

1      **Using Machine Learning Techniques for Data Quality**  
2      **Monitoring at CMS Experiment**

3      by

4      Guillermo A. Fidalgo Rodríguez

5      A thesis presented for the degree of

6      BACHELLOR'S OF SCIENCE??

7      in

8      Physics

9      UNIVERSITY OF PUERTO RICO  
10     MAYAGÜEZ CAMPUS

11     2018

12     Approved by:

13     \_\_\_\_\_  
14     Sudhir Malik, Ph.D.

15     President, Graduate Committee  
Date

16     \_\_\_\_\_  
17     Héctor Méndez, Ph.D.

18     Member, Graduate Committee  
Date

19     \_\_\_\_\_  
20     Samuel Santana Colón, Ph.D.  
21     Member, Graduate Committee

Date

22     \_\_\_\_\_  
23     Rafael A. Ramos, Ph.D.  
24     Chairperson of the Department

Date

# <sup>25</sup> Abstract

<sup>26</sup> The Data Quality Monitoring (DQM) of CMS is a key asset to deliver high-quality  
<sup>27</sup> data for physics analysis and it is used both in the online and offline environment. The cur-  
<sup>28</sup> rent paradigm of the quality assessment is labor intensive and it is based on the scrutiny of  
<sup>29</sup> a large number of histograms by detector experts comparing them with a reference. This  
<sup>30</sup> project aims at applying recent progress in Machine Learning techniques to the automa-  
<sup>31</sup> tion of the DQM scrutiny. In particular the use of convolutional neural networks to spot  
<sup>32</sup> problems in the acquired data is presented with particular attention to semi-supervised  
<sup>33</sup> models (e.g. autoencoders) to define a classification strategy that doesn't assume previ-  
<sup>34</sup>ous knowledge of failure modes. Real data from the hadron calorimeter of CMS are used  
<sup>35</sup> to demonstrate the effectiveness of the proposed approach.

<sup>36</sup> *Keywords:* [DQM, online, offline, Machine Learning ]

<sup>37</sup> **Acknowledgments**

<sup>38</sup> I wish to thank United States State Department and University of Michigan for pro-  
<sup>39</sup> viding the opportunity to work abroad at CERN during the 2018 Winter Semester. I also  
<sup>40</sup> wish to thank CERN staff, CMS Experiment , Texas Tech University, and the University  
<sup>41</sup> of Puerto Rico at Mayagüez, with special thanks to Dr. Federico de Guio for his local  
<sup>42</sup> mentorship and Dr. Nural Akchurin, and Dr. Sudhir Malik for their guidance. A very  
<sup>43</sup> special thanks to Dr. Jean Krisch for accepting me for this great research opportunity and  
<sup>44</sup> Dr. Steven Goldfarb for being a wonderful host and overall local guidance at CERN.

# <sup>45</sup> List of Figures

46	2.1 CMS Detector . . . . .	4
47	2.2 The trajectory of a particle traveling through the layers of the detector leaving behind it's signature footprint . . . . .	5
49	4.1 Occupancy maps with 5x5 affected regions . . . . .	11
50	4.2 Weights and Biases . . . . .	12
51	4.3 Gradient Descent algorithm . . . . .	13
52	4.4 Loss Function surface . . . . .	14
53	5.1 Occupancy Maps with 1x1 bad regions. A) Good image B) Dead image C) Hot image . . . . .	15
55	5.2 Two Convolutional Layers Model . . . . .	17
56	5.3 Confusion Matrix results and Learning curve for $5 \times 5$ damaged area with on the same location for all trials . . . . .	17
58	5.4 Confusion Matrix results and Learning curve for $5 \times 5$ damaged area with on the random location for all trials . . . . .	18
60	5.5 Confusion Matrix results and Learning curve for $5 \times 5$ damaged area with an extra class to identify with random location for all trials . . . . .	18
62	5.6 Confusion Matrix results and Learning curve for $1 \times 1$ damaged area with random location for all trials . . . . .	19
64	5.7 Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics. . . . .	20
66	5.8 $Z \rightarrow \nu\bar{\nu}$ background prediction for all search bins, including the break- down of the various uncertainties. . . . .	21

# 68 **Contents**

69	<b>Abstract</b>	i
70	<b>Acknowledgments</b>	ii
71	<b>List of Figures</b>	iii
72	<b>1 Introduction</b>	1
73	<b>2 The CMS Experiment</b>	3
74	<b>3 Data Collection and Data Quality Monitoring</b>	6
75	3.1 What is Data Collection for CMS? . . . . .	6
76	3.2 What is Data Quality Monitoring? . . . . .	7
77	<b>4 What is Machine Learning?</b>	9
78	4.1 Developing the Algorithm . . . . .	10
79	4.2 Teaching the Algorithm . . . . .	12
80	<b>5 Results</b>	15
81	5.1 SL Models for known anomalies in the HCAL data for DQM . . . . .	16
82	5.1.1 Two Convolutional Layers for binary classification . . . . .	16
83	5.1.2 Systematics . . . . .	20
84	5.1.3 $Z \rightarrow \nu\bar{\nu}$ Estimation for the Search Bins . . . . .	21
85	<b>6 References</b>	22

<sup>86</sup> **Chapter 1**

<sup>87</sup> **Introduction**

<sup>88</sup> The work for this thesis was performed at CERN on CMS Experiment. CERN stands  
<sup>89</sup> for European Organization for Nuclear Research. It was founded in 1954 and is located  
<sup>90</sup> at the Franco-Swiss border near Geneva. At CERN, physicists and engineers are probing  
<sup>91</sup> the fundamental structure of the universe. They use the world's largest and most complex  
<sup>92</sup> scientific instruments to study the basic constituents of matter – the fundamental parti-  
<sup>93</sup> cles. The instruments used at CERN are purpose-built particle accelerators and detectors.  
<sup>94</sup> Accelerators boost beams of particles to high energies before the beams are made to col-  
<sup>95</sup> lide with each other or with stationary targets. Detectors observe and record the results  
<sup>96</sup> of these collisions. The accelerator at CERN is called the Large Hadron Collider (LHC),  
<sup>97</sup> the largest machine ever built by humans and it collides particles (protons) at close to the  
<sup>98</sup> speed of light. The process gives the physicists clues about how the particles interact, and  
<sup>99</sup> provides insights into the fundamental laws of nature. Seven experiments at the LHC use  
<sup>100</sup> detectors to analyze particles produced by proton-proton collisions. The biggest of these  
<sup>101</sup> experiments, ATLAS and CMS, use general-purpose detectors designed to study the fun-  
<sup>102</sup> damental nature of matter and fundamental forces and to look for new physics or evidence  
<sup>103</sup> of particles that are beyond the Standard Model. Having two independently designed de-  
<sup>104</sup> tectors is vital for cross-confirmation of any new discoveries made. The other two major  
<sup>105</sup> detectors ALICE and LHCb, respectively, study a state of matter that was present just  
<sup>106</sup> moments after the Big Bang and preponderance of matter than antimatter. Each experi-

107 ment does important research that is key to understanding the universe that surrounds and  
108 makes us.

