

# Discrimination of Signal-Background Events with Supervised and Weakly Supervised Learning in the Search for New Bosons Decaying to $Z + \gamma$ Final State

---

Nkateko Baloyi

*Supervisor(s):*

Bruce Mellado

Xifeng Ruan



A research report submitted in accordance with the requirements for the degree of  
Master of Science in the field of e-Science

in the

School of Computer Science and Applied Mathematics

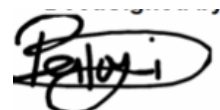
University of the Witwatersrand, Johannesburg

14 July 2020

# Declaration

I, Nkateko Baloyi, declare that the content in this research report is my own original work and It has not been submitted to any other university for a degree.

Part of this research have contributed and submitted to the High Energy Particle Physics (HEPP) 2019 and conference proceedings and the South African Institute of High Energy Physics (SAIP) 2019 conference proceedings to be published under the topics, **Application of Boosted Decision trees in the Search for a New Boson decaying to  $Z_\gamma$  Final State** and **Discrimination of Signal-Background Events with Supervised and Semi-supervised Learning in the Search for New Bosons Decaying to  $Z + \gamma$  Final State**, respectively.



Nkateko Baloyi

14 July 2020

## *Abstract*

The theory of Standard Model (SM) has successfully driven experiments and predictions since it was developed in 1975 and it has never been contradicted by experimental results. In 2012, the SM led to the discovery of Higgs boson ( $h$ ) at the Large Hadron Collider (LHC) which completed the particle spectrum of the SM. The discovery of  $h$  inspired experimental studies to further understand the  $h$  scalar properties, opening the search for Beyond the Standard Model physics (BSM). BSM physics searches for new particles that can help understand and answer phenomena that cannot be explained by the SM. The LHC collides protons at high luminosity and high energy trying to recreate particles that occurred moments after the big bang and the BSM particles. The data produced during the collision requires advanced techniques that can search for relevant information in the data for signal ( $S$ ) events to be identified. The production of the  $S$  events comes with a huge amount of background ( $B$ ) production to which the  $S$  events cannot be easily identified. advanced machine learning (ML) and statistical techniques can be used to isolate the  $S$  events from the  $B$  events. ML is a subset of artificial intelligence that provides systems the ability to automatically learn and improve from experience. This research focuses on the application of boosted decision trees (BDT) and deep neural networks (DNN) on the Monte Carlo simulated data. Supervised learning and weakly-supervised learning (WSL) approaches are implemented to discriminate the  $S$  from the  $B$  events. The supervised learning is used as a benchmark to measure the performance of the WSL approach. Pre-selection cuts are applied on the data and four different models are applied to classify the  $S$  and  $B$  events for both BDT and DNN using the supervised learning and WSL approach. The WSL approach use two samples, one sample with only  $B$  events and the other sample mixed with  $S$  and  $B$  events to train the model. The DNN model is trained on the samples and applied to classify the  $S$  and  $B$  events on the same test data used in the supervised learning approach. The performance of the WSL models are compared to the supervised learning models performance. The results show a strong bias in the WSL approach.

# Acknowledgements

I sincerely and gratefully acknowledge my supervisor Prof. Bruce Mellado and my co-supervisor Dr. Xifeng Ruan for their guidance, patience, constant support throughout this research. I am tremendously fortunate to have the opportunity to be mentored and learn from them with the depth of knowledge they have in the field. I am also truly grateful for the opportunity to travel to CERN to be associated and learn from experts in the field of Machine Learning from different parts of the world, which was made possible by my supervisor.

In addition, I gratefully acknowledge the e-Science programme for this amazing opportunity to enroll in the programme. I am truly grateful for the funding I received from DST-CSIR through the programme, and I thank SA-CERN for making it possible for me to travel to CERN.

Lastly, I would like to express my heartiest gratitude to my parents, Ntombhi Mazini and David Baloyi, for their support and prayers throughout my academic journey. Special thanks to my friends Risuna Nkolele, Zwidofhelangani Ndamulelo, Orapeleleng Mogawana for their support, and my amazing team, Joshua Choma, Theodore Gaelejwe, Ntsoko Phuti and Kehinde Tomiwa for the role they played in me executing this research.

# Contents

<b>Declaration</b>	<b>i</b>
<b>Abstract</b>	<b>ii</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>List of Figures</b>	<b>vi</b>
<b>List of Tables</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	2
1.2 Problem Statement . . . . .	3
1.3 Research Question . . . . .	4
1.4 Research Aims . . . . .	4
1.5 Objectives . . . . .	4
1.6 Overview . . . . .	4
1.7 Outline . . . . .	5
<b>2 Literature Review</b>	<b>6</b>
2.1 Introduction . . . . .	6
2.1.1 Supervised Learning . . . . .	6
2.1.2 Unsupervised Learning . . . . .	7
2.1.3 Semi-supervised Learning . . . . .	7
2.1.4 Weakly Supervised Learning . . . . .	8
2.1.5 Boosted Decision Trees (BDT) . . . . .	9
The Mathematics of Decision Trees and Gradient Boosting . .	10
2.1.6 Deep Neural Networks . . . . .	13
The Mathematics of Neural Networks . . . . .	13

2.1.7	Activation Functions . . . . .	16
2.1.8	Classifier Performance . . . . .	17
2.2	Related Studies . . . . .	19
<b>3</b>	<b>Research Methodology</b>	<b>22</b>
3.1	Data . . . . .	22
3.2	Methods . . . . .	24
3.2.1	Supervised Learning . . . . .	24
3.2.2	Weakly Supervised Learning . . . . .	24
3.2.3	Application . . . . .	25
Application of Boosted Decision Trees . . . . .	25	
Application of Deep Neural Network . . . . .	26	
3.2.4	Performance Metrics . . . . .	27
3.3	Summary . . . . .	27
<b>4</b>	<b>Results and Discussion</b>	<b>28</b>
4.1	Boosted Decision Trees Results . . . . .	28
4.1.1	Supervised Learning . . . . .	28
4.1.2	Weakly Supervised Learning . . . . .	31
4.2	Deep Neural Networks Results . . . . .	34
4.2.1	Supervised Learning . . . . .	34
4.2.2	Weakly Supervised Learning . . . . .	36
4.3	Discussion . . . . .	39
<b>5</b>	<b>Conclusions and Future Work</b>	<b>41</b>
5.1	Conclusions . . . . .	41
5.2	Future Work . . . . .	41
<b>A</b>		<b>42</b>
A.1	1D Distributions . . . . .	42
	<b>Bibliography</b>	<b>44</b>

# List of Figures

2.1	Illustration of a Decision Tree [34] . . . . .	11
2.2	Illustration of a shallow Neural Network (Left) [41] and Deep Neural Network (Right) [42] . . . . .	14
2.3	Plot of the ReLu Activation function (Left) and the Leaky ReLu activation function (Right) . . . . .	17
2.4	Plot of the Sigmoid function . . . . .	17
4.1	ROC curves illustrating the performance of the the BDT supervised learning model. . . . .	29
4.2	$S$ and $B$ Output distributions . . . . .	30
4.3	ROC curves illustrating the performance of the the BDT WSL model. . . . .	32
4.4	$S$ and $B$ Output distributions . . . . .	33
4.5	ROC curves illustrating the performance of the DNN supervised learning model. . . . .	35
4.6	$S$ and $B$ Output distributions . . . . .	36
4.7	ROC curves illustrating the performance of the DNN WSL model. . . . .	37
4.8	$S$ and $B$ Output distributions . . . . .	38
A.1	1 Dimensional distributions of the input variables . . . . .	43

# List of Tables

2.1	Confusion Matrix . . . . .	18
3.1	Input variables description . . . . .	23
3.2	BDT Hyper-parameters . . . . .	26
3.3	DNN Hyper-parameters . . . . .	26
4.1	BDT supervised learning Table of Results . . . . .	29
4.2	BDT WSL Table of Results . . . . .	31
4.3	DNN supervised learning Table of Results . . . . .	34
4.4	DNN WSL Table of Results . . . . .	37



# Chapter 1

## Introduction

The Large Hadron Collider (LHC) is the largest and most powerful particle accelerator in the world. It was built between 1998 and 2008 by the European Organization for Nuclear Research for the purpose of testing predictions of different theories in high energy physics [1, 2]. The LHC runs under the Swiss-French border near Geneva in a circular tunnel of 27 km circumference. The accelerator consists of two proton beams that travel in opposite directions at near the speed of light with very high energy before colliding [2, 3]. In 2010, the LHC was successfully commissioned for proton-proton ( $pp$ ) collision with 7 tera electron Volts ( $TeV$ ) centre of mass energy [4].

The phenomena that physicists are looking for in the collisions have a very low probability of occurrence. High luminosity is required to produce more data events in order to increase the occurrence of the phenomena. Luminosity gives the measure of the number of potential collisions that occur in a given period of time, an important indicator of the performance of the accelerator [5]. An integral of the luminosity over time known as the integrated luminosity gives a measure of the collected data size measured in inverse femtobarn ( $fb^{-1}$ ). The value of 1  $fb^{-1}$  equates to 100 million collisions [4, 5].

During the  $pp$  collision in the LHC, the protons can either collide violently and produce any new particles or they can glance off each other. The probability of production of any of the particles is different, each particle have a different cross-section. The rare particles such as the Higgs bosons have a very small cross-section [6]. The large luminosity accumulated by the LHC opens the possibility of a new range of Beyond the Standard Model (BSM) physics discoveries. The data collected

in the LHC during the collisions is complex, high dimensional, highly imbalanced with a strong overlap of the classes, which requires more advanced techniques with the potential to extract information produced in the  $pp$  collision. The information that can be obtained from the collisions includes, but not limited to the half-lives of particles, the velocity, decay rate of the particles, transverse momentum. Statistical tools and advanced techniques, such as machine learning (ML) can be used to improve data processing, classification and regression tasks in particle physics data and directly improves the probability of discovering new particles. ML is developed with the concept of allowing computers to learn by themselves and improve efficiency with experience and is applied in various professional fields and industries to improve the processing and the quality of results by finding patterns in data that lead to accurate decision making [7, 8, 9].

Various ML algorithms have been applied in binary classification tasks and have achieved accuracy above the random guess threshold. Boosting techniques and deep learning techniques have gained popularity in classification and regression tasks due to their ability to improve accuracy when applied on weak learning algorithms and trained on huge amount of data, respectively. Boosted decision trees and deep neural networks have been applied widely in different professions, including high energy physics data, and they have proven to be superior in terms of accuracy [10, 11].

This study aims to highlights machine learning approaches that can be used to discriminate signal and background events on simulated data and help develop a model for real data. We develop two machine learning models following the supervised learning and the weakly supervised learning approach. The background, research aims and objectives of this study are described in this chapter.

## 1.1 Background

The standard model (SM) of physics is a theory that describes three of the four fundamental forces, the electromagnetic, weak, and strong interactions, and how these forces and elementary particles that makes up matter as we know it, relate to each other. The SM has successfully driven experiments and predictions since it was introduced in 1975 [12, 13]. This theory is a well-tested method and it has never

been contradicted by experimental results. In 2012, the theory led to the discovery of the Higgs boson ( $h$ ) [14], which completed the particle spectrum of the SM. The particle  $h$  is a very important particle in the SM as it signals the existence of the Higgs field, an invisible energy field that is accompanied by the  $h$ , and generates mass through its interaction with other particles. The discovery of  $h$  inspired many theoretical and experimental studies conducted to understand the scalar properties of this particle and how new physics can be identified.

Although the SM has successfully explained experiments and was validated by the discovery of  $h$ , it still fails to explain some phenomena [15]. The discovery of  $h$  ended the era of theories driving the experimental results. This means that experiments must be conducted in search of new physics Beyond the SM (BSM), to help understand and explain some phenomena that cannot be explained by the SM. The LHC data has been made public to learn more about the data and analysis studies have been done on this data [16]. The studies conducted show that there are excesses observed in the LHC data release, which suggests the possibility of new physics [17, 18]. These excesses have been observed in both the ATLAS and CMS experiments and a number of final states resulting in these excesses has been identified. The Wits Institute for collider Particle Physics and their collaborators proposed a hypothesis called the Madala hypothesis, an extension of the SM, which searches for BSM physics.

