



Anomaly Detection Using Machine Learning Techniques for Beam Injections from the SPS to the LHC at CERN

MARC FERRIGGI

Supervised by Dr. Gianluca Valentino

Department of Computer Science
Faculty of ICT
University of Malta

May, 2019

*A FYP submitted in partial fulfilment of the requirements for the
degree of B.Sc. (Hons.) Computing Science AND Statistics
and Operations Research.*

Statement of Originality

I, the undersigned, declare that this is my own work unless where otherwise acknowledged and referenced.

Candidate Marc Ferriggi

Signed _____

Date April 4, 2019

Acknowledgements

Abstract

Contents

1	Introduction	1
2	Background and Literature Review	2
2.1	The LHC Machine	2
2.2	The Instruments Used to Gather Data	3
2.3	Feature Scaling and Reduction Techniques	4
2.4	Unsupervised Anomaly Detection Techniques	5
2.5	Applying Machine Learning to LHC Data	6
	References	7

List of Figures

2.1	The CERN Particle Accelerator Complex	3
-----	---	---

List of Tables

List of Abbreviations

AGM Abort Gap Monitor	4
BLM Beam Loss Monitors	3
BPM Beam Position Monitors	4
CERN European Organization for Nuclear Research	2
DBSCAN Density Based Spatial Clustering of Applications with Noise	5
Gy/s Grays per Second	4
IQC Injection Quality Check	3
LHC Large Hadron Collider	2
LOF Local Outlier Factor	5
LS Logging Service	2
MJ Mega Joule	2
MKD horizontally deflecting extraction kicker magnets	4
mm millimetres	4
PCA Principal Component Analysis	4
RF Radiofrequency	4
SPS Super Proton Synchrotron	2
TDI Beam Absorber for Injection	4
TeV teraelectronvolts	2
TL Transfer Line	4

Introduction

Background and Literature Review

2.1 The LHC Machine

The Large Hadron Collider (LHC) is a “two-ring-superconducting-hadron accelerator and collider” installed at the European Organization for Nuclear Research (CERN) between the years 1984 and 1989 [1]. The collider is 26.7km long and its purpose is to accelerate and collide two proton beams [2].

In order to fill the LHC to its required centre-of-mass energy of 14 teraelectronvolts (TeV), twelve injections from the Super Proton Synchrotron (SPS) consisting of a number of electron bunches of around 1 Mega Joule (MJ) of stored energy are required [3]. Thus, in order to fill the LHC, approximately 4 minutes per beam is required. Furthermore, the whole experiment process of filling the LHC, performing the required checks, running the tests and dumping the beam should take a theoretical minimum of 70 minutes [1]. This is expected to take around 6 times longer due to unsuccessful or anomalous proton injections [1].

Clearly, filling the LHC is a challenging task given the high energy of the beam, the very small apertures and the delivery precision’s tight tolerances. Thus, multiple sensors are installed around the CERN particle accelerator complex [4] which gather readings and data that can be used to check the quality of the injected beam.

For this particular study, data generated from the sensors around the injection from the SPS to the LHC will be of particular interest. This data is stored using CERN’s Logging Service (LS) [5]. While many studies have been made using this logged data and lots of statistical tests have been done with regards to injection quality checks for

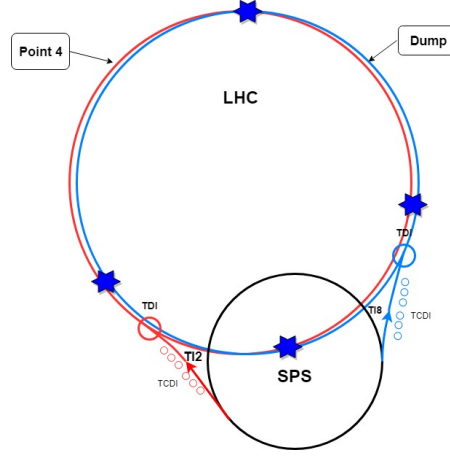


Figure 2.1: Diagram of the particular area of interest of the CERN Particle Accelerator Complex for this study

the LHC (such as [3] and [6]), no literature was uncovered where researchers used unsupervised machine learning methods to analyse this data. Figure 2.1 highlights the particular area of interest of this study.

The Injection Quality Check (IQC) software currently installed has a set of hard-coded rules for detecting anomalies in the SPS-LHC injection [3], however there are documented cases in the past where situations occurred which were outside the originally foreseen rules and were therefore not caught as anomalies.

