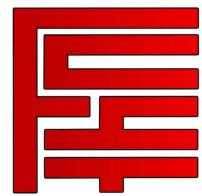




Mayor de San Simón University
Faculty of Science and Technology
Systems Engineering Degree



Driving Anomaly Detection

Degree project submitted in compliance with the requirements
to opt for the Degree in Systems Engineering

Presented by:

Evelyn CUSI LÓPEZ

Advisor:

Ph.D. Eduardo DI SANTI

COCHABAMBA - BOLIVIA

March, 2020

Agradecimientos

To my professor: Eduardo Di Santi, who gave me his valuable and selfless orientation and guidance in the preparation of this work, for his patience, support and ideas that motivated this research; to whom I owe much of my learning and my taste for the Artificial Intelligence's area.

To my family, for their constant understanding and encouragement, in addition to their unconditional support throughout my studies.

To my friends, who in one way or another supported me in carrying out this work.

To the Faculty of Science and Technology and the Mayor de San Simon University, for open me their doors and being the home of my training throughout my university career.

*To my mother, for being with me, for teaching me to grow and to get up,
for supporting and guiding me, for being the base that helped me get
here.*

Contents

Agradecimientos	iii
Abstract	xv
1 INTRODUCTION	1
1.1 Approach of the problem	1
1.2 General objective	2
1.3 Specific objectives	2
1.4 Justification	3
1.4.1 Practical justification	3
1.4.2 Methodological justification	3
1.5 Limits and scope	4
1.6 Research method	4
2 FUNDAMENTALS OF THE DETECTION OF ANOMALIES	5
2.1 Anomaly detection	5
2.2 Challenges in anomaly detection	7
2.2.1 Anomaly detection approaches	7
Supervised anomaly detection	7
Semi-supervised anomaly detection	8
Unsupervised anomaly detection	8
2.3 Related work	8
2.4 Focus on the problem	9
3 MACHINE LEARNING FOR ANOMALY DETECTION	11
3.1 Supervised Learning, Unsupervised Learning and Semi Supervised Learning	11
3.1.1 Supervised Learning	11
3.1.2 Unsupervised Learning	12
3.1.3 Semi Supervised Learning	12
3.2 Generative and Discriminative Models	13
3.3 Artificial neural networks	14
3.3.1 Neurons or nodes	14
3.3.2 Types of activation functions	16
Sigmoid activation function (logistics function)	16
Hyperbolic Tangent Function - Tanh	16
Rectified Linear Unit function (ReLU)	16
Leaky ReLU (LReLU)	18

Exponential Linear Unit Function (ELU)	19
3.3.3 Architecture of the Neural Networks	19
3.3.4 Learning process of Neural Networks	20
Backpropagation	20
3.4 Types of Neural Networks	21
3.4.1 Autoencoders	21
3.4.2 Convolutional neural networks	21
3.4.3 Recurrent neural networks	22
3.4.4 LSTM	23
3.4.5 GRU	25
3.5 Anomaly detection techniques	25
3.5.1 One-Class SVM	26
3.5.2 Isolation Forest	26
3.5.3 Autoencoders	28
3.6 Evaluation Metrics	29
3.6.1 Classification Accuracy	29
3.6.2 Logarithmic Loss	29
3.6.3 Confusion Matrix	30
3.6.4 Area Under Curve (AUC) Recall or Sensitivity or True Positive Rate (TPR)	30
Specificity or True Negative Rate (TNR)	31
ROC (Receiver Operating Characteristics)	31
3.6.5 F1 Score	31
Precision	31
Recall	32
4 DATA CAPTURE AND PREPARATION	33
4.1 Data Capture	33
4.2 Data preparation	34
4.2.1 Data selection	35
4.3 Data preprocessing	36
4.3.1 Data cleaning	36
Compensation techniques for incomplete records	36
Remove noise from data (smoothing)	37
Fusion or data integration	39
Data transformation	40
Data reduction	41
5 GENERACIÓN DEL MECANISMO DE DETECCIÓN DE ANOMALÍAS	51
5.1 Entorno de desarrollo	51
5.2 Conjunto de datos normales y anómalos	52
5.2.1 Generación de series temporales	52
5.3 Modelo de detección de anomalías	53
5.3.1 Modelo del comportamiento normal	54
Arquitectura del modelo	54

5.3.2	Método de detección de anomalías	57
Umbralización	57	
Isolation Forest	61	
One-Class SVM	63	
Evaluación del método de detección de anomalías	64	
6	RESULTADOS Y EVALUACIÓN	69
6.1	Evaluación de desempeño	69
6.1.1	Evaluación en términos de rendimiento de detección	69
6.2	Resultados	70
6.2.1	Detección de anomalías del tipo zig zag	70
6.2.2	Detección de anomalías del tipo giros a alta velocidad	72
6.2.3	Detección de anomalías del tipo frenos en seco	73
6.2.4	Detección de falsos positivos	75
7	CONCLUSIONES Y TRABAJOS FUTUROS	77
7.1	Conclusiones	77
7.2	Trabajos futuros	78
BIBLIOGRAPHY		81
References		81
A	Experimentos de diferentes arquitecturas para los autoencoders	87
A.1	Redes densas	87
A.1.1	Redes densas para 3 componentes	87
A.1.2	Evaluación redes densas	88
A.2	Redes convolucionales	88
A.2.1	Redes convolucionales para 3 componentes	88
A.2.2	Evaluación redes convolucionales	89
A.3	Redes recurrentes	89
A.3.1	Redes recurrentes para 3 componentes	89
A.3.2	Evaluación redes recurrentes	90
B	Arquitectura del Sistema de Demostración	91
B.1	Arquitectura Física	91
B.2	Arquitectura Lógica	92

List of Figures

1.1 Deaths due to traffic accidents by region depending on user's type (de Salud (OMS), n.d.).	2
1.2 Problem tree (Own elaboration).	3
2.1 Example of point anomalies in a 2-dimensional data set (Varun & Arindam, 2009).	6
2.2 Contextual anomaly t_2 in a temperature time serie (Varun & Arindam, 2009).	6
2.3 Collective anomaly corresponding to premature atrial contraction in a human electrocardiogram (Varun & Arindam, 2009).	7
2.4 Proposed anomaly detection method (Own elaboration).	10
3.1 Graphic of a biological neuron. Reproduced from (Wikipedia, n.d.).	14
3.2 Graphic of an artificial neuron. Reproduced from (Jayesh, n.d.).	15
3.3 Activation functions (Jing & Guanci, 2018).	17
a Function Sigmoid	17
b Function Tanh	17
c Function ReLU	17
d Function Leaky ReLu	17
e Function ELU	17
3.4 Architecture of an artificial neuron. Reproduced from (Michael, 2015).	19
3.5 Graphic of an Autoencoder (Own elaboration).	21
3.6 Architecture of a Convolutional Neural Network (CNN) (Muhammad, n.d.).	22
3.7 Sequential processing in a recurrent neural network (RNN) (Olah, n.d.).	23
3.8 LSTM structure. Played from Yan (Yan, n.d.).	24
3.9 GRU structure. Reproduced from (Zhang, Lipton, Li, & Smola, 2019).	25
3.10 One-Class SVM. Reproduced from (Alashwal, Bin D., & Othman, 2006)	26
3.11 Performance comparison between the One-Clas SVM and Isolation Forest algorithms. Reproduced from (0.22, n.d.).	27
3.12 Detection of anomalies with autoencoder (Own elaboration).	28
3.13 Example of an AUC-ROC curve (Özler, n.d.).	32
4.1 Data collection, with interval of one second (Own elaboration).	34
4.2 Windshield cell mount, horizontal position (Own elaboration).	34
4.3 Fragment of data set obtained (Own elaboration).	35
4.4 Graph of sensors captured in different positions (Own elaboration).	36
4.5 Histogram of data set's frequencies (Own elaboration).	38
4.6 Table of data set's statistical results (Own elaboration).	38
4.7 Rule 68-95-99.7 (Galarnyk, n.d.).	39

4.8	Result of applied Rule 68-95-99.7 on values' accelerometer sensors in Z (Own elaboration).	40
4.9	Division of the dataset (Own elaboration).	43
4.10	Visualization of the captured driving parameters (Own elaboration).	44
4.11	Graph resulting from applying different types of normalizations to data set (Own elaboration).	45
a	Min max Scaler	45
b	Standard Scaler	45
c	Max Scaler	45
d	Robust Scaler	45
4.12	Gráfico de la varianza vs. el número de componentes (Elaboración propia).	49
5.1	Gráfica resultante de diferentes tamaños de series de tiempo (Elaboración propia).	53
a	Serie de tiempo de 2 pasos.	53
b	Serie de tiempo de 3 pasos.	53
c	Serie de tiempo de 4 pasos.	53
d	Serie de tiempo de 5 pasos.	53
5.2	Resultados (Elaboración propia).	58
a	Resultados de la red NN_33	58
b	Resultados de la red CNN_33	58
c	Resultados de la red RNN_33	58
5.3	Curva de los valores de reconstrucción obtenidos con el modelo del comportamiento normal (Elaboración propia).	59
5.4	Resultados de la obtención de codos con diferentes valores de Sensibilidad, para los valores de reconstrucción obtenidos con el modelo del comportamiento normal (Elaboración propia).	60
5.5	La figura de la izquierda muestra el aislamiento de una anomalía, que requiere solo tres particiones. A la derecha, el aislamiento de un punto normal requiere seis particiones (Wolpher, n.d.).	61
5.6	Representación gráfica del modelo de comportamiento normal o autoencoder (Elaboración propia).	62
5.7	Representación gráfica del error de reconstrucción usado para el entrenamiento de los bosques de aislamiento y los SVM de una clase (Elaboración propia).	63
5.8	Mecanismo de detección de anomalías (Elaboración propia).	66
6.1	Resultados de la detección de anomalías del tipo zig zag (Elaboración propia).	71
6.2	Resultados de la detección de anomalías del tipo giros a alta velocidad (Elaboración propia).	73
6.3	Resultados de la detección de anomalías del tipo frenos en seco (Elaboración propia).	74
6.4	Resultados de la detección de falsos positivos (Elaboración propia).	75
B.1	Arquitectura física del Sistema de Demostración (Elaboración propia).	91
B.2	Arquitectura Lógica del Sistema de Demostración (Elaboración propia).	92

List of Tables

3.1	Confusion matrix, for a binary classification (Own elaboration)	30
4.1	Table of dataset's division (Own elaboration)	43
4.2	Tabla con estadísticas descriptivas de los datos escalados con diferentes técnicas (Elaboración propia)	47
5.1	Tabla del conjunto de anomalías (Elaboración propia)	52
5.2	Tabla de los métodos comparados (Elaboración propia)	54
5.3	Arquitectura densa para una secuencia de 3 pasos y 3 componentes principales (Elaboración propia)	55
5.4	Arquitectura convolucional para una secuencia de 3 pasos y 3 componentes principales (Elaboración propia)	55
5.5	Arquitectura recurrente para una secuencia de 3 pasos y 3 componentes principales (Elaboración propia)	56
5.6	Evaluación de las redes NN_33, CNN_33 y RNN_33 (Elaboración propia)	56
5.7	Evaluación de la detección de anomalías para cada codo obtenido con los diferentes valores de sensibilidad (Elaboración propia)	60
5.8	Evaluación de la detección de anomalías usando Isolation forest para valores compresos (Elaboración propia)	63
5.9	Evaluación de la detección de anomalías usando Isolation forest para errores de reconstrucción (Elaboración propia)	63
5.10	Evaluación de la detección de anomalías usando One-Class SVM para valores compresos (Elaboración propia)	64
5.11	Evaluación de la detección de anomalías usando One-Class SVM para el error de reconstrucción del autoencoder (Elaboración propia)	64
5.12	Comparación de los mejores métodos de detección de anomalías (Elaboración propia)	65
6.1	Matriz de confusión, para el mecanismo de detección de anomalías (Elaboración propia)	70
6.2	Resultados anomalías tipo Zig Zag (Elaboración propia)	72
6.3	Resultados del tipo Giros a alta Velocidad (Elaboración propia)	73
6.4	Resultados del tipo Frenos en seco (Elaboración propia)	74
A.1	Arquitectura densa para 3 componentes principales (Elaboración propia)	87
A.2	Tabla de evaluación de redes densas (Elaboración propia)	88
A.3	Arquitectura convolucional para 3 componentes principales (Elaboración propia)	88

A.4 Tabla de evaluación de redes convolucionales (Elaboración propia).	89
A.5 Arquitectura recurrente para 3 componentes principales (Elaboración propia).	89
A.6 Tabla de evaluación de redes recurrentes (Elaboración propia).	90

Resumen

This paper describes the development of a mechanism to detect driving anomalies, which is implemented using a mobile device and Machine Learning techniques.

The objective is create a tool capable to find anomalous behaviors in the driving of human or autonomous agent, having a previous knowledge of normal driving conducts of them. Likewise it presents a background of driving anomalies detection's researches around the world, the driving parameters obtained by the mobile device are analized, and a proposal is presented to identify anomalies through the use of Neural Networks and Isolation Forests, a method of Machine Learning that is commonly used for anomalies' detection.

The work has two main parts: a model adjusted to the normal driving behavior of an agent, and a method of detecting anomalies, which were iteratively trained with 30,000 samples, which correspond only to normal driving behavior.

The detection accuracy of the complete mechanism proposed in this document is 67.68%, which was evaluated with 44040 samples, of which 164 correspond to anomalous samples, thus being one of the most outstanding contributions for the detection of driving anomalies with a semi supervised approach.

Chapter 1

INTRODUCTION

This document describes the development of a method for detecting anomalies in car driving. The use of Machine Learning techniques is proposed to generate a mechanism that identifies driving anomalies, so that they can be used to promptly alert agents and thus correct their driving behaviors.

The main idea is to generate a model that learns the normal driving behavior of a specific agent, to later autonomously detect those unexpected behaviors and report them as anomalies, so that a traffic accident can be avoided or the effects of it reduced.

1.1 Approach of the problem

Due to the serious consequences they cause on people and the high economic costs associated with them, traffic accidents are classified as a global social and public health problem.

According to the World Health Organization (WHO) each year there are approximately 1.25 million deaths due to traffic accidents, adding that half of all these victims are pedestrians, cyclists and motorcyclists (See Figure 1.1 page 2). It can also be said that they are one of the most important causes of death in the world, and the main cause of death among people between the ages of 15 and 29.

On the other hand, according to the Cochabamba Transit Operating Unit, the accidents recorded in 2017 caused the death of 200 people and left approximately 2200 injured.

Figure 1.2 shows the causes for which a traffic accident is caused, it can be seen that a large part of these are due to the human factor, however there are others that involve environmental and mechanical factors, so it is impossible completely avoid them.

That is why it is necessary to have mechanisms to prevent and/or act in a timely manner against possible traffic accidents, which is why this work focuses on studying driving behaviors, in order

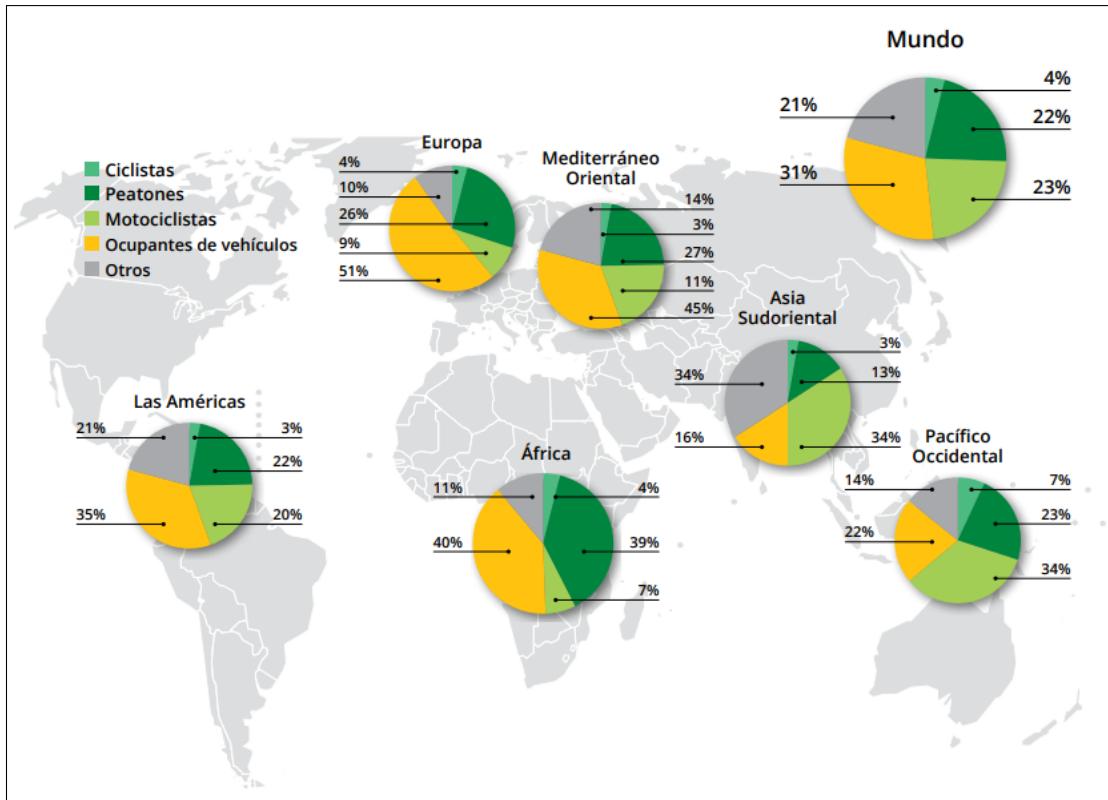


FIGURE 1.1: Deaths due to traffic accidents by region depending on user's type (de Salud (OMS), n.d.).

to generate alerts when finding a behavior anomalous in driving, so that the effects of it can be avoided or in any case minimized.

1.2 General objective

The general objective of the present work is to develop a mechanism for detecting driving anomalies, through the use of a mobile device and Machine Learning algorithms, in order to alert in a timely manner the finding of an anomalous driving pattern, such as tiredness, drunkenness, or health problems, eg. epilepsy.

1.3 Specific objectives

- Capture the driving parameters of a driver by using sensors of a mobile device.
- Scale parameters using data pre-processing techniques.
- Generate a Machine Learning model that adjusts to normal driving behavior.
- Define an anomaly detection method to generate an abnormal driving alert.

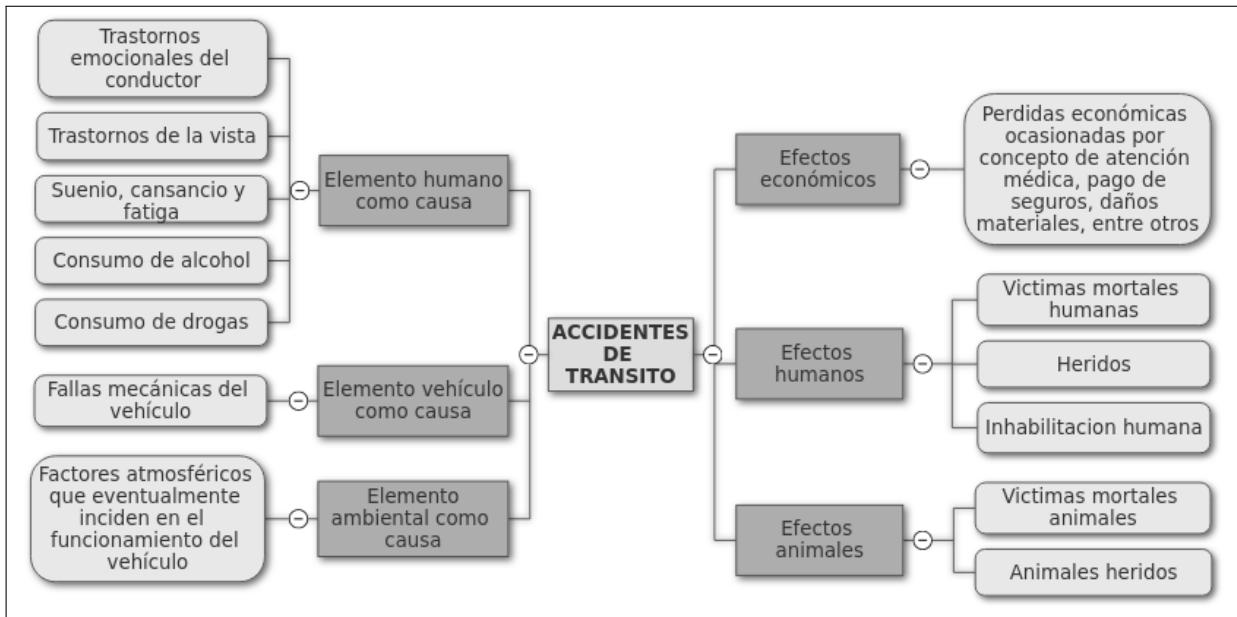


FIGURE 1.2: Problem tree (Own elaboration).

- Evaluate the anomaly detection method with new samples

1.4 Justification

Traffic accidents charge an unacceptable number of victims each year, especially in the poorest regions of the world. This is due to various aspects, but the main one lies in the low level of citizen awareness that exists, which leads to many people driving under the influence of alcohol, with excessive speed, manipulating their mobile devices, among others. For this reason, this work seeks to establish patterns of driving behaviors through the use of a mobile device and Machine Learning techniques, so that it is possible to detect timely driving anomalies.

1.4.1 Practical justification

Detect driving anomalies allows the authorities or agents to generate a timely alert so that they can correct their driving behaviors quickly. This way we can avoid traffic accidents or minimize their effects, thus reducing the amount of damage, both material and personal.

1.4.2 Methodological justification

The study carried out in the development of this research allows highlighting the efficiency of Artificial Intelligence techniques in detection of anomalies.

1.5 Limits and scope

Due conducting field tests for this investigation is quite dangerous, the examples of anomalous driving were limited to:

- Dry brakes.
- Turn right and left at high speed.
- Turn abrupt zig zag.

Thus, experiments and tests were performed only on a small set of anomalous samples, therefore it is not expected that the proposed detection model will work correctly on those examples that were not considered.

1.6 Research method

The present study was conducted with an experimental approach, with the following hypothesis:

Is it possible to detect driving anomalies by using a mobile device and Machine Learning algorithms?

Chapter 2

FUNDAMENTALS OF THE DETECTION OF ANOMALIES

This chapter discusses necessary concepts that are needed to understand the detection of anomalies, as well as different projects and investigations carried out in field of the detection of driving anomalies to date.

2.1 Anomaly detection

To understand what anomaly detection implies, it is necessary to assimilate what an anomaly is, and in what ways these can occur. Therefore, it can be said that anomalies, or outliers, are patterns in the data that do not fit a well-defined notion of normal behavior.

Anomalies can be classified into one of the following three categories:

1. **Point anomalies:** Point anomalies are simply unique and anomalous instances within a larger data set, that is, they are separated from the rest of the data. For example, in Figure 2.1, points o_1 , o_2 and region O_3 are outside the limits of normal regions (N_1 and N_2), and therefore are point anomalies because they are different from the normal data set.

This type of anomaly is considered the simplest and is the focus of most research focused on anomaly detection.

2. **Contextual (or conditional) anomalies:** These are points that are only considered anomalous in a specific context. The notion of this context is induced by the structure in the data set and must be specified as part of the problem formulation.

These types of anomalies have been explored more frequently in time series and spatial data. Figure 2.2 shows an example of a time series of the monthly temperature of an area in the last 5 years, it should be taken into account that the temperature at time t_1 is the same as at time t_2 , but occurs in a different context , therefore t_2 is considered an anomaly.