## 1.2 Problem Statement

BSM physics requires particles to be accelerated at extremely high energy and high luminosity to increase the probability of the BSM events to be created. During the collision, a large amount of background is created with a small amount of signal. For new particles to be discovered, the signal and background events must be isolated. This requires statistical tests and advanced techniques that can extract the small signal within the large background. Machine learning techniques can be developed to discriminate the signal and background events, and improve the purity of signal in the search for new physics.

### 1.3 Research Question

From highly imbalanced dataset, can we identify and discriminate signal from the background events with supervised learning, to help improve the purity of the signal in order to perform the search for BSM physics?

When working with mixed unlabeled dataset, can we develop a model from data using the weakly supervised learning approach that can isolate signal from background?

### 1.4 Research Aims

This study aims to develop machine learning models that can discriminate signal and background events in the search for new bosons decaying to  $Z + \gamma$  final state.

### 1.5 Objectives

The primary goal is to develop a supervised learning model trained on labeled data to distinguish signal and background events. The secondary goal is to develop a weakly supervised learning model that can be employed on unlabeled data by assigning weak-labels on the unlabeled data. The weakly supervised approach implemented on weak-labeled data is used to distinguish the same test data used in the supervised learning approach to compare the output in order to help develop a model that can be implemented on real data.

### 1.6 Overview

Studies show that excesses have been observed in the LHC Run1 and Run2 data release, which suggest that BSM physics is possible. This inspired the search for new BSM physics that can explain and help understand life moments after the big bang. The search for BSM physics requires particles to be accelerated at extreme high energy and luminosity to increase the BSM cross-section. This ends up with large background and small signal events produced. In order to search for new

physics, the signal must be isolated from the background to check for unknown signal. With the vast amount of data generated by the LHC during collisions, advanced machine learning and statistical techniques are the key to isolating signal from background to improve the quality of signal used in detecting possible BSM physics. In this study, we aim to develop supervised learning and weakly supervised learning models using simulated Monte Carlo data, that can be implemented on real data.

## 1.7 Outline

The remaining chapters of this research report are designed as follows. Chapter 2 describes the concepts of machine learning and the performance metrics. Boosted decision trees and deep neural networks are described in detail. The related study is presented in the form of literature review. Chapter 3 describes the data and outlines the research methodology. The data preparation, application of boosted decision trees and deep neural networks are described. The results and discussion of the boosted decision trees and deep neural networks findings are presented in Chapter 4. The conclusions and possible future research are presented in Chapter 5.

## Chapter 2

# Literature Review

### 2.1 Introduction

Machine Learning (ML) is the umbrella term for a set of statistical and computational techniques for interpreting data which evolves from computational learning and the study of pattern recognition. ML interprets data and make predictions by exploring algorithms that learn and extract patterns in data such as various regression methods, classification methods, etc., and improve automatically through experience [19]. This chapter focus on the binary classification tasks using two ML approaches, supervised and weakly supervised learning. The literature of these two ML approaches, the ML algorithms implemented and the performance metrics that are used in this study are defined in this chapter.

#### 2.1.1 Supervised Learning

Supervised learning is a type of ML that learns from labeled data in order to make predictions on new unseen data. The supervised learning algorithm learns a function based on the input-output pairs to map the input vector to the output. The training sample consist of labeled data with instance  $x$  and a corresponding label  $y$ , that is,  $\{(x_i, y_i)\}_{i=1}^n$ . For a given set of training data,  $\{(x_i, y_i)\}$  with  $x_i$  representing a feature vector and  $y_i$  representing the label of training example  $i$ , the algorithm learns a function  $f : \chi \rightarrow Y$ , where the goal of the function  $f(x)$  is to predict the correct label  $y$  when given test data  $x$  such that  $\chi$  is the domain of instances and  $Y$  is the domain of labels [20, 21].

There are two types of Supervised learning, classification and regression. Classification is a technique of identifying to which category does a set of observations belong to with discrete  $y$  labels. Regression techniques are used to predict a target continuous value  $Y$  based on the relationship of variables in the data. This study focus on the classification techniques since the goal of the study is to discriminate instances from two classes.

### 2.1.2 Unsupervised Learning

Unsupervised learning is a type of ML that trains a machine to estimate a model that can find unknown patterns in data that is not labeled, classified or categorized. Unsupervised learning is divided into the parametric unsupervised learning and the non-parametric supervised learning. The parametric unsupervised learning assumes that the data points are independently and identically drawn from a common distribution, based on a fixed set of parameters. For a given set of input examples  $X = x_1, x_2, \dots, x_n$ , the machine estimate a model representing the probability distributions for a new input instance  $x_i$ , such that  $P(x_i|x_1, \dots, x_{i-1})$ . Non-parametric unsupervised learning groups the data into clusters based on shared attributes. In clustering, the common attributes in the data are identified and the data is then grouped based on the similarities and dissimilarities of the instances [22, 23].

### 2.1.3 Semi-supervised Learning

Semi-supervised learning is a class of ML techniques that make use of unlabeled data and labeled data for training. Semi-supervised learning use unlabeled data with a small proportion of labeled instances to model a classifier that takes an advantage of using labeled and unlabeled instances. Labeled data is difficult to obtain since it requires experts to draw conclusion on the labels of the instances. Most often data collected is unlabeled and advanced techniques such as unsupervised learning techniques are used to label the data. Semi-supervised learning used to label unlabeled instances use self training which employs a large amount of unlabeled data with a proportion of labeled data to build a classification model that can label unlabeled instances. The model is built on the labeled instances, and then

tries to add labels to the unlabeled instances. For a given set of  $n$  examples containing independently identically distributions,  $(x_1 \dots x_n) \in \mathcal{X}$  with the corresponding labels  $(y_1 \dots y_n) \in \mathcal{Y}$  and another given set of unlabeled examples  $x_{n+1} \dots x_{n+u}$ , the semi-supervised learning make use of the combined given data to build a classifier that could be obtained by following the supervised learning approach where the unlabeled data is discarded or make use of the unsupervised learning by discarding the labeled data. Assuming that the data comes from a Gaussian mixture model in Equation 2.1

$$f(x|\theta) = \sum_{k=1}^n \alpha_k f(x|\theta_k), \quad (2.1)$$

given that  $\alpha_k$  is defined as the mixture coefficient,  $\sum_{k=1}^n \alpha_k = 1$  and  $\theta = \theta_i$  are the parameters. Let the target label  $y_i$  be a random variable with distribution  $P(y_i|x_i, z_i)$ , where  $z_i$  is a mixture variable that determines the random label distribution with the feature vector  $x_i$ . The maximum posterior gives:

$$g(x) = \underset{E \in (Y, N)}{\operatorname{argmax}} \sum_{k=1}^n P(y_i = E | z_i = k, x_i) P(z_i = k | x_i), \quad (2.2)$$

and

$$P(z_i = k | x_i) = \frac{\alpha_k f(x_i | \theta_k)}{\sum_{j=1}^n \alpha_j f(x_i | \theta_j)}. \quad (2.3)$$

The goal is to estimate  $P(y_i = E | z_i = k, x_i)$  and  $P(z_i = k | x_i)$  from the training data [24, 25, 26].

### 2.1.4 Weakly Supervised Learning

Weakly supervised learning is a term that covers the processes applied to build a model by learning with weak supervision. Weak supervision uses labeled data, but the labels are not quite correct. There might be noise in the labeling process of the data, or it might be that the labels represents which of the collection of classes a given observation belongs to. The goal of weak supervision is the same as the supervised learning approach, to develop a model that accurately outputs the labels, however it is not trained on correctly labeled training data. For a given unlabeled



data  $X = (x_1, \dots, x_n)$  there are weak labeled training instances from the weak supervision  $P_i(y|x), i = 1$ , such that each instance has a coverage set  $C_i$  which represents the set points  $x$ . The accuracy of the model is defined as the probability of the correct target label  $y$  over  $C_i$ , assumed to be  $< 1$ . Weak supervision is divided into three classes, the incomplete supervision, inexact supervision and the inaccurate supervision. The incomplete supervision is derived from the semi-supervised or active learning where only a proportion of instances in the training data are labeled. This method uses the active learning to label the unlabeled instances assuming that the cost of the labeling of these instances depends only by the number of queries. In the inaccurate supervision, the training data consists of instances with correct and incorrect labels. This method employs several learners to identify the unlabeled instances and examine with the training data in order to make corrections. In the inexact supervision, the training data is assigned labels but the labels are not as exact as desired. In this method, the training data consists of labeled bags which contains unlabeled instances. Each bag consists of the collection of unlabeled instances where a bag is negatively labeled if all the instances in the bag are negative. The task is to learn  $F : \chi \rightarrow Y$  from a given training data  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ . The bag is defined as  $X_i = \{x_{i1}, \dots, x_{i,n_i} \subseteq \chi\}$  where  $n_i$  is the number of instances in the bag  $X_i$  and  $y_i \in Y = \{Y, N\}$ , and each instance is defined as  $x_{ik} \subseteq \chi (k \in \{1, \dots, n_i\})$ . A bag  $X_i$  is positively labeled if there exist  $x_{ip}$  which is positive and  $P \in \{1, \dots, n_i\}$  is unknown. The goal of this method is to correctly predict the labels of the new unseen bags [27, 28, 26].

### 2.1.5 Boosted Decision Trees (BDT)

Decision tree is data mining tree like structure model that performs a top-down split greedy approach [29]. The model structure consists of a root node, branch and leaf node where each internal node represents the input attribute, the branch represents the outcome and the leaf node represents the target output. The root node is at the top and instances are sorted down from the root node to the leaf node, in which splitting is done on the entire training set. Decision trees are very popular and they are commonly used in data mining to derive strategy for decision making, and they have gained popularity and have been widely used in ML to predict target output. Although decision trees are very powerful, they are prone to bias

and variance. Studies have been done to simplify decision trees and the techniques developed by J.Quinlan (1999) [30] show that decision trees can be simplified while still retaining their accuracy. Ensemble techniques have been implemented on decision trees to improve the prediction of decision trees by constructing more than one decision tree to build a better model. The principle behind ensemble methods is to combine several weak learners in order to form one strong optimal model with improved prediction performance. Bagging and boosting are some of the popular ensemble techniques implemented to improve the performance of weak models by lowering error rates and prediction errors in classification and regression tasks, respectively. Bagging and boosting in decision trees are used to reduce variance of the model and reduce error from the prior tree, respectively. In bagging, each model is trained on samples independently with replacement from the original training dataset and the target output is the weighted sum of the predictions made by all the models. Boosting trains the initial model on the training set with replacement from the training dataset while passing through all the training patterns in order to note patterns that gives the least errors [31, 32].

### The Mathematics of Decision Trees and Gradient Boosting

Figure 2.1 is an illustration of a decision tree. Each node represents the attributes (features), each branch and leaf represents the decision and the outcome, respectively. To build a tree, the attribute with the highest information gain ( $IG$ ) is chosen. The entropy ( $H$ ), known as the information theory, is calculated first according to Equation 2.4

$$H(s) = \sum_{c \in C} -p(c) \log_2 p(c), \quad (2.4)$$

where  $H(s)$  is the entropy ( $H$ ) in a set, measuring the amount of uncertainty in the set,  $s$  is the set for which  $H$  is calculated,  $C$  is the set of classes in  $s$ , and  $p(c)$  is the proportion of the number of entries in class  $C$  to the entries in set  $s$ . The  $IG$  is then calculated according to Equation 2.5

$$IG(A, s) = H(s) - \sum_{t \in T} p(t) H(t), \quad (2.5)$$

where  $A$  is an attribute,  $IG(A)$  is the measure of a decrease in entropy  $H$  after splitting  $s$  into  $A$ ,  $T$  is a subset created when  $s$  was split to  $A$ ,  $p(t)$  is the proportion of

the total entries in subset  $t$  to the entries in  $s$  and  $H(t)$  is the entropy  $H$  of subset  $t$  [33, 34].