109

110      [Chapter 2](#) presents a basic description of the Large Hadron Collider and CMS Detector

111

112      ?? gives a brief motivation

113

114      ?? is dedicated to a study optimizing

115

116      ?? ptimated.

117

118      ?? details an improvarger production cross-section than Z+jets process used before.

119

120      The conclusions and results of each chapter are presented in the corresponding chap-  
121 ter.

122

123      This thesis work has been presented at several internal meetings of the CMS Experi-  
124 ment and at the following international meetings and conferences:

125      1. **Andrés Abreu** gave a talk “*Estimation of the Z Invisible Background for Searches*

126            *for Supersymmetry in the All-Hadronic Channel*” at “APS April 2018: American  
127            Physical Society April Meeting 2018, 14-17 Apr 2018”, Columbus, OH

128      2. **Andrés Abreu** gave a talk “*Phase-2 Pixel upgrade simulations*” at the “USLUA

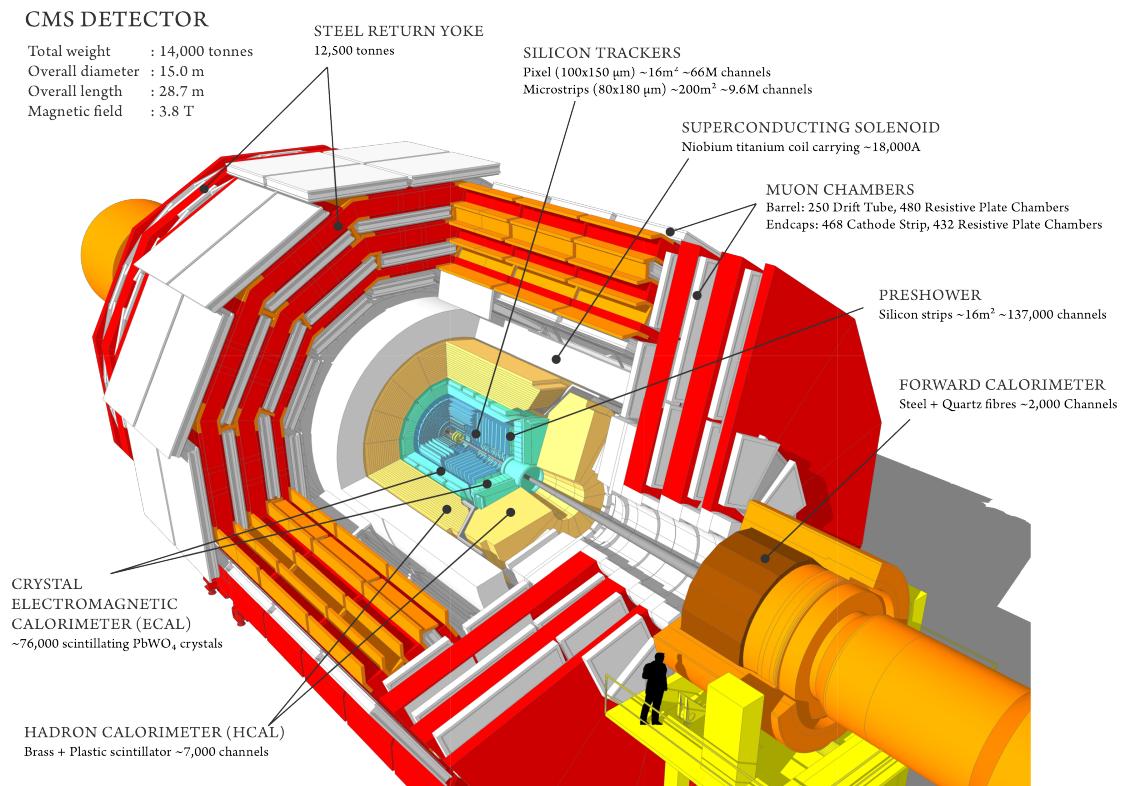
129            Annual meeting: 2017 US LHC Users Association Meeting, 1-3 Nov 2017”, Fer-  
130            milab, Batavia, IL

<sup>131</sup> **Chapter 2**

<sup>132</sup> **The CMS Experiment**

<sup>133</sup> The Compact Muon Solenoid (CMS) detector is a general purpose particle detector  
<sup>134</sup> designed to investigate various physical phenomena concerning the SM and beyond it,  
<sup>135</sup> such as Supersymmetry, Extra Dimensions and Dark Matter. As its name implies, the  
<sup>136</sup> detector is a solenoid which is constructed around a superconducting magnet capable of  
<sup>137</sup> producing a magnetic field of 3.8 T. The magnetic coil is 13m long with an inner diameter  
<sup>138</sup> of 6m, making it the largest superconducting magnet ever constructed. The CMS detector  
<sup>139</sup> itself is 21m long with a diameter of 15m and it has a weight of approximately 14,000  
<sup>140</sup> tons. The CMS experiment is one of the largest scientific collaborations in the history  
<sup>141</sup> of mankind with over 4,000 participants from 42 countries and 182 institutions. CMS is  
<sup>142</sup> located at one of these points and it essentially acts as a giant super highspeed camera  
<sup>143</sup> that makes 3D images of the collisions that are produced at a rate of 40 MHz (40 million  
<sup>144</sup> times per second). The detector has an onion-like structure to capture all the particles that  
<sup>145</sup> are produced in these high energy collisions most of them being unstable and decaying  
<sup>146</sup> further to stable particles that are detected. CMS detector was designed with the following  
<sup>147</sup> features (as shown in [Figure 2.1](#)) :

- <sup>148</sup> 1. A **magnet** with large bending power and high performance muon detector for good  
<sup>149</sup> muon identification and momentum resolution over a wide range of momenta and  
<sup>150</sup> angles.

