

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

²⁵ **Abstract**

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

³⁶ *Keywords:* [DQM, online, offline, Machine Learning]

³⁷ **Acknowledgments**

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

⁴⁵ List of Figures

| | | |
|----|---|----|
| 46 | 2.1 CMS Detector | 4 |
| 47 | 2.2 The trajectory of a particle traveling through the layers of the detector | |
| 48 | 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 **Contents**

| | | |
|----|--|-----|
| 54 | Abstract | i |
| 55 | Acknowledgments | ii |
| 56 | List of Figures | iii |
| 57 | 1 Introduction | 1 |
| 58 | 2 The CMS Experiment | 3 |
| 59 | 3 Data Collection and Data Quality Monitoring | 6 |
| 60 | 3.1 What is Data Collection for CMS? | 6 |
| 61 | 3.2 What is Data Quality Monitoring? | 7 |
| 62 | 4 What is Machine Learning? | 9 |
| 63 | 4.1 Developing the Algorithm | 10 |
| 64 | 4.2 Teaching the Algorithm | 12 |
| 65 | A Appendix Title | 15 |
| 66 | B References | 16 |

⁶⁷ **Chapter 1**

⁶⁸ **Introduction**

⁶⁹ The work for this thesis was performed at CERN on CMS Experiment. CERN stands
⁷⁰ for European Organization for Nuclear Research. It was founded in 1954 and is located
⁷¹ at the Franco-Swiss border near Geneva. At CERN, physicists and engineers are probing
⁷² the fundamental structure of the universe. They use the world's largest and most complex
⁷³ scientific instruments to study the basic constituents of matter – the fundamental parti-
⁷⁴ cles. The instruments used at CERN are purpose-built particle accelerators and detectors.
⁷⁵ Accelerators boost beams of particles to high energies before the beams are made to col-
⁷⁶ lide with each other or with stationary targets. Detectors observe and record the results
⁷⁷ of these collisions. The accelerator at CERN is called the Large Hadron Collider (LHC),
⁷⁸ the largest machine ever built by humans and it collides particles (protons) at close to the
⁷⁹ speed of light. The process gives the physicists clues about how the particles interact, and
⁸⁰ provides insights into the fundamental laws of nature. Seven experiments at the LHC use
⁸¹ detectors to analyze particles produced by proton-proton collisions. The biggest of these
⁸² experiments, ATLAS and CMS, use general-purpose detectors designed to study the fun-
⁸³ damental nature of matter and fundamental forces and to look for new physics or evidence
⁸⁴ of particles that are beyond the Standard Model. Having two independently designed de-
⁸⁵ tectors is vital for cross-confirmation of any new discoveries made. The other two major
⁸⁶ detectors ALICE and LHCb, respectively, study a state of matter that was present just
⁸⁷ moments after the Big Bang and preponderance of matter than antimatter. Each experi-

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

90

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

92

93 ?? gives a brief motivation

94

95 ?? is dedicated to a study optimizing

96

97 ?? ptimated.

98

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

100

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

103

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

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

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

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

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

¹¹² **Chapter 2**

¹¹³ **The CMS Experiment**

¹¹⁴ The Compact Muon Solenoid (CMS) detector is a general purpose particle detector
¹¹⁵ designed to investigate various physical phenomena concerning the SM and beyond it,
¹¹⁶ such as Supersymmetry, Extra Dimensions and Dark Matter. As its name implies, the
¹¹⁷ detector is a solenoid which is constructed around a superconducting magnet capable of
¹¹⁸ producing a magnetic field of 3.8 T. The magnetic coil is 13m long with an inner diameter
¹¹⁹ of 6m, making it the largest superconducting magnet ever constructed. The CMS detector
¹²⁰ itself is 21m long with a diameter of 15m and it has a weight of approximately 14,000
¹²¹ tons. The CMS experiment is one of the largest scientific collaborations in the history
¹²² of mankind with over 4,000 participants from 42 countries and 182 institutions. CMS is
¹²³ located at one of these points and it essentially acts as a giant super highspeed camera
¹²⁴ that makes 3D images of the collisions that are produced at a rate of 40 MHz (40 million
¹²⁵ times per second). The detector has an onion-like structure to capture all the particles that
¹²⁶ are produced in these high energy collisions most of them being unstable and decaying
¹²⁷ further to stable particles that are detected. CMS detector was designed with the following
¹²⁸ features (as shown in [Figure 2.1](#)) :

- ¹²⁹ 1. A **magnet** with large bending power and high performance muon detector for good
¹³⁰ muon identification and momentum resolution over a wide range of momenta and
¹³¹ angles.

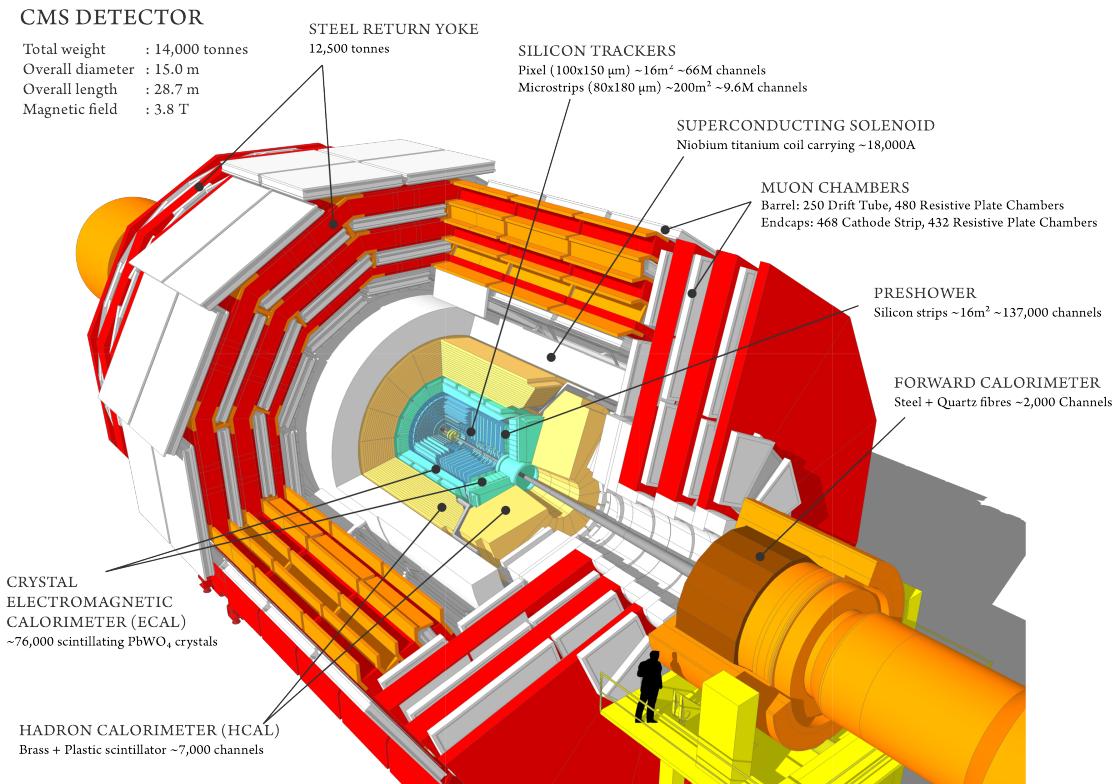


Figure 2.1: CMS Detector

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

138 A property from these particles that is exploited is their charge. Normally, particles
139 produced in collisions travel in a straight line, but in the presence of a magnetic field,
140 their paths are skewed and curved. Except the muon system, the rest of the subdetectors
141 lie inside a 3.8 Tesla magnetic field . Due to the magnetic field the trajectory of charged
142 particle produced in the collisions gets curved (as shown in [Figure 2.2](#)) and one can
143 calculate the particle's momentum and know the type of charge on the particle. The
144 Tracking devices are responsible for drawing the trajectory of the particles by using a
145 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.