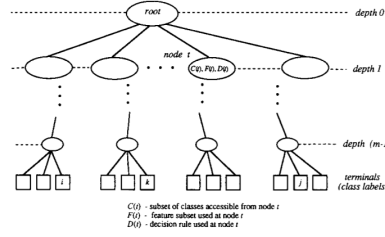


FIGURE 2.1: Illustration of a Decision Tree [34]

In gradient boosting, weak learning models are combined together to build one complex model [35, 36]. Consider the input training data  $(x_i, y_i)_{i=1}^n$ , loss function  $L(y, F(x))$  and iterations  $m$ .

An initial model consist of a function  $F_0(x)$  defined with a constant value where the function is defined according to Equation 2.6:

$$F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma) \quad (2.6)$$

$$F_m(x) = F_{m-1}(x) + \operatorname{argmin}_{h_m \in H} \left[ \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + h_m(x_i)) \right]. \quad (2.7)$$

The model predicts  $y$  by minimizing the mean squared error.

(2.8)

The differential of Equation 2.6 with respect to  $\gamma$ , minimizes to Equation 2.9

$$F_0(x) = \frac{\sum_{i=1}^n y_i}{n}, \quad (2.9)$$

where  $i$  is the indexes of the training set. The pseudo-residual of each error is calculated according to Equation 4

$$r_{im} = \left[ \frac{\delta L(y_i, F(x_i))}{\delta F(x_i)} \right]_{F(x)=F_{m-1}(x)}. \quad (2.10)$$

From the given dataset, new data is derived and modified as defined in Equation 2.11

$$D = (x_i, r_{im}). \quad (2.11)$$

A base learner  $h(x)$  which represent the direction in which the loss function is minimum with respect to  $F_{m-1}$  is trained and fitted on the modified dataset.  $F_m(x)$  is then calculated according to Equation 2.12:

$$F_m(x) = F_{m-1}(x) + \gamma h_m(x), \quad (2.12)$$

where  $h(x)$  is calculated using Equation 2.13

$$h_x = F_{m+1} - F_m(x), \quad (2.13)$$

where  $F_{m+1}$  attempts to correct the errors made by the previous model. The base learner  $h_m(x)$  is then fitted on the pseudo-residuals  $r_m$ . Weights  $W$  are moved in the direction of decreasing loss function. The learning rate hyper-parameter controls how much weights moves. The weight  $W$  are updated according to Equation 2.14

$$W \leftarrow W - \alpha \frac{\delta L}{\delta W}, \quad (2.14)$$

optimum  $\gamma$  found by solving Equation 2.5

$$\gamma_m = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)). \quad (2.15)$$

When the optimum  $\gamma_m$  is solved, Equation 2.12 is updated.

## 2.1.6 Deep Neural Networks

Neural networks are a set of algorithms inspired by the biological neural networks that recognise patterns in data and improves automatically through experience. Neural networks consists of the input layer, hidden layer with nodes and an output layer. The goal of the neural network is to take an input and learn a function that can map the input to the output [37]. The nodes are linked to each other by a connection represented by weights, where each node has an activation function. In a feed-forward network, the patterns of the input information are fed into the input layer which the neural network learns, triggering the hidden layers which apply different transformations to the information received from the previous layer, allowing the neural network to learn more complicated patterns. The hidden layer send the information extracted from the patterns to the output layer [38].

Neural network learns by a process of back propagation, where the goal of the back propagation is to optimize the weights by computing the gradient descent with regard to the weights and biases. A loss function is used to measure the difference between the neural network output to the desired output and to estimate the loss. Once the loss is estimated, the neural network propagates the loss information back-wards, from the output layer, through the nodes in the hidden layers, to the input layer, causing the neural network to adjust weights. The loss function is minimized by using the gradient descent technique which allows the weights to be updated by calculating the gradient of the loss function in order to determine the direction of the global minimum [39, 40].

### The Mathematics of Neural Networks

Figure 2.2 on the Left illustrate the structure of a neural network with dimension of 4 input, a hidden layer and 2 output layers. The input  $x_i$  is connected to the output node by weights  $w_i$ . Each node receives a vector ( $x$ ) which contains the feature values and compute a corresponding output target ( $y$ ). Each node consists of a set of parameters, weights ( $w$ ) and a bias ( $b$ ). The node calculates the weighted average of the vector  $x$  and add  $b$  such that:

$$z = \sum_{i=1}^m w_i x_i + b, \quad (2.16)$$

where  $m$  is the total number of inputs. The input  $x$  and  $w$  represented as vectors gives the dot product of the two product shown in Equation 2.17 below:

$$z = w^T \cdot x + b. \quad (2.17)$$

The output  $y$  is passed through an activation function  $g$  such that:

$$y = g(z). \quad (2.18)$$

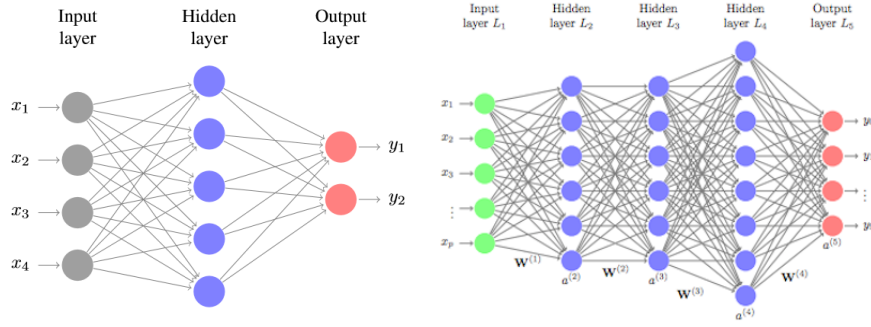


FIGURE 2.2: Illustration of a shallow Neural Network (Left) [41] and Deep Neural Network (Right) [42]

The advancement of computers and technology improved the computation, the amount of data generated and data collection processes and advanced neural networks algorithms. These advancements allowed the improvement of neural networks to run high dimensional parameter and deep neural, known as deep neural network (DNN). Figure 2.2 on the Right shows an architecture of DNN. DNNs are different from shallow neural networks since they have more than one hidden layer, whereas shallow neural networks have only one. When writing the notation of a selected layer in DNN, vector  $a$ , which is the activation function corresponding to the selected layer, is used. The layers in the DNN performs calculations following Equation 2.19:

$$z_i^l = w_i^T \cdot a^{l-1} + b_i, \quad (2.19)$$

and the activation is calculated according to Equation 2.20:

$$a_i^l = g^l(z_i^l), \quad (2.20)$$

where  $i$  represent the index of the node in a layer and  $l$  represent the layer. The  $b$  of each node in a layer and vectors of  $w$  are independently stacked together to form a vertical vector  $b$  and matrix  $W$ , respectively. The matrix  $W^l$  is represented as:

$$\begin{bmatrix} w_1^l T \\ w_2^l T \\ w_3^l T \\ \vdots \\ w_{n-1}^l T \end{bmatrix}$$

and the vertical vector  $b^l$  is represented as follows:

$$\begin{bmatrix} b_1^l \\ b_2^l \\ b_3^l \\ \vdots \\ b_{n-1}^l \end{bmatrix}$$

When considering multiple entries instead of working with one layer as we have from the equations derived above, vectorization is applied across the multiple entries. Given  $m$  entries and  $j$  features, the vertical vectors  $x, a$  and  $z$  create  $X, A$  and  $Z$ , respectively. Equation 2.19 and Equation 2.20 becomes:

$$z = W^l . A^{l-1} + b^l, \quad (2.21)$$

and

$$A^l = g^l(z^l). \quad (2.22)$$

Different activation functions can be used for each layer, which introduce non-linearity in the output of the node. The  $W$  and  $b$  parameters are adjusted using the Equation 2.23 and Equation 2.24 in the back propagation:

$$W^l = W^l - \alpha \frac{\delta L}{\delta W^l}, \quad (2.23)$$

and

$$b^l = b^l - \alpha \frac{\delta L}{\delta b^l}, \quad (2.24)$$

where  $\delta W$  and  $\delta b$  are the partial derivatives of  $W$  and  $b$ , and  $\alpha$  is the learning rate. The learning rate  $\alpha$  is a hyper-parameter which allows control of the adjustments. Setting  $\alpha$  is important since it controls how your neural network model learns. If the  $\alpha$  is set too low, the model learns too slow and if set too high, the model might miss global minimum [41, 43].

### 2.1.7 Activation Functions

Activation functions determines whether a node should be activated or not activated given an input or a set of inputs. They introduce non-linear properties in neural network, allowing the network to learn and "represents non-linear complex arbitrary functional mappings between inputs and outputs". There are several activation functions used in neural networks, rectified linear units (ReLU) and sigmoid are some of the popular and commonly used activation functions. ReLU is the most commonly used activation function in neural networks that gained its popularity in the past few years, almost all deep learning applications use ReLU. ReLU is linear for all positive values and zero for all the negative values.

$$\text{ReLU} : f(x) = \begin{cases} 0, & \text{for } x < 0 \\ 1, & \text{for } x \geq 0 \end{cases}$$

This means that ReLU does not saturate when  $x$  gets larger [44], it allows for quick training of neural networks, and it avoids and rectifies the vanishing gradient problem. X.Gloret al [45], demonstrated that ReLU allows for better training of deeper neural networks when compared with the other activation functions. Although ReLU is commonly used in all deep learning applications and has improved the training of deep learning, it has some drawbacks. One of its major drawbacks is a dying ReLU problem which causes several nodes to die and remain inactive [45, 46]. This problem can be solved with the use of Leaky ReLU which introduces a small slope to the gradient of ReLU. Figure 2.3 is an illustration of the ReLU and Leaky ReLU activation functions .



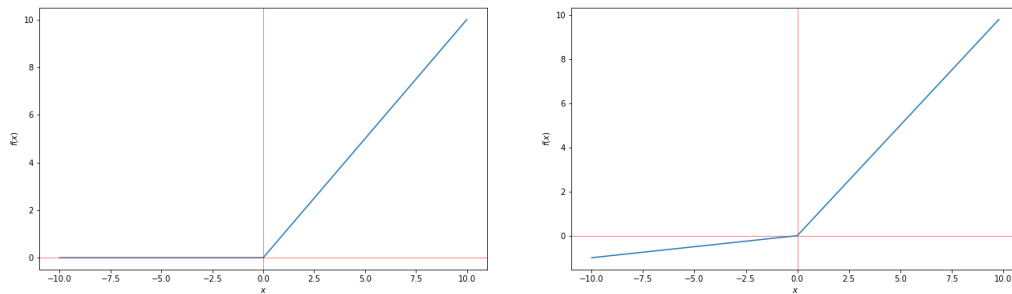


FIGURE 2.3: Plot of the ReLu Activation function (Left) and the Leaky ReLu activation function (Right)

Another popular and commonly used activation function is the sigmoid function, defined in Equation 2.25.

$$f(x) = \frac{1}{1 + e^{-x}}. \quad (2.25)$$

Sigmoid is characterised by a 'S' shaped curve, ranging between 0 and 1. Sigmoid functions have a problem of vanishing gradient which prevents the weights from changing its value, results in back-propagation errors [47].

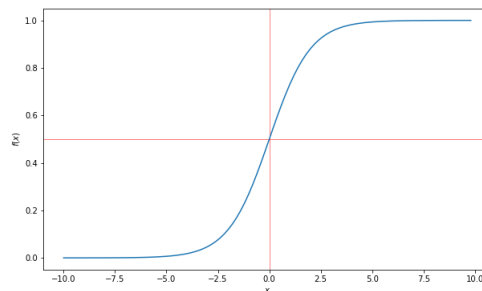


FIGURE 2.4: Plot of the Sigmoid function

## 2.1.8 Classifier Performance

In ML, the performance measurement of a model is an essential task. In classification tasks, the classifier produce a discrete class label indicating the predicted class of the instance. To measure the performance of the classifier, classification metrics are used. The area under the curve (AUC) and receiver operating characteristics

(ROC) curves analysis are one of the most important evaluation metrics to measure the performance of the classifier. The AUC and ROC analysis have been adopted in different fields of study as a standard evaluation metrics to measure how well a model is capable of distinguishing classes. To distinguish between the actual class label and predicted class label, the labels  $\{Y, N\}$  are used for the classifier's predicted class label [48]. In a binary classification task, for each instance in a given classifier, there are four possible outcomes:

True Positive ( $TP$ ) = Number of positive instances correctly classified as positive.