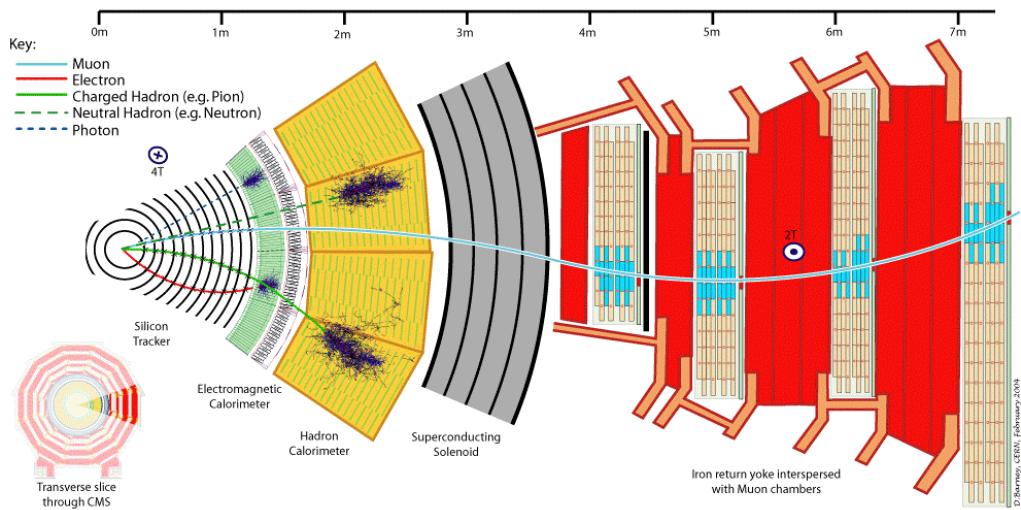


**Figure 2.1:** CMS Detector

- 151     2. An **inner tracking system** capable of high reconstruction efficiency and momen-  
152       tum resolution requiring **pixel detectors** close to the interaction region.
- 153     3. An **electromagnetic calorimeter** able to provide good electromagnetic energy res-  
154       olution and a high isolation efficiency for photons and leptons.
- 155     4. A **hadron calorimeter** capable of providing precise missing-transverse-energy and  
156       dijet-mass resolution.

157       A property from these particles that is exploited is their charge. Normally, particles  
158       produced in collisions travel in a straight line, but in the presence of a magnetic field,  
159       their paths are skewed and curved. Except the muon system, the rest of the subdetectors  
160       lie inside a 3.8 Tesla magnetic field . Due to the magnetic field the trajectory of charged  
161       particle produced in the collisions gets curved (as shown in [Figure 2.2](#) ) and one can  
162       calculate the particle's momentum and know the type of charge on the particle. The  
163       Tracking devices are responsible for drawing the trajectory of the particles by using a  
164       computer program that reconstructs the path by using electrical signals that are left by

the particle as they move. The Calorimeters measure the energy of particles that pass through them by absorbing their energy with the intent of stopping them. The particle identification detectors work by detecting radiation emitted by charged particles and using this information they can measure the speed, momentum, and mass of a particle. After the information is put together to make the “snapshot” of the collision one looks for results that do not fit the current theories or models in order to look for new physics.



**Figure 2.2:** The trajectory of a particle traveling through the layers of the detector leaving behind its signature footprint

The project focusses specifically on data collected from one of the Calorimeters, - the Hadron Calorimeter (HCAL). The HCAL, as its name indicates, is designed to detect and measure the energy of hadrons or, particles that are composed of quarks and gluons, like protons and neutrons. Additionally, it provides an indirect measurement of the presence of non-interacting, uncharged particles such as neutrinos (missing energy) . Measuring these particles is important as they can tell us if new particles such as the Higgs boson or supersymmetric particles (much heavier versions of the standard particles we know) have been formed. The layers of the HCAL are structured in a staggered fashion to prevent any gaps that a particle might pass through undetected. There are two main parts: the barrel and the end caps. There are 36 barrel wedges that form the last layer of the detector inside the magnet coil, there is another layer outside this, and on the endcaps, there are another 36 wedges to detect particles that come out at shallow angles with respect to the beam line.

# <sup>184</sup> Chapter 3

## <sup>185</sup> Data Collection and Data Quality

### <sup>186</sup> Monitoring

#### <sup>187</sup> 3.1 What is Data Collection for CMS?

<sup>188</sup> During data taking there are millions of collisions occurring in the center of the detec-  
<sup>189</sup> tor every second. The data per event is around one million bytes (1 MB), that is produced  
<sup>190</sup> at a rate of about 600 million events per second [1], that's about 600 MB/s. Keeping  
<sup>191</sup> in mind that only certain events are considered "interesting" for analysis, the task of de-  
<sup>192</sup> ciding what events to consider out of all the data collected is a two-stage process. First,  
<sup>193</sup> the events are filtered down to 100 thousand events per second for digital reconstruction  
<sup>194</sup> and then more specialized algorithms filter the data even more to around 100 200 events  
<sup>195</sup> per second that are found interesting. For CMS there is a Data Acquisition System that  
<sup>196</sup> records the raw data to what's called a High-Level Trigger farm which is a room full  
<sup>197</sup> of servers that are dedicated to processing and classify this raw data quickly. The data  
<sup>198</sup> then gets sent to what's known as the Tier-0 farm where the full processing and the first  
<sup>199</sup> reconstruction of the data are done. [2]

## 200 3.2 What is Data Quality Monitoring?

201 To operate a sophisticated and complex apparatus as CMS, a quick online feedback on  
202 the quality of the data recorded is needed to avoid taking low quality data and to guarantee  
203 a good baseline for the offline analysis. Collecting a good data sets from the collisions  
204 is an important step towards search for new physics as deluge of new data poses an extra  
205 challenge of processing and storage. This all makes it all the more important to design  
206 algorithms and special software to control the quality of the data. This is where the Data  
207 Quality Monitoring (DQM) plays a critical in the maintainability of the experiment, the  
208 operation efficiency and performs a reliable data certification. The high-level goal of  
209 the system is to discover and pinpoint errors, problems occurring in detector hardware  
210 or reconstruction software, early, with sufficient accuracy and clarity to maintain good  
211 detector and operation efficiency. The DQM workflow consists of 2 types: **Online** and  
212 **Offline**.

213 The **Online** DQM consists of receiving data taken from the event and trigger his-  
214 tograms to produce results in the form of monitoring elements like histogram references  
215 and quality reports. This live monitoring of each detector's status during data taking gives  
216 the online crew the possibility to identify problems with extremely low latency, mini-  
217 mizing the amount of data that would otherwise be unsuitable for physics analysis. The  
218 scrutiny of the Online DQM is a 24/7 job that consists of people or shifters that work at the  
219 CMS control center constantly monitoring the hundreds of different plots and histograms  
220 produced by the DQM software. This consumes a lot of manpower and is strenuous work.

221 The **Offline** DQM is more focused on the full statistics over the entire run of the  
222 experiment and works more on the data certification. In the offline environment, the  
223 system is used to review the results of the final data reconstruction on a run-by-run basis,  
224 serving as the basis for certified data used across the CMS collaboration in all physics  
225 analyses. In addition, the DQM framework is an integral part of the prompt calibration  
226 loop. This is a specialized workflow run before the data are reconstructed to compute and  
227 validate the most up-to-date set of conditions and calibrations subsequently used during

228 the prompt reconstruction.

229 This project aims to minimize the DQM scrutiny by eye and automate the process so  
230 that there is a more efficient process to monitor the detector and the quality of the data by  
231 implementing Machine Learning techniques.