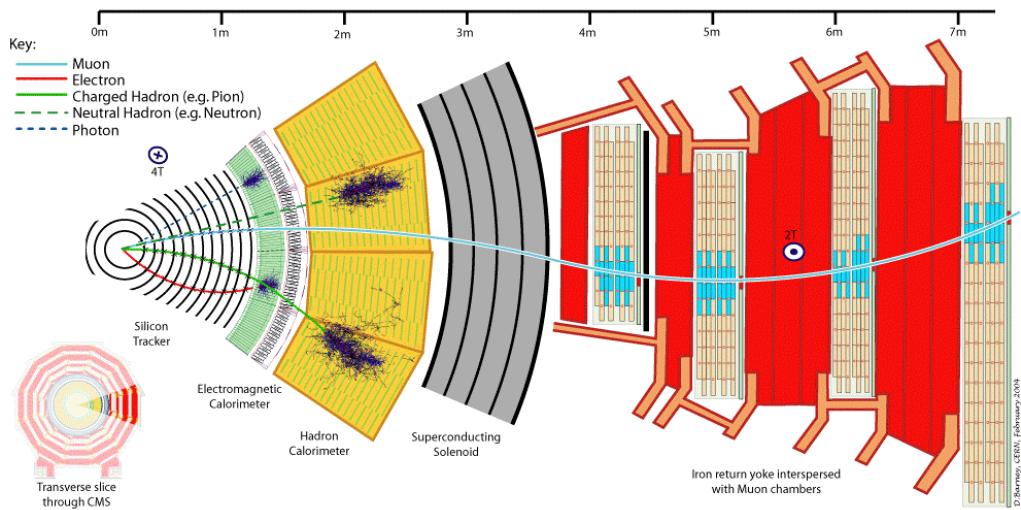


Figure 2.2: The trajectory of a particle traveling through the layers of the detector leaving behind its signature footprint

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.

¹⁶⁵ Chapter 3

¹⁶⁶ Data Collection and Data Quality

¹⁶⁷ Monitoring

¹⁶⁸ 3.1 What is Data Collection for CMS?

¹⁶⁹ During data taking there are millions of collisions occurring in the center of the de-
¹⁷⁰ tector every second. The data per event is around one million bytes (1 MB), that is produced
¹⁷¹ at a rate of about 600 million events per second [1], that's about 600 MB/s. Keeping
¹⁷² in mind that only certain events are considered "interesting" for analysis, the task of de-
¹⁷³ ciding what events to consider out of all the data collected is a two-stage process. First,
¹⁷⁴ the events are filtered down to 100 thousand events per second for digital reconstruction
¹⁷⁵ and then more specialized algorithms filter the data even more to around 100 200 events
¹⁷⁶ per second that are found interesting. For CMS there is a Data Acquisition System that
¹⁷⁷ records the raw data to what's called a High-Level Trigger farm which is a room full
¹⁷⁸ of servers that are dedicated to processing and classify this raw data quickly. The data
¹⁷⁹ then gets sent to what's known as the Tier-0 farm where the full processing and the first
¹⁸⁰ reconstruction of the data are done. [2]

¹⁸¹ 3.2 What is Data Quality Monitoring?

¹⁸² To operate a sophisticated and complex apparatus as CMS, a quick online feedback on
¹⁸³ the quality of the data recorded is needed to avoid taking low quality data and to guarantee
¹⁸⁴ a good baseline for the offline analysis. Collecting a good data sets from the collisions
¹⁸⁵ is an important step towards search for new physics as deluge of new data poses an extra
¹⁸⁶ challenge of processing and storage. This all makes it all the more important to design
¹⁸⁷ algorithms and special software to control the quality of the data. This is where the Data
¹⁸⁸ Quality Monitoring (DQM) plays a critical in the maintainability of the experiment, the
¹⁸⁹ operation efficiency and performs a reliable data certification. The high-level goal of
¹⁹⁰ the system is to discover and pinpoint errors, problems occurring in detector hardware
¹⁹¹ or reconstruction software, early, with sufficient accuracy and clarity to maintain good
¹⁹² detector and operation efficiency. The DQM workflow consists of 2 types: **Online** and
¹⁹³ **Offline**.

¹⁹⁴ The **Online** DQM consists of receiving data taken from the event and trigger his-
¹⁹⁵ tograms to produce results in the form of monitoring elements like histogram references
¹⁹⁶ and quality reports. This live monitoring of each detector's status during data taking gives
¹⁹⁷ the online crew the possibility to identify problems with extremely low latency, mini-
¹⁹⁸ mizing the amount of data that would otherwise be unsuitable for physics analysis. The
¹⁹⁹ scrutiny of the Online DQM is a 24/7 job that consists of people or shifters that work at the
²⁰⁰ CMS control center constantly monitoring the hundreds of different plots and histograms
²⁰¹ produced by the DQM software. This consumes a lot of manpower and is strenuous work.

²⁰² The **Offline** DQM is more focused on the full statistics over the entire run of the
²⁰³ experiment and works more on the data certification. In the offline environment, the
²⁰⁴ system is used to review the results of the final data reconstruction on a run-by-run basis,
²⁰⁵ serving as the basis for certified data used across the CMS collaboration in all physics
²⁰⁶ analyses. In addition, the DQM framework is an integral part of the prompt calibration
²⁰⁷ loop. This is a specialized workflow run before the data are reconstructed to compute and
²⁰⁸ validate the most up-to-date set of conditions and calibrations subsequently used during

209 the prompt reconstruction.

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

²¹³ Chapter 4

²¹⁴ What is Machine Learning?

²¹⁵ Machine Learning (ML) can be defined as an application of Artificial Intelligence that
²¹⁶ permits the computer system to learn without being told explicitly. In ML a computer
²¹⁷ program is said to learn from experience E with respect to some class of tasks T and
²¹⁸ performance measure P, if its performance at tasks in T, as measured by P, improves
²¹⁹ with experience E [3]. ML has made tremendous strides in the past decades and has
²²⁰ become very popular recently due to its multifaceted applications. It is being used on
²²¹ social media, marketing, and in the scientific community as well. Some examples of
²²² ML applications are: the algorithms used on application in smartphones to detect human
²²³ faces, self-driving cars, computer games, stock prediction, and voice recognition. An
²²⁴ interesting characteristic of ML algorithms is that the more data one inputs the better is
²²⁵ the performance. The ML application has a very wide spectrum covering almost every
²²⁶ aspect of human endeavor that involves a lot of data. Scientific analysis today generates
²²⁷ enormous data and is hence a perfect use case to apply ML techniques. In this work
²²⁸ we use enhanced ML techniques based on progress in the recent past.

²²⁹ In general, there are two main categories to classify machine learning problems: **Su-**
²³⁰ **pervised Learning** (SL) and **Unsupervised Learning** (UL). SL is the most used ML
²³¹ approach and has proven to be very effective for a wide variety of problems. Examples
²³² of common SL problems are: spam filters, predicting housing prices, identifying a ma-
²³³ lignant or benign tumor, etc. These types of problems are characterized by providing a

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

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

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

255 4.1 Developing the Algorithm

256 To develop an ML algorithm the following are taken into consideration, what is the
257 task? and what is the method to approach the task? In our case, we are looking into images
258 that have information about the activity that the channels in the HCAL are detecting.
259 These images are called ”occupancy maps” and they are a visual way of monitoring the
260 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.

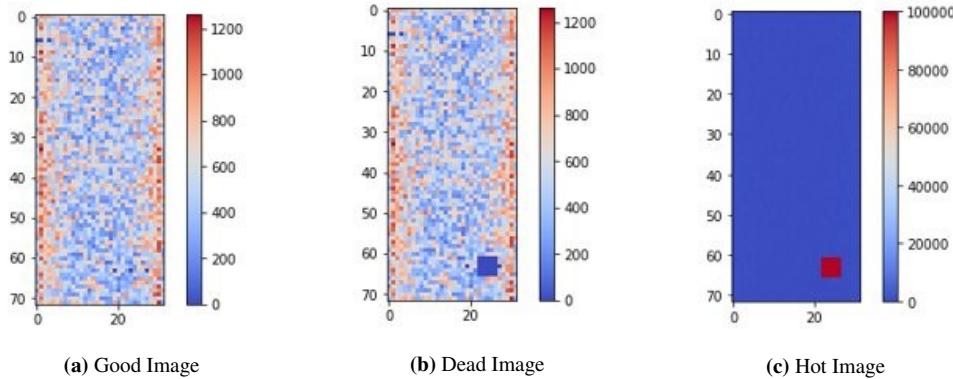


Figure 4.1: Occupancy maps with 5x5 affected regions

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.