False Positive ( $FP$ ) = Number of negative instances incorrectly classified as positive.

True Negatives ( $TN$ ) = Number of negative instances correctly classified as negative.

False Negative ( $FN$ ) = Number of positive instances incorrectly classified as positive.

Table 2.1 represents the confusion matrix. The major diagonal (TP and TN) represents the number of correctly classified instances and the second diagonal (FP and FN) represent the number of instances that are incorrectly classified.

TABLE 2.1: Confusion Matrix

	Positive (1)	Negative (0)
Positive (1)	TP	FP
Negative (0)	FN	TN

The true positive rate ( $TPR$ ) and the false positive rate ( $FPR$ ) are defined in Equation 2.26 and Equation 2.27, respectively, where  $tp$  represents the number of correctly classified positives,  $P$  represents the total number of positive instances,  $fn$  represents the number of incorrectly classified negatives and  $N$  represents the total number of negative instances.

$$TPR \approx \frac{tp}{P}, \quad (2.26)$$

$$FPR \approx \frac{fn}{N}. \quad (2.27)$$

The ROC curve is a plot of the  $TPR$  against the  $FPR$  at different classification threshold. The AUC measure the area under the ROC curve to provide the overall measure of performance ranging from 0 to 1. An excellent performing classifier has AUC closer to 1 indicating that the classifier has a good measure of separability. A poor performing classifier has AUC closer to 0, indicating that the classifier has the worst measure of separability. The higher the AUC, the better the model is at predicting instances of different classes.

Other common evaluation metrics are defined in Equation 2.28 - Equation 2.30. Precision and recall are also important evaluation metrics in machine learning. Precision indicate the ability of the classifier to predict relevant instances, that is, how well the classifier can label negative instances as negatives and not as positives. Recall indicates how well the classifier can correctly classify all the positive instances by providing a percentage of total positive results correctly classified. Accuracy provides the percentage of the correctly classified instances [49]. Accuracy is not a reliable metric for the real performance of a classifier, however, it is commonly used with other evaluation metrics.

$$precision = \frac{TP}{TP + FP}, \quad (2.28)$$

$$accuracy = \frac{TP + TN}{P + N}, \quad (2.29)$$

$$recall = \frac{TP}{P}. \quad (2.30)$$

## 2.2 Related Studies

Discrimination of signal and background instances produced during the PP collision in the LHC is a very crucial step in searching for new signal instances produced during the PP collision without losing or missing the signal instances [10, 50]. Studies have been conducted using modern techniques that can discriminate these two types of instances with maximum accuracy and also out perform the traditional physics methods.

High energy physics data is often imbalanced with the class of interest being the minority class. H.Han et al. [51] compares different techniques that can be used to handle highly imbalanced dataset and highlights metrics that can be a measure of a classifier performance when working with imbalanced dataset. There are different techniques used to correct imbalance dataset, synthetic minority Over-sampling technique (SMOTE) is one of the methods used to over-sample the minority class. In this study, SMOTE is compared with other over-sampling techniques. A confusion matrix is used to measure the performance of the classifier, and the results obtained show that the SMOTE approach achieves better results.

BDT and ANN have been widely used in high-energy physics. A comparison study of this two techniques applied on particle identification using Monte Carlo samples is discussed in this paper [11]. The results obtained show that both BDT and ANN techniques in this task are very competitive. The BDT performance improved as number of variables (features) are increased, and proved to be more efficient and superior than the ANN. The results also show that BDT also improves performance when the number of trees are increased.

L.Dery et al. [52] introduces a new machine paradigm for classification task called weakly supervised classification that can be trained on incomplete and/or unlabeled data and still match the performance of a fully supervised learning technique. The classifier is trained on unlabeled data using an approach called learning with label proportions (LLP). Data is split into bags containing unlabeled instances. "For a binary classification task with classes,  $A$  and  $B$ , the training set contains at least one instance of  $A$  classifier". The goal is to optimize a model that can at least identify one instance of  $A$  on unseen bag, given only a fraction of class  $A$  in a given bag.

E.Metodiev et al. [53] introduces another type of weak learning paradigm that can also be used in binary classification, called classification without labels (CWoLa) approach. This approach allows a classifier to be trained on real data set, which implies that this can be a better approach when dealing with the LHC data since the data is not labeled and mixed. The classifier is trained on two mixed samples,  $M1$  and  $M2$ , consisting of both signal and background instances in each sample with

different proportions. Using full supervision, the classifier is trained to discriminate  $M1$  from  $M2$  without providing the classifier with the proportions of each sample. The same classifier used to discriminate  $M1$  from  $M2$ , is then applied to discriminate signal instances from background instances without labels. This approach is found to match the performance of full supervision approach.

## Chapter 3

# Research Methodology

### 3.1 Data

This study employs the Monte Carlo dataset implementing the LHC data recorded from 2015-2017 in the  $pp$  collision at  $\sqrt{s}=13$  TeV generated using PYTHIA8 [54, 55, 56]. One sample follows the 2015 and 2016 LHC data normalized to luminosity of  $36.1 \text{ fb}^{-1}$ , and the other sample follows the 2017 LHC data normalized to luminosity of  $43.8 \text{ fb}^{-1}$ . The signal data consist of  $Z$  boson produced from gluon-gluon fusion sample, real missing transverse energy ( $E_T^{miss}$ ) from neutrinos produced by  $Z$  decay, and the heavy pseudo-scalar particle decaying to  $Z$  boson and heavy scalar sample. Missing energy is the energy that is expected to be detected during particle collision but is not detected on the particle detector since it is carried by particles that do not interact with strong forces. The background data is a mixture of the  $t\bar{t}\gamma$  and  $\gamma$  – Jet background samples and is produced. The background contains fake  $E_T^{miss}$ . The dataset for the supervised approach consist of labeled signal ( $S$ ) and background ( $B$ ) events with 11730  $S$  events and 151534  $B$  events decaying to the  $Z + \gamma$  final state. The second dataset for the weakly supervised learning approach used in this study consists of two different samples of unlabeled data. The first sample is the Monte Carlo data collected from 2015 to 2017 in the  $pp$  collision that consists of the  $Z + \gamma$  events with the invariant mass outside 240-280 GeV. The second sample consist of different statistics of the  $Z + \gamma$  events of the first sample and the BSM monte carlo data. The first sample is made up of only  $B$  events, while the second sample consist of mixed  $S$  and  $B$  events. Pre-selection cuts are applied on the data for both the supervised and weakly supervised approach into three categories with respect to the missing transverse energy significance ( $S_{E_T^{miss}}$ ). The cuts are applied based on

the following:

- low (Low) category :  $2.5 \leq S_{E_T^{miss}} \leq 3.5 \sqrt{GeV}$
- The intermediate (Int) category :  $3.5 \leq S_{E_T^{miss}} \leq 5.5 \sqrt{GeV}$
- The high (High) category :  $S_{E_T^{miss}} \geq 5.5 \sqrt{GeV}$ .

The pre-selection cuts are applied to maximize  $S$  and minimize  $B$ , where the cut points are scanned regions of  $S$  to  $B$  efficiency. The training data consists of 114284 events with only 8211  $S$  events and 106073  $B$ , and the test data contains 48980 events with 3515  $S$  events and 45461  $B$ , before applying the pre-selection cuts. The Low category consists of 1336  $S$  and 7875  $B$  training events, and 529  $S$  and 3369  $B$  test events. The Int category consists of 3486  $S$  and 3996  $B$  events in the training data and only 1744  $B$  and 1524  $S$  events in the test data. The High category consist of 2017  $S$  and 1393  $B$  events in the training data and 874  $S$  and 588  $B$  in the test data. Table 3.1 below is the description of the features used in this study. The features used include two leptons, a photon and multiple jets.

TABLE 3.1: Input variables description

Variable	Description
mu	Average bunch crossing
dphifjmet	The angular distance $\Delta(\phi)$ between forward jets and MET
dphisjmet	The angular distance $\Delta(\phi)$ between soft jets and MET
ssumpt2	Subleading vertex sum $P_t$ squared
djpt	Scalar difference between the vectorial sum $P_t$ of all the jets and leading vectorial sum $P_t$
dsumpt2	Difference between leading and subleading vertex sum $P_t$ squared
dphirefjetmet	The angular distance $\Delta(\phi)$ between vectorial sum $P_t$ of all jets and MET

The 1-dimensional distributions of these input variables are provided in Appendix A. The distributions for the  $S$  and  $B$  have been normalised in order for the distributions to be comparable. All the variables except for the metsig variable,  $E_{miss}^T$ , strongly overlap. This demonstrate that the  $B$  and  $S$  events share similar characteristics that cannot easily be distinguished .

## 3.2 Methods

This section describes the methods applied in this study. The method is divided into two sections, the supervised learning and the weakly supervised learning approach.

### 3.2.1 Supervised Learning

The ROOT data format is converted into a python readable format using uproot I/O library, primarily intended to stream data into ML libraries in python. The data is split into 70:30 train and test for each category before pre-processing. The MinMax scaler is applied to transform the feature events by scaling to a range of 0-1. Min-max scaler is used instead of standard scaler since we are not concerned with the standardisation along the variance axes. The data described in Section 3.1 shows that the data highly imbalanced, with 1:13  $B$  to  $S$  ratio. The training data for each category is split into 80 % training and 20 % validation before oversampling the minority class. Mixmax scaler is then used to scale the data points between 0 and 1 on the training data, and is then fitted on the validation and test data independently. Random oversampling technique is applied to improve the  $S$  to  $B$  ratio by supplementing multiple copies of some of the  $S$  events. Oversampling is only applied on the training data of each category after splitting the train, validation and test data to avoid having multiple copies of events instances of the training data also in the validation and test data.

### 3.2.2 Weakly Supervised Learning

The weakly supervised learning method employed in this study is a new method derived from the inaccurate weak supervision. This approach uses the unlabeled data samples of only background and the mixed sample. The two samples are assigned weak labels, where the background only events in sample 1 ( $M1$ ) are labeled 0 and the mixed  $S$  and  $B$  events in sample 2 ( $M2$ ) are labeled 1. The samples are split into 70:30 training and testing data and a model is trained to distinguish events from  $M1$  and  $M2$ . The same model trained to distinguish  $M1$  and  $M2$  events, is used to discriminate the  $S$  and  $B$  events in the test data used in the supervised learning approach. E. Metodiev et al. [53] shows that a weakly learning model trained



to distinguish events from different samples can be employed to discriminate the  $S$  and  $B$  and still performs just as well as the supervised learning model. This study takes a different approach by using the same model trained to distinguish events from  $M1$  and  $M2$ , to discriminate  $S$  and  $B$  from the supervised learning approach in order to measure the performance of the weakly supervised learning with the supervised learning results to help develop a model for real data. The performance of the model is measured by comparing the ROC curves of the supervised learning categories and the output distributions.

### 3.2.3 Application

#### Application of Boosted Decision Trees

The K-fold cross validation is used to first find the optimal maximum depth. The optimal maximum depth is determined by evaluating the boosted decision tree on a held-out data set using the 10-fold cross validation for the maximum depth range of 1 to 25. The data is re-sampled, randomly split into 10 folds and fit the trees using K-1 folds and validate on the remaining folds. The scores are noted and the process is repeated until all the K-folds serves as the validation set. The classification accuracy is measured on the validation folds. The optimal maximum depth is chosen based on the best bias-variance trade-off. The optimal maximum depth is then fixed, and a GridSearchCV hyper-parameter tuning is used to apply all the possible combinations of parameters provided in the grid through a list of dictionaries to find the optimal parameters to build a BDT model. The BDT is trained for each value in the grid and the performance score of each iteration is provided. The optimal parameter selected is the one that gave the best score. The validation data is used to evaluate the performance of the model while fine tuning the parameters to ensure there is little overtraining observed from the training and validation accuracy before applying it to the test data. The parameters used to build a BDT model for each category are defined in Table 3.2.