<sup>232</sup> **Chapter 4**

<sup>233</sup> **What is Machine Learning?**

<sup>234</sup> Machine Learning (ML) can be defined as an application of Artificial Intelligence that  
<sup>235</sup> permits the computer system to learn without being told explicitly. In ML a computer  
<sup>236</sup> program is said to learn from experience E with respect to some class of tasks T and  
<sup>237</sup> performance measure P, if its performance at tasks in T, as measured by P, improves  
<sup>238</sup> with experience E [3]. ML has made tremendous strides in the past decades and has  
<sup>239</sup> become very popular recently due to its multifaceted applications. It is being used on  
<sup>240</sup> social media, marketing, and in the scientific community as well. Some examples of  
<sup>241</sup> ML applications are: the algorithms used on application in smartphones to detect human  
<sup>242</sup> faces, self-driving cars, computer games, stock prediction, and voice recognition. An  
<sup>243</sup> interesting characteristic of ML algorithms is that the more data one inputs the better is  
<sup>244</sup> the performance. The ML application has a very wide spectrum covering almost every  
<sup>245</sup> aspect of human endeavor that involves a lot of data. Scientific analysis today generates  
<sup>246</sup> enormous data and is hence a perfect use case to apply ML techniques. In this work  
<sup>247</sup> we use enhanced ML techniques based on progress in the recent past.

<sup>248</sup> In general, there are two main categories to classify machine learning problems: **Su-**  
<sup>249</sup> **pervised Learning** (SL) and **Unsupervised Learning** (UL). SL is the most used ML  
<sup>250</sup> approach and has proven to be very effective for a wide variety of problems. Examples  
<sup>251</sup> of common SL problems are: spam filters, predicting housing prices, identifying a ma-  
<sup>252</sup> lignant or benign tumor, etc. These types of problems are characterized by providing a

253 “right answer” as a reference. For example, spam filter algorithms identify emails that  
254 are spams by training on a dataset that has examples of such emails. In case of predicting  
255 house prices, the algorithm is trained on a dataset of houses involving features like the  
256 area of the house, number of rooms, and the selling price of the house.

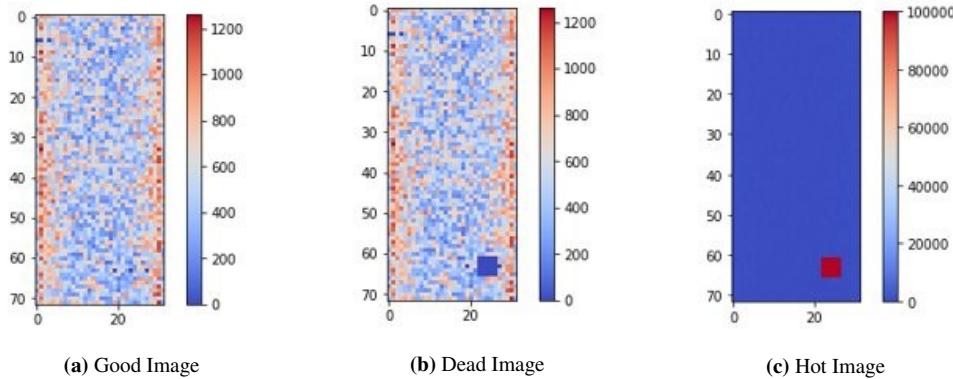
257 UL algorithms are different in the sense that they do not have the “right answers”  
258 given to the machine. Instead, UL algorithms are used for finding patterns and make  
259 clusters from the given data. That is what also forms the basis of a search engine (e.g.  
260 Google news). Clicking on a link to a news article, one gets many different stories of  
261 different journals that have some correlation with the article searched. This happens be-  
262 cause the ML algorithm is capable of learning features and shared patterns from a bunch  
263 of data without being given any specifics. Another interesting UL problem is the so-called  
264 “cocktail party” that involves distinguishing the voice of two people recording on two mi-  
265 crophones located at different places. The ML algorithm is able to separate the sources of  
266 the voices in the recordings by learning the voice features that correspond to each person,  
267 showing the power of the UL algorithm.

268 In this study, I have focused on an SL approach and a variant of the UL approach,  
269 called the **Semi-Supervised Learning** approach (SSL). The SSL is named so because  
270 the data involves looking at images that are already known to be “Good” but one doesn’t  
271 necessarily know every possible situation that produces a “Bad” image. The purpose is to  
272 define a metric for a “good” image and subsequently decide if an image is “bad” in case  
273 it deviates too much from an acceptable value.

## 274 4.1 Developing the Algorithm

275 To develop an ML algorithm the following are taken into consideration, what is the  
276 task? and what is the method to approach the task? In our case, we are looking into images  
277 that have information about the activity that the channels in the HCAL are detecting.  
278 These images are called ”occupancy maps” and they are a visual way of monitoring the  
279 health of the detector itself (see [Figure 4.1](#)). There are two common problems that can be

identified by viewing occupancy maps which are called "dead channels" and "hot towers". They are referred to as "**dead**" and "**hot**" respectively in the rest of this document. Dead channels mean that on a certain place in the occupancy map there is not any readout from the channels on the HCAL and hot channels mean that there are channels that are being triggered by noise or are damaged in a way that makes them readout too much activity.



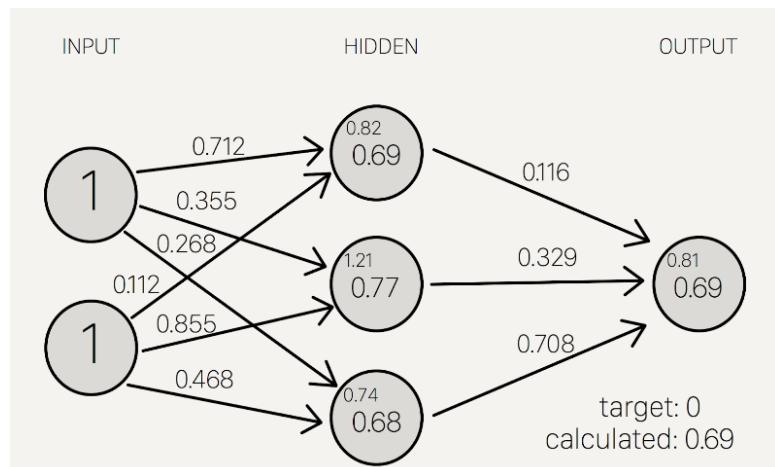
**Figure 4.1:** Occupancy maps with 5x5 affected regions

The problem is the following, to create a model that can detect and classify what type of scenario is occurring on each occupancy map. For this, we want to go with a SL approach which means that we will give the model the images as the input and it will train on these images by learning to identify patterns or features in the image and try to do a "fit" from the images to their corresponding labels. After the training, the algorithm will be given a testing set for us to evaluate the model's ability to correctly detect if there is a problem with the image and what type of problem is being detected. The output of the model will be the predicted class of the test image. The predictions are based on the labels and their corresponding images that were given to the model during training. This means that if the model was trained with 3 different types of images with their corresponding label the model will only work well for images that present similar patterns or characteristics to those presented in the training. For example, if we only train the model to distinguish between "good" and "hot" then when the model encounters images that aren't either of these two, like an image labeled "dead", then the model will not know what to do with this image and will give it an incorrect label. After the SL model has been tested the next step is trying an SSL model. The term semi-supervised

simply means that there isn't a ground truth label that is being given to the model during training because either there isn't necessarily a ground truth, or we don't know what the ground truth is. What we do know, is what is considered as a "good" image and what this approach hopes to accomplish is to use the error in the reconstruction of the input image and use that information to discriminate between the "good" vs the "bad" images.