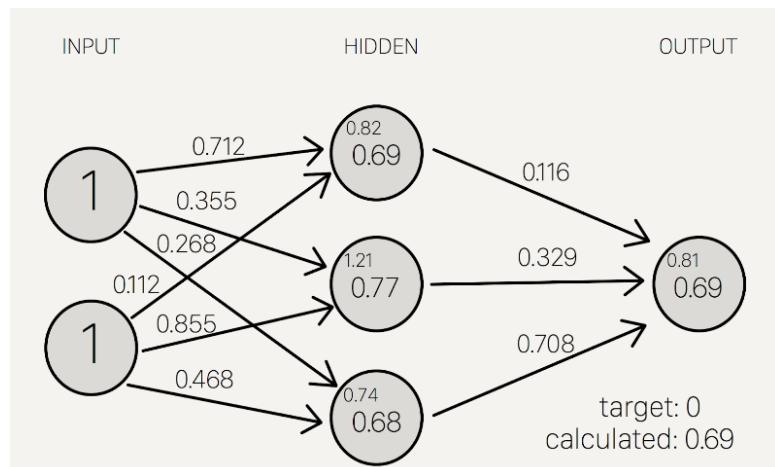


Figure 4.2: Weights and Biases

"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

η : Learning Rate (typically 0.01 to 0.001)

η :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

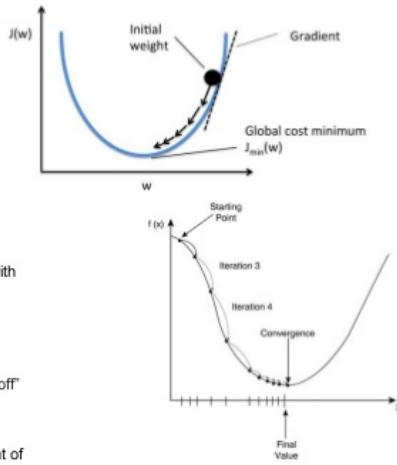


Figure 4.3: Gradient Descent algorithm

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

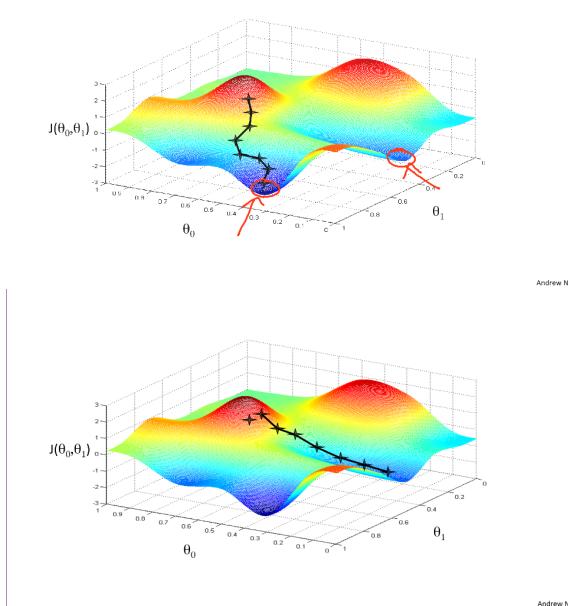


Figure 4.4: Loss Function surface

307 **Chapter 5**

308 **Results**

309 **5.1 Introduction**

310 A detailed explanation for the estimation of the $Z \rightarrow \nu\bar{\nu} + \text{jets}$ background is presented
311 in this chapter. The following estimation procedure builds upon the 2016 SUSY analysis
312 summarized in ?? and aims to refine the overall background calculation as well as reduce
313 the uncertainties associated with the previous method. To accomplish this, an additional
314 $\gamma + \text{jets}$ CS is used in conjunction with the tight $Z \rightarrow \mu^+ \mu^-$ control region used in the
315 2016 analysis estimation. The new $\gamma + \text{jets}$ CS provides a more data-driven estimation
316 procedure with the added benefit of a substantially larger production cross-section than
317 $Z + \text{jets}$ processes [?].

318 **5.2 The Irreducible $Z \rightarrow \nu\bar{\nu}$ Background**

319 An important source of background in searches for SUSY in the 0-lepton final state
320 comes from events in which a Z boson, accompanied by jets, decays into a pair of neu-
321 trinos ($Z \rightarrow \nu\bar{\nu} + \text{jets}$). The resulting final state is comprised of a large p_T^{miss} (from the
322 neutrino pair) and multiple hadron jets, closely mimicking the SUSY signal. For similar
323 searches, the $Z \rightarrow \nu\bar{\nu}$ contribution can make up a large portion of the background in many
324 of the search bin regions (higher than 50% in some regions) [?, ?], and can compose up

325 to about a third of the total SM background. For this particular analysis, the $Z \rightarrow \nu\bar{\nu}$ ac-
 326 counts for about 17% of the total SM background and owes its low value to the dedicated
 327 top-tagging algorithm outlined in ??.

328

329 There are several different methods that have been developed to estimate the $Z \rightarrow$
 330 $\nu\bar{\nu}+\text{jets}$ background [?]. Two of the commonly used methods involve the use of a control
 331 region dominated by $Z \rightarrow ll+\text{jets}$, where the l stands for lepton (either a muon or an
 332 electron, in this case) or $\gamma+\text{jets}$ events. The $Z \rightarrow ll$ channel has the advantage of having
 333 very similar kinematics to the $Z \rightarrow \nu\bar{\nu}$ region but suffers from low statistics (due in part
 334 to its small branching ratio), specially in the tight search regions used in typical SUSY
 335 searches. On the other hand, the $\gamma+\text{jets}$ region has a much higher production cross-section
 336 but involves a completely different process. The hybrid method described in this chapter
 337 makes use of both control regions in order to estimate the $Z \rightarrow \nu\bar{\nu}$ background corrections,
 338 and aims to improve on the results of the 2016 method described in ??.



Figure 5.1: Leading-order Feynman diagrams for $Z+\text{jets}$ and $\gamma+\text{jets}$ processes. The ‘V’ in the figure can represent either Z or γ .

339 5.3 The Loose $\gamma+$ jets Control Region

340 The loose $\gamma+\text{jets}$ control region, as well as the photon ID/isolation selection, is de-
 341 scribed in this section. The $\gamma+\text{jets}$ control region is chosen with the purpose of substitut-
 342 ing the loose muon control region used in the 2016 version of this analysis for the calcu-
 343 lation of the shape correction scale factors applied to the final estimation of the $Z \rightarrow \nu\bar{\nu}$
 344 background. The $\gamma+\text{jets}$ control region is presumed to be better suited for this estima-
 345 tion due to having a much higher cross-section and therefore, an expected reduction in
 346 the statistical and systematic uncertainties associated to the shape correction factors are

347 expected.

348 **5.3.1 Photon ID and Isolation**

349 Three different working points are provided by the CMS EGM physics object group
350 (POG) for simple cut-based photon identification [?]. The three working points, called
351 loose, medium and tight, are chosen according to the requirements of the particular analy-
352 sis and differ on the amount of background rejection they offer, as well as on their average
353 photon selection efficiency. The higher the efficiency of a given ID, the lower the amount
354 of background that is rejected, and vice-versa. ?? shows the cut values that are applied
355 to photons that are found within both the ECAL barrel and endcap range. The associated
356 values to the efficiency and the background rejection rate are shown for each of the three
357 different photon ID selections.

358

359 In order to obtain the high efficiency and background rejection rates shown, a robust
360 set of identification and isolation criteria are selected. A total of five parameters are
361 used for this simple cut based method. For photon identification, the H/E and the $\sigma I\eta I\eta$
362 variables are found to provide the best results. The H/E parameter is defined as the ratio
363 of the HCAL tower energy over the ECAL seed cluster energy. A threshold value is
364 selected on H/E to remove background from electrons that are detected in both the ECAL
365 and HCAL but have no reconstructed track [?]. The $\sigma I\eta I\eta$ variable is known as the
366 photon shower shape variable and is defined in the ECAL as the energy weighted standard
367 deviation of a single crystal within the 5×5 crystal η range, centered around the crystal
368 with maximum energy [?]. This variable is a key component in the identification of both
369 electrons and photons since it provides a measure of the shower width where most of the
370 energy has been deposited in a given ECAL crystal.

Table 5.1: Identification and isolation cut values for photons provided by the CMS EGM POG. Values are provided for the three working points described (loose, medium and tight) for both the ECAL barrel and endcaps. The photon selection efficiency of the three working points, as well as their associated background rejection rate, are provided.

| Barrel | Loose (90.06%) | Medium (80.19%) | Tight (70.01%) |
|---|---|---|---|
| Background Rejection | Loose (85.73%) | Medium (88.87%) | Tight (90.66%) |
| H/E | 0.105 | 0.035 | 0.020 |
| $\sigma I\eta I\eta$ | 0.0103 | 0.0103 | 0.0103 |
| ρ -corrected PF charged hadron isolation | 2.839 | 1.416 | 1.158 |
| ρ -corrected PF neutral hadron isolation | $9.188 + 0.0126 \cdot p_T^\gamma + 0.000026 \cdot p_T^{\gamma^2}$ | $2.491 + 0.0126 \cdot p_T^\gamma + 0.000026 \cdot p_T^{\gamma^2}$ | $1.267 + 0.0126 \cdot p_T^\gamma + 0.000026 \cdot p_T^{\gamma^2}$ |
| ρ -corrected PF photon isolation | $2.956 + 0.0035 \cdot p_T^\gamma$ | $2.952 + 0.0040 \cdot p_T^\gamma$ | $2.065 + 0.0035 \cdot p_T^\gamma$ |

| End Cap | Loose (90.81%) | Medium (80.06%) | Tight (70.11%) |
|---|--|---|---|
| Background Rejection | Loose (76.90%) | Medium (81.50%) | Tight (84.34%) |
| H/E | 0.029 | 0.027 | 0.025 |
| $\sigma I\eta I\eta$ | 0.0276 | 0.0271 | 0.0271 |
| ρ -corrected PF charged hadron isolation | 2.150 | 1.012 | 0.575 |
| ρ -corrected PF neutral hadron isolation | $10.471 + 0.0119 \cdot p_T^\gamma + 0.000025 \cdot p_T^{\gamma^2}$ | $9.131 + 0.0119 \cdot p_T^\gamma + 0.000025 \cdot p_T^{\gamma^2}$ | $8.916 + 0.0119 \cdot p_T^\gamma + 0.000025 \cdot p_T^{\gamma^2}$ |
| ρ -corrected PF photon isolation | $4.895 + 0.0040 \cdot p_T^\gamma$ | $4.095 + 0.0040 \cdot p_T^\gamma$ | $3.272 + 0.0040 \cdot p_T^\gamma$ |