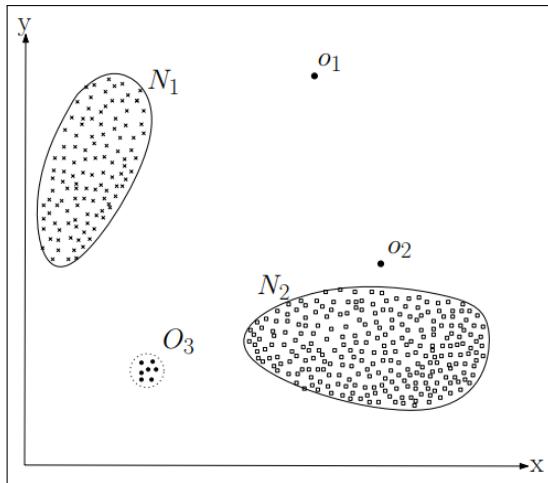


FIGURE 2.1: Example of point anomalies in a 2-dimensional data set (Varun & Arindam, 2009).

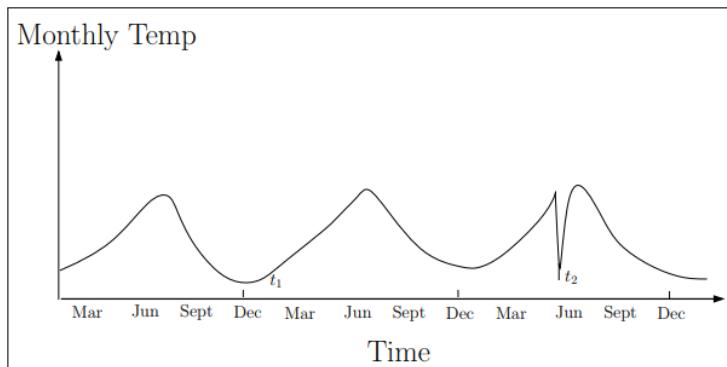


FIGURE 2.2: Contextual anomaly t_2 in a temperature time serie (Varun & Arindam, 2009).

3. **Collective anomalies:** If a collection of related data instances is anomalous with respect to the entire data set, it is called a collective anomaly. The instances of individual data in a collective anomaly may not be anomalies by themselves, but their joint appearance as a collection is anomalous.

Figure 2.3 illustrates an example showing a human electrocardiogram output, it can be noted that the region highlighted in red denotes an anomaly because there is the same low value for an abnormally long time (corresponding to a premature atrial contraction). It should be taken into account that this low value by itself is not considered an anomaly.

In relation to the above, it can be defined as detection of anomalies, or outliers, to the identification of data points, elements, observations or events that do not fit the expected pattern of a particular group.

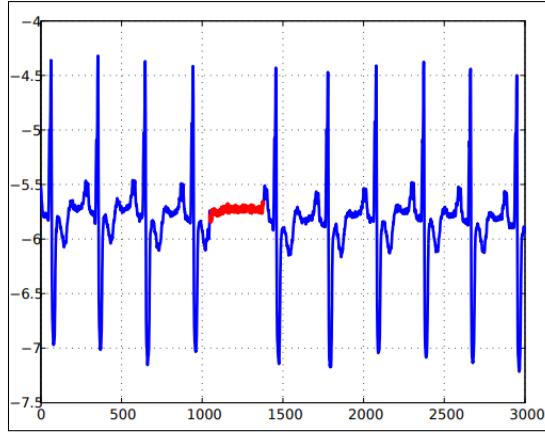


FIGURE 2.3: Collective anomaly corresponding to premature atrial contraction in a human electrocardiogram (Varun & Arindam, 2009).

Anomaly detection is used in different application domains, for example: image processing, card fraud detection, network intrusion detection systems, etc.

2.2 Challenges in anomaly detection

On an abstract level, anomaly detection may seem like a simple task. However, it can be a very challenging task. Below are some of these challenges.

- The definition of normal regions is quite difficult. In many cases, the boundaries between anomalies and normal data are not accurate. Therefore, normal observations could be considered anomalies and vice versa.
- Therefore, normal observations could be considered anomalies and vice versa.
- Most of the time the approaches for detecting anomalies in a specific field cannot be used in another field.
- The low availability of positive samples (anomalies) for training and validation of anomaly detection model.

2.2.1 Anomaly detection approaches

The approaches that can be used for this purpose are classified into following categories:

Supervised anomaly detection

The use of supervised learning techniques requires the availability of a set of labeled training data, for both normal and anomalous classes. The main focus is to build a predictive model for normal classes vs. anomalies, then take any instance of unseen data, compare with model and determine to which class it belongs.

There are two main drawbacks that arise with the use of this technique.

- The number of anomalous instances is much lower than that of the normal instances, which creates an imbalance of class distribution during training.
- Obtaining accurate and representative labels, particularly for anomaly class, is a challenge.

Semi-supervised anomaly detection

These techniques require a training set with labeled instances, but only for the normal class, this makes its use more applicable than the supervised techniques, since no labels are required for the anomaly class.

The typical approach used in these techniques is to build a model for the class corresponding to normal behavior, and use model to identify anomalies in the test data.

Unsupervised anomaly detection

Techniques that operate in an unsupervised manner do not require training data, which is why they are the most widely used. These techniques assume that normal instances are much more frequent than anomalies in test data, in case this assumption is not true, such techniques suffer from a high rate of false alarms.

2.3 Related work

The identification of abnormal driving behaviors is an indispensable part of improving driving safety, however, as previously described, this is not a simple task. In recent years, several techniques have been proposed to detect driving behaviors. This section is dedicated to review them.

Dang-Nhac et al. (2018), propose a combined system that consists of two modules: one to detect the type of vehicle of users and the other to detect events of instant driving, regardless of the orientation and position of smartphones, this system It achieves an average accuracy of 98.33% in detection of vehicles's type (car, motorcycle, bicycle, among others) and an average accuracy of 98.95% in recognition of motorcyclist driving events when using Random Forest as a classifier.

On the other hand Ferreira et al. (2017) present a quantitative evaluation of 4 Machine Learning algorithms (Bayesian Network BN, Artificial Neural Network ANN, Random Forest RF and Support Vector Machine SVM) with different configurations, applied in the detection of 7 types of driving events, between events normal and aggressive, using data collected from 4 Android smartphone sensors (accelerometer, linear acceleration, magnetometer and gyroscope); resulting in the gyroscope and accelerometer being the best sensors to detect driving events and that

Random Forest (RF) is by far the best-performing Machine Learning Algorithm, followed by the simplest form of ANN the Multi Layer Perceptron (MLP).

Johnson and Trivedi (2011) propose the MIROAD system which shows that Dynamic Time Warping (DTW) is a valid algorithm to detect potentially aggressive driving maneuvers, where almost all aggressive events (97%) were correctly identified, using the set of T sensors (accelerometer, gyroscope and tone of voice). Likewise, in the work of Kridalukmana, Yan-Lu, and Naderpour (2017), a system focused on developing driver awareness through notifications in critical situations that can trigger unsafe driving maneuvers is proposed, using a model to detect dangerous situations based on Object-Oriented Bayesian Network (OOBN).

As well as the works presented previously there is a large number of works (Bhoyar, Lata, Katkar, Patil, & Javale, 2013; Chen, Yu, Zhu, Chen, & Li, 2015; Eren, Makinist, Akin, & Yilmaz, 2012; Boonmee & Tangamchit, 2009; Koh & Kang, 2015) that use smart phone sensors (accelerometer and gyroscope) for aggressive driving detection, due of advantage of not buying or installing any device and besides being highly portable, however it depends a lot on the performance of GPS receiver and is not applicable in areas not available for GPS.

There are also other approaches to detection of aggressive driving of a driver, for instance in the work of Who-Lee, Sik-Yoon, Min-Song, and Ryoung-Park (2018), a method based on Convolutional Neural Network (CNN) is proposed to detect the emotion of aggressive driving, by using a driver's facial images obtained with a NIR light camera and a thermal camera.

2.4 Focus on the problem

It is clear that this topic was extensively researched and that it has a wide variety of solution proposals, however most of these are based on detection by supervised learning techniques, which presents the great disadvantage of requiring data labeled to generate the model detection; In addition, much of related work proposes generalized models for detection and not specific models for each agent, which is crucial because each agent has individual driving behaviors and different conditions driving, that is, the driving of an agent that circulates through paved avenues will be different from driving of an agent that circulates through cobble streets or driving of an agent that circulates through avenues or busy streets will be different from that of the agents that circulate through relatively decongested streets.

The proposal made in this work is intended to provide a prototype of a tool that helps analyze the sequences of motion sensors of a mobile device and allows detecting anomalies from information obtained from this analysis.

To achieve this goal, the data set of motion sensors of a smartphone is captured, using a mobile application, then data is divided and prepared with data pre-processing techniques.

Once these phases are completed, a model is trained with the training set and validated with the development set. Finally, an optimal model and a technique are chosen to classify outliers, in order to cover a full range of normal behaviors and exclude anomalies as accurately as possible. This method can best be seen in Figure 2.4.

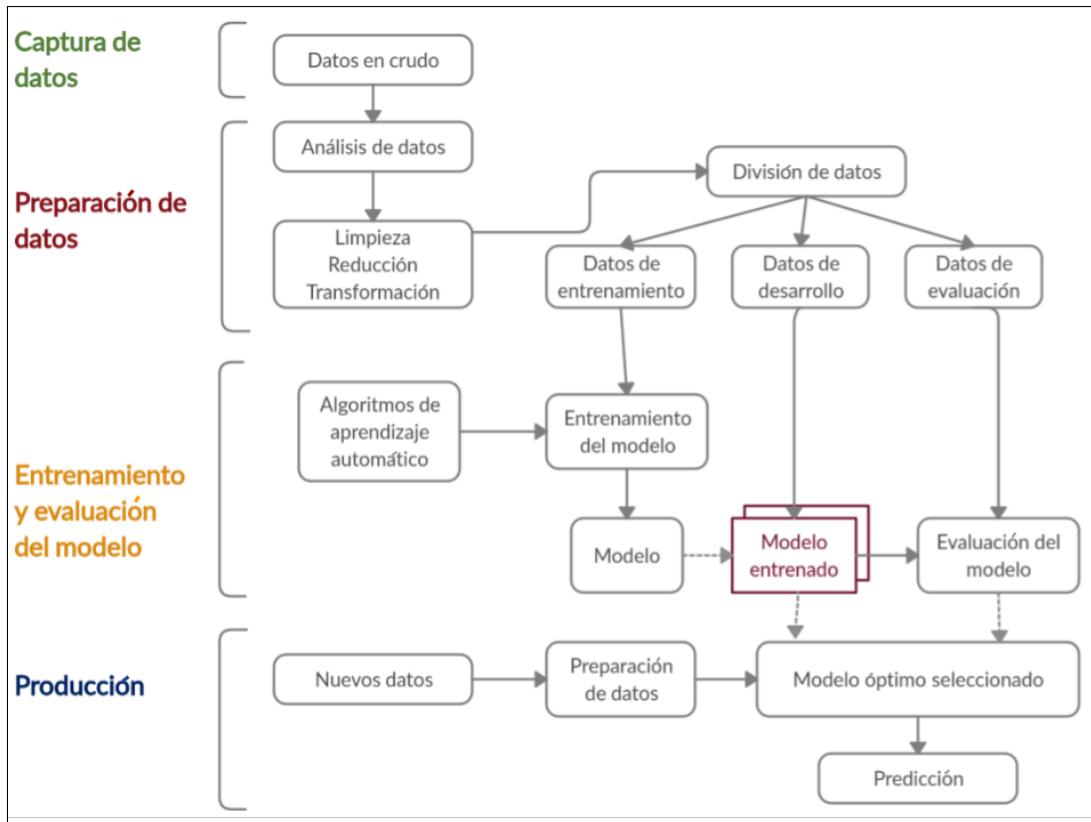


FIGURE 2.4: Proposed anomaly detection method (Own elaboration).

Chapter 3

MACHINE LEARNING FOR ANOMALY DETECTION

The term Machine Learning refers to the automatic detection of significant patterns within a data set (Shai & Shai, 2014). In recent decades it has become a common tool in almost any task that requires the extraction of information from a large amount of data, which is why it has become one of the fastest growing areas of information technology.

Although Machine Learning can solve some problems that are solved with traditional algorithms, it has overcome these problems such as image recognition, voice, language, writing, games, robotics, data analysis, time series analysis, etc. From this perspective, it is expected that through the application of Machine Learning a model can be generated that fits to a normal behavior expected for the agent.

Therefore, this chapter details the theoretical bases necessary to address the development of the method for driving anomaly detection. First, the different learning paradigms that exist and the different model approaches are described, then, a description of the functioning of neural networks is made and the different types of networks are shown, as well as different techniques of anomaly detection that exist and finally shows different evaluation metrics presented by machine learning models.

3.1 Supervised Learning, Unsupervised Learning and Semi Supervised Learning

There are several ways to classify learning paradigms that exist, however in this work only supervised, unsupervised and semi-supervised will be treated.

3.1.1 Supervised Learning

Supervised Learning is one that has input variables (X) and an output variable (Y), this type of learning uses an algorithm to learn the mapping function from input to output.

$$Y = f(X) \quad (3.1)$$

The objective of this type of learning is to approximate the mapping function so that when you have new input data (X) you can predict the output variables (Y) for that data.

This type of learning addresses two types of problems: classification and regression. Classification problems are those where the output variable is a category, such as: "Red", "Blue", or "Healthy", "Sick", on the other hand in Regression problems the output variable is a real value, such as: "price" or "height". Some of the most common types of problems built on classification and regression include recommendation and prediction of time series.

3.1.2 Unsupervised Learning

On the other hand, Non-Supervised Learning is one where there is only input data (X) and there are no corresponding output variables, its main objective is to model the structure or underlying distribution in data to learn more about them.

As for learning problems without supervision, they can be grouped into two: grouping and association. **Grouping** is one where you want to discover the groupings inherent in data set, such as grouping customers by purchasing behavior. On the other hand, **Association** is one that wishes to discover rules that describe large portions of its data, for example, people who buy X also tend to buy Y. Some of the most popular unsupervised learning algorithms are: Kmeans (for clustering problems) and Apriori algorithm (for learning problems of association rules).

3.1.3 Semi Supervised Learning

Finally, there is Semi-supervised Learning, which covers those problems where there is a large amount of input data (X) and only some of data is labeled (Y). These types of problems are between supervised and unsupervised learning, it is also important to point out that many of Machine Learning's problems in the real world are in this area, this is because it is expensive or it may take a long time to label the data set, while unlabeled data is cheap, in addition to being easy to collect and store. These types of problems can use a combination of supervised and unsupervised techniques to be solved.

Since Supervised Learning methods require a large amount of labeled training data, it is important to clarify that the collection of negative samples (abnormal driving) is difficult and risky for this particular study; In addition, the supervised approach has a potential limitation, which is: the detection of new atypical patterns, this because the resulting model is only trained to recognize a limited set of anomalous patterns, so at the time a new one is presented pattern this model will be unable to recognize it.

On the other hand, unsupervised approach has advantage of not requiring tagged information, however it often suffers from high false alarm rates and low detection rates (Xue, Shang, & Feng, 2010).

In many applications, including the one of present study, normal samples are easy to obtain, while anomalous ones are quite difficult to obtain, consequently, for implementation of this study, the application of the Semisupervised approach has been chosen. Thus, as mentioned in Chapter 2, Semi-supervised anomaly detection approach only has normal samples in the training set; that is, information about anomalies cannot be obtained, therefore unknown samples are classified as outliers, as long as their behavior is very different from that of normal samples already known.

As mentioned in this section, all of these learning approaches are based on generating a Model capable of helping either classification, grouping, etc; However, there is more than one type of model. The following section will detail in detail the different types of models that exist.

3.2 Generative and Discriminative Models

When using Machine Learning there are two main approaches to understand (model) the real world and make decisions. These two approaches are discriminative and generative models. More formally the generative and discriminative models represent two different strategies to estimate the probability that a particular object belongs to a category (Hsu & Griffiths, 2010).

Discriminative models are based on conditioned probability $P(Y|X)$, that is, they learn a direct map of a set of characteristics X to labels of Y classes. These types of models try to model simply depending on observed data (set of data), they also make less assumptions about distributions; however, they depend largely on quality of data. Some examples of discriminative models are: Logistic Regression, SVM (Support Vector Machine), Neural Networks, Random Forest, among others.

On the other hand the **generative models** point to a complete probabilistic description of data set, its objective is to develop joint probability distribution $P(X, Y)$, either directly or by calculating $P(Y|X)$ and $P(X)$, then infer conditional probabilities required to classify new data. These models help to specify the uncertainty of a model, some examples of generative models are: Gaussian Mixture Model, Hidden Markov Model, Restricted Boltzmann Machine, Generative Adversarial Networks (GAN), among others.

Discriminative models have been at the forefront of Machine Learning's success in recent years, since these models make predictions that depend on a given input, although they cannot generate new samples or data, so in the present study preference will be given to use of discriminative models.

Next, we will review the fundamental theoretical bases of some Semi-Supervised Learning techniques with a discriminative approach, to later detail what type of algorithms will be applied in the method proposed in this research work.

3.3 Artificial neural networks

The Artificial Neural Network or ANN¹ is a paradigm of information processing inspired by the way in which biological nervous system processes information. It consists of a large number of highly interconnected processing elements (neurons) that work in unison to solve a specific problem.

3.3.1 Neurons or nodes

Biological neurons (nerve cells) are fundamental units of brain and nervous system. Neurons are the cells responsible for receiving sensory information from external world through dendrites, processing it, and exiting through the axon (See Figure 3.1).

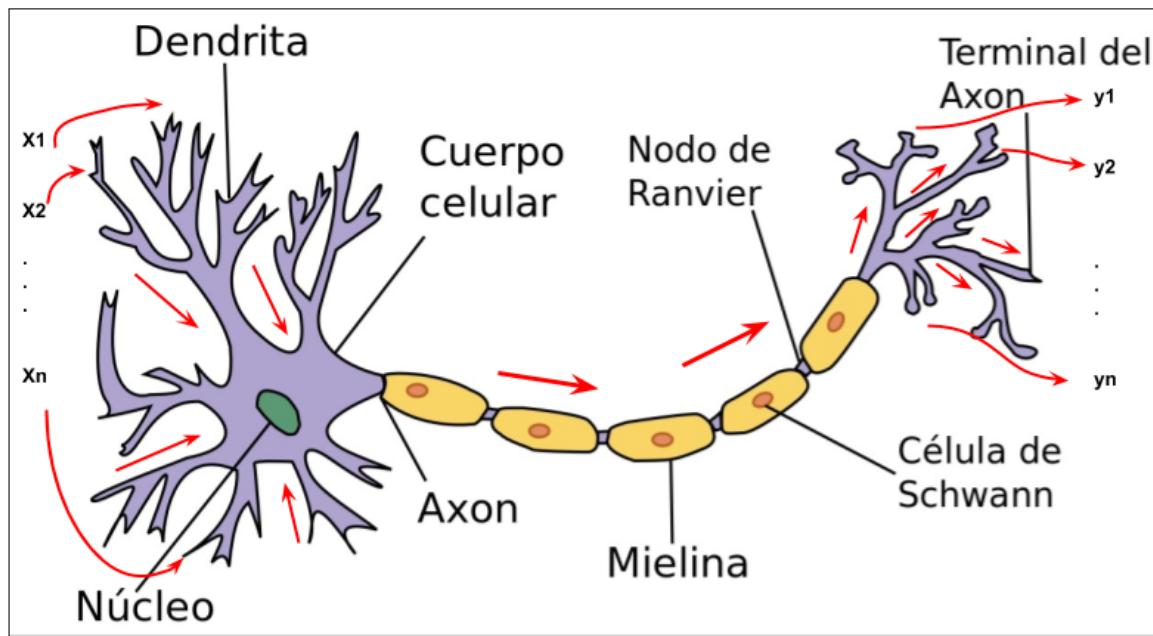


FIGURE 3.1: Graphic of a biological neuron. Reproduced from (Wikipedia, n.d.).

A brain neuron can receive about 10,000 entries and in turn send its output to hundreds of neurons.

The connection between neurons is called a **synapse**, this is not a physical connection, due there is 2 mm. of separation between neurons. These connections are unidirectional, where

¹ANN, Artificial Neural Network

information's transmission is done electrically inside the neuron and chemically between neurons, thanks to neurotransmitters.

An **artificial neuron** is an elementary processor, because it processes a vector $x(x_1, x_2, \dots, x_n)$ of inputs and produces a unique response or output. The main elements of an artificial neuron are the following:

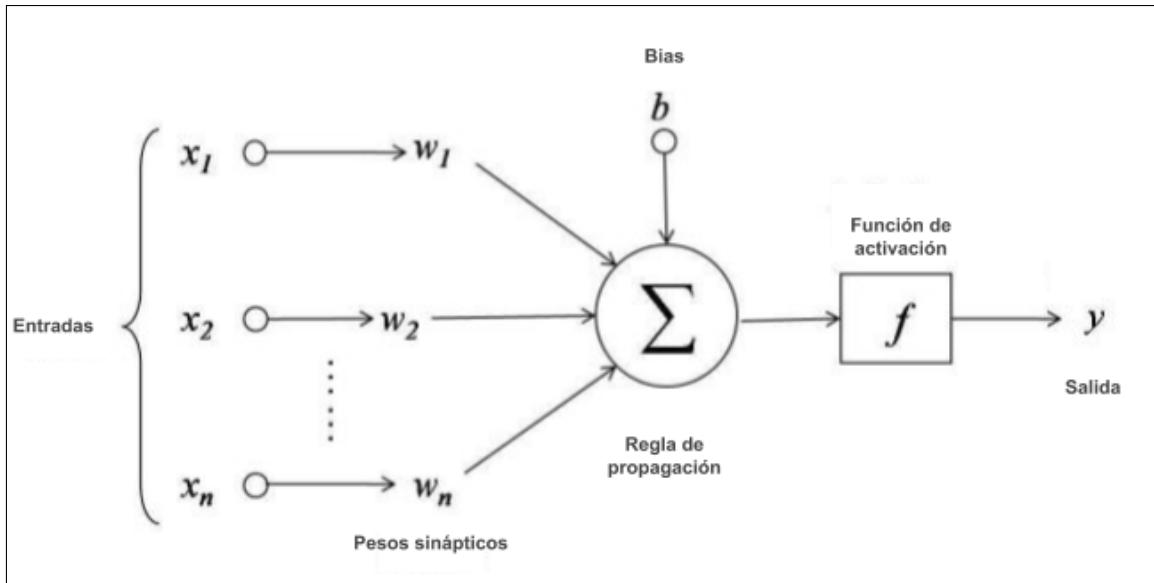


FIGURE 3.2: Graphic of an artificial neuron. Reproduced from (Jayesh, n.d.).

- **Inputs** that receive data from other neurons, these inputs would be dendrites of a biological neuron.
- **Synaptic weights** w_{ij} . In an artificial neuron, those inputs that come from other neurons are assigned a weight (importance factor). This weight is a numerical value that is modified during the training process of a neural network, and therefore it is here that information that makes the network serve one purpose or another is stored.
- **Propagation Rule**. With inputs and synaptic weights, some type of operation is usually done to obtain the potential postsynaptic value; one of the most common operations is to add up all entries, but taking into account importance (synaptic weight) of each one; This operation is called *weighted sum* 3.2, however other operations are also possible. Another propagation rule that is usual is Euclidean distance.

$$h_i(t) = \sum_j w_{ij}x_j \quad (3.2)$$

- **Activation Function**. The value obtained with the propagation rule is filtered through a function known as the *activation function* and is what gives the output of neuron. The activation function is important because it is the one that decides whether a neuron should

be activated or not, and if this function is not applied, the output signal of neuron would simply be a linear function.