## 4.2 Teaching the Algorithm

The way an ML algorithm learns is by an iterative process called an optimization algorithm in which the predicted output value of the model is compared to the desired output (See [Figure 4.2](#)) and the weights and biases of the model are adjusted such that the predicted output is closer to the desired output.



**Figure 4.2:** Weights and Biases

"Optimization algorithms helps us to **minimize** (or **maximize**) an **Objective function** (*another name for Error function*)  $E(x)$  which is simply a mathematical function dependent on the Model's internal **learnable parameters** which are used in computing the target values(Y) from the set of *predictors*(X) used in the model. For example - we call the **Weights(W)** and the **Bias(b)** values of the neural network as its internal learnable *parameters* which are used in computing the output values and are learned and updated in the direction of optimal solution i.e. minimizing the **Loss** by the network's training process and also play a major role in the **training** process of the Neural Network Model." [4].

## Gradient Descent

The “Learning” in Machine Learning.

Update the values of X (punish) it when it is wrong.

$$X = X - \eta \nabla(X)$$

X: weights or biases

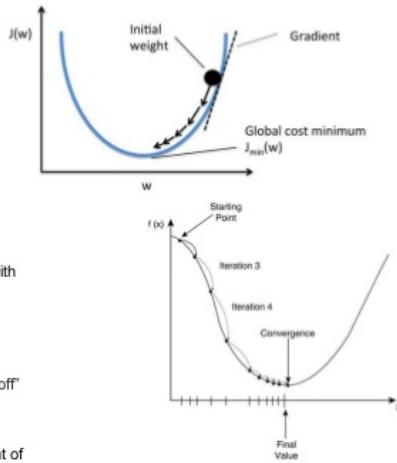
$\eta$ : Learning Rate (typically 0.01 to 0.001)

$\eta$  :The rate at which our network learns. This can change over time with methods such as Adam, Adagrad etc.  (hyperparameter)

$\nabla(X)$ : Gradient of X

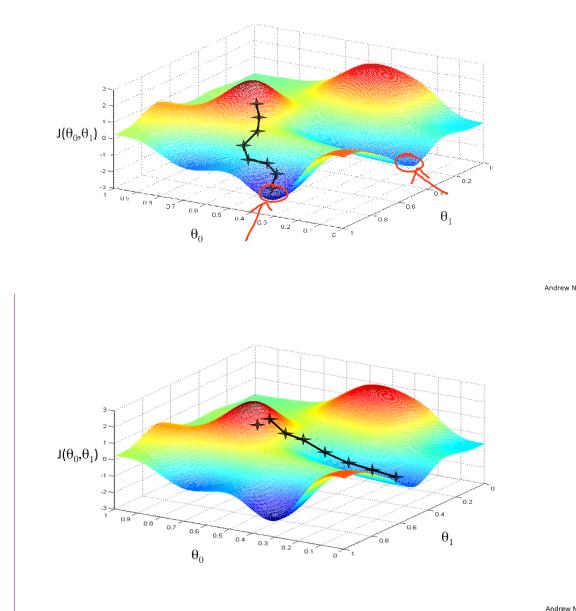
We seek to update the weights and biases by a value indicating how “off” they were from their target.

Gradients naturally have increasing slope, so we put a negative in front of it to go downwards



**Figure 4.3:** Gradient Descent algorithm

- 319 The most basic and probably the most used optimizer is called Gradient Descent (GD).
- 320 GD is based on the concept of using the gradient of a loss or cost function and moving
- 321 the weights and biases of the ML model so that the predicted value is taking a step in the
- 322 decreasing direction of this error function (See [Figure 4.3](#)). In general, the “terrain” of the
- 323 loss function is not a smooth bowl-shaped surface like the one present in the image. The
- 324 most general form of the surface is more similar to a rocky mountain (See [Figure 4.4](#)),
- 325 which presents a problem when using simple optimizers like GD.



**Figure 4.4:** Loss Function surface

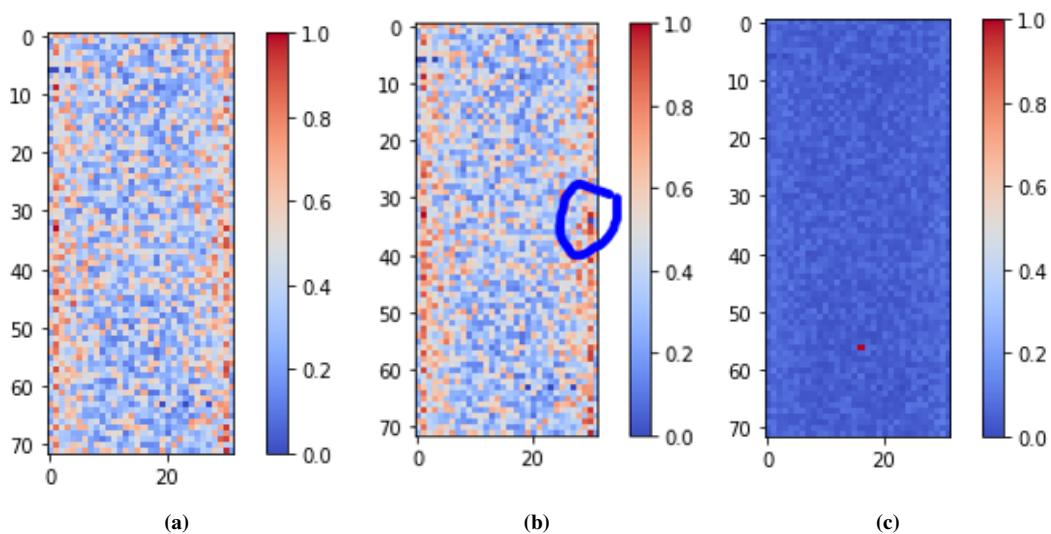
326 

# Chapter 5

327 

## Results

328 Here first the limitations of Scikit-learn predefined ML models - Logistic Regres-  
329 sion(LR) and Multi-Layer-Perceptron(MLP), are described. The Logistic Regression  
330 Model seems to work almost perfectly with all 3 classes when the bad region size is  
331  $5 \times 5$  (as in [Figure 4.1](#)) with either the same or randomized location. When the bad region  
332 size is  $1 \times 1$  like in [Figure 5.1](#) the LR Model performs poorly with an accuracy of approxi-  
333 mately 20%. The MLP does not seem to work in any of the used cases that are studied as  
334 it always performs poorly with an accuracy of  $\approx 40\%$ .