³⁷¹ The other three parameters considered, comprise the isolation portion of the photon
³⁷² selection cuts. These are the ρ -corrected particle flow (PF) charged hadron, neutral hadron
³⁷³ and photon isolation parameters. As can be seen from ??, two of these parameters (the
³⁷⁴ neutral hadron and photon isolation) have a dependence on p_T^γ . These cuts are used to
³⁷⁵ ensure that the identified photon is well isolated within its own cone and by rejecting
³⁷⁶ photons that are identified within close proximity to either a charged or a neutral hadron
³⁷⁷ [?]. The value ρ included in the name of each of these parameters refers to the total

378 pileup density [?]. Therefore, the term “ ρ -corrected” implies that these values, which are
 379 sensitive to the residual contamination that arises from pile-up, have been corrected to
 380 include these contributions.

381 5.3.2 Photon Selection

382 The event selection process for the γ +jets control region starts with photon candidates
 383 that have a $p_T > 200$ GeV and are within the acceptance range of the CMS ECAL (given
 384 by $|\eta| < 1.4442$ for the barrel and $1.566 < |\eta| < 2.5$ for the endcaps). The photons are
 385 subjected to pass the loose ID/isolation cuts described in ?? in order to remove $\sim 85\%$
 386 of the background processes and obtain a prompt photon sample that is $\sim 90\%$ pure, on
 387 average. Additional restrictions include some of the same requirements imposed on the
 388 signal baseline selection discussed in ?? These include $N_j \geq 4$, an $H_T > 300$ GeV,
 389 the $\Delta\phi$ requirements for leading jets and the lepton vetoes described in ?? The lepton
 390 veto, in particular, greatly improves the prompt photon selection by removing many of
 391 the events in the simulated samples where a lepton gets misidentified as a photon.

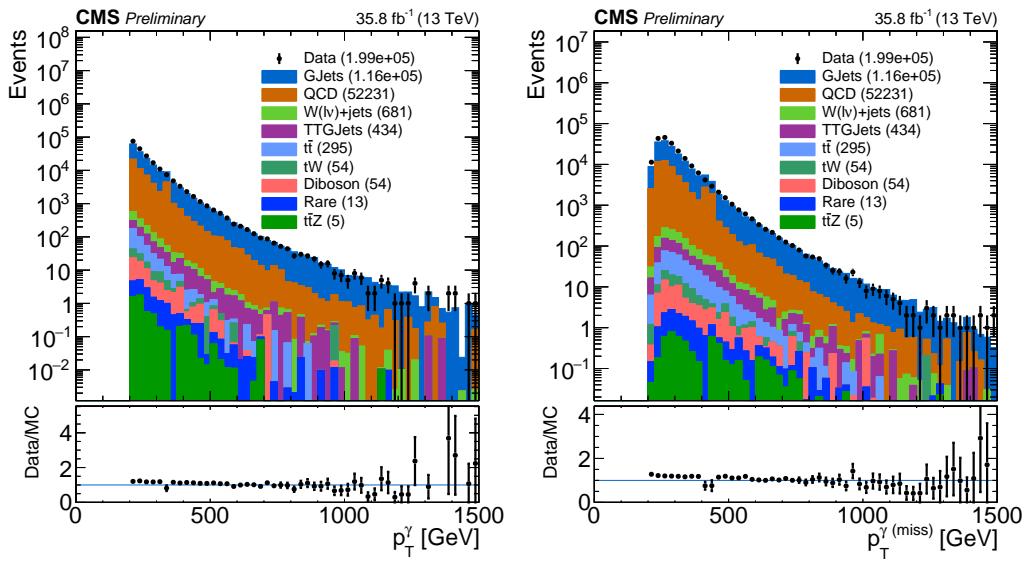


Figure 5.2: Shown are both the p_T^γ (left) and $p_T^{\gamma(\text{miss})}$ (right) distributions before applying any corrections. $p_T^{\gamma(\text{miss})}$ is obtained by adding the p_T^γ to the total p_T^{miss} in every event.

392 To further emulate the $Z \rightarrow \nu\bar{\nu}$ +jets background, a variable in which the photons are
 393 treated as p_T^{miss} is defined. We call this variable $p_T^{\gamma(\text{miss})}$ and we obtain it by adding the p_T^γ

394 for every event to the total p_T^{miss} in the event. Both the p_T^γ and the resulting $p_T^{\gamma^{(miss)}}$ distri-
 395 butions are shown in ?? as data/MC comparison plots, where the simulated backgrounds
 396 are stacked in order of ascending contribution.

397

398 The main contributions from simulation arise from the γ +jets, QCD and to a lesser
 399 extent, $t\bar{t}\gamma$. Other non-dominant backgrounds in the control region include contributions
 400 from $W(l\nu)+\text{jets}$, $t\bar{t}$, Diboson, tW , $t\bar{t}Z$ and rare processes. Most of these lesser back-
 401 grounds are nearly negligible (several orders of magnitude lower than the dominant back-
 402 grounds) and are considered to be mostly composed of fake photons. In addition to the
 403 cuts described, all of the simulation samples are subjected to weights that apply correc-
 404 tions to pileup as well as the b-tagging efficiency. Data, on the other hand, is obtained
 405 from a sample that contains events with at least one identified photon. Photons in this
 406 sample are also subjected to the high-level trigger HLT_Photon175, which restricts the
 407 selection to photons that have a $p_T > 175$ GeV. Both simulation and data are subjected to
 408 the same selection criteria established in this section.

409 5.3.3 Photon Purity and Fake Rate

410 Three different types of photons make up the γ +jets CS: prompt photons, produced
 411 either directly or through fragmentation, and fake photons. Prompt photons are defined
 412 as photons which are formed shortly after the proton-proton collision (i.e. before the pro-
 413 duced quarks and gluons have had enough time to form hadrons). Two types of photons
 414 fit in this category. The first type, which we designate as direct photons, are photons that
 415 are produced directly from the proton-proton interaction [?]. A secondary type of prompt
 416 photon, that is virtually indistinguishable from the direct photons at the detector level,
 417 originates from the decay of π^0 mesons and are called fragmentation photons. The final
 418 type of photon found in the CS corresponds to fake (or non-prompt) photons. The fake
 419 photon contribution typically arises from leptons (mostly electrons) whose tracks are not
 420 properly reconstructed, yet leave energy measurements in the ECAL.

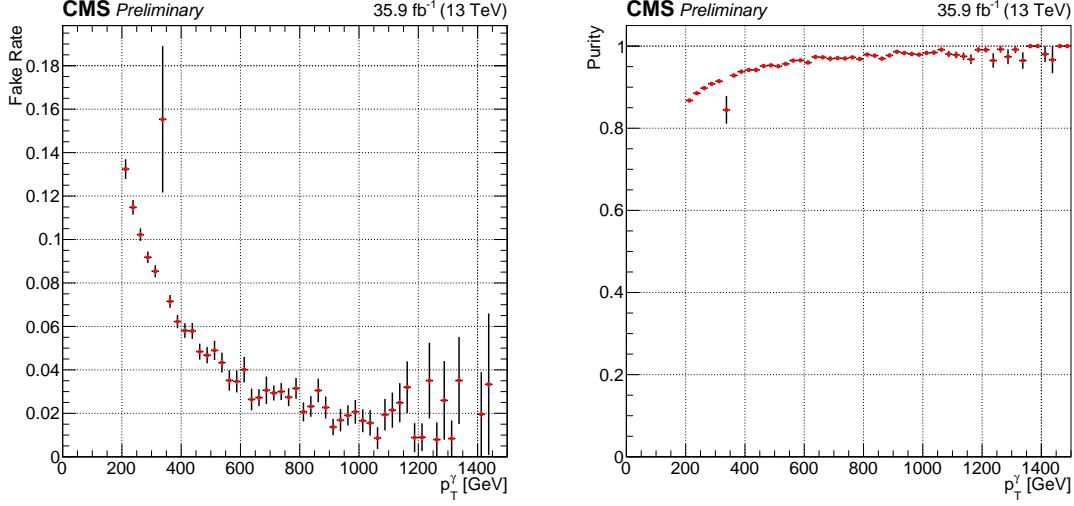


Figure 5.3: Plots for Fake Rate (left) and Purity (right) as a function of the photon p_T are shown. The events are selected are required to have a $p_T > 200$ GeV, be within the ECAL acceptance range, and pass the loose ID selection cuts. This selection was produced in order to verify the values given by the E/ γ POG. As can be seen, the efficiency (purity) is seen to agree with the values of the loose photon ID/isolation selection.