2.2 The Instruments Used to Gather Data

Throughout this study, different data recorded as the beam leaves the SPS and enters the LHC will be used as input parameters to the chosen anomaly detection algorithms. This data was recorded using different sensors located in different parts of the injection life cycle. This section describes the different types of sensors that were used to collect the data, highlighting the particular details which need to be considered when analysing this data.

The Beam Loss Monitors (BLM) are some of the most safety critical modules of the LHC because a loss of a very small fraction of this beam may damage parts of the machine or cause a quench in the superconducting magnets [7]. A high beam loss reading could also indicate over-injection. In fact, an injection of a high intensity beam into the LHC is only allowed if there is a low intensity bunch circulating the LHC in order to

avoid settings errors [6]. The BLM module is the mostly used module in the current IQC checks [3]. The BLMs must be reliable; the probability of not detecting a dangerous loss was found to be 5×10^{-6} per channel and they are only expected to generate 20 false dumps per year [7]. The BLMs are extensively logged to a database for offline analysis [7].

For this particular study, the readings logged for the Beam Absorber for Injection (TDI) BLMs and the Transfer Line (TL) BLMs in TI2 and TI8 will be used (refer to Figure 2.1). These readings come in 10 second windows around the injection of a bunch in Grays per Second (Gy/s).

The Beam Position Monitors (BPM) were installed as a system for fast monitoring of the beam's position with respect to its orbit drift [8]. The trajectory offsets recorded by the BLMs in the transfer lines must be minimised in order to reduce losses [3]. In fact, if the change in orbit substantially exceeds its provided boundary values then the beam should be dumped [8] so as to not cause any damage to the equipment. Unlike the TDI BLMs, the BPM system is independent to the collimator system. For this study, the readings from the transfer line BPMs around TI2 and TI8 will be used (refer to Figure 2.1). Raw values for these readings are stored by the LS in millimetres (mm) and are logged every 1 - 5 seconds on average.

When filling the LHC, it is necessary to keep an abort gap of at least $3\mu\text{s}$ in order to accommodate for the horizontally deflecting extraction kicker magnets (MKD) rise time [9]. As the LHC is filling to nominal intensity, this gap will be populated with un-trapped particles and particles leaking out of their Radiofrequency (RF) buckets [9]. The Abort Gap Monitor (AGM) was hence specifically designed to measure this particle population in the abort gap [10]. This monitor can be found in Point 4 (refer to Figure 2.1) in the LHC [10]. The raw values extracted for this study are stored in number of particles and come in 10 second groups around the moment of injection.

2.3 Feature Scaling and Reduction Techniques

Feature Scaling and Feature Reduction are two important pre-processing steps that should be considered when using machine learning in the data science process. Standard Scaling in particular will be used in this study as a pre-processing step to Principal Component Analysis (PCA). Standard Scaling ensures that all the features have the

properties of a standard normal distribution [11], which is especially important since PCA involves finding the components that maximise the variance [12].

Apart from scaling, another challenge for outlier detection algorithms is data involving high dimensions since the contrast between different points diminishes as the number of dimensions increases [13]. This phenomenon is known as ‘The Curse of Dimensionality’ and a technique to reduce the effect of this phenomenon is to use a dimension reduction technique and run the outlier detection algorithm on this new lower-dimensionalised dataset. In this study, PCA will be used as a dimension reduction technique.

PCA uses statistical and mathematical techniques to reduce the dimension of large data sets, thus allowing a large data set to be interpreted in less variables called principal components [14]. This technique works with the hope that the variance explained by an acceptably small number of principal components is large enough to explain the underlying structure of the dataset reasonably [12]. In fact, this non-parametric method has been used as a means of revealing the simplified structures’ underlying complex datasets with minimal effort. The fact that this technique is non-parametric gives it the advantage that each result is unique and only dependent on the provided data set since no parameter tweaking is required [12] however, this is also a weakness of this technique as there is no way of exploiting prior expert knowledge on the data set.

2.4 Unsupervised Anomaly Detection Techniques

Unsupervised machine learning algorithms refer to the class of machine learning algorithms where only the input features are available to the learner as there is no access to output labels corresponding to each input feature vector, or the aim of the algorithm is simply to observe or detect patterns in the available data. A. Hyvärinen states in [15] that some of the goals of unsupervised learning include data visualisation, noise reduction, feature extraction and finding interesting components; all of which are of particular interest in this study.

