# One-class classification – From theory to practice: A case-study in radioactive threat detection ☆

Colin Bellinger [a,*], Shiven Sharma [b], Nathalie Japkowicz [c]

[a] *Computing Science, University of Alberta, Edmonton, Canada*
[b] *Fluent Solutions Inc., Ottawa, Canada*
[c] *American University, Washington, D.C., USA*

A B S T R A C T

Over the years, the acceptance of machine learning as a valuable tool in the real-world has caused much interest in the research community; this is particularly the case as the field of Big Data is coming into prominence. However, real-world data comes with a myriad of challenges, amongst the most prominent of which is the fact that it can exhibit a high level of imbalance. This can come in the form of both within- and between-class imbalance. While a significant amount of research has been devoted to the impact of within-class imbalance over binary classifiers, very little attention has been given to their impact on one-class classifiers, which are typically used in situations of extreme between-class imbalance. During our collaboration with Health Canada into the identification of anomalous gamma-ray spectra, the issue of within-class imbalance in a one-class classification setting was highly significant. In this setting, the imbalance comes from the fact that the background data that we wish to model is composed of two concepts (background no-rain and rain); the rain sub-concept is rare and corresponds to spectra affected by the presence of water in the environment. In this article, we present our work into developing systems for detecting anomalous gamma-rays that are able to handle both the inherent between-class and within-class imbalance present in the domain. We test and validate our system over data provided to us by Health Canada from three sites across Canada. Our results indicated that oversampling the sub-concept improves the performance of the baseline classifiers and multiple classifier system when measured by the geometric mean of the per-class accuracy.

© 2018 Elsevier Ltd. All rights reserved.

## 1. Introduction

This paper summarizes and advances our previous research into the one-class classification of gamma-ray spectra. In particular, we explore the difficulties caused by the existence of potentially rare regions in the target class; these regions are more formally known as sub-concepts (Weiss, 2003). They make learning a good model difficult because they may be sparse and separated by areas of low-density or the other class. We consider the benefits of both multiple classifier systems and preprocessoing the data to balance the sub-concepts priors in order to mitigate the negative effects of the less frequent sub-concepts.

The application of this research involved a collaboration with physicists at the Radiation Protection Bureau of Health Canada. The overarching objective is to aid government monitors in predicting the occurrence of a pending radioactive event. Depending on the setting, the specific event may involve a malfunction at a nuclear facility or the nefarious use of nuclear material. Thus, the classification objective is to design and develop a system capable of accurately identifying rare spectra signifying potential dangers.

Our initial attempts with standard one-class classification did not satisfy our objective of high accuracy on the target and outlier classes. By examining the misclassifications using principle component analysis, we were able to identify that the errors were occurring in a common area of the data space that forms a sub-concept associated with heavy rainfall.

To address these false positives, our previous work proposed a multiple classifier system-based approach with a cascade architecture (Sharma, Bellinger, Japkowicz, Berg, & Ungar, 2012). We referred to this as a two-tiered system; the key objective was to simplify the target distribution by separating it into its rain and no-rain sub-concepts. We have subsequently been provided with new

---

gamma ray spectra datasets and have identified a limitation in the multiple classifier system approach related to the degree of imbalance in the rain sub-concept. This has inspired us to consider the problem from the perspective of within-class imbalance (Japkowicz & Stephen, 2002), and enabled us to utilize the wealth of research available in the binary class imbalance literature to shape our solution.

The main contributions of this work are to:

- Summarize and extend our previous work on two new radiation monitoring datasets that come from unique domains;
- Explore the problem of gamma-ray spectral classification from the perspective of within-class imbalance, and demonstrate how it can impact the performance of one-class classifiers;
- Propose a solutions to the within-class imbalance problem by extending sampling methods from the binary classification literature to balance the sub-concepts; and
- Show that synthetically oversampling the imbalanced subconcept improves performance more than the multi-classifier approach previously proposed.

## 2. Related work

Class imbalance has been studied in both extreme cases and moderate cases. The former refers to the situations where no, or almost no, training instances from a concept of interest are available. Alternatively, in the moderate cases, the relative training balance is skewed enough to negatively impact performance. Nonetheless, binary methods with some pre-processing or weighting may still be applied. Both of these veins of research have influenced our solution to the gamma-ray spectra classification problem. As such, we discuss the relevant work below. This commences with the moderate case and proceeds to the extreme setting.