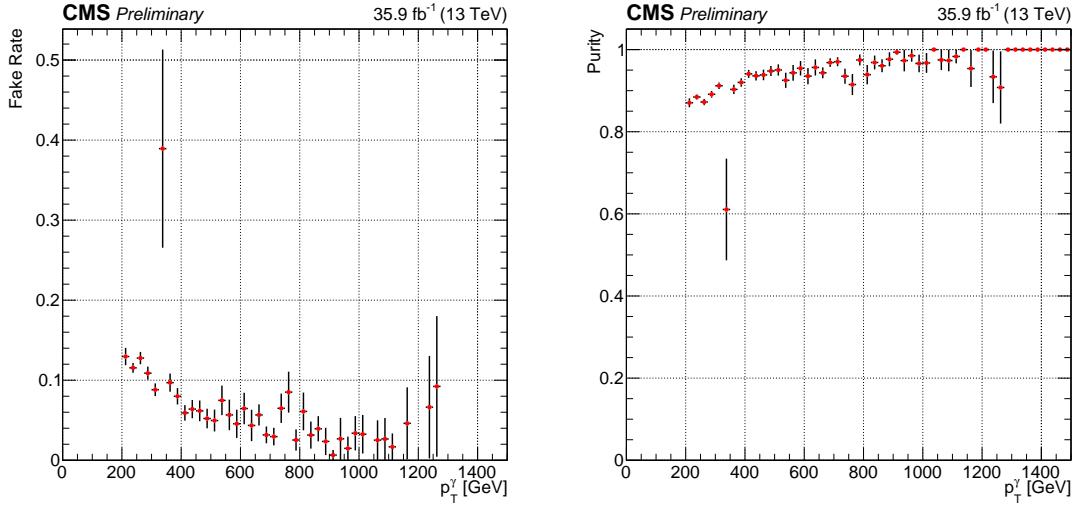


Figure 5.4: Plots for Fake Rate (left) and Purity (right) as a function of the photon p_T are shown. These plots include photons with the full control region selection. Aside from exhibiting lower statistics, the plots seem to agree with the fake rate and purity before all the control region cuts are applied.

421 In order to identify prompt photons, reconstructed photons from the $\gamma + \text{jets}$ and QCD
 422 samples are matched to generator-level photons in space and momentum by requiring
 423 $\Delta R(\gamma_{\text{gen}}, \gamma_{\text{reco}}) < 0.4$ and $0.5 < p_T^{\text{gen}}/p_T^{\text{reco}} < 2.0$, respectively. Any reconstructed photon
 424 which fails to get matched to a generator level photon is labeled as a fake/non-prompt
 425 photon. Direct photons are identified by further requiring that the reconstructed photons
 426 be matched to a parton (a gluon or quark) in space as $\Delta R(\gamma, \text{parton}) > 0.4$. This require-

ment is intended to distinguish the reconstructed photons from highly boosted π^0 's, which compose a large portion of the experimentally indistinguishable fragmentation photons. Finally, fragmentation photons are obtained exclusively from QCD simulation and are required to have $\Delta R(\gamma, \text{parton}) < 0.4$ in order to avoid double counting photons from the γ +jets sample.

432

With all three types of photons defined, a study can be carried out from simulation to estimate their respective contributions to the defined control region. The study takes into account that any reconstructed photon in the γ +jets or QCD samples can only be categorized as prompt (through direct production or fragmentation) or non-prompt (fake). The purity and fake rate can then be defined in terms of the relative proportions of prompt or non-prompt photons with respect to the sum of the contributions from all three types of photons. Identified direct photons are taken from the γ +jets sample exclusively. Meanwhile, the fragmentation and fake photon contributions are taken from the QCD sample. The three quantities are then added together and their respective contributions are determined in terms of the photon p_T (??).

443

The photon purity (?? and ??, right) is defined in terms of the prompt and non-prompt photons as:

446

$$p_\gamma = \frac{\text{prompt}}{\text{prompt} + \text{fake}} ,$$

where the prompt photon portion comes from the sum of the direct photons (extracted from the γ + jets sample) and the fragmentation photons (extracted from the QCD sample). The remaining non-prompt (or fake) photons all come from photons in the QCD sample that were not matched to truth-level photons in space and momentum with the specified required conditions. Meanwhile, the photon fake rate (?? and ??, left) is defined from this same combination of samples as:

454

$$455 \quad f = \frac{fake}{prompt+fake},$$

456 ?? shows the purity and fakerate for photons that pass the loose ID/selection, have a
 457 $p_T > 200$ GeV and are within the ECAL acceptance range. A sample is obtained in which
 458 77% of the photons are direct, 12% are fragmentation and 11% are fakes. This implies an
 459 average purity of $\sim 89\%$ for this sample, well within the value that is expected. ?? shows
 460 the same ratios for the loose $\gamma +$ jets control region described in ?? . Although the amount
 461 of statistics has decreased due to the additional cuts, a similar trend can be observed.

462 5.4 The $Z \rightarrow \mu^+ \mu^-$ Control Region

463 The $Z \rightarrow \mu^+ \mu^-$ control region defined in this section is in every respect identical to
 464 the one applied in the 2016 analysis (??). The only difference between the 2016 method
 465 and the one discussed in this chapter is that the Drell-Yan (DY) sample is only used for
 466 the normalization correction of the $Z \rightarrow \nu \bar{\nu}$ background. Therefore, the loose $\mu \mu$ control
 467 region is not used or applied in the calculation of the scale factors. In the following
 468 subsections only the tight $\mu \mu$ control region, and its usage to obtain the normalization
 469 scale factor R_{norm} , is discussed.

470 5.4.1 Muon ID and Isolation

471 The muons are selected using the “medium muon” selection [?], per the recommen-
 472 dation of the muon POG. The muon candidates in this selection satisfy $p_T > 10$ GeV and
 473 $|\eta| < 2.4$. Other additional cuts are applied to aid in the muon candidate selection, such as
 474 an impact parameter cut. Muons are also subjected to a PF relative-isolation (also referred
 475 to as mini-isolation) in which the cone size is inversely proportional to the muon p_T . This
 476 requirement enforces the p_T within the isolation cone to be at most 20% of the muon p_T in
 477 order to eliminate events with an isolated muon. Details of the medium photon selection

478 are included in ?? and ??, while details of the impact parameter cut are summarized in
 479 ??.

| | |
|---|--------------|
| Muon Medium ID | |
| Loose muon ID | Yes |
| Fraction of valid tracker hits > | 0.80 |
| Good Global muon OR Tight segment compatibility > | Yes OR 0.451 |

Table 5.2: Muon Medium ID 2016 HIP Safe

| | |
|-------------------------------------|-------|
| Good Global muon | |
| Global muon | Yes |
| Normalized global-track $\chi^2 <$ | 3 |
| Tracker-Standalone position match < | 12 |
| Kick finder < | 20 |
| Segment compatibility > | 0.303 |

Table 5.3: Muon Medium ID HIP Safe Good Global Muon

| | |
|-----------------------|-----|
| Muon Impact Parameter | |
| d0 < | 0.2 |
| dz < | 0.5 |

Table 5.4: Additional Impact Parameter cut on Muons

480 5.4.2 Muon Selection in the Tight Control Region

481 Events are selected from data samples that contain exactly two oppositely charged
 482 muons ($\mu^+\mu^-$), which fall within the invariant mass $81 < m_{ll} < 101$ GeV window for
 483 the Z boson. Additional cuts for the tight muon selection include baseline requirements
 484 such as an $H_T > 300$ GeV, $N_j \geq 4$, the $\Delta\phi$ baseline cut on leading jets, a $p_T^{miss} > 250$
 485 GeV, an $m_{T2} > 200$ GeV and at least 1 top-tagged jet $N_t \geq 1$. In addition, the p_T of
 486 the two muons are required to be $p_T > 50$ GeV for the leading muon and $p_T > 20$ GeV
 487 for the sub-leading one. The only difference, when compared to the signal region is the
 488 missing lepton veto, in addition to the dimuon events being treated as p_T^{miss} . This makes
 489 for a region that exhibits very similar kinematics to the Z $\rightarrow\nu\bar{\nu}$ signal region, yet suffers
 490 from a lack of statistics.