3.3.2 Types of activation functions

There are different activation functions, then only the most used in field of neural networks will be presented.

Sigmoid activation function (logistics function)

A sigmoid function is a mathematical function that has a characteristic "S" shaped curve or a sigmoid curve that ranges between 0 and 1 (See Figure 3.3a), so this function is usually used in models where you need to predict a Probability as an exit. This function is defined by the following formula:

$$f(x) = \frac{1}{1 + e^{-x}} \quad (3.3)$$

The sigmoid function was successfully applied in problems of binary classification, modeling of logistic regression tasks, as well as other neural network domains, however, it suffers significant drawbacks that include acute wet gradients during backward propagation from deeper hidden layers to input layers, gradient saturation, slow convergence and non-zero centered output, which causes gradient updates to propagate in different directions.

Hyperbolic Tangent Function - Tanh

It is quite similar to Sigmoid but has a much better performance compared to multilayer neural network training, its nature is nonlinear. This function is centered at 0 and its range is between -1 and 1 (See Figure 3.3b), therefore, its output is defined by:

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (3.4)$$

Although this function has a better performance than the sigmoid, it could not solve the leakage gradient problem that sigmoid functions have. One of the main advantages of tangential function is that it produces a zero-centered output, which helps the backward propagation process.

Tangent functions have been used primarily in recurrent neural networks for natural language processing (Dauphin, Fan, Auli, & Grangiera, 2017) and speech recognition tasks (Mass, Hannun, & Ng, 2013).

Rectified Linear Unit function (ReLU)

The ReLU function was proposed by Nair and Hinton in 2010, and since then it has been the most widely used activation function for machine learning applications with neural networks.

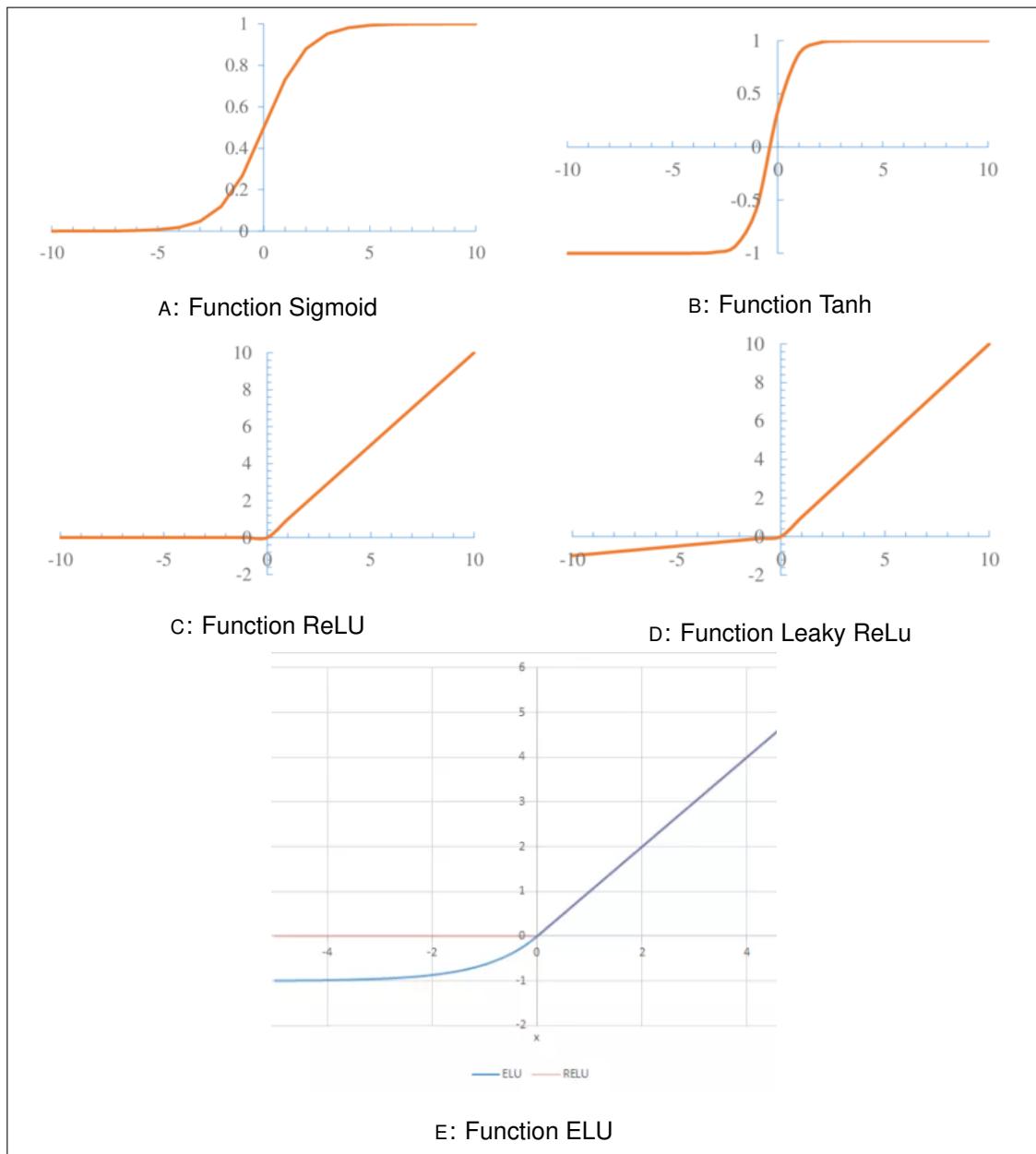


FIGURE 3.3: Activation functions (Jing & Guanci, 2018).

ReLU is a faster learning activation function (Lecun, Bengio, & Hinton, 2015), so it proved to be the most successful and most used function. This function offers better performance and generalization than sigmoid and tangent functions in learning with neural networks.

ReLU represents an almost linear function and, therefore, retains the properties of linear models that makes it easy to optimize, with gradient descent methods.

The ReLU activation function performs a threshold operation for each input element where values below zero are set to zero (See Figure 3.3c), so ReLU is defined by:

$$f(x) = \max(0, x) = \begin{cases} \text{si } x_i \geq 0 & x_i \\ \text{si } x_i < 0 & 0 \end{cases} \quad (3.5)$$

This function rectifies the values of inputs below zero, forcing them to become zero, thereby eliminating leakage gradient problem observed in previous types of activation function. The ReLU function has been used within the hidden units of neural networks.

The main advantage of using ReLU is that it guarantees a faster calculation, since it does not calculate exponentials and divisions, with an improved general calculation speed (Zeiler et al., 2013). Another property of ReLU is that it introduces the shortage in hidden units, since it reduces values between zero and maximum. However, ReLU has the limitation that it is easily overfitted compared to sigmoid function, although abandonment technique has been adopted to reduce overfitted effect of ReLU and rectified networks improved performance of neural networks.

ReLU has a significant limitation that it is sometimes fragile during training, causing the death of some gradients. This makes some neurons also dead, to solve problems of dead neurons, the Leaky ReLU activation function was proposed.

Leaky ReLU (LReLU)

The year 2013 was proposed as an activation function, this function introduces a small negative slope to ReLU to keep and keep weight updates alive during the propagation process (Mass et al., 2013). The parameter α was introduced as a solution to the problems of dead neurons of ReLU. This function calculates gradient with a very small constant value for negative gradient α in the range of 0.01, so LReLU (See Figure 3.3d) is calculated as:

$$f(x) = \alpha x + x = \begin{cases} \text{si } x_i > 0 & x_i \\ \text{si } x_i \leq 0 & \alpha x_i \end{cases} \quad (3.6)$$

Exponential Linear Unit Function (ELU)

The ELU² function tends to converge the cost to zero faster and produces more accurate results. Unlike other activation functions ELU has an additional alpha constant that should be a positive number.

It is very similar to ReLU since both have an identity function for positive inputs, however in ELU negative inputs it softens slowly until its output is equal $-\alpha$ while in ReLU it softens sharply. ELU function is calculated according to equation 3.7.

$$f(x) = \begin{cases} \text{si } x_i > 0 & x_i \\ \text{si } x_i \leq 0 & \alpha * (e^{x_i} - 1) \end{cases} \quad (3.7)$$

3.3.3 Architecture of the Neural Networks

A regular neural network consists of a chain of interconnected layers of neurons, these layers are: an input layer, one or several hidden layers and an output layer.

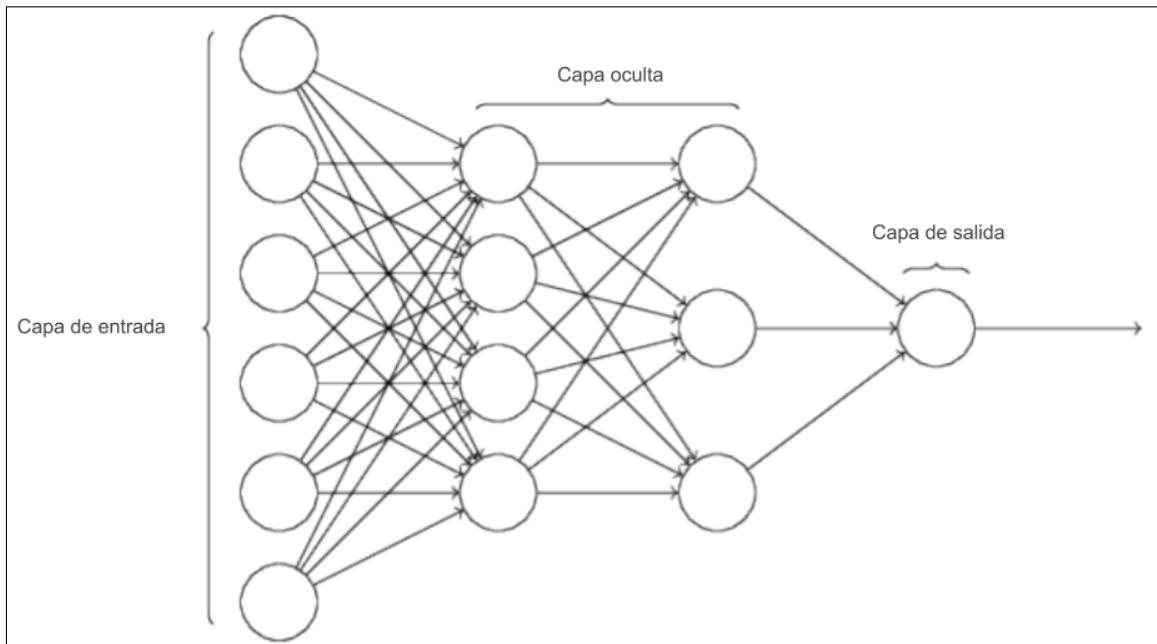


FIGURE 3.4: Architecture of an artificial neuron. Reproduced from (Michael, 2015).

Figure 3.4 is an example of a very simple neural network; In this figure the leftmost layer is called **input layer**, and the neurons within this layer are called input neurons. The rightmost or **output layer** contains output neurons. And finally the two intermediate layers are **hidden**

²ELU, Exponential Lineal Unit

layers of neural network, these are called that because neurons of this layer are not input or output; A neural network can have one or more hidden layers.

Unlike the human brain, an Artificial Neural Network has a fairly strict predefined structure, connections between neurons are always forward (feedforward): connections range from the neurons of a given layer to neurons of next layer, that is, There are no side connections or back connections. This means that a neuron that was activated in layer 3 cannot activate a neuron in layer 2 or earlier.

3.3.4 Learning process of Neural Networks

A key feature of neural networks is their iterative learning process, that is, each sample of training set is presented to the network, one at a time, so that the weights associated with input values are adjusted each time. During this learning phase, the network trains by adjusting the weights to predict a correct output for input samples.

Neural networks have the advantage of having a high tolerance for noisy data, as well as a high capacity to classify patterns with which they have not been trained. The most popular neural network training technique is the backpropagation algorithm.

Once the structure of a network is defined for a particular application, it is ready to be trained. To begin this process, initial weights are chosen at random, and then proceed with training (learning).

Backpropagation

A neural network propagates the signal of input data forward through its parameters at the time of decision; and then spread back the information about the error, so that parameters can be altered. This happens by following steps below:

- The network guesses output data, using its parameters.
- The network measures its accuracy with a loss function.
- The error is propagated backwards to adjust the wrong parameters.

Therefore it can be said that Backpropagation algorithm takes the error associated with an erroneous assumption by neural network, and uses that error to adjust parameters of neural network in the direction that generates the least error.

3.4 Types of Neural Networks

3.4.1 Autoencoders

An Autoencoder is an Artificial Neural Network used for unsupervised machine learning, it is trained to reconstruct its own inputs, that is, predict the value of output \hat{x} given an input x via a hidden layer h , see Figure 3.5. Autoencoders are usually used for dimensionality reduction and feature learning.

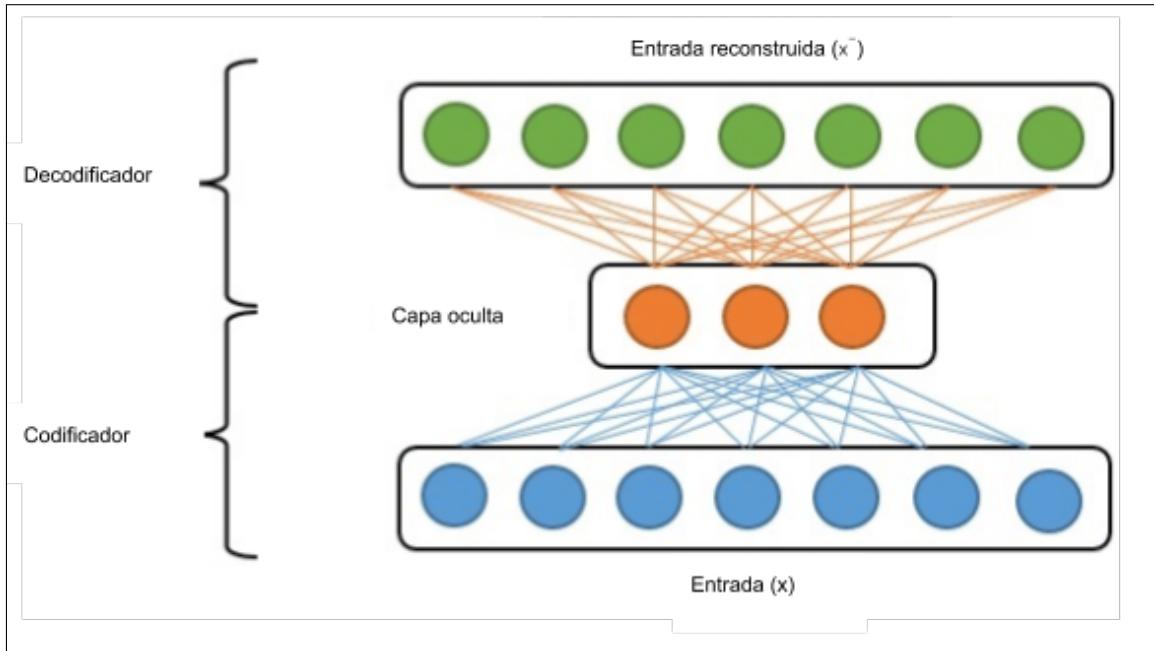


FIGURE 3.5: Graphic of an Autoencoder (Own elaboration).

Autoencoders are composed of two parts: the encoder and decoder. Encoder learns a compressed representation of input data, this can be defined with the coding function $h = \text{encoder}(x)$, which is defined by a linear or nonlinear function. If function of encoder is non-linear, auto-encoder will be able to learn more features than a linear PCA. The purpose of decoder is to reconstruct its own input via the decoding function, $\hat{x} = \text{decoder}(h)$.

The difference between input and reconstructed input is the reconstruction error. During training, autoencoder minimizes the reconstruction error as an objective function. Autoencoders are often used for data generation as generative models. Decoder of an autoencoder can generate an output given an artificially assigned compressed representation.

3.4.2 Convolutional neural networks

For some types of data, specifically for images, conventional neural networks are not well adapted; which implies that in the study Lecun et al. (1998) propose convolutional neural

networks (CNN³) to solve this problem. CNNs have revolutionized image processing and eliminated manual feature extraction. A CNN acts directly on matrices, or even on tensors for images with three RGB color channels; so currently, CNNs are widely used for image classification, object recognition, face recognition, among others.

CNNs not only provide better performance compared to other detection algorithms; but even in some cases they outnumber humans, such as in classification of objects in specific categories such as particular breed of a dog or a species of bird (Russakovsky, 2014).

By stacking multiple and different layers in a CNN, complex architectures are created for classification problems. The four types of layers that are most common are: convolution layer, grouping/subsampling layer, nonlinear layer and fully connected layer. An example of a CNN can be seen in Figure 3.6, where the first and third layers are convolutional layers, the second and fourth are subsampling layers and finally the fifth and sixth layers with completely connected layers.

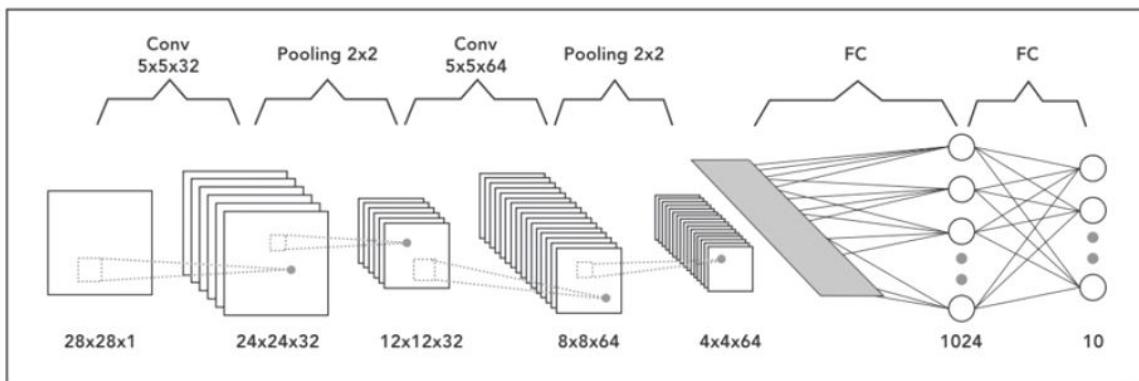


FIGURE 3.6: Architecture of a Convolutional Neural Network (CNN) (Muhammad, n.d.).

3.4.3 Recurrent neural networks

To understand the importance of time series, the following analogy can be taken, human beings do not start to think from scratch every second, so when reading a document each word is understood based on understanding of the previous words, that is to say , everything is not eliminated and you start thinking of zero every time, given this statement you can say that thoughts of human beings have persistence.

Traditional neural networks do not have data persistence, which for some specific problems, including the one of this work, is a major deficiency. In order to solve these types of problems, Recurrent Neural Networks (RNN⁴) appear, which are a type of artificial neural network proposed in the 80s (Rumelhart, Hinton, & Williams, 1986; Elman, 1990; Werbos, 1988) designed to

³CNN, Convolutional Neural Network

⁴RNN, Recurrent Neural Network

recognize patterns in data streams, such as text, genomes, handwriting, numerical time series data emanating from sensors, among others.

RNNs are a particular family of neural networks where the network contains one or more feedback connections, so that the activation of a group of neurons can flow in a loop. This property allows model to retain information about past, which allows it to discover temporal correlations between events that are far from each other in data.

RNN have a certain memory of what happened previously in a sequence of data, this helps system to gain context of data. Theoretically it is said that RNN has infinite memory, that is, these types of networks have ability to look back indefinitely; however, in practice you can only look back a few last steps.

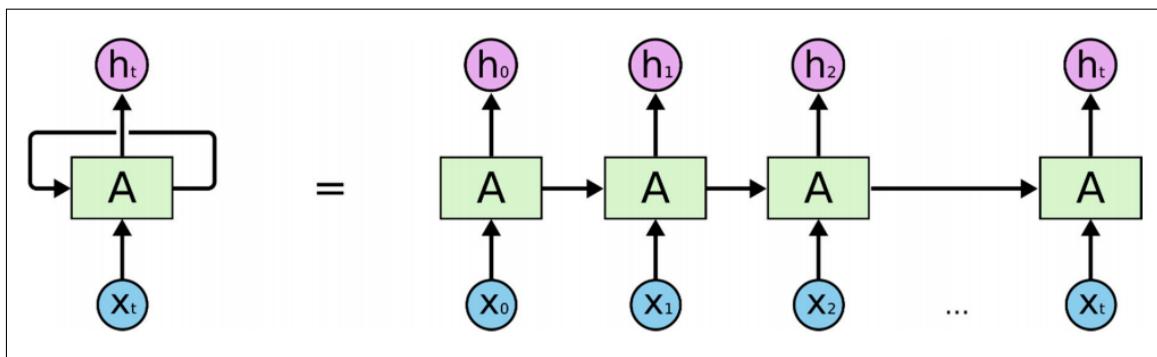


FIGURE 3.7: Sequential processing in a recurrent neural network (RNN) (Olah, n.d.).

Figure 3.7 illustrates a simple RNN with an input unit, an output unit and a recurring hidden unit expanded in a complete network, where x_t is input in the time step t and y_t is output in the time step t . On the other hand, in training process, RNNs use the Backpropagation algorithm over time (BPTT⁵). The BPTT process uses a backward, layer by layer, work approach of final output of network, adjusting the weights of each unit according to calculated unit portion of total output error. The repetition of information loops results in large updates of neural network model weights and leads to an unstable network due to accumulation of error gradients during the update process. Therefore, BPTT is not efficient enough to learn a pattern of long-term dependence due to disappearance of gradient and the explosion of gradient problems (Bengio, Simard, & Frasconi, 1994). To overcome disappearance and explosion gradient problems that standard RNNs have, LSTM and GRU can be used (Pascanu, Mikolov, & Bengio, 2013).

3.4.4 LSTM

LSTM⁶ is an evolution of RNN, it was introduced by Hochreiter and Schmidhuber in (1997) to board the problems of standard RNNs that were mentioned before, adding additional interactions

⁵BPTT, Backpropagation Through The Time

⁶LSTM, Long Short-Term Memory

per module (or cell). LSTMs are a special type of RNN, capable of learning long-term dependencies and remembering information for extended periods by default.

The LSTM model is organized in form of a chain structure. However, the repetition module has a different structure. Instead of a single neural network like a standard RNN, it has four interactive layers with a unique method of communication. The LSTM's structure is shown in Figure 3.8.