TABLE 3.2: BDT Hyper-parameters

Hyper-parameters	Low Category	Int category	High Category	No-cut
NTrees	700	800	100	1500
Max-Depth	3	3	1	3
Learning rate	0.01	0.01	0.01	0.1
BoostType	Gradient	Gradient	Gradient	Gradient

### Application of Deep Neural Network

The k-fold cross validation is used to first determine the activation functions. The held-out approach 5-fold cross validation is applied for ReLu, Sigmoid and Leaky ReLu activation function. The activation functions for the layers were then chosen based on the accuracy and the highest AUC score. The activation functions are fixed during the optimal hyper-parameter search. GridSearchCV is used to find the optimal batch size, number of hidden layers, drop out rates and the epochs. A dictionary of hyper-parameters with the parameter name and an array of values are provided to evaluate a model in the parameter grid. The gridsearch construct and evaluate the model for each combination of parameters. The combination of parameters which gives the best score are used. The optimal hyper-parameters are added to the sequential model. Adam Optimizer is used to find the individual learning rate of each parameter. Table 3.3 show the parameters used to build the DNN model. ReLu activation function is used to activate the hidden layers, and sigmoid function is used to activate the output layer.

TABLE 3.3: DNN Hyper-parameters

Parameters	Low Category	Int category	High Category	No-cut
No. Hidden Layers	3	3	3	4
Batch Size	100	100	60	20
Epochs	100	100	100	10
Dropout	0.1	0.1	0.1	0.2

### 3.2.4 Performance Metrics

To measure the performance of the models, ROC curves, AUC score, confusion matrix and output distributions are used. The ROC curve is a summary of the confusion matrix, however, we use both the ROC AUC score and the confusion matrix for analysis. The output distributions are used for easy read and interpretation of the results. For a model that can distinguish  $S$  and  $B$  events with no misclassifications, we expect a ROC curve with an AUC score of 1, and an output distribution that shows an ideal separability with no overlapping of the  $S$  and  $B$  classes.

## 3.3 Summary

Simulated Monte Carlo dataset is used in this study. Optimal BDT and DNN models used in this study are developed following the supervised and weakly supervised learning approaches. The weakly supervised learning approach implemented is called weakly supervised learning where two samples of unlabeled data are used, with one sample containing mixed  $S$  and  $B$  events. A model is trained to first distinguish events from each sample, and then used to discriminate  $S$  and  $B$  events from new dataset.

## Chapter 4

# Results and Discussion

In this Chapter the supervised learning and the weakly supervised learning (WSL) results achieved with BDT and DNN are discussed. The output distributions show how well the models can distinguish  $S$  and  $B$  events. The ROC curve is the plot of signal efficiency against background rejection, a summary of the confusion matrix. Using the output distributions, table of results and the ROC curves, the supervised learning results the WSL and both the BDT and DNN overall performance are presented. The instances will be referred to as events in this chapter.

### 4.1 Boosted Decision Trees Results

The results obtained using BDT are presented in this section. The BDT supervised learning approach results are presented first in Subsection 4.1.1. The WSL approach results are presented in Subsection 4.1.2.

#### 4.1.1 Supervised Learning

Table 4.1 shows the performance of the supervised learning BDT. The Low category correctly classified 75% of  $S$  and 71%  $B$ , misclassifying 25% of  $S$  and 29% of  $B$ . The Int category model correctly classified 68%  $S$  and 74%  $B$ . The High category model correctly classified 68% of  $B$  and 52% of  $S$  events. The No-cut category model correctly classified 85% of  $S$  and 89% of  $B$ .

TABLE 4.1: BDT supervised learning Table of Results

	Correctly Classified		Misclassified	
	$S$	$B$	$S$	$B$
Low Category	0.75	0.71	0.25	0.29
Int Category	0.68	0.74	0.32	0.26
High Category	0.52	0.68	0.48	0.32
No-cut Category	0.85	0.89	0.15	0.11

Figure 4.1 represent the ROC curves of the BDT supervised learning. The Low category model achieved 78% AUC represented by the blue curve showing that the is 22% likelihood of misclassifying events. The Int category model achieved 79% AUC as shown by the green curve and the High category represented by the maroon curve achieved 69% AUC. The is 21% and 31% likelihood of misclassifying events with the Int and Hihg categories, respectively. The red curve represents the No-cut category with 94% AUC. The No-cut illustrates a low likelihood of misclassifying events with the model.

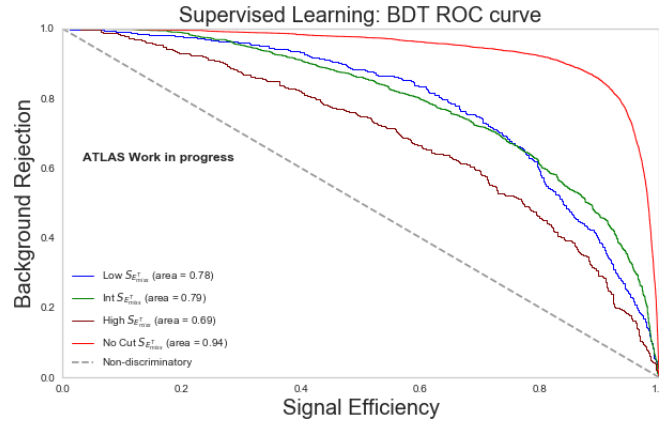


FIGURE 4.1: ROC curves illustrating the performance of the the BDT supervised learning model.

The BDT output in Figure 4.2 shows how well the BDT can distinguish between the  $S$  and  $B$  events. The red stepfilled distribution represent the training  $B$  events and the the blue stepfilled distribution represent the  $S$  train events. The purple and green step distribution represent the  $B$  and  $S$  test. The Low BDT output shows that the model distinguished the  $S$  and  $B$  with less uncertainty. A separation of  $S$  and  $B$  is observed with a small amount of  $S$  and  $B$  overlapping between  $-2$  and  $2$  regions. The output distribution shows that the is a close fit between the train and test output, which demonstrate that the model is slightly overtraining. The  $S$  and  $B$  separation is also observed on the Int category distribution with a slight moderate amount of  $S$  and  $B$  overlapping. The distribution demonstrate a close fit of the test to the train events, with more overtraining observed in the  $S$  region. The High category model performed poorly. The is a huge overlap of  $S$  and  $B$  with little separation observed. The No-cut category BDT output demonstrate a clear separation of  $S$  and  $B$  with little overlapping observed. The is little overtraining observed on the  $S$  region and an almost perfect fit on the  $B$  region.

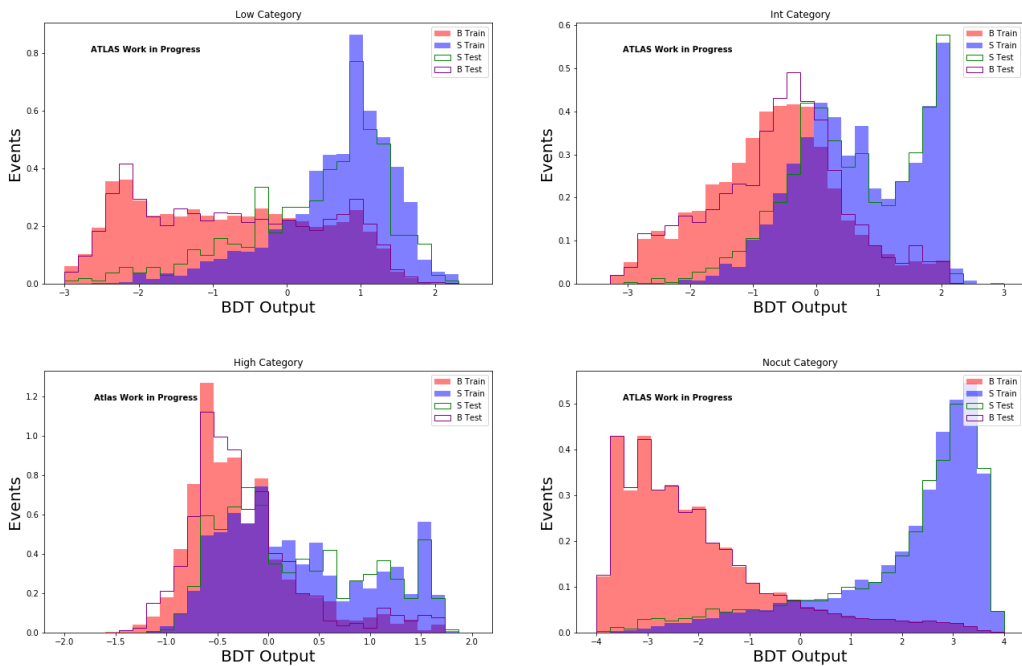


FIGURE 4.2:  $S$  and  $B$  Output distributions

### 4.1.2 Weakly Supervised Learning

Table 4.2 shows the performance of the models when predicting  $S$  and  $B$ . The Low category model is bias towards the  $B$ . The model correctly classified 97% of  $B$  and misclassified 84% of  $S$  events. The Int category correctly classified 69% of  $S$  and 74% of  $B$  events, misclassifying 31% and 26% of  $S$  and  $B$ , respectively. The High category model is bias towards the  $S$ . The model correctly classified 99% of  $S$  and misclassified 94% of  $B$ . The No-cut category is also bias towards the  $B$  with 97% of the  $B$  correctly classified and only 61% of  $S$  events classified correctly.

TABLE 4.2: BDT WSL Table of Results

	Correctly Classified		Misclassified	
	$S$	$B$	$S$	$B$
Low Category	0.16	0.97	0.84	0.03
Int Category	0.69	0.74	0.31	0.26
High Category	0.99	0.06	0.01	0.94
No-cut Category	0.61	0.97	0.39	0.3

Figure 4.3 shows the WSL ROC curves, a summary of Table 4.2. The Low category represented by the blue curve shows that the model achieved 77% AUC, which illustrates that the model have 23% likelihood of misclassifying events. The Int and High category curves in green and maroon, respectively, achieved 79% and 64% AUC. The No-cut category model achieved 93% AUC. The likelihood of the No-cut model misclassifying events is low. The high AUC scores achieved by the models are due to the high ratio of  $B$  over  $S$ .

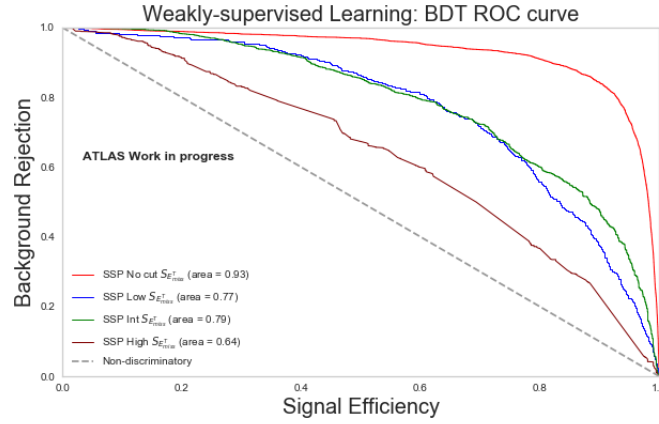


FIGURE 4.3: ROC curves illustrating the performance of the the BDT WSL model.