491 **5.5 Analysis**

492 In this section a detailed explanation of the calculation of the scale factors for both
493 shape and normalization corrections is provided. The following methods make use of
494 the loose γ +jets and the tight $\mu\mu$ control regions defined in the previous sections. The
495 procedure involves extracting the shape corrections S_γ from the γ +jets control region and
496 afterwards obtain a single normalization correction factor R_{norm} from the tight $\mu\mu$ control
497 region. Both factors will then be applied to the final prediction of the $Z \rightarrow \nu\bar{\nu}$ background
498 in each of the required search bins.

499 **5.5.1 Shape Correction Using the γ + jets Control Sample**

500 In this section the validation of the γ +jets simulation is discussed in terms of the
501 shape of the loose photon control region. As it was shown in ??, this control region
502 has high purity for γ +jets events, particularly in regions of high p_T ($\gtrsim 300$). In order
503 to apply this correction factor it is assumed that the shape differences between data and
504 simulation are similar between $Z \rightarrow \nu\bar{\nu}$ and γ +jets events. This assumption is validated in
505 studies which compare the cross-section ratio of Z +jets to γ +jets events [?]. ?? shows the
506 results of this study, conducted in 2014, for both data and MadGraph simulation with an
507 integrated luminosity of 19.7fb^{-1} and a center-of-mass energy of 8 TeV. It can be seen that
508 for values of $p_T^{Z/\gamma} \gtrsim 300$ GeV, the ratio of the cross-section of both processes becomes
509 nearly constant. It is then a matter of applying a factor to account for the difference in the
510 amount of events between the Z +jets and γ +jets events in order to obtain the total amount
511 of $Z \rightarrow \nu\bar{\nu}$ +jets events. This factor is obtained from the tight $\mu\mu$ control region, as shown
512 in ??.

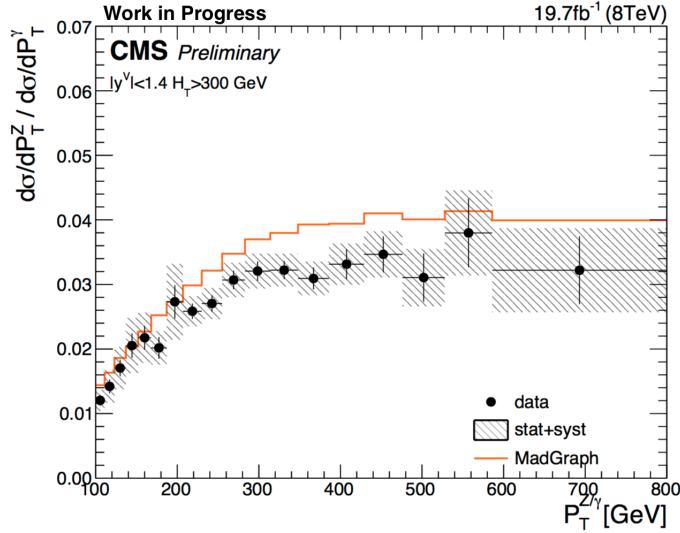


Figure 5.5: Results of study of the Z+jets to γ +jets cross-section ratio for both data and MadGraph simulation. For high values of the vector boson transverse momentum, the ratio between these processes is observed to be nearly constant.

513 In order to obtain the shape corrections, the ratio between data and simulation of the
 514 jet multiplicity distribution is used (??). This is due to it exhibiting the highest difference
 515 between the observed data and MC. The re-weight for the γ +jets simulation sample is
 516 then accomplished by applying the N_{jet} dependent factor $S_\gamma(N_j)$. This scale factor is
 517 determined by taking the ratio of the data and simulation, after subtracting all other MC
 518 samples from data events:

$$519 \quad S_\gamma^i = \frac{\text{Data}^i - \text{MC}_{\text{other}}^i}{\text{MC}_{\gamma+\text{jets}}^i},$$

520 where i denotes any given bin in the N_j distribution. The shape correction factors S_γ^i are
 521 displayed graphically in ?? (right) for each N_j bin. These factors correct for differences
 522 in the jet multiplicity shape, while the overall normalization is estimated from the tight
 523 $\mu\mu$ control region. ?? shows the N_j distribution in the tight $\mu\mu$ control region after the
 524 calculated scale factors have been applied. The S_γ correction will be applied to the $Z \rightarrow \nu\bar{\nu}$
 525 simulation final prediction for each of the analysis search bins. The uncertainty associated
 526 with the scale factor is estimated from the event yields in the loose photon control region.
 527 This uncertainty will form part of the total systematic uncertainty in the final prediction.

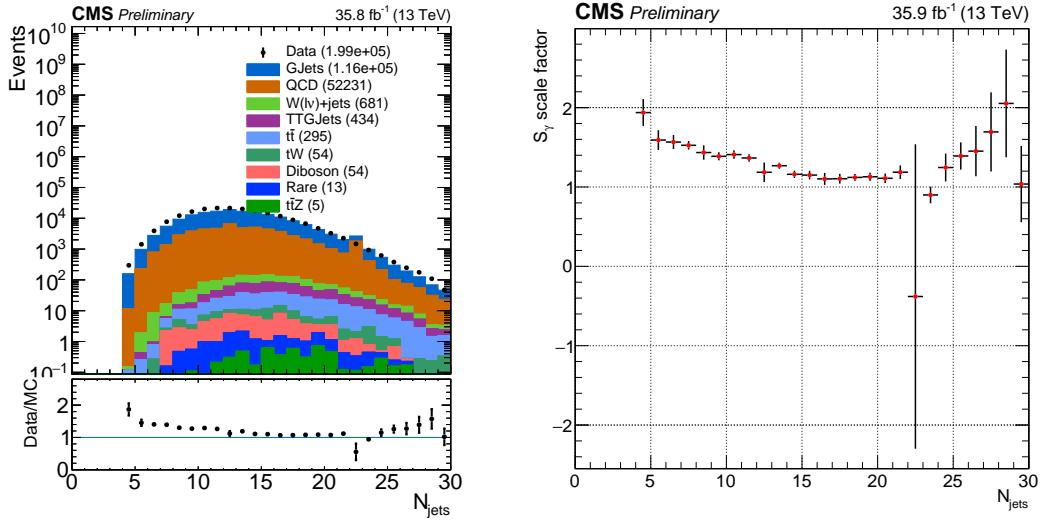


Figure 5.6: Jet multiplicity and the associated S_γ scale factor in the loose photon control region before any corrections are applied.

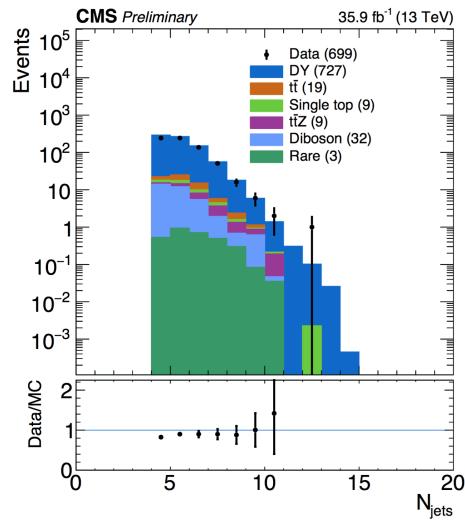


Figure 5.7: N_{jet} distribution in the tight $\mu\mu$ control region after S_γ corrections.

528 The effect of the $S_\gamma(N_j)$ scale factor is shown for various distributions. These results
 529 show that the overall agreement between data and simulation improves after applying the
 530 corresponding shape corrections.

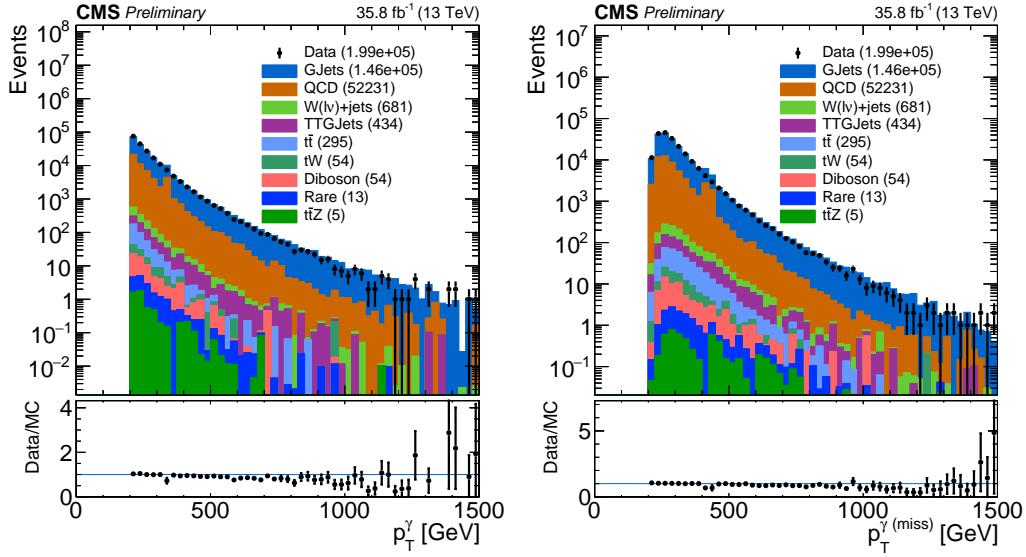


Figure 5.8: p_T^γ (left) and $p_T^{\gamma(\text{miss})}$ (right) distributions after applying the $S_\gamma(N_j)$ scale factor. Comparing to ??, an improvement in the agreement between data/MC can be observed.