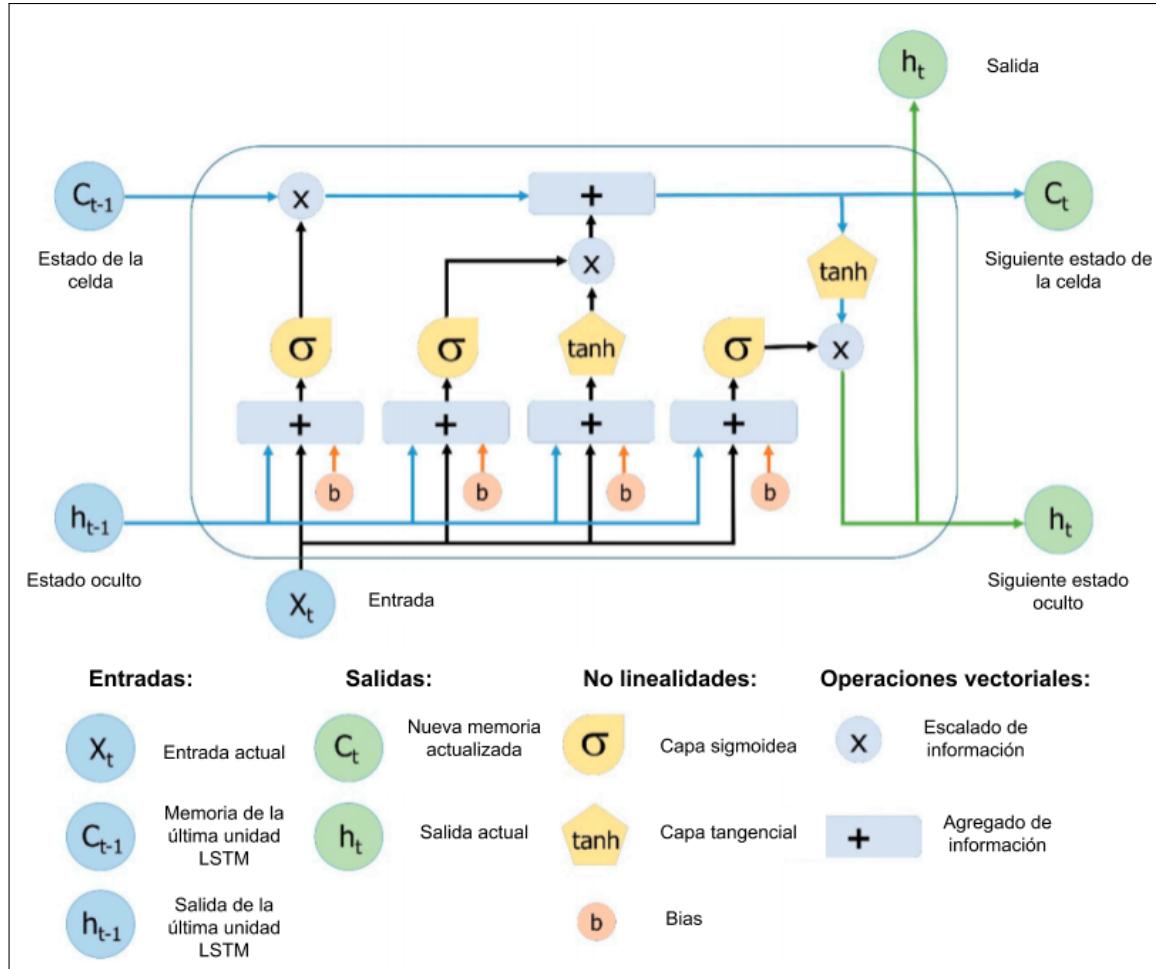


FIGURE 3.8: LSTM structure. Played from Yan (Yan, n.d.).

A typical LSTM network is made up of memory blocks called cells. Two states are transferred to next cell, the state of cell and hidden state. The **cell's state** is the main chain of data flow, which allows data to flow forward, essentially, without changes. However, some linear transformations may occur. Data can add or remove the cell's state through sigmoid gates.

A door is similar to a layer or a series of matrix operations, which contains different individual weights. LSTMs are designed to avoid the problem of long-term dependence because it uses doors to control memorization process.

3.4.5 GRU

GRU⁷ was first designed by Kyunghyun Cho in (2014a). This RNN structure only contains two doors. The update door controls information that flows into memory, while the reset door controls information that flows out of memory. Similar to LSTM, GRU has gate units that modulate flow of information within the unit; however, without having a separate memory cell. One could say that GRU is a slightly more simplified variant of RNN that often offers comparable performance to LSTM and is significantly faster to calculate.

In summary, GRUs have following two distinctive characteristics:

- **Reset doors** help capture short-term dependencies in time series.
- **Update doors** help capture long-term dependencies in time series.

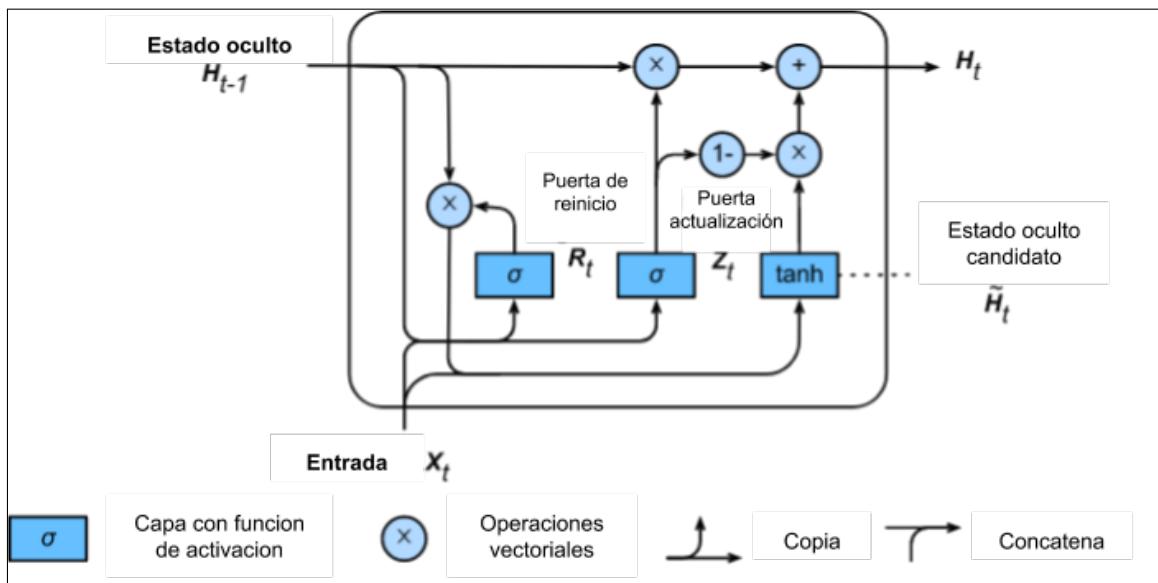


FIGURE 3.9: GRU structure. Reproduced from (Zhang et al., 2019).

3.5 Anomaly detection techniques

There are several approaches available for anomaly detection, in this section some of the most used algorithms will be described.

⁷GRU, Gated Recurrent Unit

3.5.1 One-Class SVM

One-Class SVM⁸ (OC-SVM) is a widely used approach to discover anomalies in an unsupervised manner (Schölkopf and Smola, 2002). OC-SVMs are one of the most widely used semi-supervised learning techniques, because it gives good results even for large data sets. However, this algorithm has the disadvantage that it requires a lot of time and memory in practice and its complexity grows quadratically with number of records.

This algorithm is only trained with positive examples (normal classes). The general idea of this algorithm is to transform the attribute space and draw a divisional hyperplane so that observations are as far as possible from origin, as can be seen in Figure 3.10.

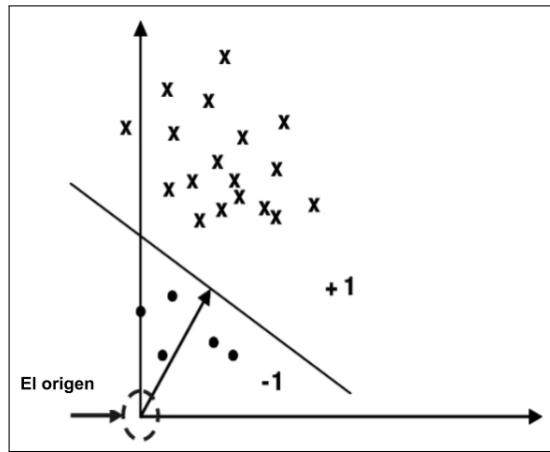


FIGURE 3.10: One-Class SVM. Reproduced from (Alashwal et al., 2006)

As a result, a margin is obtained, on one side of which observations of training sample are grouped as densely as possible (normal observations $\hat{y}_i = 1$), and on the other, abnormal values are found ($\hat{y}_i = -1$), not similar to what the algorithm saw during training.

3.5.2 Isolation Forest

This model is used in a scenario similar to One Class SVM, specifically in an unsupervised environment. The isolation forest takes a different approach to the OC SVM, since instead of grouping normal data, it tries to isolate the anomalous data.

The basic component of Isolation Forest is the isolation tree, which is a simple binary tree where in each T_i node both characteristic and threshold for division rule are randomly selected. An existing node stops generating children if and only if there is only one example following the division rule for that specific route (which means that the example has been isolated) or a maximum height has been reached. This means that at the end of the training process we

⁸SVM, Support Vector Machine

will have a completely over-adjusted random classification tree, which can be used for anomaly detection purposes. The main intuition of this algorithm is that if an example is anomalous, it will be isolated after some cuts in features space, which translates into having a low height in isolation tree.

Unlike other methods such as grouping or classification, isolation forests do not learn a profile of what is normal, but instead directly attack anomalies. No distance metric is used in this algorithm which saves time in calculations; so isolation forests have a linear temporal complexity.

A comparison of the ability of Isolation Forest and One-Class SVM algorithms to cope with different two-dimensional data sets is presented in Figure 3.11, with the aim of giving some intuition about behavior of these algorithms.

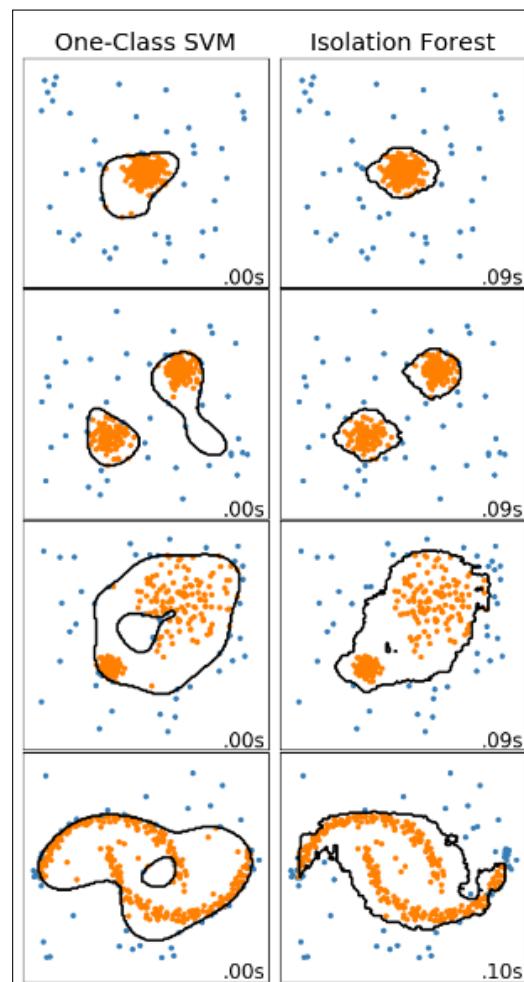


FIGURE 3.11: Performance comparison between the One-Class SVM and Isolation Forest algorithms. Reproduced from (0.22, n.d.).

3.5.3 Autoencoders

Today, autoencoders have been widely used in image classification, machine translation and voice processing; This is due to its ability to compress data without supervision. As far as is known, Hawkins et al. (2002) and Williams and Baxter (2002) were the first to propose autoencoders for anomaly detection. Since then, the ability of auto-encoders to detect outliers was demonstrated in different domains such as the detection of anomalies in X-rays.

The traditional method of autoencoder-based anomaly detection is mainly based on reconstruction error, considering as anomalies those samples that present a high reconstruction error. In training phase, only normal data is used to train the auto-encoder, in order to minimize the reconstruction error, so that auto-encoder can recognize characteristics of normal data. In test phase, the trained auto-encoder will be able to reconstruct normal data with small reconstruction errors, but they will fail with anomalous data that auto-encoder has not encountered before and, therefore, have relatively higher reconstruction errors compared to normal data. Therefore, when comparing whether the reconstruction score of an anomaly is above a predefined threshold, auto-encoder will determine if the data presented for test is anomalous (Guo et al., 2018) (See Figure 3.12) . Equation 3.8 shows how this technique determines what an anomaly is and what is not; where S_z represents the reconstruction.

$$C(z) = \begin{cases} \text{if } S_z \leq \text{Threshold} & \text{Normal behavior} \\ \text{if } S_z > \text{Threshold} & \text{Anomaly} \end{cases} \quad (3.8)$$

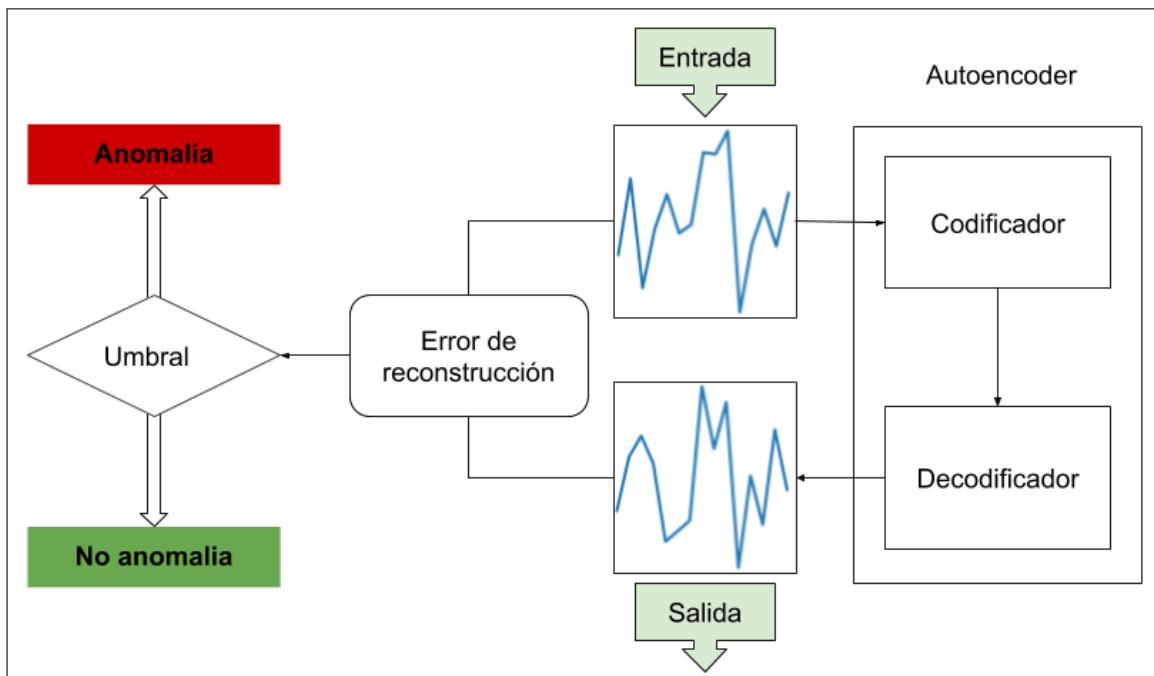


FIGURE 3.12: Detection of anomalies with autoencoder (Own elaboration).

3.6 Evaluation Metrics

Evaluating machine learning algorithms is an essential part of any project, because a model can provide satisfactory results when evaluated with a metric; but it can give a poor result when another metric is used. Classification accuracy is commonly used to measure the performance of a model, however, it is not enough to judge a model. Different types of evaluation metrics will be covered below.

3.6.1 Classification Accuracy

The classification accuracy is the relationship between the number of correct predictions and the total number of input samples.

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions made}} \quad (3.9)$$

This metric works well if there are the same number of samples that belong to each class; for example, if you consider that you have a data set that contains 98% of samples that belong to class A and 2% that belong to class B, the model could easily obtain 98% training accuracy by simply predicting each training sample as class A.

This implies that precision can give the false feeling of achieving high precision, which becomes a real problem if it involves problems that involve the detection of high-risk things; for example, a rare but deadly disease since cost of not diagnosing a sick person's disease is much higher than cost of sending a healthy person for further analysis.

3.6.2 Logarithmic Loss

Logarithmic loss, also known as Log Loss, works by penalizing false classifications, and also has a good performance for the classification of various classes. When working with Log Loss, the classifier must assign a probability to each class of all samples; assuming that there are N samples that belong to M classes, Log Loss is calculated according to the following equation:

$$\text{LogarithmicLoss} = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} * \log(p_{ij}) \quad (3.10)$$

where:

y_{ij} , indicates whether sample i belongs to class j or not.

p_{ij} , indicates the probability that the sample i belongs to class j .

Log Loss has no upper limit and exists in the range $[0, \infty)$. When a Log Loss close to 0 is obtained, it indicates greater precision, while if it is far from 0 it indicates less precision. In general, it can be said that minimizing Log Loss provides greater accuracy for classifier.

3.6.3 Confusion Matrix

Confusion matrix, also called error matrix, is the most common method to evaluate accuracy of a classification result (Smits, Dellepiane, & Schowengerdt, 1999). This matrix is a cross-tabulation of expected data and results of classification model. The number of columns and rows is equal to number of categories in classification and different statistical measures are derived from them.

		Prediction	
		Positive Class	Negative Class
Real	Positive Class	True Positive (TP)	False Negative (FN)
	Negative Class	False Positive (FP)	True Negative (TN)

TABLE 3.1: Confusion matrix, for a binary classification (Own elaboration).

In Table 3.1 the rows of matrix represent the real values, while columns are associated with the data classified by the model (predictions). The main diagonal, which appears in green, indicates the successes or True Positives (TP) and True Negatives (TN), which are all those data where model obtains the same result that was expected to be obtained. As for all the other values of the matrix, they belong to those data that were classified incorrectly, these are classified into two classes: False Positives (FP), which in the matrix are presented in red and False Negatives (FN) that in the matrix they were represented in orange.

The overall accuracy of matrix is calculated by dividing the sum of correctly classified samples by the total number of samples taken 3.11. This value is a measure of classification as a whole, since it indicates the probability that a sample is correctly classified.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (3.11)$$

Accuracy is not an adequate measure for the evaluation of atypical value detection algorithms, because most of the time a false negative is much more expensive than a false positive, in addition to training sets presenting a large amount of data normal compared to amount of anomalous data. There are two common metrics to evaluate outlier detection algorithms, AUC and F1 Score; These two metrics will be deepened in detail.

3.6.4 Area Under Curve (AUC)

Area Under Curve (AUC) is one of the most used metrics for evaluation. It is used for binary classification problems.

The AUC of a classifier is equal to the probability that the classifier classifies a positive example chosen at random higher than a negative example chosen at random. Before defining AUC, following terms must be understood:

Recall or Sensitivity or True Positive Rate (TPR)

True positive rate, also known as sensitivity, corresponds to the proportion of positive data points that are correctly considered positive, with respect to all positive data points. It is defined according to equation 3.12.

$$Sensitivity = \frac{TP}{FN + TP} \quad (3.12)$$

Specificity or True Negative Rate (TNR)

True negative rate, also known as specificity, corresponds to the proportion of negative data points that are correctly considered negative, with respect to all negative data points. It is defined according to equation 3.13.

$$Specificity = \frac{TN}{FP + TN} \quad (3.13)$$

ROC (Receiver Operating Characteristics)

ROC is the curve drawn by connecting the points on the X axis = FPR (False positive rate) and the y axis = TPR (True positive rate) for different values of a discrimination threshold (decision limit to determine if a value corresponds to a class or not) for a model, that is, different thresholds are chosen for a model, the TPR and FPR are calculated for each threshold, then it draws the ROC curve and finally the AUC is calculated, which is the area under the curve ROC . An example of the AUC-ROC curve is shown in Figure 3.13.

There are two main reasons why this curve is needed, first is that it reflects how good the model is to separate two classes and second is that it helps to choose the best threshold; for example, an AUC equal to 0.5 means that model separates two possible random results and an AUC of 1 (maximum value) implies a perfect separation; therefore, it can be said that the higher the value of AUC, better will be the performance of evaluated model.

3.6.5 F1 Score

F1 Score defines how accurate is a model, that is, how many instances it classifies correctly, as well as indicating how robust is model. This metric is necessary when you want to find a balance between accuracy and recovery, since it gives a fair evaluation even when the data set is unbalanced.

Precision

Precision is the number of correct positive results divided by the number of positive results predicted by classifier.

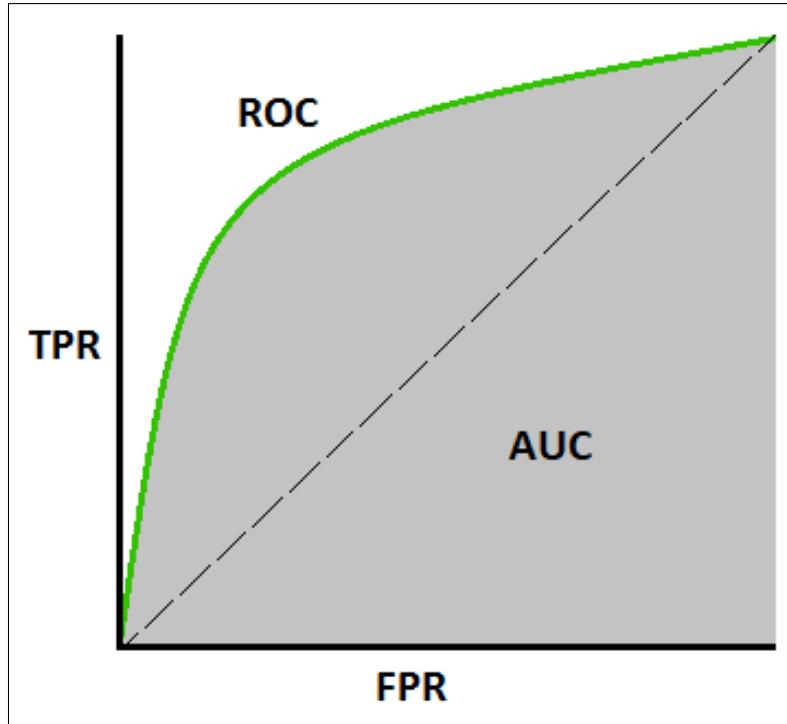


FIGURE 3.13: Example of an AUC-ROC curve (Özler, n.d.).

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3.14)$$

Recall

It is the number of correct positive results divided by the number of all samples that should have been classified as positive.

$$\text{Recall} = TPR = \frac{TP}{TP + FN} \quad (3.15)$$

Therefore, F1 Score is expressed mathematically as:

$$F1 = 2 * \frac{1}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \quad (3.16)$$

Next chapter details the process of capturing and preparing the data set, an important stage, because this set will be the one with which proposed anomaly detection mechanism will be trained.

Chapter 4

DATA CAPTURE AND PREPARATION

Having a large amount of data in any anomaly detection problem is what makes it possible to generate more accurate models, because you never know what characteristics can indicate an anomaly, having multiple types of data is what allows you to go more beyond a mere detection of specific anomalies and being able to identify more sophisticated contextual or collective anomalies. However, obtaining this data is not always a simple task, so you often have to find a way to generate it.

This chapter will detail the method of data collection that was performed for this research and the different data analysis techniques that were applied.

4.1 Data Capture

Currently there are several approaches to access the information of driver (agent) and vehicle. In the first approach, a set of sensors and additional hardware are previously implemented in the vehicle, for example, telematic boxes (black boxes provided by car insurance companies), on-board diagnostic adapters (OBD-II) plugged into the controller of network area vehicle (CAN) (Zaldivar, Calafate, Cano, & Manzoni, 2011; Araujo, Igreja, R., & Araujo, 2012), the information recorded by these devices can be retrieved or sent via Internet. However, this strategy requires that vehicles install additional devices, which implies a higher cost. To overcome these inconveniences, there is an alternative approach which is to use smartphones to collect data through a set of integrated sensors, such as inertial sensors (accelerometers and gyroscopes), global positioning systems (GPS), magnetometers, microphones, sensors Image (cameras), light sensors, proximity sensors, direction sensors (compass), among others.