### 2.1. Binary class imbalance

Class imbalance appears in a wide variety of important and challenging binary classification tasks. Some prominent examples of imbalanced classification problems are: oil spill classification, gene function annotation, and medical and text classification (Akbani, Kwek, & Japkowicz, 2004; Blondel, Seki, & Uehara, 2011; Kubat, Holte, & Matwin, 1998; Nguwi & Cho, 2009). Applications in such areas have demonstrated that it can pose a major challenge for classification systems (He & Garcia, 2009; Japkowicz & Stephen, 2002). In the literature, two forms of imbalance have been identified, namely between-class and within-class imbalance. It has been found that, in many cases, data complexity, such as class overlap, noise and sub-concepts, contribute much of the difficulty to imbalanced problems (Batista, Prati, & Monard, 2004; Denil & Trappenberg, 2010; García, Sánchez, & Mollineda, 2007; Japkowicz, 2001; Prati, Batista, & Monard, 2004).

The issue of sub-concepts is highly relevant in this work. As Stefanowski (2016) highlight, apart from imbalance, the performance of classifiers can be impacted by the presence of sub-concepts (i.e., small dusjuncts). Research in binary classification has established that sub-concepts, particularly rare sub-concepts, can lead to a degradation in classification performance (Japkowicz, 2003). Our work here shows that this is also an issue in one-class classification. For binary classification, Jo and Japkowicz (2004) propose a method for dealing with both within and between class imbalance by clustering and random oversampling. Napierała, Stefanowski, and Wilk (2010), examined further means of managing the affect of noise, overlap and sub-concepts with data cleaning and oversampling based on the local characteristics of the data. They empower rare, but relevant, sub-concepts, whilst removing noise and borderline instances.

The Synthetic Minority Oversampling TEchniques (SMOTE) is the standard method applied for synthetic oversampling in the literature (Chawla, Lazarevic, Hall, & Bowyer, 2003). SMOTE generates new instances of the minority class by interpolating them at random points on the edges connecting nearest neighbors in the minority class. This results in samples created within the convex-hull formed by the minority class. The manifold-based synthetic oversampling sampling method was recently proposed (Bellinger, Drummond, & Japkowicz, 2017); it's approach of modeling data as low-dimensional manifolds is particularly helpful on sparse, high-dimensional domains, as we have here. Samples are generated from the induced manifold, which leads to a better representation of the probability density.

### 2.2. One class classification

The goal in one-class classification is to induce a binary class predictor, $f: x \rightarrow y$, that learns a functional mapping from the feature vector $x$ to the corresponding label $y$, where $y \in \{0, 1\}$. Learning takes place on a given set of training examples $X$ sampled from the target class $y = 0$. This is a challenging learning problem because a classifier must be induced, and the model must be selected without seeing examples of the other class $y = 1$.

One-class classifiers typically induce their decision boundaries using one of three modeling paradigms: density-based, recognition-based and boundary-based. Each of these paradigms have been widely applied. Density-based methods have been applied to one-class classification problems, such as diseases and infection detection, and to monitor content (Cohen, Sax, Geissbuhler et al., 2008; Tarassenko, Hayton, Cerneaz, & Brady, 1995; Zuo, Wu, Hu, & Xu, 2008). Reconstruction-base classifiers have been applied to predict failures in helicopter gear boxes, classify documents and to detect nuclear tests (Japkowicz, 1999; Manevitz & Yousef, 2001). One-class support vector machines (SVM) and support vector data description are the standard boundary based methods for one-class classification. These have had a significant amount of success in applications of text, image retrieval and human health (Chen, Zhou, & Huang, 2001; Erfani, Rajasegarar, Karunasekera, & Leckie, 2016; Manevitz & Yousef, 2001; Zhang, Wang, Xu, & Liu, 2006).

Whilst much research has been undertaken to understand the data properties that impact the performance of binary classifiers induced over imbalanced datasets, the relationship between one-class classifiers and data properties has not been as thoroughly considered. A recent study on the impact of data complexity on one-class classifiers is conducted in Bellinger, Sharma, Zaiane, and Japkowicz (2017); the authors highlight that multi-modality resulting from the presence of sub-concepts, as well as class overlap, can cause a significant degradation in the performance of both binary and one-class classifiers as imbalance increases.

In order to make the one-class classifiers more robust to the presence of sub-concepts, Sharma (2016) and Sharma, Bellinger, and Nathalie (2012) demonstrated that by isolating and learning over each sub-concept, better one-class classifier systems can be produced. Isolation is performed by clustering as in Jo and Japkowicz (2004), and a separate one-class classification model is induced for each cluster. This is the motivation for our multi-classifier system, and corresponds to a general approach for building ensembles of one-class classifiers (Jackowski, Krawczyk, & Woniak, 2014; Krawczyk, 2015; Lipka, Stein, & Anderka, 2012); by grouping multiple classifiers into as single system, their collective strengths can be harnessed.

Finally, we note that whilst each of these multiple classifier methods can help to deal with sub-concepts, they implicitly assume that the sub-concepts are well represented. Our results suggest that the relative frequency of the target class sub-concepts have implications on the performance of mutli-classifier systems.