The output distributions in Figure 4.4 represent the BDT training output of  $M1$  and  $M2$  (Left) and the performance of the BDT model when predicting the test  $S$  and  $B$  (Right). The Low category train output shows that there is less certainty in the results. There is more  $M1$  and  $M2$  events overlapping and a slight separation achieved. The Low category  $S$  and  $B$  test output distribution shows an improvement when the model predicts  $S$  and  $B$ . There is less overlapping of  $S$  and  $B$  observed as compared to the  $M1$  and  $M2$  distribution. The Int category distributions show that the performance of the model on the  $M1$  and  $M2$  train is close to the performance of the  $S$  and  $B$  test. There is  $S$  and  $B$  separation observed on the distribution with a moderate amount of events overlapping. The High category BDT output distributions show that the model performed poorly on classification of both  $M1$  and  $M2$  events, and also the  $S$  and  $B$  events. There is a strong overlap of events on both the train and test distributions. The No-cut category BDT output distribution of the  $M1$  and  $M2$  shows that the model failed to distinguish  $M2$  events from  $M1$  events. The distributions show that the  $M2$  events are misclassified as  $M1$  events. The  $S$  and  $B$  output distribution shows a separation of the  $S$  and  $B$  events, with the  $B$  events correctly classified and fewer  $S$  events correctly classified. The No-cut model is biased towards the  $B$  events.



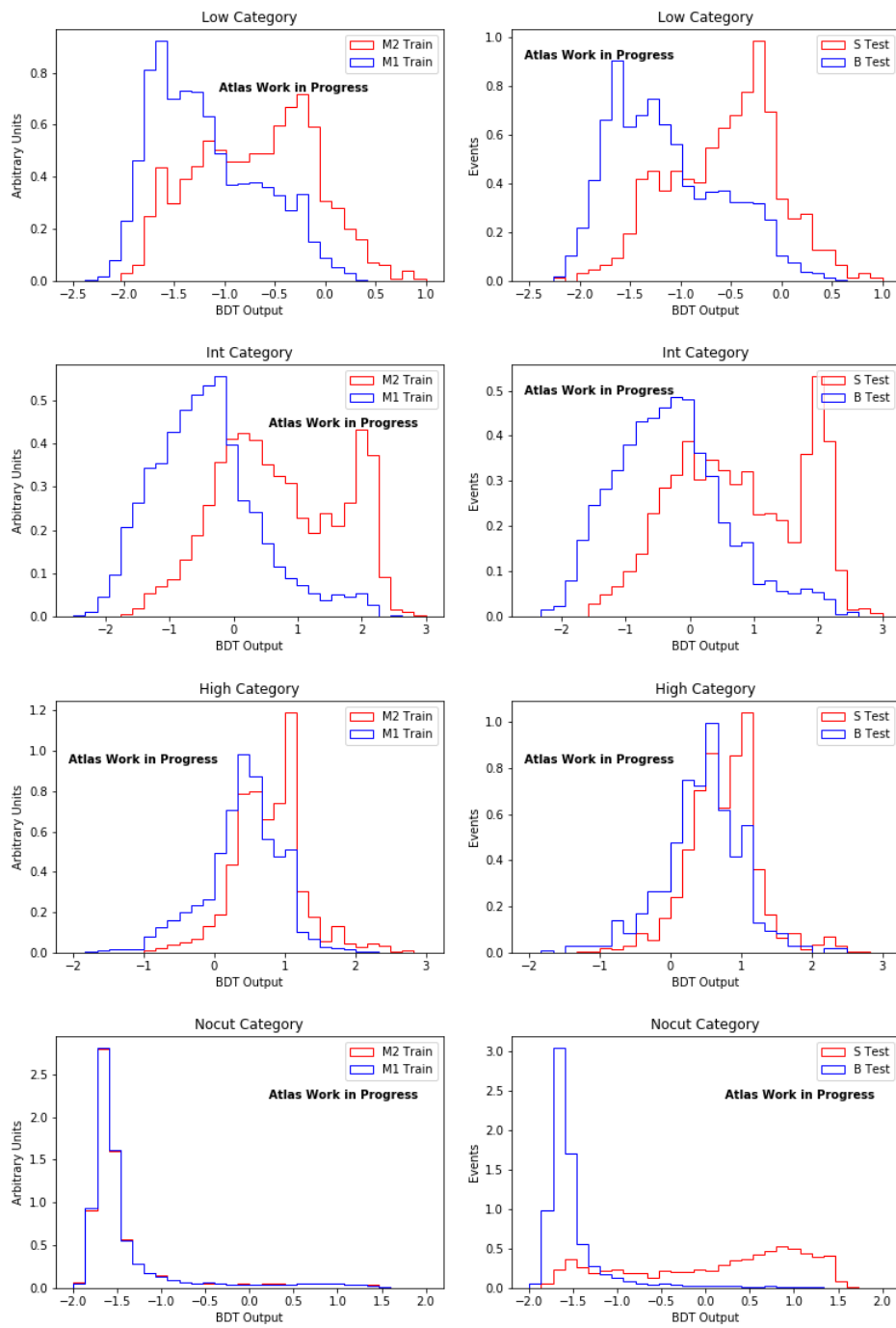


FIGURE 4.4: S and B Output distributions

## 4.2 Deep Neural Networks Results

The DNN results presented in this Section illustrate the performance of the DNN models implemented following the supervised approach and WSL approach. Subsection 4.2.1 presents the results of the supervised DNN models and the WSL results are presented in Subsection 4.2.2.

### 4.2.1 Supervised Learning

Table 4.3 shows a summarised classification results of the DNN supervised learning approach. There is a slight bias towards *S* observed in the Low category where 80% of *S* are correctly classified and 37% of *B* are misclassified as *S*. 63% of *B* are classified correctly and only 20% of *S* events are misclassified as *B*. The Int category achieved 77% and 62% of correctly classified *S* and *B* events, respectively. The High category shows bias towards *S*, where 82% of *S* are correctly classified and 66% of *B* events are misclassified as *S*. The No-cut category achieved 84% and 90% *S* and *B* events classified correctly, respectively, with only 16% of *S* misclassified as *B* and 10% of *B* misclassified as *S*.

TABLE 4.3: DNN supervised learning Table of Results

	Correctly Classified		Misclassified	
	<i>S</i>	<i>B</i>	<i>S</i>	<i>B</i>
Low Category	0.80	0.63	0.20	0.37
Int Category	0.77	0.62	0.23	0.38
High Category	0.82	0.34	0.18	0.66
No-cut Category	0.84	0.90	0.16	0.1

The ROC curves in Figure 4.5 summarises Table 4.3. The Low category represented by the blue curve shows that the model achieved 78% AUC. This means there is a 22% likelihood of misclassifying events with the Low category model. The Int category represented by the green curve shows that the model achieved 75% AUC. The High category model achieved 61% AUC which illustrates that there is a 39% likelihood of misclassifying events and only 61% likelihood of correctly classifying the events.

The No-cut category represented by the red curve shows a high likelihood of correctly classifying the events. The model achieved 94% AUC.

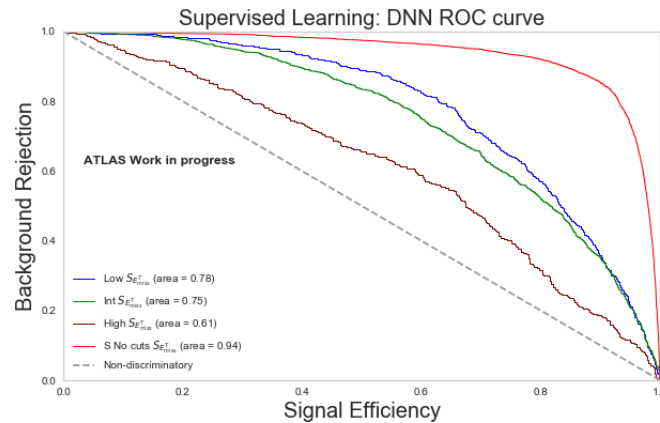
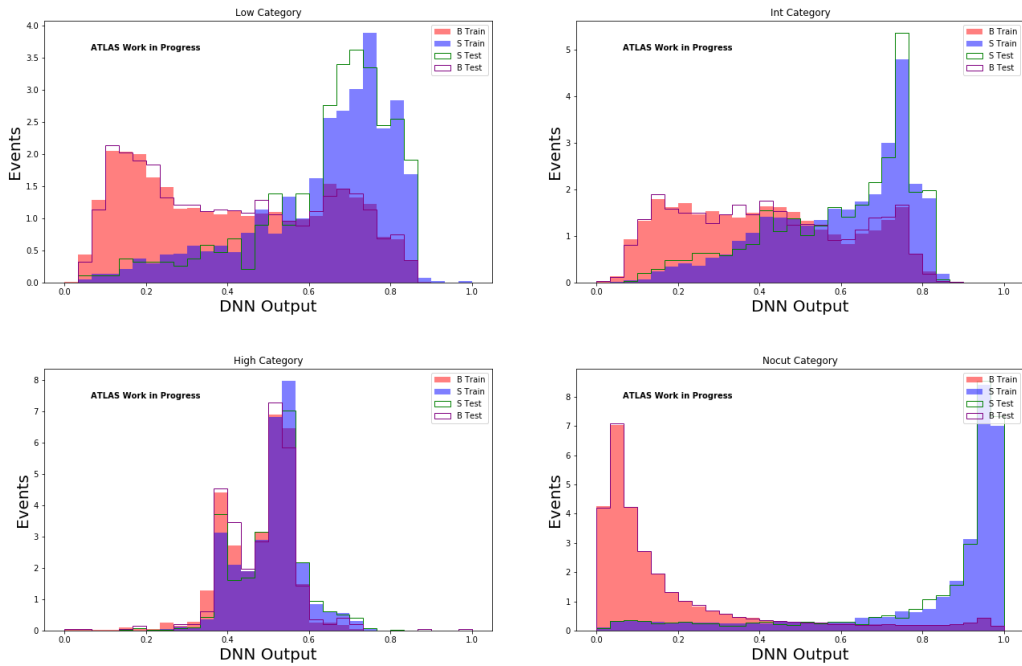


FIGURE 4.5: ROC curves illustrating the performance of the DNN supervised learning model.

The DNN output distribution in Figure 4.6 shows the  $S$  and  $B$  separation achieved by the DNN models. The Low and Int category distributions shows that the models could separate the  $S$  and  $B$  events uncertainty, demonstrated by the overlapping  $S$  and  $B$  events. The is close fit of the test distribution to the train, with overtraining observed on both the Low and Int Categories. The High category model performed poorly, with the  $S$  and  $B$  strongly overlapping. The No-cut category distribution shows a clear separation of the  $S$  and  $B$  and an almost perfect fit on the  $B$  region with slightly-to-no overtraining observed on the  $S$  region.

FIGURE 4.6:  $S$  and  $B$  Output distributions

## 4.2.2 Weakly Supervised Learning

The WSL tabulated results in Table 4.4 shows the performance of the models on the Low, Int, High and No-cut category. The Low category model demonstrate a strong bias towards the  $B$  where 97% of  $B$  is classified correctly as  $B$  and only 3% is misclassified as  $S$ . The model misclassified 82% of  $S$  as  $B$  and only 18% of  $S$  is correctly classified. The Int category shows that the model achieved 67% and 71% of  $S$  and  $B$  events correctly classified. The High category correctly classified 99% of  $S$  with 96% of misclassified  $B$ , demonstrating a strong bias towards  $S$ . The No-cut model is also bias towards the  $B$ . The is 97% of  $B$  events classified correctly and only 58% of  $S$  is correctly classified. The is a strong bias demonstrated by the WSL models.

TABLE 4.4: DNN WSL Table of Results

	Correctly Classified		Misclassified	
	$S$	$B$	$S$	$B$
Low Category	0.18	0.97	0.82	0.03
Int Category	0.67	0.71	0.33	0.29
High Category	0.99	0.04	0.01	0.96
No-cut Category	0.58	0.97	0.42	0.03

Figure 4.7 shows that ROC curves and AUC of the WSL models. The Low and Int categories achieved 77% and 72% AUC represented by the blue and green curves, respectively. The Low and Int models have 23% and 28% likelihood of misclassifying events. The High category achieved 60% AUC and the No-cut category achieved 93% AUC. The No-cut category represented by the red curve shows that there is a low likelihood of misclassifying events. The high AUC score achieved by the models shows that there is a high ratio of  $B$  over  $S$  as Table 4.4 demonstrated a strong bias in the Low, Int and No-cut categories.