For this research work, the use of smartphones was chosen to access the type of driving information, for reasons presented above, with this approach an Android-based mobile application was developed to collect sensor data: accelerometer and gyroscope, in intervals of 1 second, which in the first instance will be stored internally in the mobile device. (Ver Figura 4.1)

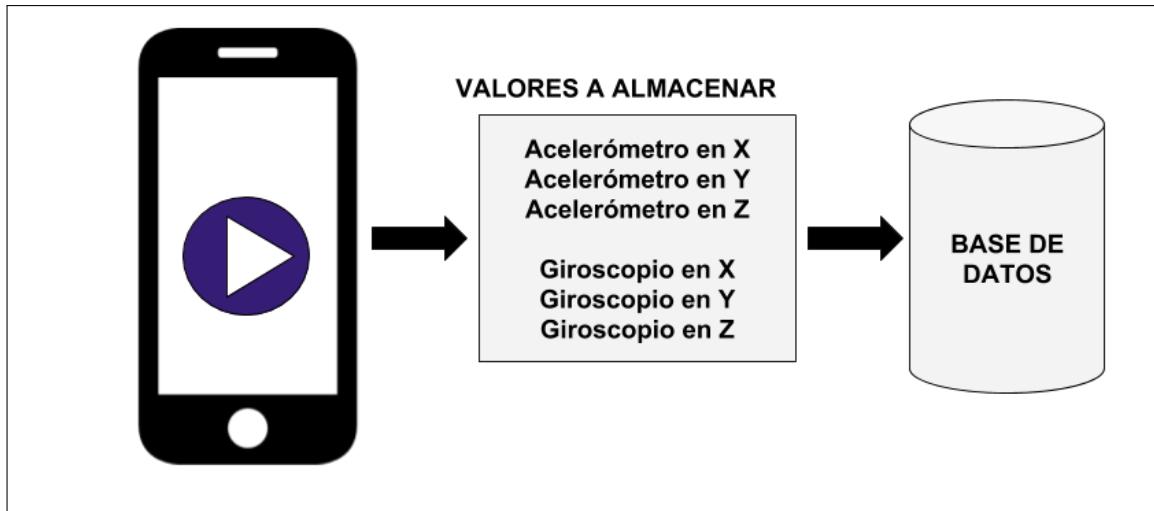


FIGURE 4.1: Data collection, with interval of one second (Own elaboration).

For data collect, a windshield cell phone holder was used as seen in Figure 4.2; the capture was made in two different positions (vertical and horizontal).



FIGURE 4.2: Windshield cell mount, horizontal position (Own elaboration).

Each capture, regardless of position in which it was made, resulted in a dataset (dataset), where for each time T (1 sec.) There are six variables: accelerometer in X (acc x), accelerometer in Y (acc y), accelerometer in Z (acc z), gyroscope in X (gyr x), gyroscope in Y (gyr y) and gyroscope in Z (gyr z). A fragment of data set that was obtained in a capture is shown in Figure 4.3

4.2 Data preparation

Machine Learning depends heavily on the data. They are the most crucial aspect that makes algorithm training possible and explains why machine learning has become so popular in recent

<u>_id</u>	<u>acc_x</u>	<u>acc_y</u>	<u>acc_z</u>	<u>gyr_x</u>	<u>gyr_y</u>	<u>gyr_z</u>	<u>fecha</u>
	Filter	Filter	Filter	Filter	Filter	Filter	Filter
1	3.69699096...	-0.4218902...	9.07658386...	0.00050354...	7.62939453...	0.00115966...	2019-10-19
2	3.68769836...	-0.3159942...	9.03315734...	0.00183105...	-0.0009918...	9.15527343...	2019-10-19
3	3.69015502...	-0.3061065...	9.10673522...	0.00210571...	-0.0009918...	0.00035095...	2019-10-19
4	3.68812561...	-0.3355712...	9.43148803...	0.00343322...	-0.0135345...	-0.0023040...	2019-10-19
5	3.69512939...	-0.3061065...	9.13452148...	0.00263977...	-0.0009918...	-0.0009765...	2019-10-19

FIGURE 4.3: Fragment of data set obtained (Own elaboration).

years. The main problem is that all data sets have failures, which makes data preparation a very important step in the Machine Learning process.

The main purpose of data preparation is to manipulate and transform raw data, so that data can be exposed or made more easily accessible (Pyle, 1999), to achieve this purpose a process that involves selection must be followed, Pre-processing and data transformation.

4.2.1 Data selection

Data selection involves the following steps:

- Select only a subset of available data.
- Derive or simulate some data from available data, if necessary.
- Exclude data that is not relevant to the problem.

For present work, only the first step of this phase will be emphasized, because the data available is limited since they were captured for investigation and as indicated in section 4.1, this capture was made, both vertically and horizontally, the differences between them will be analyzed below.

Fragments of obtained captures by mobile device of the same user (agent) from different positions are shown in Figure 4.4.

Although captured values are very similar to each other, data that was captured with mobile device in a horizontal position, presents less noise, this because this position favors the inertia of device when vehicle is in motion, which is a great advantage over data that were captured vertically, since these were more susceptible to shaking while vehicle was moving causing the values captured in this position to present movement values not only of vehicle but also of mobile device, which is not what is sought in the present work.

For reasons presented in previous paragraph, it was decided to work with the data captured with mobile device in a horizontal position, thus discarding those data captured vertically.

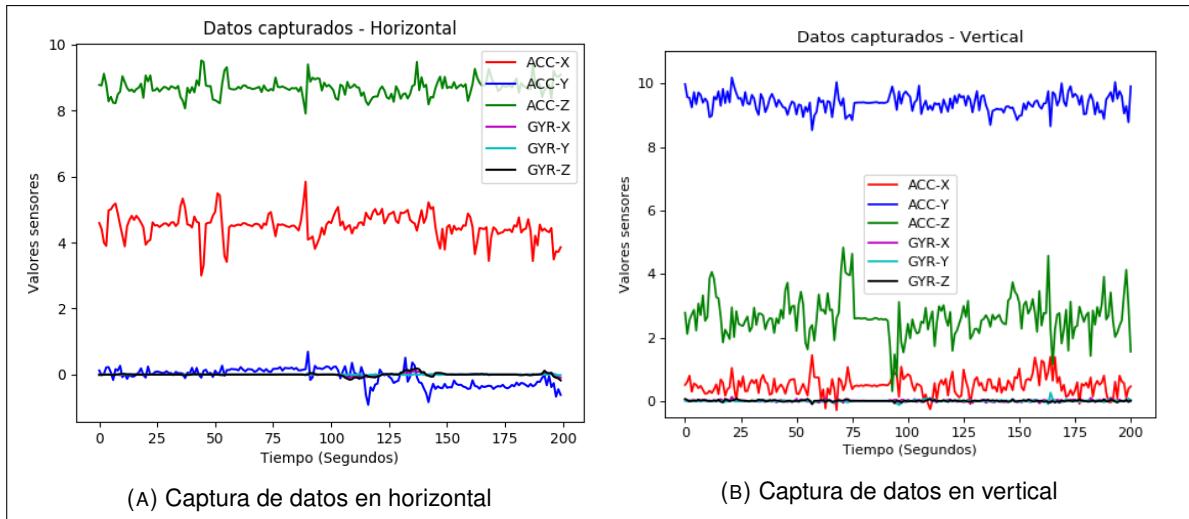


FIGURE 4.4: Graph of sensors captured in different positions (Own elaboration).

4.3 Data preprocessing

Once data with which you will work will be selected, you must proceed to preprocess them, thus entering the Pre-processing phase of data, the objective of this phase is to reduce amount of data, find relationships between them, normalize them, remove outliers and extract features of data. This phase includes several techniques such as cleaning, integration, transformation and data reduction (Suad & Wesam, 2017).

4.3.1 Data cleaning

Row data may have incomplete records, noisy values, outliers and inconsistent data. Data cleaning in most cases is the first step of data pre-processing. This technique is used to compensate for missing values, soften the noise in the data, recognize outliers and correct inconsistencies.

Compensation techniques for incomplete records

Many real-world data sets may contain missing values for different reasons; which can drastically affect the quality of machine learning model.

Below are four compensation techniques for data sets, with the aim of coping with the problem of missing values.

- **Ignore / Delete:** In some cases it is better to ignore or eliminate the tuple that contains missing values instead of filling it. Generally this technique is practiced in data sets that are very large, where deleting some data will not affect the information transmitted by data set. However, when working with a small data set, removing tuples that contain missing values could cause you to lose important information.

- **Fill out missing values manually:** Another option is also to complete missing values if you understand the nature of them, this is usually done in a small data set, since in large sets this task would require a lot of time.
- **Fill the missing values with central values (Medium / Medium):** This technique is much better than those presented above. In this technique, the mean or median of respective attribute is inserted to missing values.
- **Interpolation:** It is a reliable, accurate and scientific way to complete missing values. To use this technique, a relationship between attributes must first be developed and then the most probable and precise value for missing places is predicted, this can be achieved through regression, bayesian formulation and induction by decision trees.

Remove noise from data (smoothing)

To understand this technique you must first define what is noise in data. Noise in data is any type of random error or variation in measured attributes; on the other hand outliers present in data can also be considered as noise.

It is important to note that noise present in data set can greatly affect result of Automatic Learning algorithms, therefore those data that contain noise are not considered good data and should be eliminated as much as possible. However before removing these, you should be able to detect them; To achieve this goal there are many techniques that can be used, one of these techniques is Visualization of data, which consists in obtaining a visual representation of data, to be able to show if there is noise in data and/or outliers.

In this work, histograms of frequencies of each characteristic of data set were plotted and thus visualized of a better way if there was noise or outliers in it. In Figure 4.5 it can be shown that values of each sensor have negative and positive asymmetries, in addition to having many values far from the average, which gives an indication of possible outliers.

A table with descriptive statistics of data set is presented in Figure 4.6, this table presents: the amount of data, its mean, its standard deviation, the minimum and maximum value of data set and the 50, 25 and 75 percentiles, it should be noted that the percentile 50 is same as median.

With the information provided in Figure 4.6, it is now easier to perform an analysis of data set, first it is evident that the values obtained during the capture, have very small standard deviations between 0.05 and 0.81 and yet difference between the values minimum and maximum are really large, for example for the Accelerometer in X, difference is approximately 10 (Minimum value: -3.40 and maximum value: 7.17), this means that set of data with which it works presents a lot of noise, considering as noise or outliers, those values far from the average.

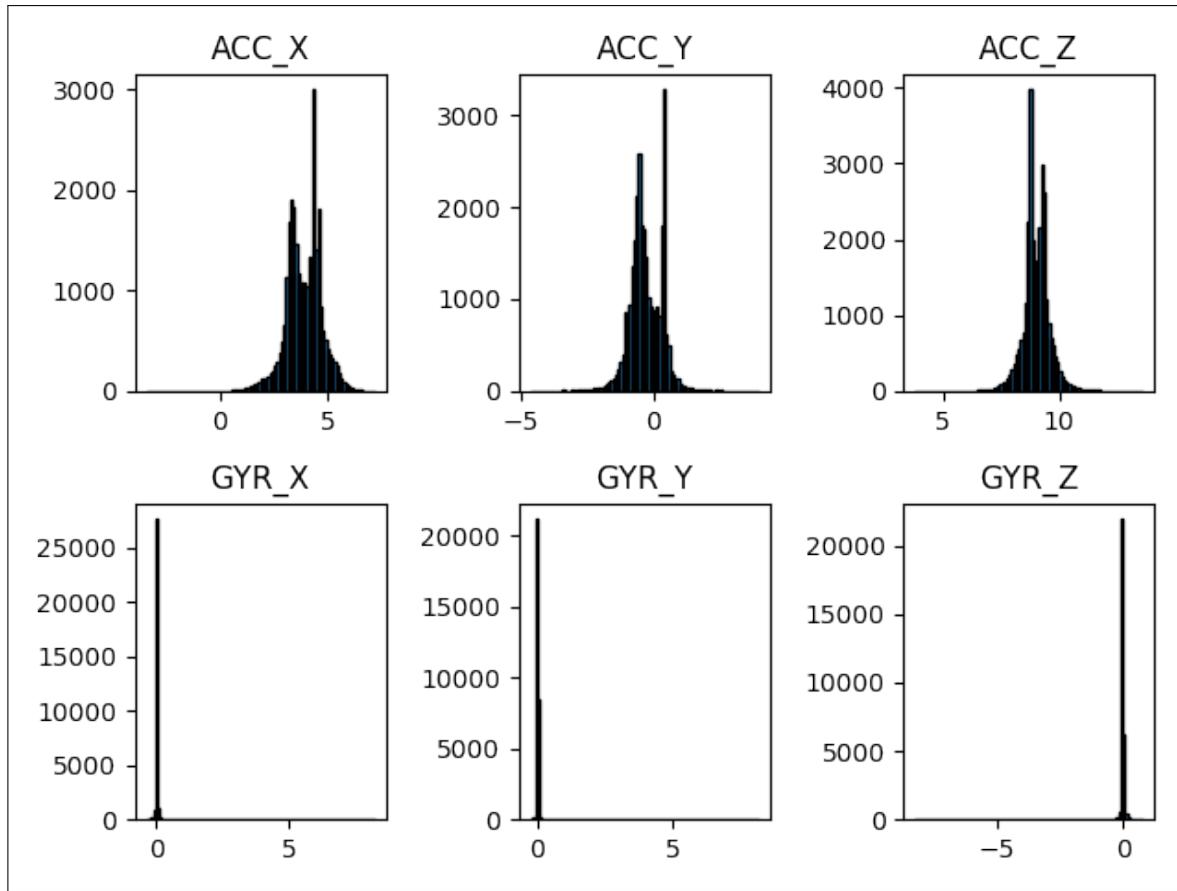


FIGURE 4.5: Histogram of data set's frequencies (Own elaboration).

	acc_x	acc_y	acc_z	gyr_x	gyr_y	gyr_z
count	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000
mean	3.874731	-0.302470	8.997619	0.000086	0.000281	-0.000389
std	0.809801	0.614389	0.563649	0.055470	0.052041	0.074062
min	-3.403488	-4.633270	3.801849	-0.388611	-0.251450	-8.214310
25%	3.311012	-0.684528	8.707027	-0.005192	-0.004990	-0.004837
50%	3.907730	-0.381134	8.988884	-0.000015	-0.000015	-0.000031
75%	4.412102	0.235142	9.307396	0.004745	0.004807	0.004578
max	7.164932	3.947861	13.609085	8.208740	8.212128	0.818634

FIGURE 4.6: Table of data set's statistical results (Own elaboration).

To eliminate the noise presented by data set, Rule 68-95-99.7, also known as Empirical Rule, was applied, assuming that data set with which it works has a normal distribution, standard deviation can be used to determine the proportion of values that fall within a particular range

of average value. For such distributions, it is always the case that 68% of the values are less than a standard deviation (1SD) of average value, that 95% of values are less than two standard deviations (2SD) of average and that the 99% of values are less than three standard deviations (3SD) from average. Figure 4.7 shows this concept schematically.

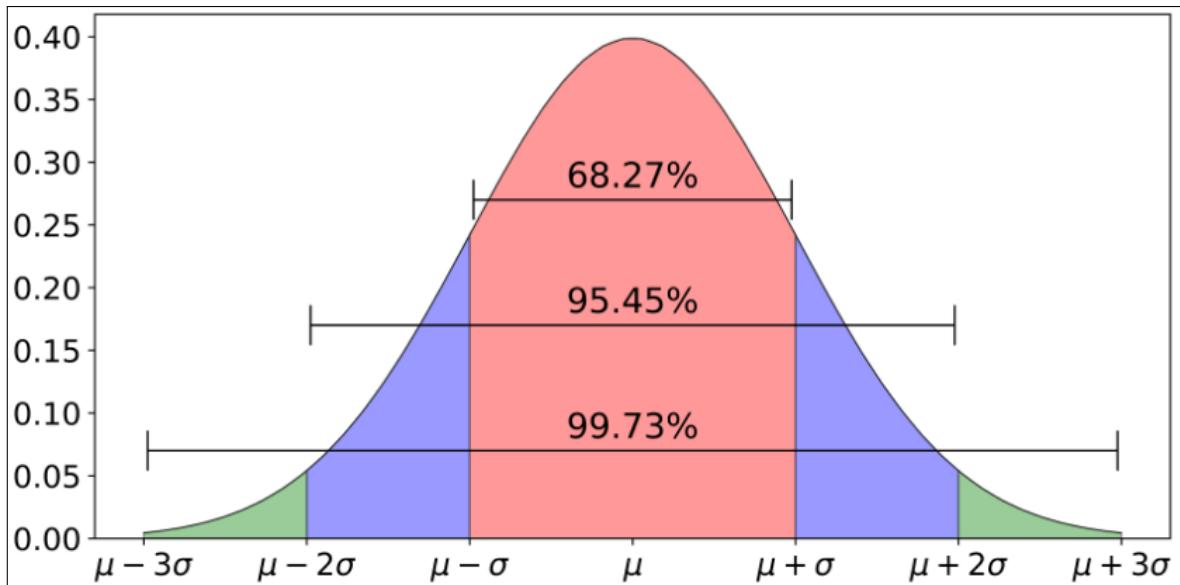


FIGURE 4.7: Rule 68-95-99.7 (Galarnyk, n.d.).

In Figure 4.8, it can be observed how rule 68-95-99.7 behaves on the values of accelerometer sensor in Z, so to eliminate noise from data set there are two options, eliminate values after three standard deviations from the mean, in case of considering that data set presents few anomalies or outliers, or eliminating values after two standard deviations in case of being sure that data set presents a great amount of noise in data, in this case most of data belong to normal driving behaviors so that only values found after three standard deviations from mean will be eliminated.

Fusion or data integration

When working with real-world data, it is possible that the required data is not found in same data set, in these cases, it is necessary to collect data from different sources and merge them into a single data set; this process is called Fusion or Data Integration. One of the most common problems of this process is redundancy.

This process was not applied in investigation because data was only captured through mobile device, so there is no problem of having more than one data source.

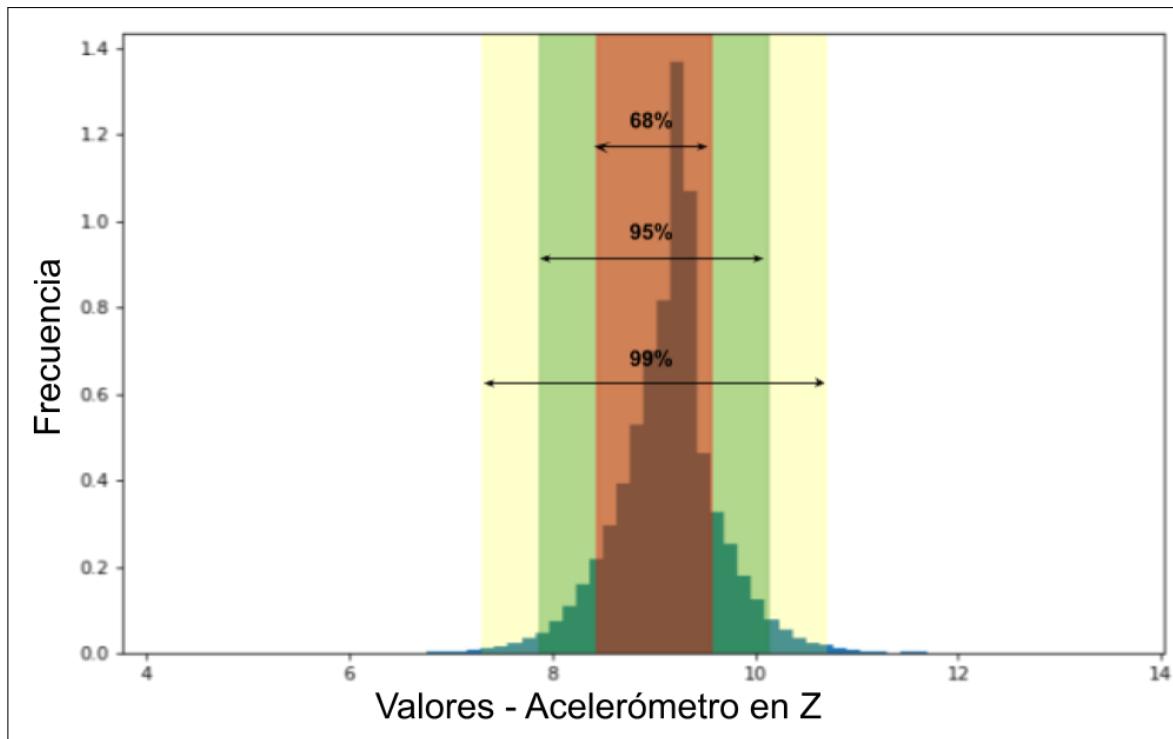


FIGURE 4.8: Result of applied Rule 68-95-99.7 on values' accelerometer sensors in Z (Own elaboration).

Data transformation

Data transformation is the process where the nature of data is changed, this process uses some strategies to be able to extract important information from data set; some of strategies for data transformation are:

- **Aggregation:** In this technique, the summary or aggregation operation is applied to data. For example: daily sales data can be used to calculate monthly and annual sales amount, and then add these calculated data to data set.
- **Discretization:** With this technique, raw values of a numerical attribute are constructed and replaced by interval values.
- **Attribute Construction/Feature Engineering:** This technique is useful for generating additional information from those data that are not representative enough by themselves, and may also be adequate when there are fewer features but still contain hidden information to extract.
- **Normalization / Standardization:** Normalization or standardization is defined as process of rescaling original data without changing its behavior or nature. A new limit is defined (generally between 0 and 1) and data is converted accordingly. This technique is useful in classification algorithms that involve neural networks or distance-based algorithms (for

example, KNN¹, K-Means², which is used for clustering. This algorithm is able to gradually learn how to group unlabeled values into groups by an analysis of average distance of these values.). Some standardization techniques are:

- *Normalization Min-Max*: In this method, each input is normalized between defined limits:

$$x_{normalize} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (4.1)$$

It presents the problem that it compresses the input data between fixed limits, which are usually 0 and 1. This means that if there is noise, it will be expanded, which makes this method not suitable for stable signals.

- *Standard Scaler*: It is an alternative method to the variables' scaling, it consists in subtracting from each data the variable's average and dividing it by standard deviation.

$$x_{normalize} = \frac{x - x_{average}}{x_{std}} \quad (4.2)$$

This method is suitable for normalizing stable signals, however, both mean and standard deviation are very sensitive to anomalous values. An alternative solution to this is the elimination of anomalies before normalization.

- *Scaling over the maximum value*: This method presents the idea of scaling data by dividing it by its maximum value.
- *Robust scaler*: Robust scaling involves eliminating the median and scaling data according to interquartile range (IQR). This method is robust for outliers.

Data reduction

This process is based on adoption of some strategies, such that analysis of reduced data produces the same information produced by original data. Some of strategies include: principal component analysis (PCA), selection of a attributes' subset, grouping and sampling among others.

Principal Component Analysis (PCA)

¹KNN, is a classification algorithm (or regression) that, to determine the classification of a point, combines the classification of the nearest K points.

² textbf{K-Means}, is a clustering algorithm that attempts to divide a set of points into K groups; so the points in each group tend to be close to each other.

Principal Component Analysis (PCA) is a technique that is widely used for applications such as dimensionality reduction, lossy data compression, feature extraction and data visualization (Bishop, 2006).

PCA is an unsupervised and non-parametric statistical technique, which is used for dimensionality reduction. This technique is an important step because high dimensionality in the field of machine learning can lead to model's overfitting, thus reducing its ability to generalize, Bellman (2003) describes this phenomenon as the "Curse of Dimensionality". In addition, the use of this technique can directly improve performance of Machine Learning models.

PCA combines the input variables in a specific way, then gets rid of the "less important" variables and at the same time preserves the most valuable parts (or principal components³) of all variables.