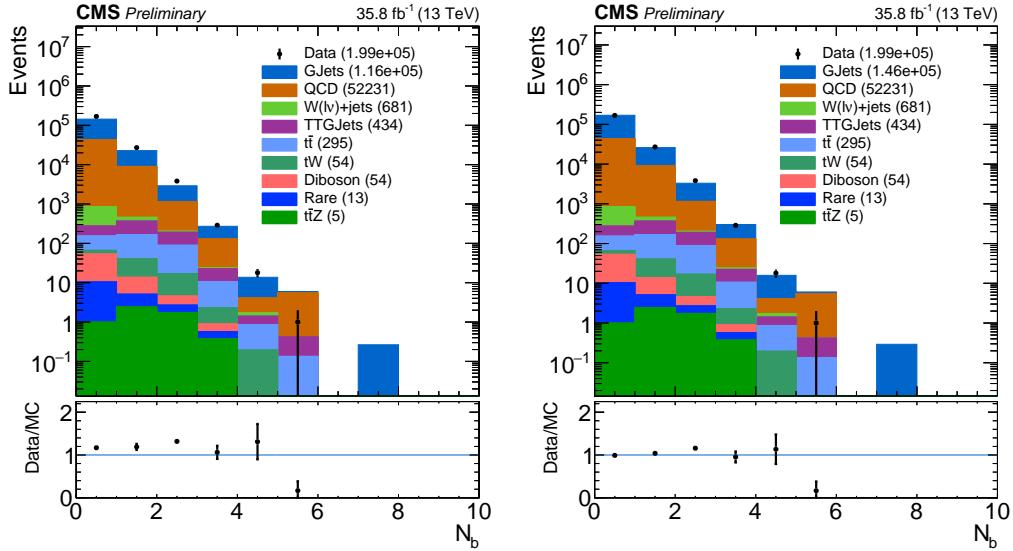


Figure 5.9: N_b distribution before (left) and after (right) applying the $S_\gamma(N_j)$ scale factor.

531 **5.5.2 Normalization Correction Using the tight $Z \rightarrow \mu^+ \mu^-$ Control**
 532 **Sample**

533 In order to constrain the normalization of the $Z \rightarrow \nu\bar{\nu}$ simulation sample, a normal-
 534 ization correction factor R_{norm} is calculated from the tight $\mu\mu$ control region defined in
 535 ???. Two categories are considered: the zero b-tagged jet category ($N_b = 0$), and the ≥ 1
 536 b-tagged jet category ($N_b \geq 1$). Both of these categories are statistically consistent with

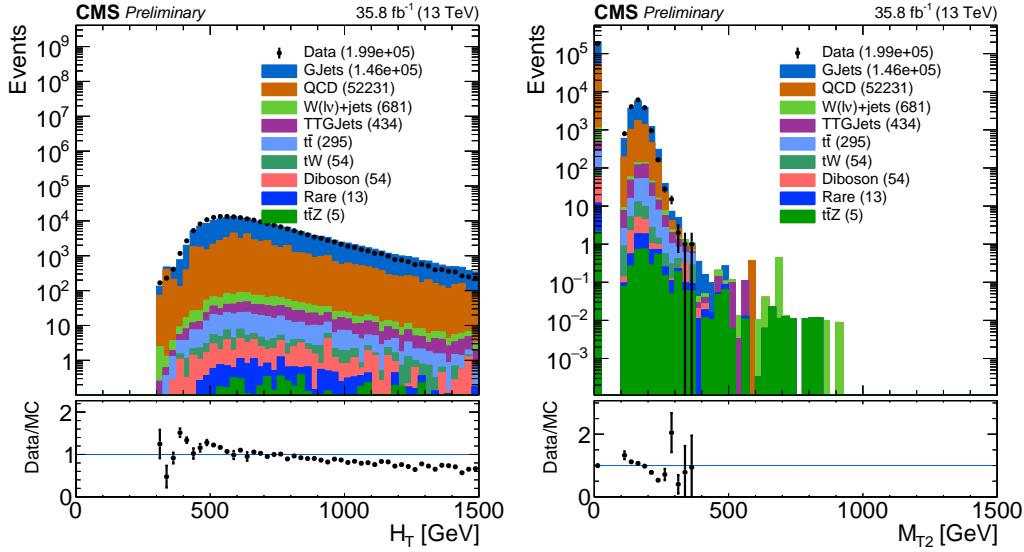


Figure 5.10: H_T and m_{T2} distributions applying the $S_\gamma(N_j)$ scale factor.

537 each other but the inclusive region ($N_b \geq 0$) has a lower overall uncertainty. The method
 538 used to calculate the normalization scale factor requires that the N_j -dependent shape cor-
 539 rection factors already be applied. Then, the R_{norm} factor can be extracted from the ratio
 540 of the total event yield in data to that in the simulation. This factor is found to be:

541

$$R_{norm} = 1.070 \pm 0.085,$$

542 where the uncertainty includes only the associated statistical uncertainties on data and
 543 simulation. This uncertainty is found to be propagated to the final background prediction,
 544 see ??.

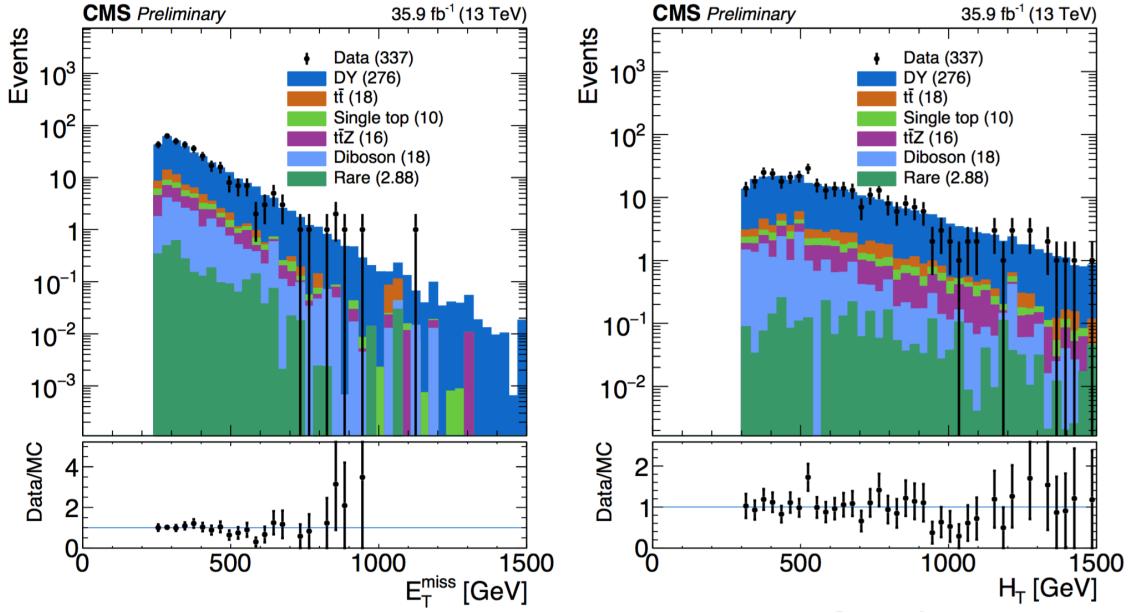


Figure 5.11: Shown are data/MC comparisons for the p_T^{miss} (left) and H_T (right) distributions after applying both the N_j -dependent shape corrections (S_γ) and the global normalization scale factor (R_{norm}).

546 Data/MC comparisons are shown in ?? and ?? after applying R_{norm} for several distri-
 547 butions in the study. With this final global scale factor all the required ingredients for the
 548 central value of the $Z \rightarrow \nu\bar{\nu}$ background prediction are obtained.