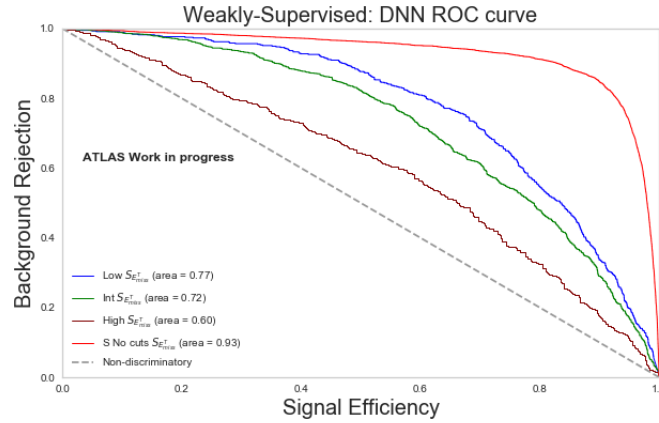


FIGURE 4.7: ROC curves illustrating the performance of the DNN WSL model.

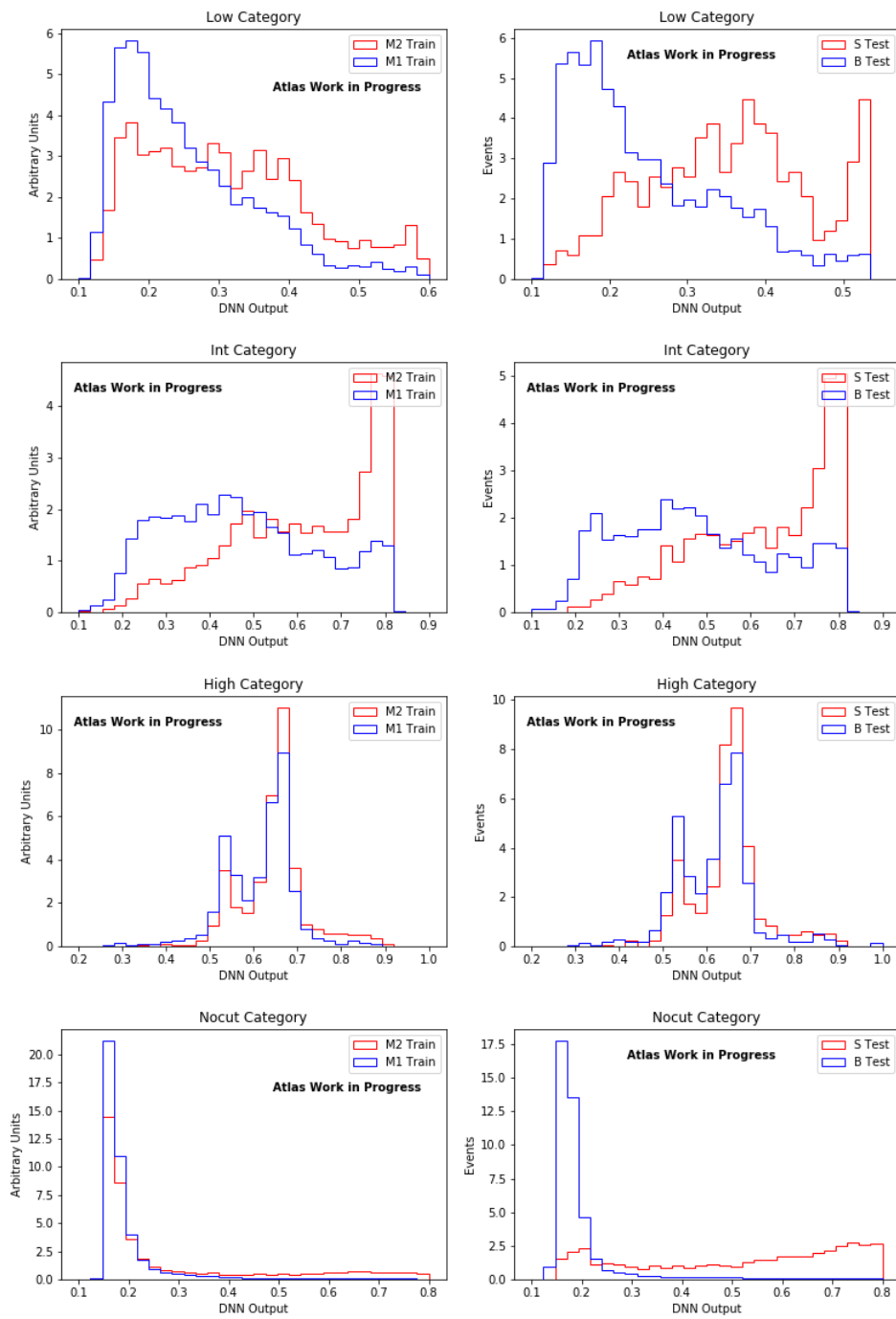


FIGURE 4.8: S and B Output distributions

Figure 4.8 represent the WSL DNN output distributions. The Low category  $M1$  and  $M2$  train shows a strong overlap on the  $M1$  and  $M2$ , demonstrating that the model cannot distinguish the events. The  $S$  and  $B$  test distribution shows an improvement from the train distribution. There is an ideal separation of  $S$  and  $B$  observed with a slight moderate of events overlapping. The Int category shows that the model could distinguish the  $M1$  and  $M2$ ,  $S$  and  $B$  events with less uncertainty. The High category performed poorly with no separation achieved and a strong overlap of the  $M1$  and  $M2$  events and also the  $S$  and  $B$  events. The No-cut category performed poorly on the  $M1$  and  $M2$  train. The model failed to distinguish  $M2$  events from  $M1$  events. The  $S$  and  $B$  test shows a separation of the events with the model bias towards the  $B$ .

### 4.3 Discussion

The performance achieved by the supervised BDT and DNN models shows that the models can distinguish  $S$  and  $B$  with less uncertainty in the results. The output distributions shows that the Low and Int categories could distinguish the  $S$  and  $B$ , observed from separation achieved, with a closely moderate amount of events still overlapping. A threshold is demonstrated by the distributions where a cut can be applied to minimize  $B$  and maximize  $S$ . The No-cut category for both BDT and DNN shows a clear separation of the  $S$  and  $B$  with high certainty. A cut can be applied to separate the  $B$  from  $S$ . The BDT and DNN results of the No-cut category shows an ideal measure of separability and that the models are able to discriminate  $S$  and  $B$  events to almost perfection when no pre-selection cuts are applied and the  $E_T^{miss}$  variable is used as a feature. The WSL approach results demonstrate a strong bias in the models for both BDT and DNN. There is a high likelihood of misclassifying  $S$  as  $B$  with the WSL models which shows there is high uncertainty in the results. The low category models misclassified 84% and 80% of BDT and DNN  $S$  as  $B$ , respectively. The strong bias in the No-cut category is visually demonstrated by the BDT and DNN output distributions where  $B$  events are well distributed and separated from the  $S$  region, while the  $S$  distribution is almost flat and overlapping with  $B$ .

The BDT and DNN results discussed in Section 4.1 and Section 4.2, respectively, illustrates that the BDT models performs better than the DNN model with a higher AUC, however, more overtraining is observed on the BDT supervised learning results than on the DNN as BDT is prone to overtraining. This makes DNN a better model for the WSL approach since it can be optimized to achieve better results with little to no overtraining, assuring certainty in the results. The output distributions and ROC curves on the Low and Int categories show that the model can discriminate  $S$  and  $B$  events with more than 70% discrimination capacity. This means that a cut can be applied where there is maximum  $B$  and less  $S$  to disregard the  $B$  without loosing a lot of  $S$  events. An ideal discrimination capacity means that when a cut is applied to maximize  $S$  and minimize  $B$ , only insignificant amount of  $S$  statistics will be lost and less  $B$  events overlapping the  $S$  region, thereby improving the purity of the  $S$ . DNN performed poorly on the High category, with 61% and 60% AUC on the supervised and WSL, respectively, showing that both models are slightly doing better than a random guess. This poor performance may have been due to fewer events present on the High category as DNN performance improves with more statistics.

The performance achieved with the WSL using BDT and DNN shows that the WSL approach can achieve results close to the supervised learning on the Int category. There is a strong bias observed on the Low, High and No-cut categories due to the high imbalance in the  $S$  and  $B$  events.



## Chapter 5

# Conclusions and Future Work

### 5.1 Conclusions

The supervised learning approach results discussed in Chapter 4 show that the BDT and DNN models can achieve an ideal discrimination capacity, with the No-cut category model outperforming the Low, Int and High categories. There is a strong bias observed in the Low, High and No-cut categories when implementing the WSL approach models for both the BDT and DNN models. Although the AUC and ROC curves achieved by the WSL approach models are close to the AUC and ROC curves achieved by the supervised learning approach, the output distributions and the table of results show a strong bias towards the  $B$  on the WSL approach. This shows that there is a high ratio of  $B$  over  $S$  and the WSL approach cannot discriminate  $S$  from  $B$  events when trained on weak-labeled highly imbalanced dataset.

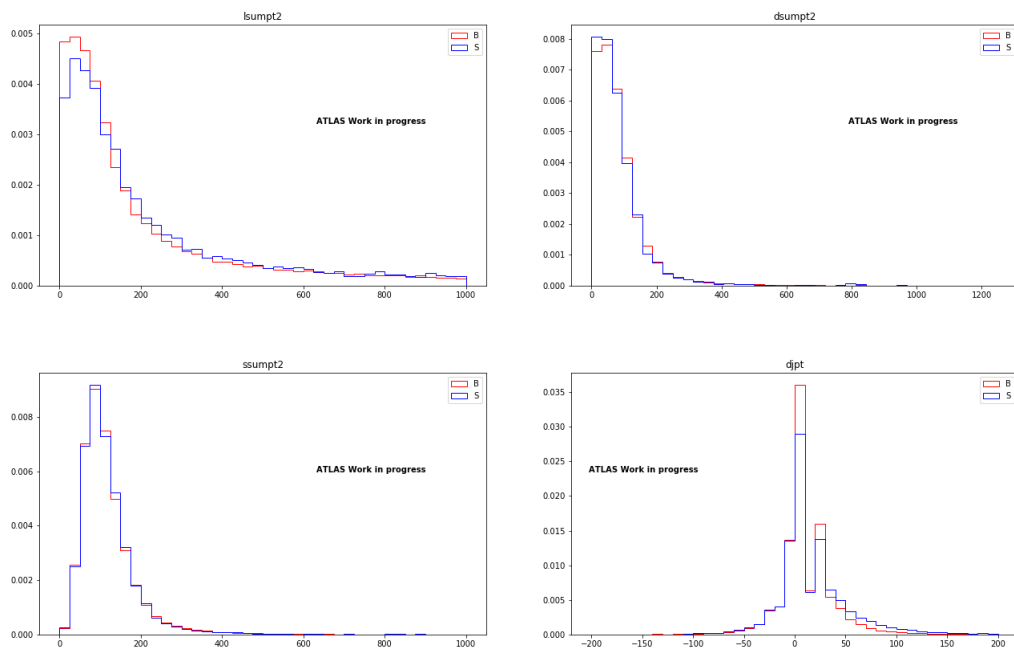
### 5.2 Future Work

Further tests still need to be performed to optimize the WSL approach.  $S$  events will be injected on the training samples in small increment to observe how the model performs as  $S$  increases in the training samples, and optimize the method. The results will be compared with a well defined anomaly detection approach to validate the results obtained with the WSL approach in order to have a well defined model that can be implemented on real data.

# Appendix A

Appendix A shows the figures of the 1-Dimensional distributions variables used in this study, described in Chapter 3.