When using PCA with a machine learning approach, we follow these steps:

1. Divide d -dimensional dataset into a training, development and test set.
2. Standardize / Normalize the dataset according to training set.
3. Construct the covariance matrix.
4. Decompose the covariance matrix into its vectors and eigenvalues.
5. Order the eigenvalues by decreasing order to classify corresponding eigenvectors.
6. Select k eigenvectors that correspond to k largest eigenvalues, where k is the dimensionality of new entity subspace ($k \leq d$).
7. Construct a projection matrix \mathbf{W} from the "top" k eigenvectors.
8. Transform d -dimensional input training set \mathbf{X} using the projection matrix \mathbf{W} to obtain the new k -dimensional feature subspace.

Dataset division

Given that there is a dataset to generate Machine Learning model, it is divided into three parts: training⁴, development⁵ and testing set⁶, however this is not a trivial task; because if it is not done correctly the result can be disastrous.

³**Principal Component**, it is a normalized linear combination of the original features in dataset.

⁴**Training set**, is the data sample used to fit the Machine Learning model.

⁵**Development set**, is the data sample used to provide an unbiased evaluation of a fitted model with the training set while adjusting the model hyperparameters.

⁶**Test set**, is the data sample used to provide an unbiased evaluation of a final fit of model in the training set.



FIGURE 4.9: Division of the dataset (Own elaboration).

Training set	70%	21000
Development set	15%	4500
Test set	15%	4500
Data set	100%	30000

TABLE 4.1: Table of dataset's division (Own elaboration).

The work of Moindrot and Genthal (2018) says that division of dataset has a great impact on productivity, so it is important that when choosing subsets they must have the **same distribution** and must be chosen randomly from dataset.

On the other hand, the size of development and test set must be large enough so that development and test results are representative for the performance of model. For large data sets (greater than one million), the development and test set can have around 10000 examples each, that is, 1% of the total data.

Other considerations to be taken into account in practice are:

- The division of training/development/test set must always be the same for all experiments, therefore a reproducible script must be available to create the training/development/test division.
- Must be checked that development and test sets came from the same distribution.

In the present investigation, the data set has 30,000 examples, so the division of data set will be as shown in Table 4.1.

Standardize / Normalize the dataset

In section 4.3.1 we have already described what is the standardization or normalization of data and some of scaling's types that exist; therefore this section will only be limited to the elaboration of an analysis to decide which scaling technique is the most suitable for data set, in Figure 4.10

a fraction of captured data set can be seen, that is the basis with which the comparative analysis will be carried out with the different scaling's types.

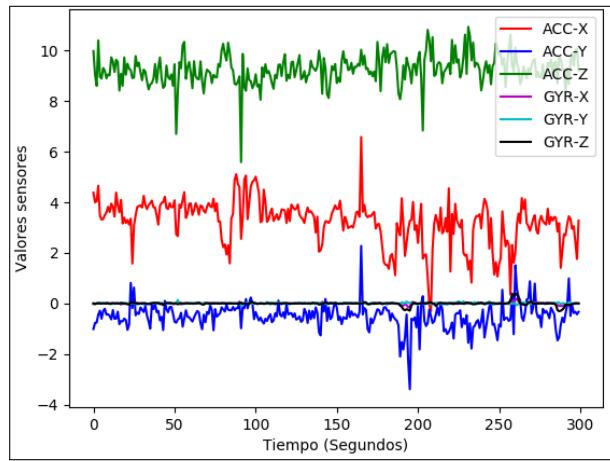


FIGURE 4.10: Visualization of the captured driving parameters (Own elaboration).

To test the different scaling's types, we fit them with data that no longer has noise or outliers from training set.

The first type of scaling that was performed on captured data set was **Min-Max Normalization**, which was performed between the limits 0 and 1, the results obtained are shown in Figure 4.11a, where it can be seen that the values of Accelerometer, in its three axes, are not deformed after being scaled with this technique and gyroscope values, which are more stable, become deformed, considering as stable data those data that are presented as a line in zero with few fluctuations; This deformation can be advantageous by making small curves that were previously imperceptible more visible, however it carries the danger that it could amplify existing noise in data that we could not eliminate in previous step to this task.

The second normalization technique that was applied to data was **standard scaling**, the result can be seen in Figure 4.11b, observing in detail the results can be seen to be very similar to those obtained with Min-Max normalization, the only differences that can be seen are that the new range of the data is broader with a mean of zero and values that oscillate mainly between 2 and -2, and that some of fluctuations presented by gyroscopes are expanded a little more.

For the third data normalization, technique applied was the **scaling on maximum value**, the results are totally different from those obtained previously, however it is clearly observed that accelerometer values do not show much change, with the difference that average of these values are different for each one, and gyroscope values behave as in the previous ones, this can be seen in figure 4.11c. This technique presents the worst results, since it does not leave data set in the same range, which complicates the work with them.

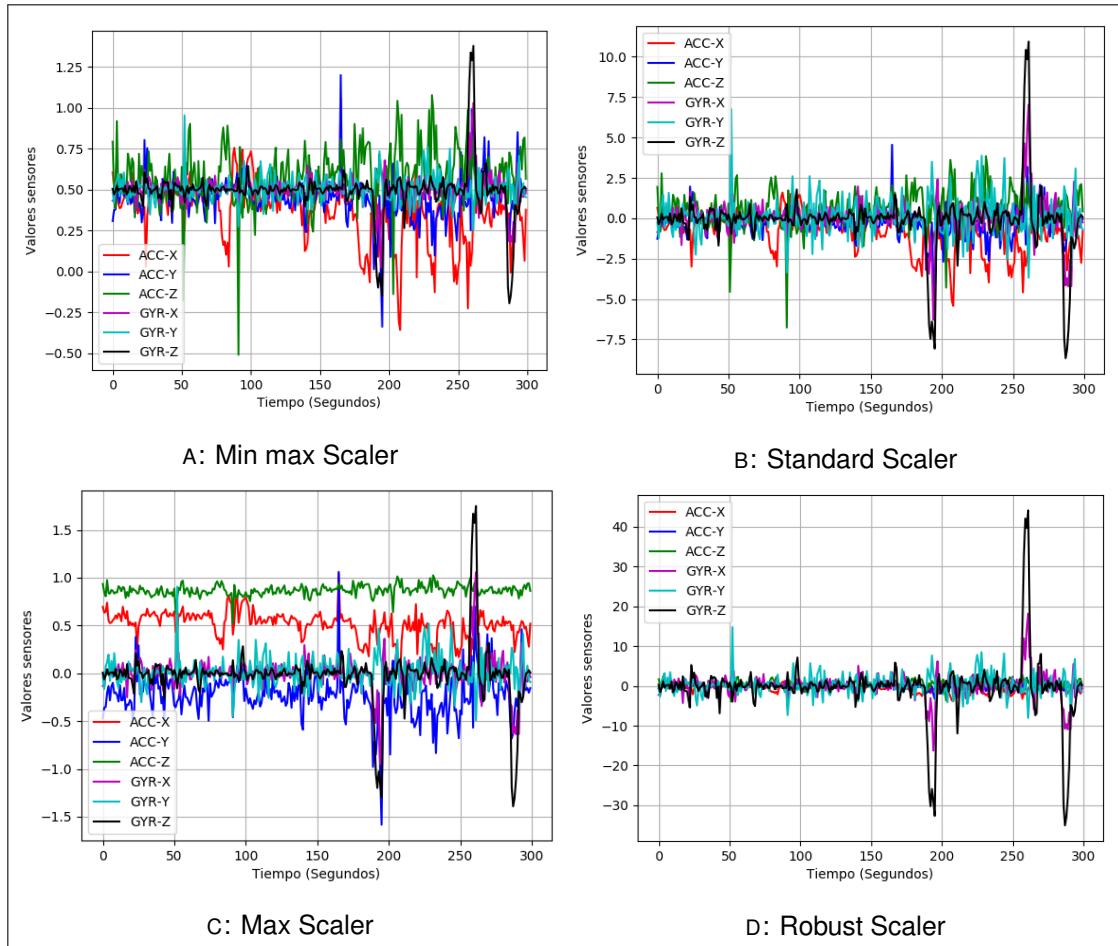


FIGURE 4.11: Graph resulting from applying different types of normalizations to data set (Own elaboration).

The last technique applied is **robust scaling**, its result can be seen in figure 4.11d, which can be considered quite similar to the first two techniques, with the difference that this scaling pronounces more the fluctuations of gyroscope's values, converting them even to a higher scale than that of accelerometer's values.

To decide the best normalization method that can be applied to the captured data, the **scaling method on the maximum value** will be completely discarded first because the results it presented were totally discouraging since the values after scaling did not share the same limits, on the other hand, the remaining three types of scaling are very similar, which complicates the choice of correct scaling, however, when using PCA, you must be very careful, because if you have a variable with a high standard deviation, this will have a greater weight in the calculation of the axis than a variable with a low standard deviation, so this can be an important decision parameter to choose the most appropriate type of scaling.

Table 4.2 shows the mean and standard deviation of each variable after the different types of scaling have been applied to the data. As can be seen in the scaling of the maximum and standard value, the standard deviations are very different, as for the Min-Max scaling, the standard deviations of each variable are almost the same since they all range between 0.1666814 and 0.169386, which It suits very well for the analysis of main components, also mentioning that this technique is the one that best preserves the relevant information in the data set, so this technique was selected because it is the most convenient for this stage.

		ACC X	ACC Y	ACC Z	GYR X	GYR Y	GYR Z
Min-Max	Average	0.500369	0.500051	0.501112	0.497923	0.499352	0.500396
	St. Deviation	0.166847	0.166814	0.167300	0.168245	0.169386	0.166852
	Minimum	-0.999198	-0.675814	-1.041074	-0.681029	-0.319990	-18.004503
	Maximum	1.178265	1.654068	1.869867	25.395520	27.227514	2.345548
Standard	Media	-0.011266	-0.009209	-0.00570	0.013265	0.009458	0.005644
	St. Deviation	1.050384	1.086118	1.117567	2.224617	2.518566	2.076506
	Minimum	-9.451768	-7.665204	-10.307538	-15.575414	-12.173164	-230.291158
	Maximum	4.256420	7.504531	9.137616	329.221266	397.424938	22.968890
Maximum Abs. Value	Average	0.615065	-0.141064	0.842598	0.000519	0.001827	-0.001747
	St. Deviation	0.128546	0.286536	0.052784	0.334925	0.337715	0.332858
	Minimum	-0.540261	-2.160843	0.356031	-2.346416	-1.631746	-36.917638
	Maximum	1.137343	1.841185	1.274447	49.564032	53.291415	3.679194
Robust	Average	-0.028625	0.083278	0.008976	0.010491	0.031468	-0.040602
	St. Deviation	0.739033	0.672602	0.956915	5.752011	5.518751	8.397460
	Minimum	-6.670812	-4.657858	-8.811955	-40.295886	-26.663430	-931.368512
	Maximum	2.974050	4.736319	7.837925	851.216772	870.859223	92.823529

TABLE 4.2: Table with descriptive statistics of scaled data with different techniques (Own elaboration).

Apply PCA to dataset

Although some of steps in the internal workings of PCA are explained in detail at the beginning of this sub-section, there is a class called PCA implemented in SCIKIT-LEARN, which automates all following steps; however, it is necessary to determine how many characteristics (principal components) to maintain and how many to eliminate; for this task there are three commonly used methods, which will be described below.

- Arbitrarily select how many dimensions you want to keep, depending of use case. For instance, in a viewing approach you could choose 2 or 3 features.
- Calculate the variance ratio for each feature, choose a threshold, and add features until the chosen threshold is reached or exceeded.
- This method is closely related to previous one, because the variance ratio for each feature is calculated, the features are ordered by variance ratio and the cumulative variance ratio explained is plotted while it maintains the features (this diagram is known as a screen diagram). You can choose how many features to include by identifying the point at which a new feature is added that has a significant drop in variance from previous feature, and choose the number of features that exist up to that point. This method is usually known as "Find elbow".

For the development of this investigation, both the first and last method of those listed above were discarded, this because the first method does not have the certainty that the quantity of features chosen is descriptive enough for dataset; while the last method does not have a mathematically precise definition, since it only finds "elbow", so it also removes control of amount of total variability in data that is finally obtained.

Once the methods that do not conform to the requirement of this work have been discarded, the only remaining option is the second method, thus being the one that will be applied for present work. Therefore, the total variability threshold to be preserved in data set was defined as 90%, later using the PCA class of SCIKIT-LEARN, the variance of each characteristic is calculated and then the results are plotted (See Figure 4.12).

And finally, it must be defined how many principal components retain at least 90% of data's variance. Figure 4.12 indicates that selecting 4 components can preserve about 97.7% of data's total variance, selecting 3 components conserves around 92.3% and selecting 2 conserves 85% of variance; therefore, the use of 3 features was decided, which will conserve 92.3% of dataset's total variance.

In various researches (Zenon, 2011; Klos & Waszczyszyn, 2011) it was proven that the application of PCA as pre-processing of the data set is an important stage, because it not only serves to compress the input data, but which also provides a satisfactory improvement in accuracy of Machine Learning models; reason why this technique is part of pre-processing stage of this work.

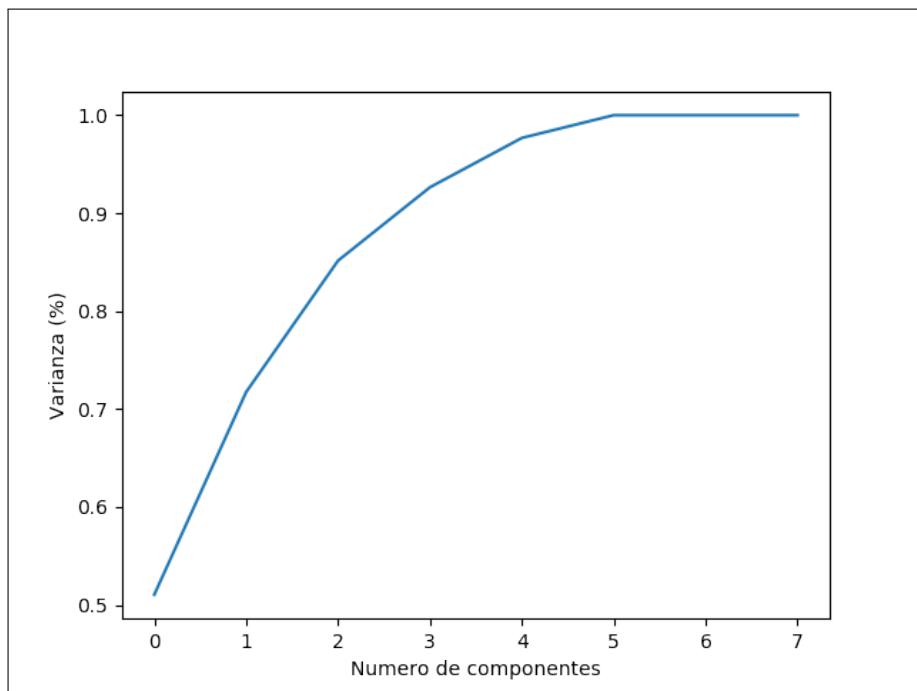


FIGURE 4.12: Variance's graph vs.components' number (Own elaboration).

Chapter 5

GENERATION OF THE ANOMALY DETECTION MECHANISM

This chapter describes the process that was followed to generate the conduction anomaly detection mechanism. As reviewed in Chapter 3, in this document, an outlier detector with a semi-supervised approach is proposed; however, before delving into the proposed method, a description of the development environment with which the experiment was used should be carried out, in addition to reviewing dataset available in this study.

5.1 Development environment

The experiment of this study was developed on a laptop with following characteristics:

- Intel Core i5-5200U 2.2GHz processor (c / TB 2.7 GHz).
- 8 GB of RAM.
- ArchLinux Operating System version 4.15.15-1-ARCH (64 bits).

It is important to clarify that methods' types used in this study are usually much more efficient in computers that have a graphics processing unit (GPU), since this unit allows parallel processing. The code developed in this work was written in PYTHON, an interpreted programming language that emphasizes code simplicity and readability, as well as being powered with support of powerful scientific libraries such as NUMPY, SCIPY, OPENCV, KERAS, MATPLOTLIB, SEABORN, etc. as for the experiments' development and the generation of detection mechanism, Jupyter Notebook was used, which is a local Python-based web application that allows you to view and execute documents that contain source code and equations. The versions of tools used are detailed below.

- **Python:** 3.6.5
- **Jupyter:** 4.3
- **Keras:** 2.2.2
- **Tensorflow:** 1.11.0

- **Scikit-learn:** 0.19.1
- **Matplotlib:** 2.0.2

5.2 Normal and anomalous dataset

En el Capítulo 4 se describió el proceso de captura y preparación del conjunto de datos, así como también su división en conjunto de entrenamiento/desarrollo/prueba; sin embargo cabe aclarar que aquel capítulo sólo se enfocó en el **conjunto de datos normales**.

A pesar de que se cuenta con una gran cantidad de datos normales, es necesario recolectar muestras que corresponden a anomalías, con el objetivo de poder validar el método que se propone en este proyecto. Por lo tanto, se realizó la captura de un conjunto de **datos anómalos**, el cual está conformado según el Cuadro 5.1.

Tipo de anomalía	No. anomalías	No. datos
Giros en Zig Zag	5	105
Giros a la derecha e izquierda a alta velocidad	7	35
Frenos en seco	6	24

TABLE 5.1: Tabla del conjunto de anomalías (Elaboración propia).

Como se mencionó en el párrafo anterior el conjunto de anomalías fue capturado para validar el método propuesto, por lo tanto este conjunto se etiquetó como positivo (con la etiqueta 1) y el conjunto de datos normales como negativo (con la etiqueta 0).

5.2.1 Generación de series temporales

Para la generación del modelo detector de anomalías se decidió ir más allá de una simple detección de anomalías puntuales y así poder detectar anomalías contextuales o colectivas; debido a ello se requiere el uso de datos en series de tiempo.

Los datos capturados por el dispositivo móvil, son dependientes del tiempo cronométrico en el que fueron capturados (un dato por segundo); por lo cual el primer paso a realizar es la generación de pequeñas fracciones de series temporales. En la Figura 5.1 se presenta los resultados de diferentes tamaños de series de tiempo, observando estos resultados en primera instancia se descarta la serie de tiempo que cuenta con dos pasos; debido a que no es lo suficientemente descriptiva. En cuanto a las series de tiempo restantes no es posible definir aún cual es la cantidad correcta de pasos, por lo cual, será un parámetro a optimizar en los diferentes experimentos que se realizará en las siguientes secciones. Cabe recalcar que el dominio de ésta variable está entre 3 y 5 pasos.

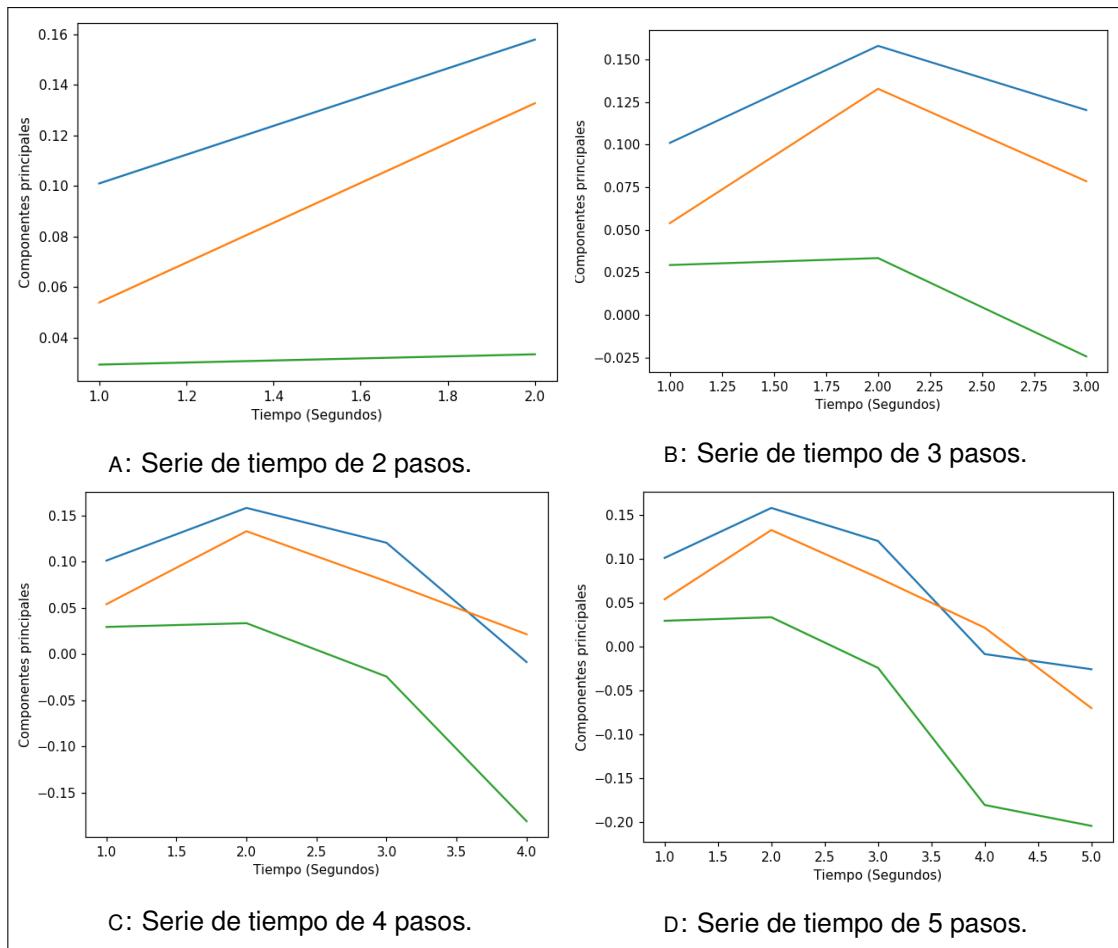


FIGURE 5.1: Gráfica resultante de diferentes tamaños de series de tiempo
(Elaboración propia).

5.3 **Modelo de detección de anomalías**

La presente investigación propone un método de detección de anomalías de conducción siguiendo un enfoque semi-supervisado, el cual consta de dos componentes: un **modelo del comportamiento normal** y un **método para la detección de valores atípicos**.

Por lo tanto, se realizó la comparación entre 3 diferentes métodos de detección, y según el rendimiento de cada uno se eligió la mejor opción. En el Cuadro 5.2 se presenta los tres diferentes métodos que fueron comparados; donde se puede observar que en todos los casos se usa un autoencoder como modelo del comportamiento normal, de esta manera, las siguientes secciones describirán la elección del autoencoder y la elección de uno de los tres diferentes métodos de detección de anomalías propuestos.

Método	Descripción
AE_T	Método de detección basado en autoencoders y umbralización (Thresholding).
AE_IF	Método de detección basado en autoencoders y aplicación de Isolation Forest.
AE_OC-SVM	Método de detección basado en autoencoders y aplicación de One-Class SVM.

TABLE 5.2: Tabla de los métodos comparados (Elaboración propia).

5.3.1 Modelo del comportamiento normal

Esta etapa es una de las partes más importantes de éste trabajo, debido a que el rendimiento del modelo de detección de anomalías depende en gran parte de la precisión de esta etapa.

Arquitectura del modelo

Como se mencionó en la anterior sección en esta etapa se utilizará un autoencoder como modelo ajustado al comportamiento normal de conducción. Por lo cual el autoencoder se entrenó con el conjunto de datos normales, de manera que el modelo aprenda a generar sólo las clases que se consideran normales y, con suerte, tendrá problemas para reconstruir anomalías, debido a que estas muestras no fueron presentadas durante el entrenamiento.