**Figure 5.1:** Occupancy Maps with  $1 \times 1$  bad regions. A) Good image B) Dead image C) Hot image

335 Also, the use of Scikit-learn's library is limited in comparison with the Keras module  
336 since one cannot customize the structure of the ML model with detail. Moreover, Keras is

337 an ML library designed for developing deep neural networks. Hence it was decided to use  
338 Keras primarily for the creation of the model. With the Keras library, numerous models  
339 were designed with both, SL method and SSL learning method. Using SL method, we are  
340 interested in detecting anomalies and classifying what type of anomaly is seen. With SSL  
341 method, we are interested in looking at the error of the reconstruction of an image to give  
342 an idea that the image given can be considered good or that it might have some unseen  
343 anomalies

344 **5.1 SL Models for known anomalies in the HCAL data  
345 for DQM**

346 We considered three SL Models for classification of known anomalies in the HCAL  
347 data for DQM. These models are based on Convolutional Neural Networks and differ  
348 in the number of layers utilized, their ordering and number of units in each layer. The  
349 Models and the corresponding results are described below.

350 **5.1.1 Two Convolutional Layers for binary classification**

351 Several variations of the two Convolutional Layers Model were tested and optimized  
352 on the DQM data. This led to an optimal value of 8 units/neurons in the Convolutional  
353 layers. The detail of selecting the number of units per layer is of great importance to find  
354 a balance between efficiency and complexity of a model. More complex models (more  
355 layers and connections) are “heavy” to train in terms of computational cost, provide better  
356 results and are prone to “overfitting” to the training data. Simpler models (fewer layers  
357 and connections) are quicker to train, efficient and computationally economic. However,  
358 simpler models are more likely to “underfit” to the data. The [Figure 5.2](#) below shows a  
359 code snippet with this model. [Figure 5.3](#) below shows the learning curve for this model  
360 trained with Good and Hot images for fixed  $5 \times 5$  location and the corresponding Confu-  
361 sion Matrix.

```

model = Sequential([
    BatchNormalization(input_shape=input_shape),
    Conv2D(8, kernel_size=(3, 3), strides=(2, 2), activation='relu'),
    Conv2D(8, kernel_size=(3, 3), strides=(2, 2), activation='relu')

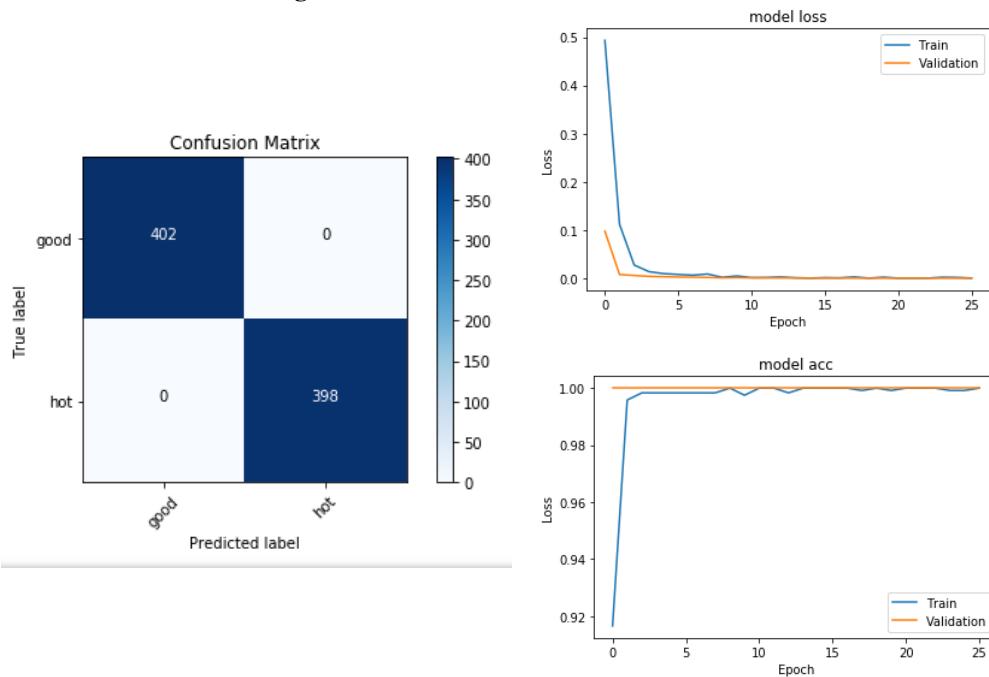
    Dropout(0.25)

    Flatten()

    Dense(2, activation='softmax')
])

```

**Figure 5.2:** Two Convolutional Layers Model



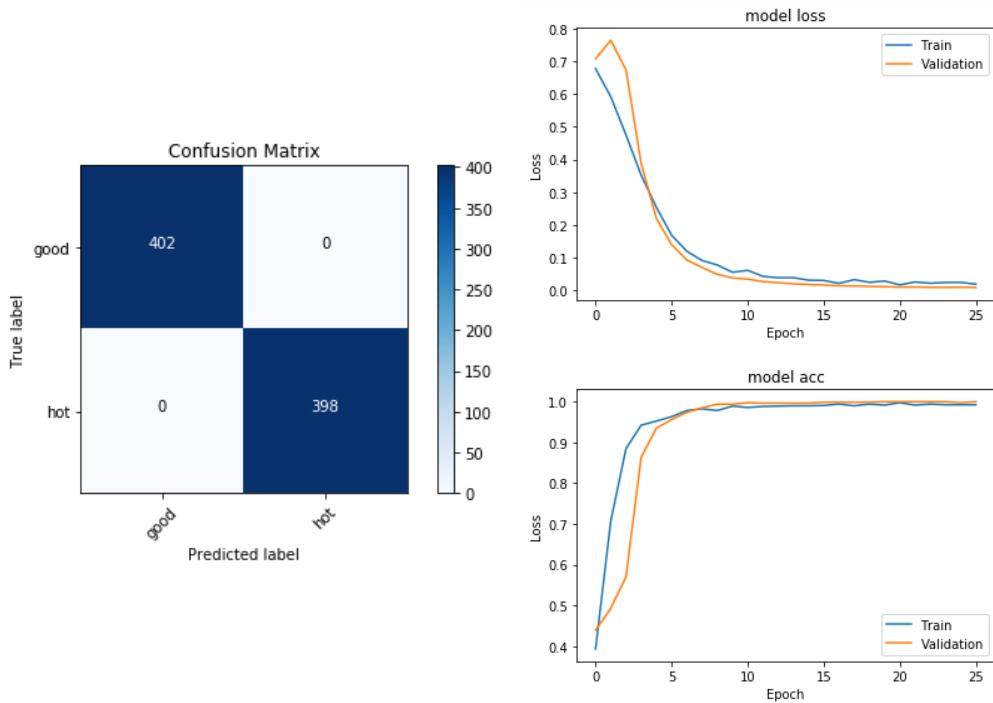
**Figure 5.3:** Confusion Matrix results and Learning curve for  $5 \times 5$  damaged area with on the same location for all trials

362     Figure 5.4 shows the learning curve for this model trained with Good and Hot images  
 363     for fixed  $5 \times 5$  location and the corresponding Confusion Matrix.