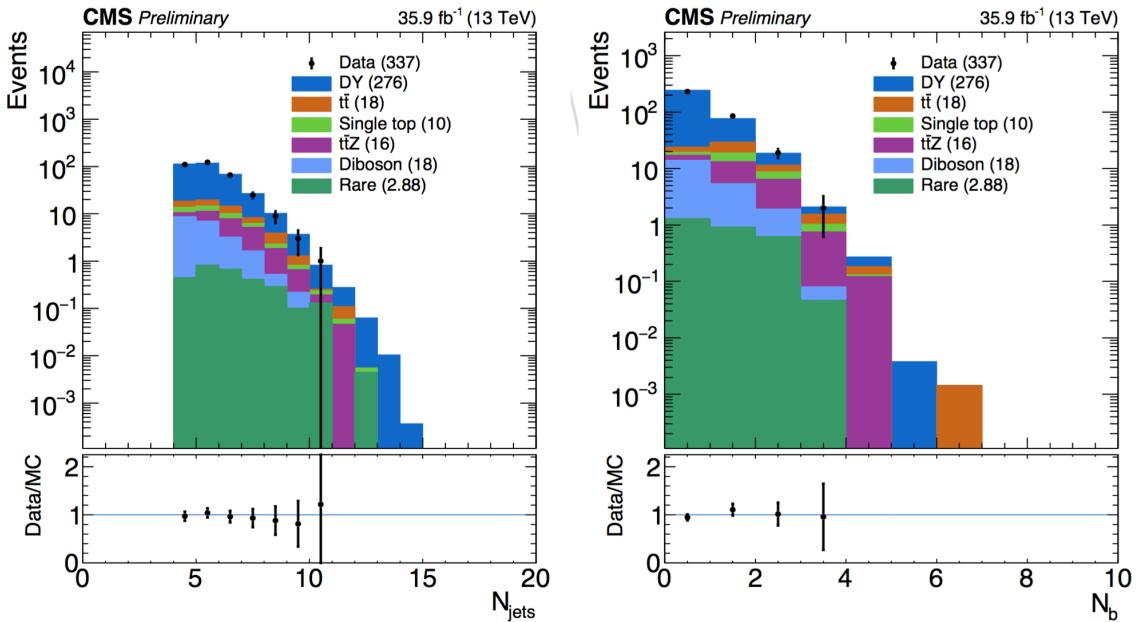


Figure 5.12: Shown are data/MC comparisons for the N_j (left) and N_b (right) distributions after applying both the N_j -dependent shape corrections (S_γ) and the global normalization scale factor (R_{norm}).

549 **5.6 Results**

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

557 **5.6.1 Systematics**

558 Two categories of uncertainties for the $Z \rightarrow \nu\bar{\nu}$ prediction are considered: uncertain-
559 ties that are associated to the use of MC simulation and the uncertainties specifically
560 associated to the background prediction method. Several sources are acknowledged in the
561 first category mentioned such as PDF and renormalization/factorization scale choices, jet
562 and p_T^{miss} energy scale uncertainties b-tag scale factor uncertainties, and trigger efficiency
563 uncertainties. Given that the simulation sample is normalized to data in the tight control
564 region, uncertainties associated with the luminosity and cross-section are excluded. In
565 addition, the overall $Z \rightarrow \nu\bar{\nu}$ statistical uncertainty from MC simulation is also taken into
566 account.

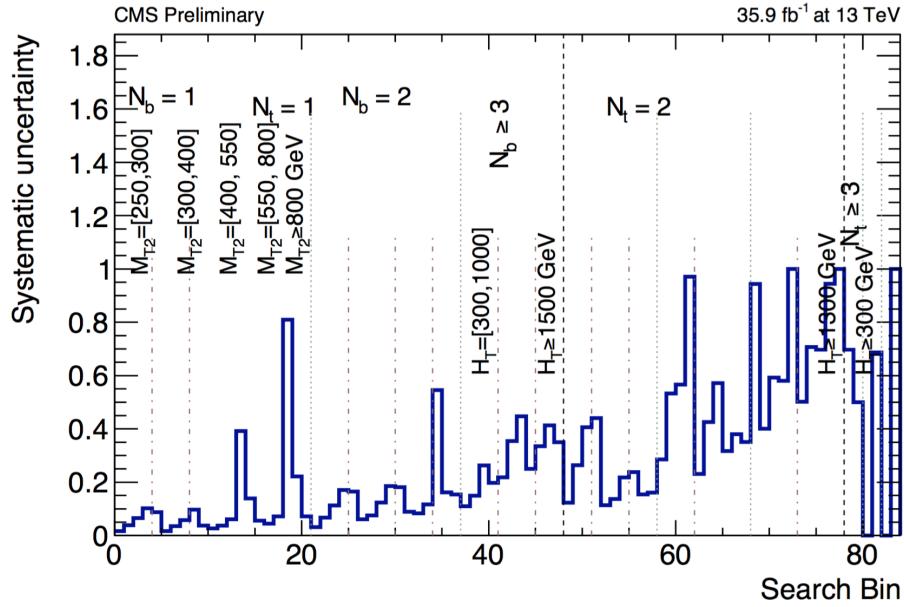


Figure 5.13: Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics.

568 The statistical uncertainty associated with each bin in the MC is propagated as a sys-
 569 tematic uncertainty. The relative uncertainty per bin can be see in ???. It shows that the
 570 uncertainties for the MC vary from as low as 1% up to 81% and even 100% in some re-
 571 gions. Since the final estimation is scaled using the global normalization factor from the
 572 tight $\mu\mu$ control region (R_{norm}), the total uncertainty, due to limited amounts of events in
 573 data, is propagated in the final prediction. This is also true for the $S_\gamma(N_j)$ scale factor, in
 574 which the residual differences in search variables other than N_j are evaluated in the loose
 575 photon control region. Both the uncertainty arising from the N_j re-weighting as well as
 576 the residual differences are evaluated together. The uncertainty from R_{norm} is propagated
 577 as a flat value of 7.9% uncertainty per each search bin.

578 5.6.2 $Z \rightarrow \nu\bar{\nu}$ Estimation for the Search Bins

579 The final estimation for the $Z \rightarrow \nu\bar{\nu}$ background calculated for all 84 search bins is
 580 shown in ???. The statistical uncertainty in bins that have zero events is treated as the
 581 average weight (the sum of the weights squared over the weight) times the poisson error
 582 on 0 which is 1.8. This average weight is calculated on the basis of a relaxed cut in which
 583 $N_b \geq 2$ is required. For comparison, a cut in which $N_t > 2$ where two tops are fake for

584 the $Z \rightarrow \nu\bar{\nu}$ is used.

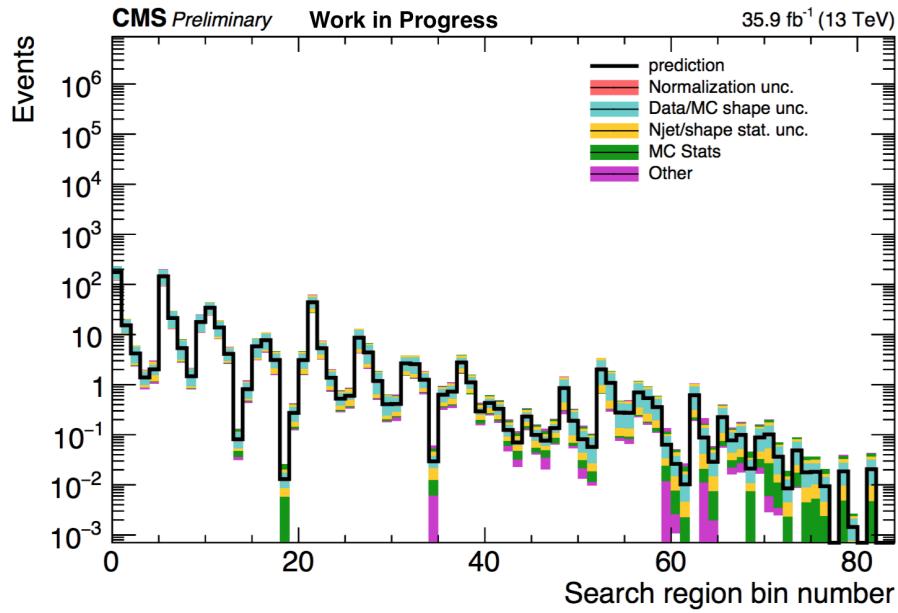


Figure 5.14: $Z \rightarrow \nu\bar{\nu}$ background prediction for all search bins, including the breakdown of the various uncertainties.

⁵⁸⁵ **Appendix A**

⁵⁸⁶ **Appendix Title**

⁵⁸⁷ **Appendix B**

⁵⁸⁸ **References**

⁵⁸⁹ [1] CERN, “Processing what to record?,” 2018.

⁵⁹⁰ [2] CMS, “The CMS Computing Project,” tech. rep., CERN, 2005.

⁵⁹¹ [3] Coursera, “Machine learning,” 2018.

⁵⁹² [4] A. S. Walia, “Types of optimization algorithms used in neural networks and ways to
⁵⁹³ optimize gradient descent,” 2018.