Para ello se probó con diferentes arquitecturas, primero la forma más simple que sólo se basa el uso de capas densas (completamente conectadas), luego se hizo pruebas con redes convolucionales y por último con redes recurrentes haciendo uso específico de capas LSTM. Por cada tipo de red se hizo la prueba con 3 diferentes tipos de entrada, es decir, se probó una diferente cantidad de pasos (entre 3 y 5) en las series temporales. Por lo tanto se realizaron 9 diferentes experimentos, de los cuales por cada tipo de red sobresalió una (usando la precisión de las redes como tipo de evaluación para desarrollar las comparaciones).

En el Cuadro 5.3 se presenta la red que obtuvo el mejor resultado de todas las Redes Densas que se probaron, esta red corresponde a la red que fue alimentada con secuencias de 3 pasos. Esta red cuenta con una capa de entrada (Input), una capa de aplanamiento (Flatten) esto debido a que la capa de entrada recibe una entrada bidimensional, un conjunto de capas densas (Dense) que van comprimiendo la información de los datos de entrada para posteriormente reconstruirlos, y por último la capa de salida es solo una capa para modificar la forma de la salida (Reshape); otro punto importante a resaltar es que las capas internas usan *elu* como función de activación y la última capa densa utilizan una función de activación tangencial (*tanh*), esto se debe a que el conjunto de datos, posterior a la obtención de componentes principales, se encuentra en el rango $(1, -1)$.

Arquitectura Densa				
NN_33				
	Tipo	Salida	Activación	# Parámetros
PCA 3	Input	(3,3)		0
	Flatten	9		0
	Dense	8	elu	80
	Dense	5	elu	45
	Dense	8	elu	88
	Dense	9	tanh	81
	Reshape	(3,3)		0

TABLE 5.3: Arquitectura densa para una secuencia de 3 pasos y 3 componentes principales (Elaboración propia).

Por otra parte el Cuadro 5.4 se presenta la red que obtiene la mejor precisión de todas las Redes Convolucionales que fueron probadas, esta red al igual que la anterior corresponde a la red que fue alimentada con secuencias de 3 pasos. Su arquitectura consta de una capa de entrada (Input), una combinación de capas de convolución de una dimensión (Conv1D) y agrupación (MaxPooling1D) hasta comprimir los datos a una dimensión de (2,4), luego un conjunto de capas convolucionales y de muestra ascendente (Upsampling1D) para decodificar la información compresa. Cabe recalcar que esta red también usa la función de activación tangencial en su última capa por las razones que se explicaron en el párrafo anterior.

Arquitectura Convolucional				
CNN_33				
	Tipo	Salida	Activación	# Parámetros
PCA 3	Input	(3,3)		0
	Conv1D	(3,2)	elu	20
	MaxPooling1D	(2,2)		0
	Conv1D	(2,4)	elu	28
	MaxPooling1D	(1,4)		0
	Conv1D	(1,6)	elu	54
	UpSampling1D	(3,6)		0
	Conv1D	(3,3)	tanh	57

TABLE 5.4: Arquitectura convolucional para una secuencia de 3 pasos y 3 componentes principales (Elaboración propia).

En el Cuadro 5.5 se muestra la red que obtuvo la mejor precisión de todas las Redes Recurrentes probadas, como en los anteriores casos ésta red es alimentada con secuencias de 3 pasos. Dicha red cuenta con una capa de entrada (Input), dos capas LSTM una que retorna sus secuencias y una que no, luego viene una capa de redimensionado, posteriormente dos capas LSTM, y finalmente un contenedor (TimeDistributed) de una capa densa.

Arquitectura Recurrente				
RNN_33				
	Tipo	Salida	Activación	# Parámetros
PCA 3	Input	(3,3)		0
	LSTM	(3,9)	elu	468
	LSTM	6	elu	384
	Reshape	(3,2)		0
	LSTM	(3,3)	elu	72
	LSTM	(3,9)	elu	468
	TimeDistributed(Dense)	(3,3)	tanh	30

TABLE 5.5: Arquitectura recurrente para una secuencia de 3 pasos y 3 componentes principales (Elaboración propia).

Es importante recalcar que las capas de redimensionamiento, agrupación, muestra ascendente y contenedores sólo fueron usadas para controlar la correcta compresión y descompresión de los autoencoders, es por ello que no se detalla a profundidad su funcionamiento.

Evaluación de autoencoders

Anteriormente se presentó los mejores representantes por tipo de red; ahora se procederá a la evaluación y comparación de estos 3 tipos de autoencoders, con el objetivo de elegir la arquitectura que se ajusta mejor al comportamiento normal de conducción.

En el Capítulo 3 se presentó los diversos tipos de evaluación que existen, en esta etapa el tipo de evaluación más apropiado es la **precisión** del modelo, debido a que se tiene un gran conjunto de datos balanceado (debido a que sólo se cuenta con comportamientos normales de conducción que corresponden a una sola clase, la "Normal"). Los resultados de la evaluación de los tres tipos de redes son mostrados en el Cuadro 5.6; dicho cuadro presenta la precisión, pérdida logarítmica y tiempo de ejecución de cada autoencoder según el conjunto de prueba; observando estos resultados se puede apreciar que las dos mejores redes son la red densa **NN_33** y la red recurrente **RNN_33** con precisiones de 90% aproximadamente, además de presentar un valor de pérdida relativamente bajo en comparación a la red **CNN_33**.

Red	Precisión	Loss	Tiempo ejecución
NN_33	0.9000740711953905	0.003956934471097257	26us/step
CNN_33	0.843777761353387	0.006740666443275081	31us/step
RNN_33	0.8899259290695191	0.003611267575787173	101us/step

TABLE 5.6: Evaluación de las redes NN_33, CNN_33 y RNN_33 (Elaboración propia).

Por otra parte la diferencia más grande entre las dos mejores redes (**NN_33** y **RNN_33**) es el tiempo de ejecución ya que de la primera es de tan solo 26 segundos/paso y de la segunda es de 101 segundos/paso, debido a estas similitudes entre ambas redes es necesario verificar visualmente los resultados de reconstrucción de cada tipo de red, de tal manera que se pueda elegir la red más adecuada para este problema.

En la Figura 5.2 se muestra los resultados de los autoencoders de siete secuencias tomadas aleatoriamente del conjunto de prueba, en la parte superior de cada figura se encuentra la secuencia de entrada y en la inferior la reconstrucción del modelo, como ya se podía esperar la red **CNN_33** presenta los peores resultados, lo cual hace que dicha red sea descartada; en cuanto a las dos redes restantes, la red **NN_33** presenta reconstrucciones muy similares a las secuencias de entrada, con algunos pequeños errores; por otra parte **RNN_33** presenta errores un poco más notorios que los obtenidos por **NN_33**. Por lo tanto se llegó a la conclusión de que la red **NN_33** se ajusta mejor al comportamiento normal de conducción, además de tener la gran ventaja de tener un tiempo de ejecución mucho menor que el de **RNN_33**, lo cual es realmente importante para los sistemas en tiempo real así como también de aquellos que cuentan con recursos de ejecución limitados, como es el caso del presente trabajo.

Una vez definido como está constituido el modelo del comportamiento normal se puede proceder con la elección del método de detección de valores atípicos.

5.3.2 Método de detección de anomalías

Al inicio de esta sección se definió tres diferentes enfoques para la detección de anomalías: la umbralización, la aplicación de bosques de aislamiento y finalmente la aplicación de SVM para una clase.

Umbralización

Esta técnica se basa en la definición de un umbral para determinar si el error de reconstrucción que obtiene el autoencoder (modelo del comportamiento normal) es lo suficientemente alto como para considerarse un valor atípico. Por lo tanto primero se debe definir la ecuación del error de reconstrucción para el modelo. En el presente trabajo el error de reconstrucción se define según la ecuación 5.1, donde x_i representa el valor real (entrada del autoencoder) y \hat{x}_i representa el valor obtenido por el autoencoder (salida del autoencoder).

$$\text{Error de reconstrucción} = S_z = |x_i - \hat{x}_i|^2 \quad (5.1)$$

En la Figura 5.3 se muestra la curva de los errores de reconstrucción obtenidos con el modelo del comportamiento normal para el conjunto de muestras normales.

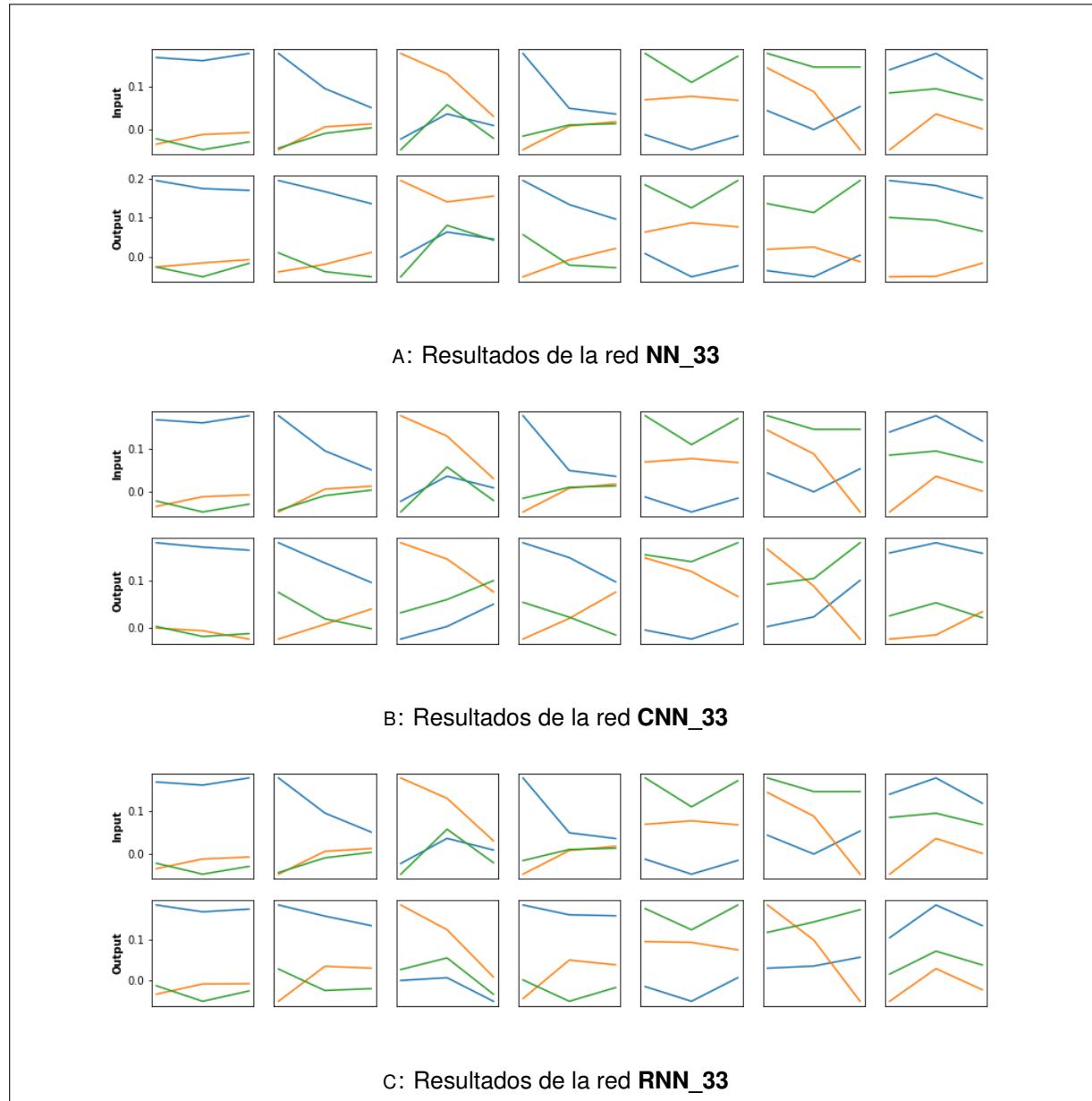


FIGURE 5.2: Resultados (Elaboración propia).

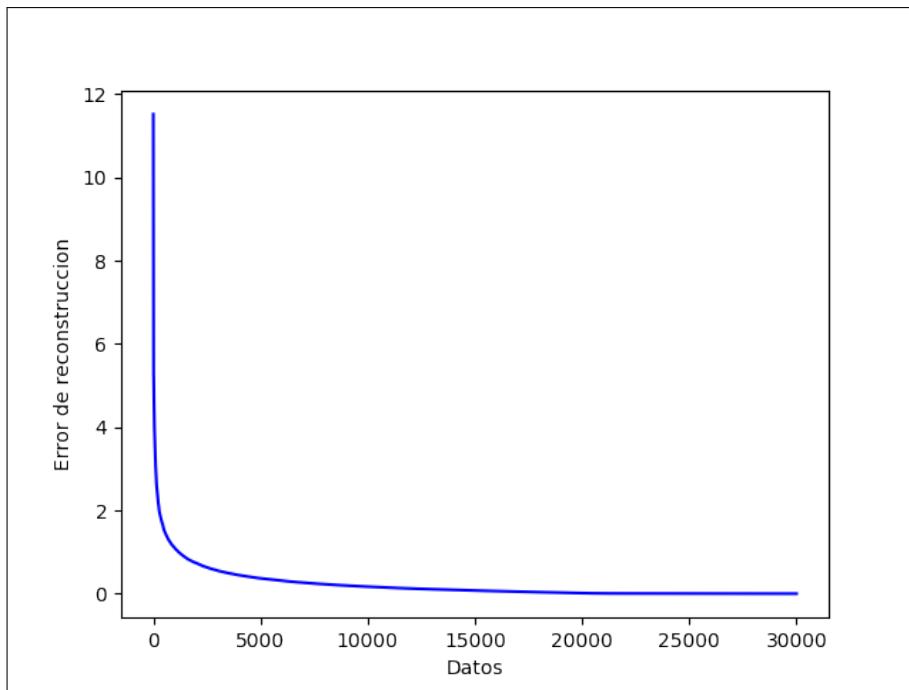


FIGURE 5.3: Curva de los valores de reconstrucción obtenidos con el modelo del comportamiento normal (Elaboración propia).

Una vez definido la ecuación de reconstrucción, se debe definir un umbral capaz de poder detectar la mayor cantidad de anomalías posibles. Esta tarea puede tornarse simple en un entorno de aprendizaje supervisado, sin embargo automatizar esta tarea en un contexto de aprendizaje no supervisado es un desafío que puede ser difícil de sobrelyear. En el presente trabajo se usó una técnica basada en encontrar un *Punto de codo* de una curva, que en este caso la curva está construida en base a los errores de reconstrucción del autoencoder.

Existen diferentes formas de hallar el punto de codo, sin embargo en este trabajo se utilizó una herramienta de Python, que automatiza esta tarea, llamada Kneedle. Esta herramienta devuelve el punto de inflexión de la función de la curva obtenida por el conjunto de valores proporcionado x y y , cabe recalcar que el punto de codo es el punto de máxima curvatura, por otra parte esta herramienta cuenta con un parámetro de sensibilidad (S), este parámetro permite ajustar qué tan agresivo se desea ser al detectar codos, los valores más pequeños para S detectan los codos más rápido, mientras que los más grandes son más conservadores, es decir, S es una medida de cuántos puntos "planos" se espera ver en la curva de datos sin modificar antes de declarar un codo.

De esta manera en el presente proyecto se realizó experimentos con diferentes valores de sensibilidad para encontrar el codo más adecuado para el conjunto de datos con el que se trabaja. En la Figura 5.4, se muestra los diferentes codos hallados para los valores de sensibilidad proporcionados (valores entre 0 y 2).

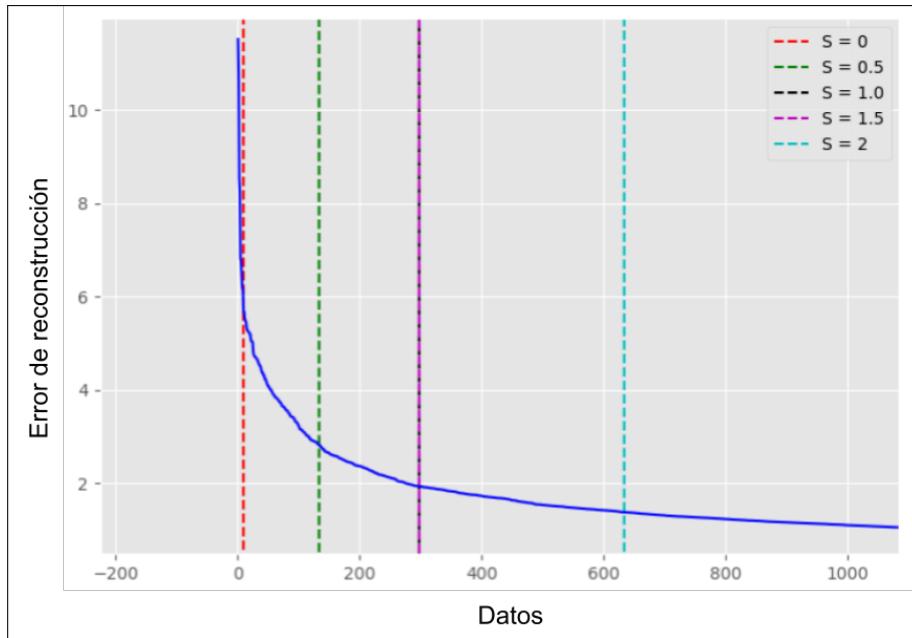


FIGURE 5.4: Resultados de la obtención de codos con diferentes valores de Sensibilidad, para los valores de reconstrucción obtenidos con el modelo del comportamiento normal (Elaboración propia).

Una vez obtenidos los codos se realizó la evaluación de cada uno de ellos, en el Cuadro 5.7 se presentan el umbral, los valores de la matriz de confusión, la sensibilidad y especificidad para cada codo. Los valores de la matriz de confusión son el resultado de aplicar el umbral de cada codo a los errores de reconstrucción obtenidos del conjunto de datos total (conjunto de datos normal y anormal equivalente a 44204 datos).

S	Umbral	VP	VN	FN	FP	Sensibilidad	Especificidad
0.0	5.665	92	43993	72	47	0.5610	0.9989
0.5	2.806	111	43814	53	226	0.6768	0.9949
1.0	1.920	120	43562	44	478	0.7317	0.9891
2.0	1.369	128	43100	36	940	0.7805	0.9787

TABLE 5.7: Evaluación de la detección de anomalías para cada codo obtenido con los diferentes valores de sensibilidad (Elaboración propia).

Según los resultados que se muestran en el Cuadro 5.7 se puede decir que mientras más pequeño es el umbral la sensibilidad (proporción de anomalías detectadas correctamente como anomalías) incrementa, sin embargo, a su vez reduce la especificidad (proporción de valores normales correctamente detectados como valores normales). Por lo tanto se debe hallar un punto intermedio, donde se pueda detectar la mayor cantidad de anomalías posibles y reducir en lo posible la cantidad de falsos positivos (datos normales que son detectados como anomalías). De esta forma el umbral más adecuado para el objetivo planteado fue 2.806, ya que

con este umbral se detecta 111 anomalías de 164 y los falsos positivos son aproximadamente el doble de los valores atípicos detectados.

De ello se deduce que, para detectar automáticamente el umbral el uso de 0.5 como parámetro S es el más adecuado, sin embargo si se desea incrementar el porcentaje de detección de anomalías a costa de incrementar el número de falsos positivos se puede usar un valor mayor a 0.5 para S y en caso de querer la menor cantidad de falsos positivos posibles se debe usar un valor menor a 0.5.

Isolation Forest

Antes de presentar como se llevará a cabo los experimentos con este algoritmo, es necesario ilustrar más detalladamente el funcionamiento del mismo. Por lo tanto la Figura 5.5 representa cómo se espera que un punto de datos anómalo se aísle rápidamente con el uso de este algoritmo, mientras que un punto de datos normal necesita más particiones para poder ser aislado.

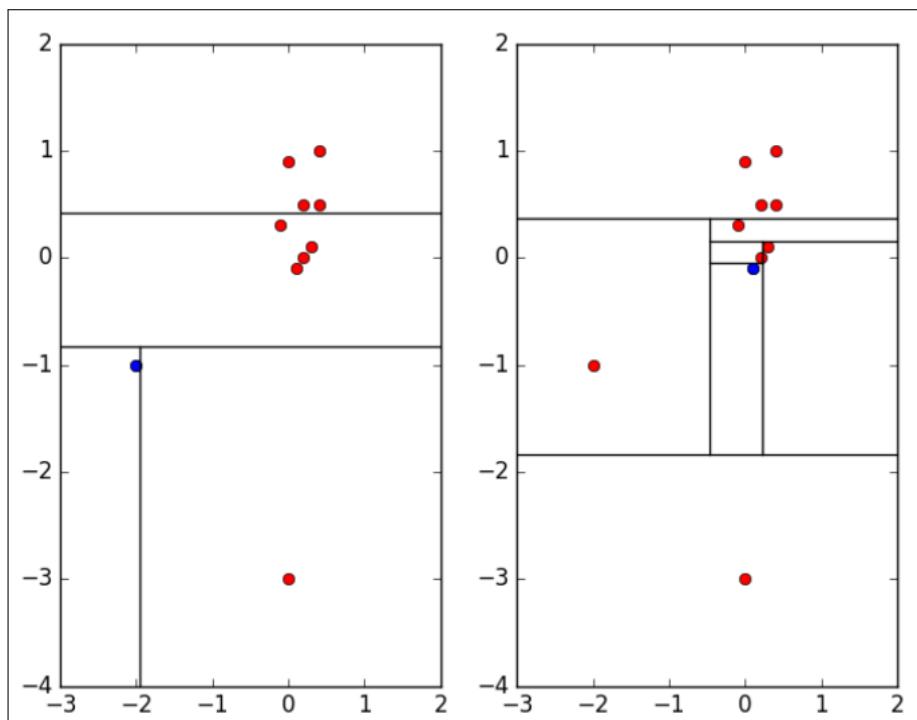


FIGURE 5.5: La figura de la izquierda muestra el aislamiento de una anomalía, que requiere solo tres particiones. A la derecha, el aislamiento de un punto normal requiere seis particiones (Wolpher, n.d.).

Una vez detallado resumidamente el funcionamiento de los bosques de aislamiento se puede proseguir con los diferentes enfoques de los experimentos que se realizará con Isolation Forest.

Existen dos enfoques que pueden realizarse con esta técnica; el primero entrena el modelo con los valores compresos del codificador del autoencoder y el segundo se entrena con los errores de reconstrucción del autoencoder. A continuación se presenta una gráfica (Ver Figura 5.6) del autoencoder (modelo del comportamiento normal) con el fin de tener un mejor entendimiento de cómo se realizarán los experimentos tanto para los bosques de aislamiento como para los SVM de una clase.