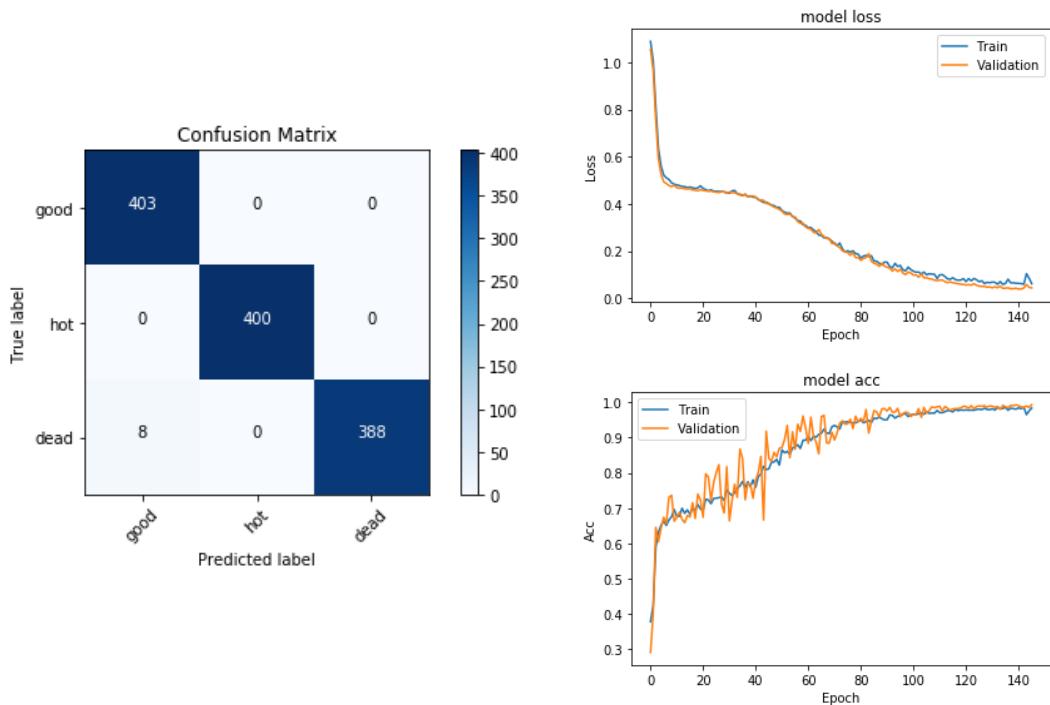
364     Figure 5.5 shows the learning curve for this model trained with Good, Hot and Dead  
 365     images for random  $5 \times 5$  location and the corresponding Confusion Matrix

366     Figure 5.6 shows the learning curve for this model trained with Good, Hot and Dead  
 367     images for random  $1 \times 1$  location and the corresponding Confusion Matrix. The corre-  
 368     sponding learning curves and confusion matrix for a fixed location for 3-class (Good,  
 369     Hot, Dead) configuration give the same behavior as 2-labels (Good, Hot) images

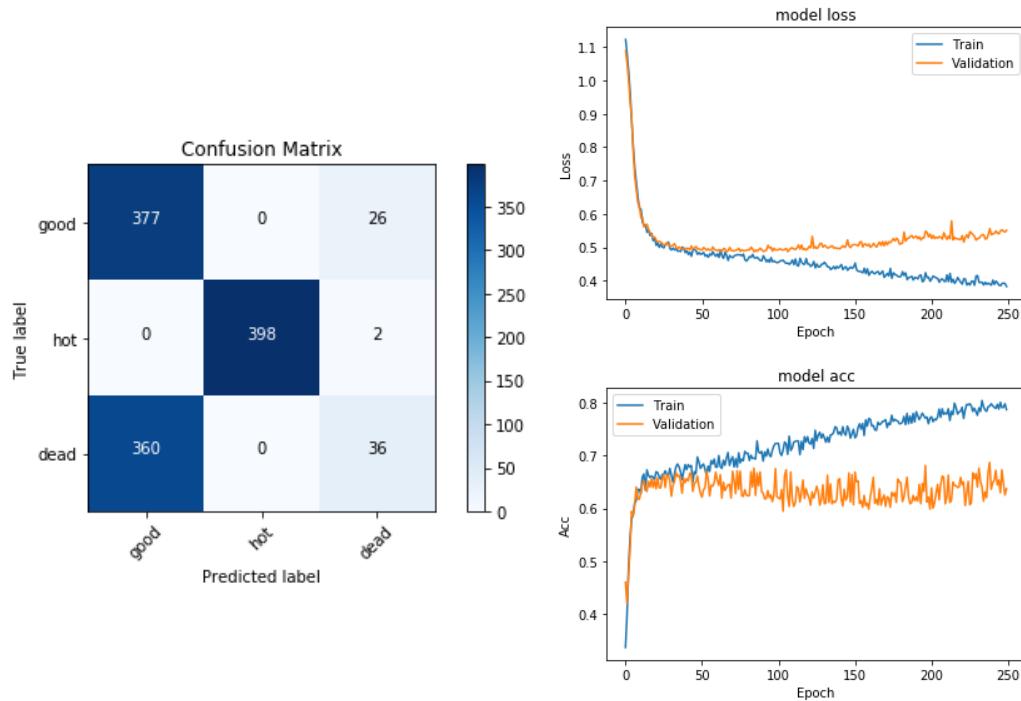
370     In a more realistic scenario, the problems with HCAL DQM would be more granular



**Figure 5.4:** Confusion Matrix results and Learning curve for  $5 \times 5$  damaged area with on the random location for all trials



**Figure 5.5:** Confusion Matrix results and Learning curve for  $5 \times 5$  damaged area with an extra class to identify with random location for all trials



**Figure 5.6:** Confusion Matrix results and Learning curve for  $1 \times 1$  damaged area with random location for all trials