In this study, Density Based Spatial Clustering of Applications with Noise (DBSCAN) and Local Outlier Factor (LOF) will both be used as unsupervised anomaly detection algorithms to detect and classify anomalous injections of the past year. Furthermore when working in 3 dimensions or less, these points can also be visualised to help the reader understand better the cause of these anomalies.

DBSCAN was created out of the necessity of having a clustering algorithm with the following requirements:

1. “Minimal requirements of domain knowledge to determine the input parameters,”
2. “Discovery of clusters with arbitrary shape,” and
3. “Good efficiency on large databases” [16]

DBSCAN manages to attain these requirements by viewing clusters as “areas of high density separated by areas of low density” [17]. The points with a lower density will thus be considered as anomalies when compared to the regular clusters which have a higher density. This algorithm also introduces the concept of *core samples* which was then used in the design of other unsupervised anomaly detection algorithms such as LOF.

The word *factor* in LOF refers to a “degree of outlier-ness” that this algorithm considers for each point in the data rather than using the concept that “being an outlier is binary” [18]. This algorithm uses a clustering technique which takes concepts from DBSCAN to measure the LOF of each point where a LOF value greater than 1 implies that the point has a lower density than its neighbours and is thus probably an outlier.

2.5 Applying Machine Learning to LHC Data

In 2017, Valentino *et. al.* released a paper on using anomaly detection techniques “to detect minor changes in the loss maps over time due to collimator settings errors or orbit variations” [2]. The authors used PCA as a dimension reduction technique and then applied LOF on the resulting 2 dimensional data. Their proposed method was shown to positively identify these anomalous loss maps based solely on BPM and BLM readings. Furthermore, they proposed using this technique to monitor losses during fills of the LHC.

References

- [1] L. Evans and P. Bryant, "LHC Machine," *JINST* 3, pp. 1–7, 2008.
- [2] G. Valentino, R. Brruce *et al.*, "Anomaly Detection for Beam Loss Maps in the Large Hadron Collider," presented at the 8th Int. Particle Accelerator Conference, 2017.
- [3] L. N. Drosdal, B. Goddard *et al.*, "Automatic Injection Quality Checks for the LHC," in *Proc. ICALEPCS*, 2011, pp. 1077–1080.
- [4] C. Lefevre, "The CERN Accelerator Complex," CERN, Tech. Rep., 2008.
- [5] C. Roderick, L. Burdzanowski, and G. Kruk, "The CERN Accelerator Logging Service - 10 Years in Operation: A Look at the Past, Present and Future," presented at the 14th Int. Conf. Accelerator & Large Experimental Physics Control Systems, 2013.
- [6] V. Kain, V. Baggiolini *et al.*, "Injection Beam Loss and Beam Quality Checks for the LHC," in *Proc. of IPAC*, 2010, pp. 1671–1673.
- [7] E. Holzer, B. Dehning *et al.*, "Beam Loss Monitoring System for the LHC," presented at the IEEE NSS, 2006.
- [8] *Protection of the CERN Large Hadron Collider*, ser. New Journal of Physics, vol. 8, no. 290, CERN, 2006. [Online]. Available: <http://www.njp.org/>
- [9] M. Meddahi, S. Bart Pedersen *et al.*, "LHC Abort Gap Monitoring and Cleaning," presented at the IPAC, 2010.
- [10] T. Lefevre, S. Bert Pedersen *et al.*, "First Operation of the Abort Gap Monitors for LHC," CERN, Tech. Rep., 2010.
- [11] (2019, April) Importance of Feature Scaling. [Online]. Available: scikit-learn.org
- [12] J. Shlens, "A Tutorial on Principal Component Analysis," April 2014.
- [13] A. Zimek, E. Schubert, and H.-P. Kriegel, "A Survey on Unsupervised Outlier Detection in High-Dimensional Numerical Data," *Statistical Analysis and Data Mining* 5, 2012.

- [14] M. Richardson, "Principal Component Analysis," May 2009, Class Lecture.
- [15] A. Hyvärinen, "Unsupervised Machine Learning," Lecture Notes, University of Helsinki.
- [16] M. Ester, H. Kriegel *et al.*, "A Density Based Algorithm for Discovering Clusters," in *Proc. KDD-96*, 1996, pp. 226–231.
- [17] (2018, November) Clustering. [Online]. Available: scikit-learn.org
- [18] M. Breunig, H. Kriegel, and J. Sander, "LOF: Identifying Density-Based Local Outliers," in *Proc. ACM SIGMOD 2000 Int. Conf. On Management of Data*, 2000.