods, while helpful, often fail to address the complexities of real-world data with uneven class distributions. The first contribution of this study is the creation of artificially imbalanced datasets from CIFAR10 and SVHN datasets, designed to systematically evaluate classifier performance under varying degrees of class disparity. The second contribution is the introduction of the Clustering and Nearest Centroid Neighbour-based Synthetic Minority Oversampling (CLNCN-SMOTE) algorithm to resolve multi-class imbalance. The algorithm is an enhancement of traditional K-means SMOTE that incorporates a nearest centroid neighbour strategy. This method effectively generates more representative synthetic samples for the minority class, thereby reducing noise and mitigating overfitting issues. Building on this, the third contribution is the development of an enhanced framework that is different from traditional methods, as it integrates oversampling with self-supervised contrastive learning and attention mechanism to tackle multi-class imbalance. The proposed framework integrated the Convolutional Block Attention Module (CBAM) within the Momentum Contrast (MoCo) framework's encoder. The enhanced encoder was initially pre-trained using unlabelled images to refine its capability. Subsequently, the encoder was utilised to extract features from the training dataset and augment them using the CLNCN-SMOTE approach in the feature space. Finally, a Multi-layer Perceptron (MLP) classifier assessed the effectiveness of the entire framework. The proposed framework leverages contrastive learning to distinguish more effectively between features of different categories and employs an attention mechanism to optimise feature selection. This improves the classifier's ability to accurately recognise features of the minority class without degrading the majority class's performance. Experiments on several benchmark datasets, including CIFAR10, SVHN, Caltech-101, ImageNet-LT, and iNaturalist 2018, demonstrate significant improvements. The proposed framework achieves a macro-averaged F1-score (FM) of 73.17% and a macro-averaged geometric mean (GM) of 83.04% on CIFAR10, with FM reaching 79.61% and GM 88.28% on various SVHN datasets. Notably, it surpasses the MoCo-v2 framework by 5% in both FM and GM on the imbalanced SVHN and CIFAR10 datasets, achieving Top-1 accuracy of 68.67% on ImageNet-LT and 73.5% on iNaturalist 2018. In conclusion, the results underscore the framework's effectiveness in tackling multi-class image imbalance and highlight its significant practical applications in fields such as medical imaging, surveillance, and anomaly detection. Future studies could explore the integration of prototype-based contrastive learning methods for further enhancement.

# ACKNOWLEDGEMENT

# TABLE OF CONTENTS

# CHAPTER 1

# INTRODUCTION

## 1.1    Research Background

Class imbalance poses a significant challenge in machine learning and deep learning, markedly influencing the performance of classifiers. These classifiers tend to favour classes with more instances, often neglecting those with fewer examples, which may represent rare but significant events. This can lead to poor classification outcomes, particularly in scenarios with limited data, where the need for extensive and costly expert labelling hinders training (Cocos et al., 2017; Hou et al., 2023). Despite its advanced capabilities, deep learning is limited by imbalanced data distributions (Chen et al., 2024). This hampers its full potential (Sampath et al., 2021; Ghosh et al., 2024).

Class imbalance in datasets occurs as two main types: inter-class and intra-class imbalance. Inter-class imbalance occurs when the minority class has fewer instances than the majority class. This often leads classifiers to misidentify the minority class as rare or even treat it as outliers or noise, which results in misclassification (Ali et al., 2015; Rezvani & Wang, 2023). Conversely, intra-class imbalance refers to attribute biases within a class, such as variation in dog breeds, colours, and poses within a dog class. While both types of imbalance affect classifier performance, this research focuses primarily on inter-class imbalance, a common challenge in many classification tasks. In particular, image recognition in natural environments presents persistent difficulties because of such imbalances, especially in long-tailed distributions that exemplify severe multi-class imbalances (Yadav & Bhole, 2020). This research aims to address multi-class imbalances, which include long-tailed imbalance.

Research addressing data imbalances has been robust, categorising them into data-level, algorithm-level, and hybrid approaches. The data-level method, which adjusts class distributions by augmenting minority classes or reducing majority classes, remains popular because of its flexibility and broad applicability (Huda et al., 2018; Zhang et al., 2023). In contrast, algorithm-level strategies require specific classifiers and offer less adaptability (Havaei et al., 2017; Ochal et al., 2023). Hybrid methods blend data with algorithm-level tactics, often utilising classifier ensembles for enhanced effectiveness. This research concentrates on data-level strategies, which have shown