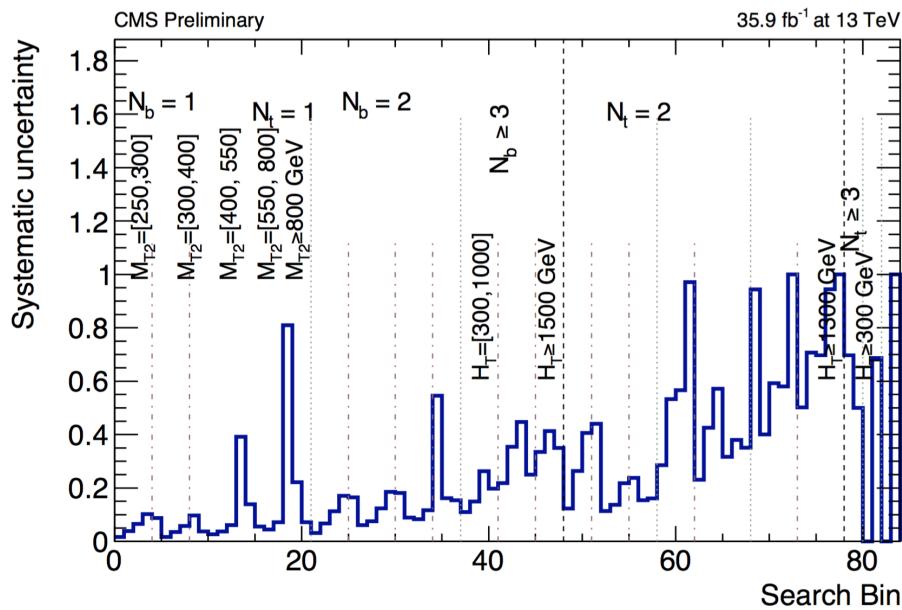
i.e.  $1 \times 1$  type. When this model is tested for problematic channels in  $1 \times 1$  configuration the learning curves for the training (blue) and validation (orange) sets depart after few epochs as shown in Figure 12 (right part). From the left part of the figure, dividing the sum of numbers along the diagonal ( $377+398+36$ ) by the sum of all the numbers in the matrix gives  $1/3$ . This demonstrates that the model is “overfitting” to the training set and misclassifies images  $\approx 33\%$  of times. Hence, we consider adding a Convolutional layer to gain more prediction accuracy as shown in the next section.

In this section the results for the final estimation of the of the  $Z \rightarrow \nu\bar{\nu}$  are presented. The current study includes preliminary results using only data obtained at the CMS detector during 2016. The results for this study are intended to confirm the assumption that the additional  $\gamma +$  jets control region introduced in this analysis reduce the overall uncertainties obtained in the 2016 analyses (described in ??). Furthermore, this study is intended as a benchmark for future analyses of the SUSY stop group based in Fermilab and will be the method used for the 2017 CMS data.

385 **5.1.2 Systematics**

386 Two categories of uncertainties for the  $Z \rightarrow \nu\bar{\nu}$  prediction are considered: uncertainties  
 387 that are associated to the use of MC simulation and the uncertainties specifically  
 388 associated to the background prediction method. Several sources are acknowledged in the  
 389 first category mentioned such as PDF and renormalization/factorization scale choices, jet  
 390 and  $p_T^{miss}$  energy scale uncertainties b-tag scale factor uncertainties, and trigger efficiency  
 391 uncertainties. Given that the simulation sample is normalized to data in the tight control  
 392 region, uncertainties associated with the luminosity and cross-section are excluded. In  
 393 addition, the overall  $Z \rightarrow \nu\bar{\nu}$  statistical uncertainty from MC simulation is also taken into  
 394 account.

395



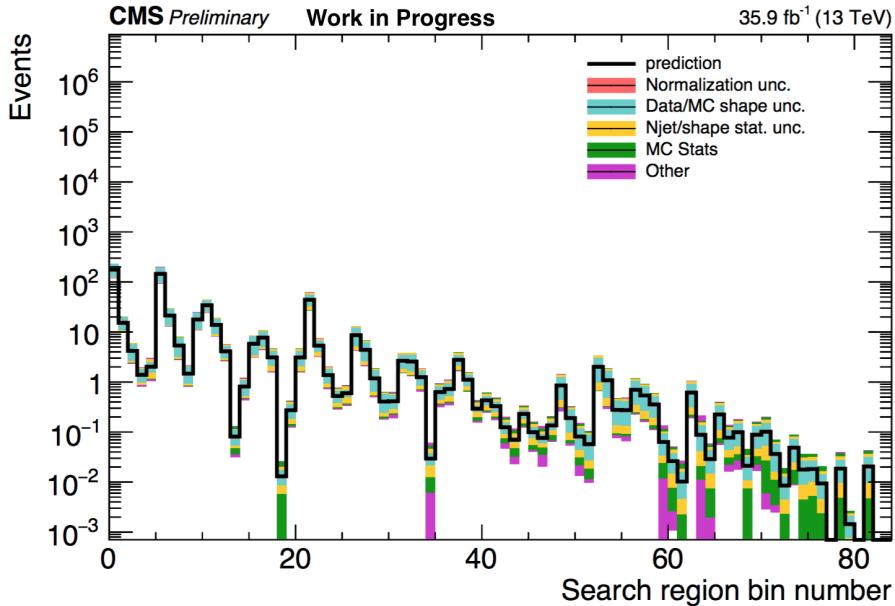
**Figure 5.7:** Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics.

396 The statistical uncertainty associated with each bin in the MC is propagated as a sys-  
 397 tematic uncertainty. The relative uncertainty per bin can be see in [Figure 5.7](#). It shows  
 398 that the uncertainties for the MC vary from as low as 1% up to 81% and even 100% in  
 399 some regions. Since the final estimation is scaled using the global normalization factor  
 400 from the tight  $\mu\mu$  control region ( $R_{norm}$ ), the total uncertainty, due to limited amounts of  
 401 events in data, is propagated in the final prediction. This is also true for the  $S_\gamma(N_j)$  scale

factor, in which the residual differences in search variables other than  $N_j$  are evaluated in the loose photon control region. Both the uncertainty arising from the  $N_j$  re-weighting as well as the residual differences are evaluated together. The uncertainty from  $R_{norm}$  is propagated as a flat value of 7.9% uncertainty per each search bin.

### 5.1.3 $Z \rightarrow \nu\bar{\nu}$ Estimation for the Search Bins

The final estimation for the  $Z \rightarrow \nu\bar{\nu}$  background calculated for all 84 search bins is shown in Figure 5.8. The statistical uncertainty in bins that have zero events is treated as the average weight (the sum of the weights squared over the weight) times the poisson error on 0 which is 1.8. This average weight is calculated on the basis of a relaxed cut in which  $N_b \geq 2$  is required. For comparison, a cut in which  $N_t > 2$  where two tops are fake for the  $Z \rightarrow \nu\bar{\nu}$  is used.



**Figure 5.8:**  $Z \rightarrow \nu\bar{\nu}$  background prediction for all search bins, including the breakdown of the various uncertainties.

<sup>413</sup> **Chapter 6**

<sup>414</sup> **References**

<sup>415</sup> [1] CERN, “Processing what to record?,” 2018.

<sup>416</sup> [2] CMS, “The CMS Computing Project,” tech. rep., CERN, 2005.

<sup>417</sup> [3] Coursera, “Machine learning,” 2018.

<sup>418</sup> [4] A. S. Walia, “Types of optimization algorithms used in neural networks and ways to  
<sup>419</sup> optimize gradient descent,” 2018.