## A.1 1D Distributions



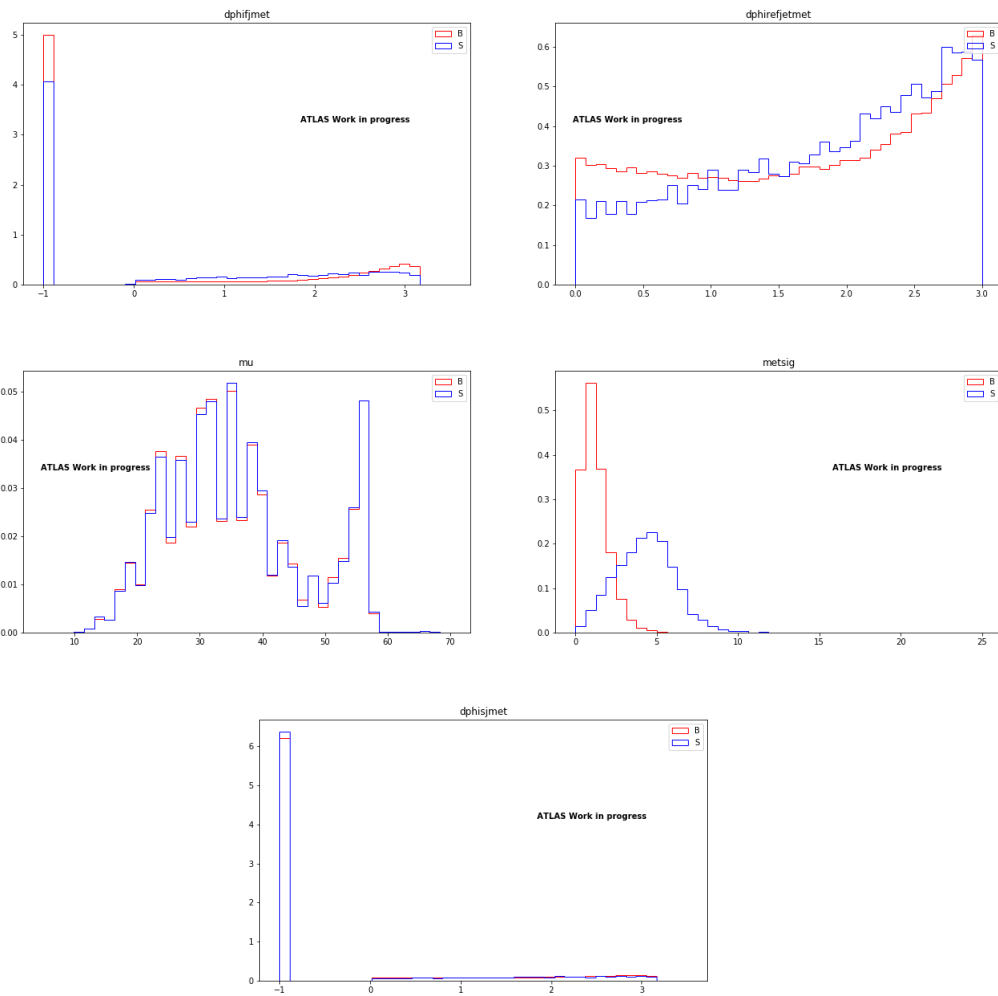


FIGURE A.1: 1 Dimensional distributions of the input variables

# Bibliography

- [1] Georges Aad et al. “The ATLAS experiment at the CERN large hadron collider”. In: *Jinst* 3 (2008), S08003.
- [2] Christopher Rhodes. “Large Hadron Collider (LHC)”. In: *Science progress* 96 (Mar. 2013), pp. 95–109. DOI: [10.3184/003685013X13623370524107](https://doi.org/10.3184/003685013X13623370524107).
- [3] Asmaa Abada et al. “HE-LHC: The high-energy large hadron collider”. In: *The European Physical Journal Special Topics* 228.5 (2019), pp. 1109–1382.
- [4] Giorgio Apollinari et al. *High-luminosity large hadron collider (HL-LHC): Preliminary design report*. Tech. rep. Fermi National Accelerator Lab.(FNAL), Batavia, IL (United States), 2015.
- [5] “High Luminosity LHC”. In: ().
- [6] CMS collaboration et al. “Measurement of differential top-quark pair production cross sections in pp collisions at  $\sqrt{s}=7$  TeV”. In: *arXiv preprint arXiv:1211.2220* (2012).
- [7] Avita Katal, Mohammad Wazid, and RH Goudar. “Big data: issues, challenges, tools and good practices”. In: *2013 Sixth international conference on contemporary computing (IC3)*. IEEE. 2013, pp. 404–409.
- [8] Donald Michie, David J Spiegelhalter, CC Taylor, et al. “Machine learning”. In: *Neural and Statistical Classification* 13 (1994).
- [9] Ian H Witten et al. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, 2016.
- [10] Pierre Baldi, Peter Sadowski, and Daniel Whiteson. “Searching for exotic particles in high-energy physics with deep learning”. In: *Nature communications* 5 (2014), p. 4308.

- [11] Byron P Roe et al. "Boosted decision trees as an alternative to artificial neural networks for particle identification". In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 543.2-3 (2005), pp. 577–584.
- [12] W Noel Cottingham and Derek A Greenwood. *An introduction to the standard model of particle physics*. Cambridge university press, 2007.
- [13] CERN Collaborators. *Feedforward Deep Learning Models*. 2019. URL: <https://home.cern/science/physics/standard-model> (visited on 11/06/2019).
- [14] Georges Aad et al. "Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC". In: *Physics Letters B* 716.1 (2012), pp. 1–29.
- [15] I Sackmann, Arnold I Boothroyd, William A Fowler, et al. "Our sun. I-The standard model: Successes and failures". In: *The Astrophysical Journal* 360 (1990), pp. 727–736.
- [16] Stefan Von Buddenbrock et al. "The compatibility of LHC Run 1 data with a heavy scalar of mass around 270\, GeV". In: *arXiv preprint arXiv:1506.00612* (2015).
- [17] Stefan Von Buddenbrock et al. "Phenomenological signatures of additional scalar bosons at the LHC". In: *The European Physical Journal C* 76.10 (2016), p. 580.
- [18] Mukesh Kumar et al. "The impact of additional scalar bosons at the LHC". In: *Journal of Physics: Conference Series*. Vol. 802. 1. IOP Publishing. 2017, p. 012007.
- [19] Željko Ivezić et al. *Statistics, data mining, and machine learning in astronomy: a practical Python guide for the analysis of survey data*. Princeton University Press, 2019.
- [20] "Introduction to statistical Machine learning". In: ().
- [21] Taiwo Oladipupo Ayodele. "Types of machine learning algorithms". In: *New advances in machine learning*. IntechOpen, 2010.
- [22] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. *Semi-supervised learning (chapelle, o. et al., eds.; 2006)*. Vol. 20. 3. IEEE, 2009, pp. 1–3.

- [23] Zoubin Ghahramani. "Unsupervised learning". In: *Summer School on Machine Learning*. Springer. 2003, pp. 72–112.
- [24] Xiaojin Jerry Zhu. *Semi-supervised learning literature survey*. Tech. rep. University of Wisconsin-Madison Department of Computer Sciences, 2005.
- [25] Xiaojin Zhu and Andrew B Goldberg. "Introduction to semi-supervised learning". In: *Synthesis lectures on artificial intelligence and machine learning* 3.1 (2009), pp. 1–130.
- [26] Zhi-Hua Zhou. "A brief introduction to weakly supervised learning". In: *National Science Review* 5.1 (2017), pp. 44–53.
- [27] Mostafa Dehghani et al. "Neural ranking models with weak supervision". In: *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM. 2017, pp. 65–74.
- [28] Jérôme Rony et al. "Deep weakly-supervised learning methods for classification and localization in histology images: a survey". In: *arXiv preprint arXiv:1909.03354* (2019).
- [29] Earl B Hunt, Janet Marin, and Philip J Stone. "Experiments in induction". In: (1966).
- [30] J Ross Quinlan. "Simplifying decision trees". In: *International Journal of Human-Computer Studies* 51.2 (1999), pp. 497–510.
- [31] Harris Drucker. "Improving regressors using boosting techniques". In: *ICML*. Vol. 97. 1997, pp. 107–115.
- [32] Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- [33] Neeraj Bhargava et al. "Decision tree analysis on j48 algorithm for data mining". In: *Proceedings of International Journal of Advanced Research in Computer Science and Software Engineering* 3.6 (2013).
- [34] S Rasoul Safavian and David Landgrebe. "A survey of decision tree classifier methodology". In: *IEEE transactions on systems, man, and cybernetics* 21.3 (1991), pp. 660–674.
- [35] Jerome H Friedman. "Greedy function approximation: a gradient boosting machine". In: *Annals of statistics* (2001), pp. 1189–1232.

- [36] Wikipedia. *Gradient Boosting*. 2019. URL: [https://en.wikipedia.org/wiki/Gradient\\_boosting](https://en.wikipedia.org/wiki/Gradient_boosting) (visited on 11/04/2019).
- [37] Patrick T Komiske, Eric M Metodiev, and Matthew D Schwartz. "Deep learning in color: towards automated quark/gluon jet discrimination". In: *Journal of High Energy Physics* 2017.1 (2017), p. 110.
- [38] Michael A Nielsen. *Neural networks and deep learning*. Vol. 25. Determination press San Francisco, CA, USA: 2015.
- [39] Warren S Sarle. "Neural networks and statistical models". In: (1994).
- [40] Jürgen Schmidhuber. "Deep learning in neural networks: An overview". In: *Neural networks* 61 (2015), pp. 85–117.
- [41] Rene Vidal et al. "Mathematics of deep learning". In: *arXiv preprint arXiv:1712.04741* (2017).
- [42] UC Business Analytics R Programming Guide. *The Standard Model*. 2019. URL: [http://uc-r.github.io/feedforward\\_DNN](http://uc-r.github.io/feedforward_DNN) (visited on 11/04/2019).
- [43] S Haykin et al. "Neural networks and learning machines. vol. 3 Pearson". In: *Upper Saddle River, NJ, USA* (2009).
- [44] Forest Agostinelli et al. "Learning activation functions to improve deep neural networks". In: *arXiv preprint arXiv:1412.6830* (2014).
- [45] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. "Deep sparse rectifier neural networks". In: *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. 2011, pp. 315–323.
- [46] Vinod Nair and Geoffrey E Hinton. "Rectified linear units improve restricted boltzmann machines". In: *Proceedings of the 27th international conference on machine learning (ICML-10)*. 2010, pp. 807–814.
- [47] wikipedia. *Vanishing Gradient Problem*. 2018. URL: [https://en.wikipedia.org/wiki/Vanishing\\_gradient\\_problem](https://en.wikipedia.org/wiki/Vanishing_gradient_problem) (visited on 11/06/2019).
- [48] Tom Fawcett. "An introduction to ROC analysis". In: *Pattern recognition letters* 27.8 (2006), pp. 861–874.
- [49] David Martin Powers. "Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation". In: (2011).

- [50] Jack H Collins, Kiel Howe, and Benjamin Nachman. “Extending the search for new resonances with machine learning”. In: *Physical Review D* 99.1 (2019), p. 014038.
- [51] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. “Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning”. In: *International conference on intelligent computing*. Springer. 2005, pp. 878–887.
- [52] Lucio Mwinmaarong Dery et al. “Weakly supervised classification in high energy physics”. In: *Journal of High Energy Physics* 2017.5 (2017), p. 145.
- [53] Eric M Metodiev, Benjamin Nachman, and Jesse Thaler. “Classification without labels: Learning from mixed samples in high energy physics”. In: *Journal of High Energy Physics* 2017.10 (2017), p. 174.
- [54] Johan Alwall et al. “The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations”. In: *Journal of High Energy Physics* 2014.7 (2014), p. 79.
- [55] Torbjörn Sjöstrand, Stephen Mrenna, and Peter Skands. “A brief introduction to PYTHIA 8.1”. In: *Computer Physics Communications* 178.11 (2008), pp. 852–867.
- [56] David J Lange. “The EvtGen particle decay simulation package”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 462.1-2 (2001), pp. 152–155.