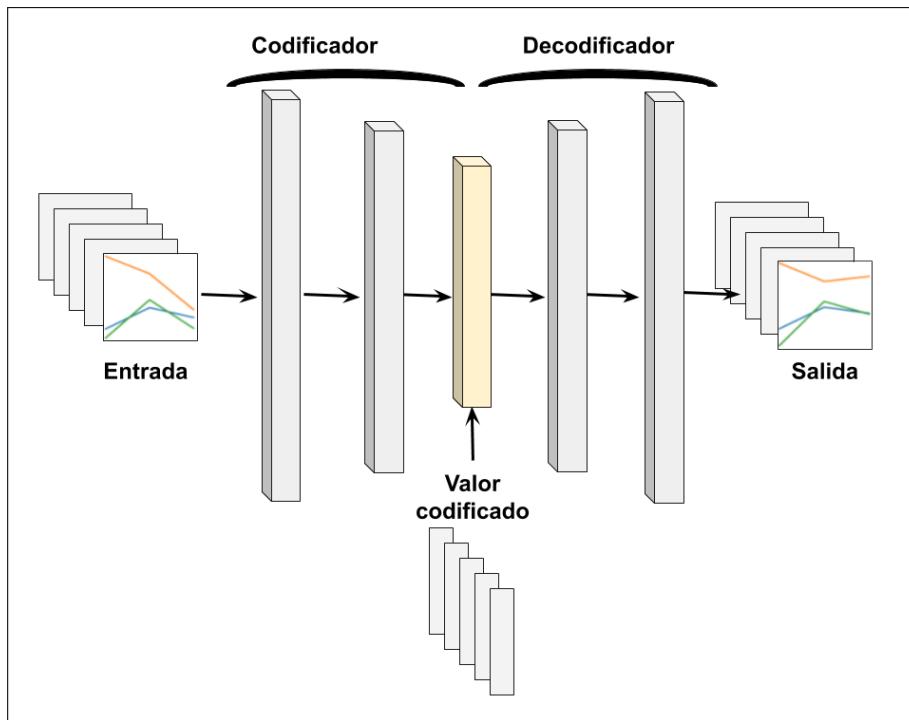


FIGURE 5.6: Representación gráfica del modelo de comportamiento normal o autoencoder (Elaboración propia).

- ***Isolation forest para valores codificados:*** Esta técnica entrena un modelo de bosque de aislamiento con los valores codificados (mediante el codificador del autoencoder, ver Figura 5.6) del conjunto de entrenamiento normal. Para los experimentos se utilizó la clase IsolationForest de SCIKIT-LEARN, esta clase tiene un parámetro llamado CONTAMINACIÓN el cual sirve para definir que cantidad del conjunto de datos está contaminado, es decir, define que cantidad de los datos de entrenamiento pueden ser valores atípicos; en la presente investigación se realizó varias pruebas con diferentes valores para el parámetro contaminación. En el Cuadro 5.8 se presenta los resultados, donde se evidencia que ninguno de los resultados es alentador, ya que la cantidad de anomalías detectadas es muy baja para los tres casos con los que se experimentó.

Contaminación (C)	VP	VN	FN	FP	Sensibilidad	Especificidad
0.0025	3	43944	161	96	0.0183	0.9978
0.0050	17	43817	147	223	0.1037	0.9949
0.0075	17	43738	147	302	0.1037	0.9931

TABLE 5.8: Evaluación de la detección de anomalías usando Isolation forest para valores compresos (Elaboración propia).

- **Isolation forest para errores de reconstrucción:** Para esta técnica se realizó el entrenamiento del bosque de aislamiento con la diferencia de los valores de entrada con los valores obtenidos por el autoencoder (Ver Figura 5.7), cabe aclarar que la diferencia mencionada anteriormente también será llamada *Error de reconstrucción* tanto en esta como en la siguiente subsección. Los resultados de esta técnica para diferentes valores de contaminación se presentan en el Cuadro 5.9, donde estos resultados se pueden considerar como óptimos, debido a que oscilan entre 62 y 67% de detecciones correctas de anomalías, además de presentar una especificidad realmente alta, del 99% aproximadamente, lo cual quiere decir que estos modelos presentan una baja tasa de falsos positivos.

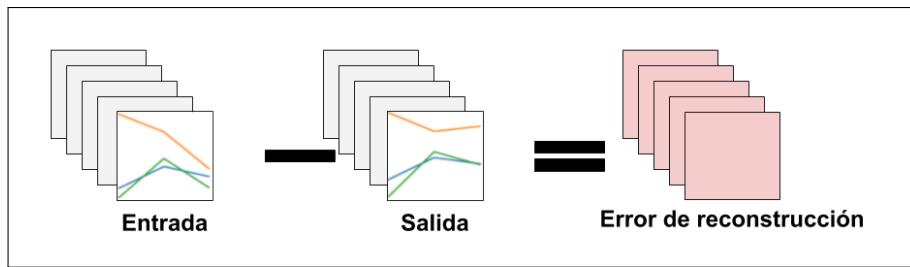


FIGURE 5.7: Representación gráfica del error de reconstrucción usado para el entrenamiento de los bosques de aislamiento y los SVM de una clase (Elaboración propia).

Contaminación (C)	VP	VN	FN	FP	Sensibilidad	Especificidad
0.0025	103	43898	61	142	0.6280	0.9968
0.0050	108	43792	56	248	0.6585	0.9944
0.0075	111	43688	53	352	0.6768	0.9920

TABLE 5.9: Evaluación de la detección de anomalías usando Isolation forest para errores de reconstrucción (Elaboración propia).

One-Class SVM

De la misma forma que en los bosques de aislamiento, se realizó dos diferentes tipos de experimentos con One-Class SVM, a continuación se detalla cada uno de ellos.

- **One-Class SVM para valores codificados:** Se debe entrenar un modelo SVM de una clase para los valores compresos obtenidos por el autoencoder; los experimentos fueron

realizados usando la clase OneClassSVM de SCIKIT-LEARN, donde se tiene diferentes parámetros que pueden ser personalizados, para la presente investigación se probó diferentes kernels, obteniendo así los resultados que se muestran en el Cuadro 5.10, donde claramente ninguno de los resultados obtenidos podría ser tomado en cuenta para ser el método de detección de anomalías de conducción ya que la sensibilidad en ninguno de los casos es superior a 50%.

Kernel	VP	VN	FN	FP	Sensibilidad	Especificidad
rbf	49	41746	115	2294	0.2988	0.9479
poly	56	22532	108	21508	0.3415	0.5116
sigmoid	13	42378	151	1662	0.0793	0.9623

TABLE 5.10: Evaluación de la detección de anomalías usando One-Class SVM para valores compresos (Elaboración propia).

- **One-Class SVM para los errores de reconstrucción:** Al igual que uno de los experimentos que se realizó con Isolation Forest, en esta técnica se usa los errores de reconstrucción (Ver Figura 5.7) para realizar el entrenamiento del modelo SVM de una clase. Como en los experimentos realizados en la anterior técnica se realizó diferentes pruebas con distintos tipos de kernel, a continuación en el Cuadro 5.11 se presenta los resultados obtenidos en los experimentos.

Kernel	VP	VN	FN	FP	Sensibilidad	Especificidad
rbf	134	41887	30	2153	0.8170	0.9511
poly	97	1559	67	42481	0.5915	0.0354
sigmoid	123	1683	41	42357	0.7500	0.0382

TABLE 5.11: Evaluación de la detección de anomalías usando One-Class SVM para el error de reconstrucción del autoencoder (Elaboración propia).

Observando los resultados del Cuadro 5.11 se puede notar que se aumentó notablemente la sensibilidad, o cantidad de anomalías detectadas correctamente, sin embargo, redujo drásticamente la especificidad ya que en algunos casos tan solo llega a un 3.5%, lo cual es muy alejado al objetivo que se persigue en el presente trabajo.

Evaluación del método de detección de anomalías

Una vez realizado los diferentes tipos de experimentos, se evaluó el mejor exponente de cada tipo, con el fin de elegir el más adecuado para la investigación. A continuación se presenta un Cuadro 5.12 con los resultados de los mejores representantes por cada tipo de técnica.

Nombre método		VP	VN	FN	FP	Sensibilidad	Especificidad
Umbralización con S=0.5		111	43814	53	226	0.6768	0.9949
Isolation Forest para errores de reconstrucción con C=0.0075		111	43688	53	352	0.6768	0.9920
One-Class SVM para errores de reconstrucción con kernel RBF		134	41887	30	2153	0.8170	0.9511

TABLE 5.12: Comparación de los mejores métodos de detección de anomalías
(Elaboración propia).

Evidentemente el mejor resultado de detección de anomalías es el que se obtuvo por el modelo SVM para una clase con una sensibilidad del 81.7%, sin embargo, este método presenta la desventaja de tener una alta cantidad de falsos positivos, es decir, por cada anomalía detectada se tendrá aproximadamente 13 alertas por falsos positivos, lo cual es una valor muy alto; y es la principal razón por la que se descarta este método.

Debido a esto sólo quedan dos métodos a comparar, donde ambos resultados son muy similares; ya que estos cuentan con una sensibilidad de 67.68%, por otra parte la especificidad tiene una pequeña variación entre ambas técnicas dando un resultado levemente mejor para la técnica de umbralización con 99.49% frente a un 99.20%.

En este punto se puede elegir cualquiera de estos dos métodos debido a las similitudes que ambos presentan. Por razones de simplificación en este estudio se eligió el método de bosque de aislamiento ya que este método hace más sencilla la detección de anomalías debido a que uno puede especificar la cantidad de contaminación que se espera del conjunto de datos, esto es mucho más ventajoso que la búsqueda de codos con el método de umbralización ya que presenta la desventaja de que es realmente complejo definir el umbral cuando se trata esta técnica en un enfoque no supervisado, como es el caso de esta etapa, además que la definición del umbral depende mucho de cuán limpio o contaminado se encuentra el conjunto de datos, haciendo más complejo el correcto tratamiento al aplicar este método.

Una vez elegido los mejores métodos para conformar el mecanismo de detección de valores atípicos, se procede a formalizar este mecanismo por medio de una gráfica (Ver Figura 5.8), la cual proporciona una representación visual del flujo del detector de anomalías propuesto; ya que es importante conocer como se compone y como funciona, especialmente antes de realizar su respectiva evaluación.

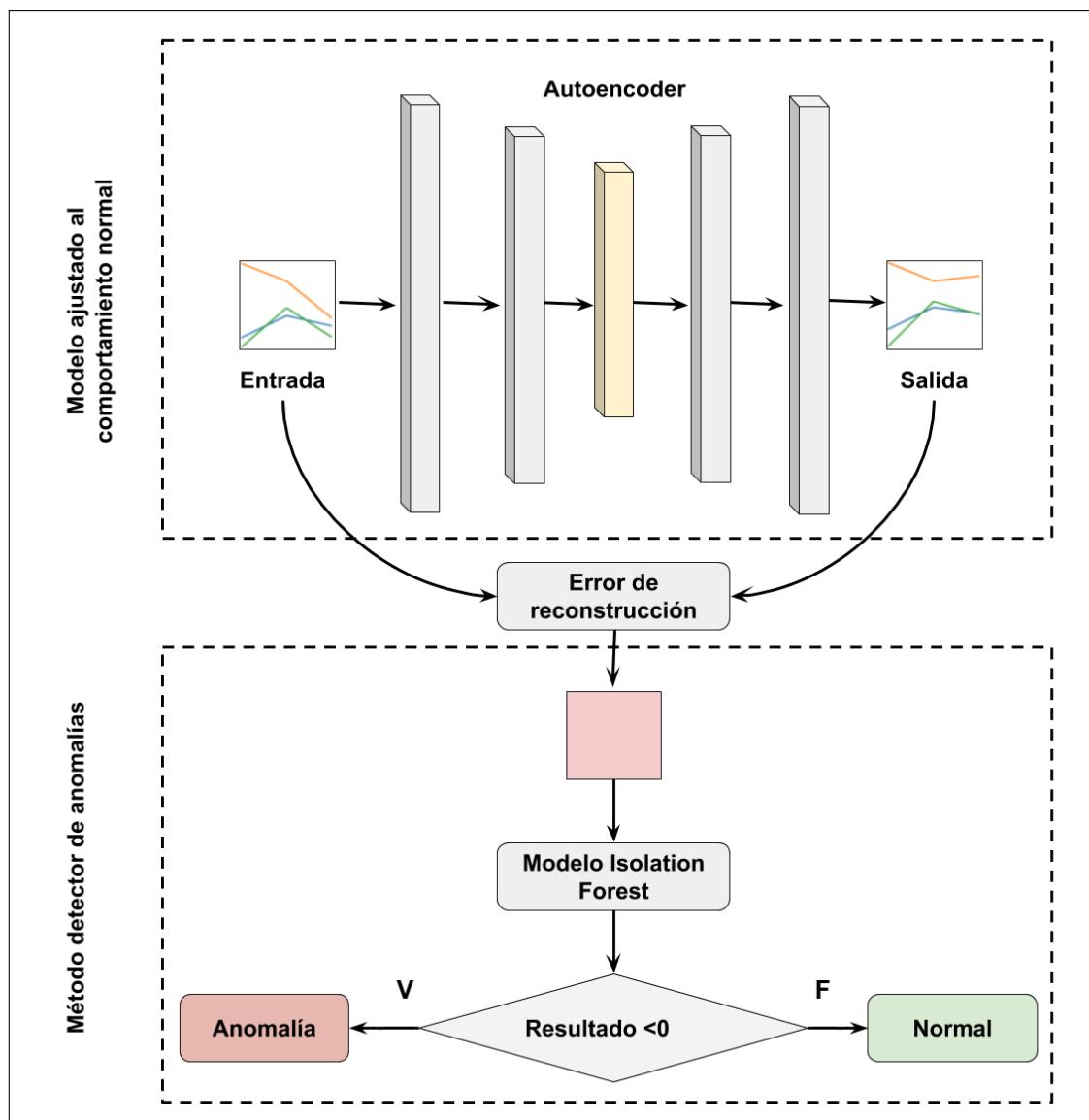


FIGURE 5.8: Mecanismo de detección de anomalías (Elaboración propia).

Observando la Figura 5.8, se puede notar claramente los componentes que conforman el mecanismo de detección de anomalías: el **modelo del comportamiento normal** y el **método detector de anomalías**. Este mecanismo funciona de una forma muy sencilla, en primer lugar se proporciona al autoencoder una secuencia de entrada con 3 pasos para 3 componentes principales (cabe aclarar que los datos de entrada han sido previamente pre-procesados), este autoencoder devuelve como salida la reconstrucción de la entrada, con la cual se obtiene el error de reconstrucción (diferencia entre la entrada real y el valor reconstruido), dicho valor es a su vez la entrada del modelo de bosque de aislamiento, el cual puede retornar dos tipos de valores (-1, 1); cuando este modelo retorna el valor 1 quiere decir que la entrada proporcionada corresponde a un valor considerado como normal y en caso de retornar -1 significa que dicha

entrada es una anomalía, terminando así el flujo del mecanismo de detección propuesto en el presente trabajo.

Chapter 6

RESULTADOS Y EVALUACIÓN

El propósito de este capítulo es presentar los resultados de la evaluación del mecanismo de detección de anomalías propuesto, para posteriormente mostrar algunos de los resultados obtenidos con el mismo.

6.1 Evaluación de desempeño

En este estudio se evaluará la efectividad de las técnicas de detección de anomalías desde las dos siguientes perspectivas:

- La capacidad del enfoque para distinguir entre datos normales y anómalos.
- La eficiencia del método de acuerdo con el tiempo requerido para entrenar el modelo y el tiempo empleado durante el proceso de detección.

6.1.1 Evaluación en términos de rendimiento de detección

Antes de evaluar el mecanismo propuesto en este estudio es importante destacar qué:

- El **modelo de comportamiento normal** fue entrenado con 21000 muestras, durante 50 iteraciones, con 4500 muestras que se usaron para validar el modelo durante la etapa de entrenamiento, y por último el conjunto de prueba con el que se realizó la evaluación final de este modelo esta conformado por 4500 muestras.
- Por otra parte el **método detector de anomalías** fue entrenado con la totalidad de los datos que se usaron en el desarrollo de la generación del modelo del comportamiento normal, es decir, con 30000 muestras.

Para evaluar el mecanismo de detección de anomalías propuesto en el estudio, se utilizó los siguientes criterios: la tasa de detección y la tasa de falsos positivos. La tasa de detección se define como el número de anomalías detectadas dividido por el número total de anomalías. La tasa de falsos positivos se define como el número de series "normales" que se clasifican como anomalías divididos por el número total de series "normales". Es importante aclarar que el conjunto de valores atípicos, con el que se cuenta en esta investigación, no fue usado para el entrenamiento del método propuesto; sin embargo este conjunto sí se usó para validar su

precisión, por lo tanto el conjunto de datos con el que se valida este mecanismo cuenta con 44204 datos.

En la Tabla 6.1 se presenta la matriz de confusión obtenida por el mecanismo propuesto, de donde se pueden obtener las siguientes afirmaciones:

- La entrada superior izquierda de la matriz muestra que 111 anomalías de 164 fueron correctamente etiquetas, es decir, que el 67.68% de las muestras de anomalías se reconocieron correctamente.
- En la fila inferior se muestra que 43688 de 44040 datos fueron etiquetadas correctamente como valores normales, es decir, el 99.20%. Por lo tanto la tasa de falsos positivos para la clase normal es $100 - 99.20\% = 0.80\%$.

		Predicción	
		Anomalía	Clase Normal
Reales	Anomalía	111	53
	Clase Normal	352	43688

TABLE 6.1: Matriz de confusión, para el mecanismo de detección de anomalías
(Elaboración propia).

Estos resultados son un gran avance para la detección de anomalías de conducción con un enfoque semi supervisado, ya que al no contar con muestras de valores atípicos en el entrenamiento es difícil tener una precisión más alta; considerando además, que uno de los valores agregados más importantes que presenta este trabajo de investigación, es el poder generar un modelo personalizado por cada tipo de agente, lo cual es realmente sobresaliente, debido a que el trabajo relacionado que se revisó, previamente a la elaboración de esta investigación, no cuenta con un ejemplar que contemple un enfoque semi-supervisado y mucho menos con modelos que se ajusten y personalicen para cada agente.

6.2 Resultados

Los resultados de este estudio proporcionan una contribución esencial en el campo de la automatización de detección temprana de conductas anómalias en la conducción de automóviles; sin embargo, éstos presentan una visión general del comportamiento del modelo propuesto, por lo cual es necesario realizar un análisis más específico de dicho comportamiento con cada tipo de anomalía presentada en el conjunto de datos, así como también el análisis sobre aquellos valores que fueron detectados erróneamente como anomalías (falsos positivos). A continuación se presentan los resultados de los análisis previamente mencionados.

6.2.1 Detección de anomalías del tipo zig zag

Esta anomalía corresponde a un comportamiento común que suelen realizar agentes que conducen bajo los efectos del alcohol; consiste en una conducción que presenta movimientos

en zig zag de forma brusca, es decir, cambios de dirección constante y a una velocidad relativamente alta. A continuación se presenta algunos de los resultados que se obtuvo con el mecanismo de detección de anomalías propuesto con este trabajo de investigación.

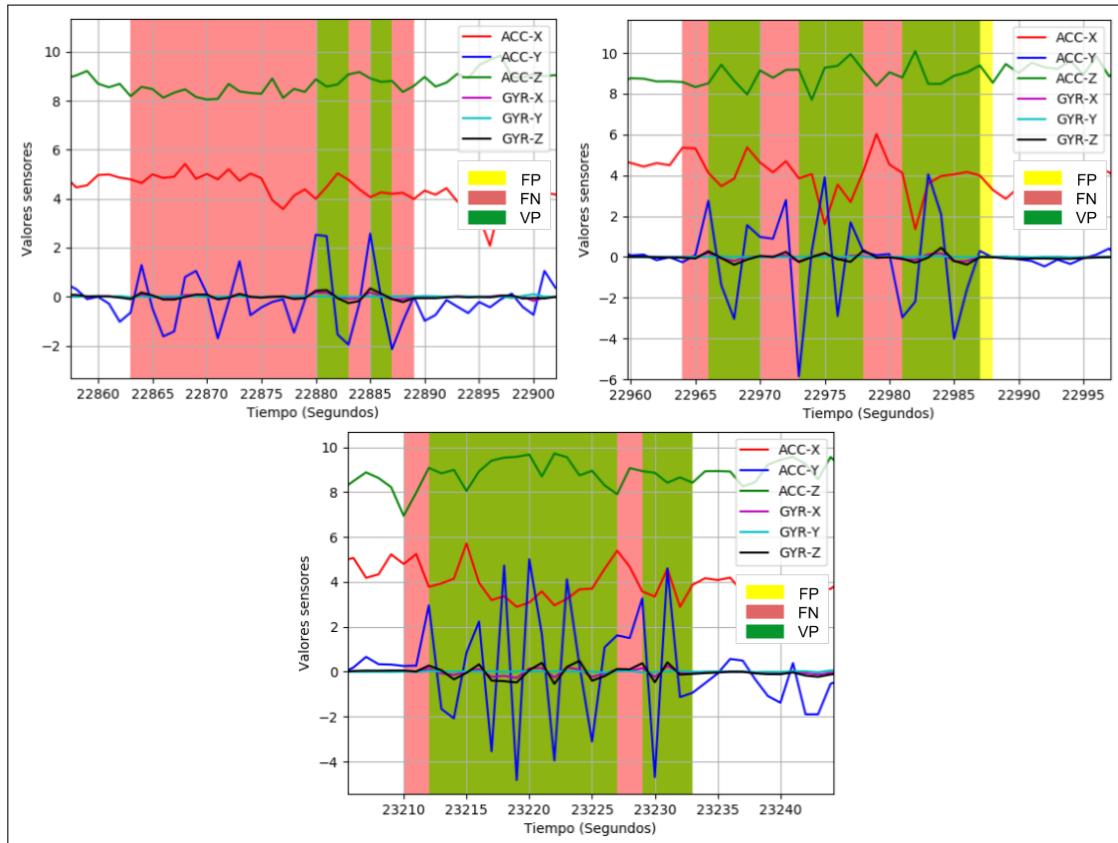


FIGURE 6.1: Resultados de la detección de anomalías del tipo zig zag
(Elaboración propia).

Antes de realizar el análisis de los resultados obtenidos se debe aclarar que aquellas secciones de las siguientes gráficas que se presentan en color rojo son los valores que pertenecen al conjunto de anomalías que no fueron correctamente detectados (Falsos negativos), las secciones en amarillo corresponden a los falsos positivos y por último las secciones verdes son los verdaderos positivos, es decir, aquellos valores que fueron detectados correctamente como anomalías.

Como se observa en la Gráfica 6.1, la imagen superior izquierda presenta una gran cantidad de falsos negativos, esto se debe a que las oscilaciones de los movimientos en Zig Zag no fueron lo suficientemente bruscos, en comparación a los demás, por otro lado la imagen superior derecha presenta una cantidad moderada de falsos negativos y un ejemplar de falso positivo, aunque el resultado no parezca del todo bueno realmente si lo es, ya que muchos de los falsos negativos se encuentran entre valores detectados correctamente, lo cual conllevaría a

una correcta generación de alarma de anomalías a pesar de no detectar como valor atípico la totalidad de los datos anómalos, en la imagen inferior se presenta un ejemplo similar, aunque en este caso no se detectan falsos positivos.

Con el fin de formalizar los resultados para las anomalías del tipo Zig Zag, en la tabla 6.2 se puede observar que 69 anomalías de 105 fueron correctamente detectadas, es decir el 65.71%.

Giros en Zig Zag		
VP	FN	Total
69	36	105

TABLE 6.2: Resultados anomalías tipo Zig Zag (Elaboración propia).

6.2.2 Detección de anomalías del tipo giros a alta velocidad

Este tipo de anomalías suelen ser comunes en agentes que conducen bajo los efectos del alcohol, drogas o con un estado emocional alterado, dichos datos se consideran anomalías ya que los giros normalmente se realizan bajando la velocidad del vehículo, y al realizar este tipo de actos un agente es propenso a ser el causante de un accidente de tránsito.

La Figura 6.2 muestra los resultados obtenidos para las anomalías del tipo giros a alta velocidad, las tres imágenes presentan resultados muy similares, todas tienen una sección en la parte inicial que se presenta como falso negativo, es decir tienen una proporción de datos que no son detectadas correctamente, posteriormente cuentan con un bloque de verdaderos positivos, y por último, dos de las tres imágenes cuentan con un ejemplar de falso positivo posterior a la anomalía. A pesar de que este tipo de anomalías no son detectadas completamente, todas presentan una sección que sí es detectado como anomalía, lo cual es suficiente para generar una alarma oportunamente.

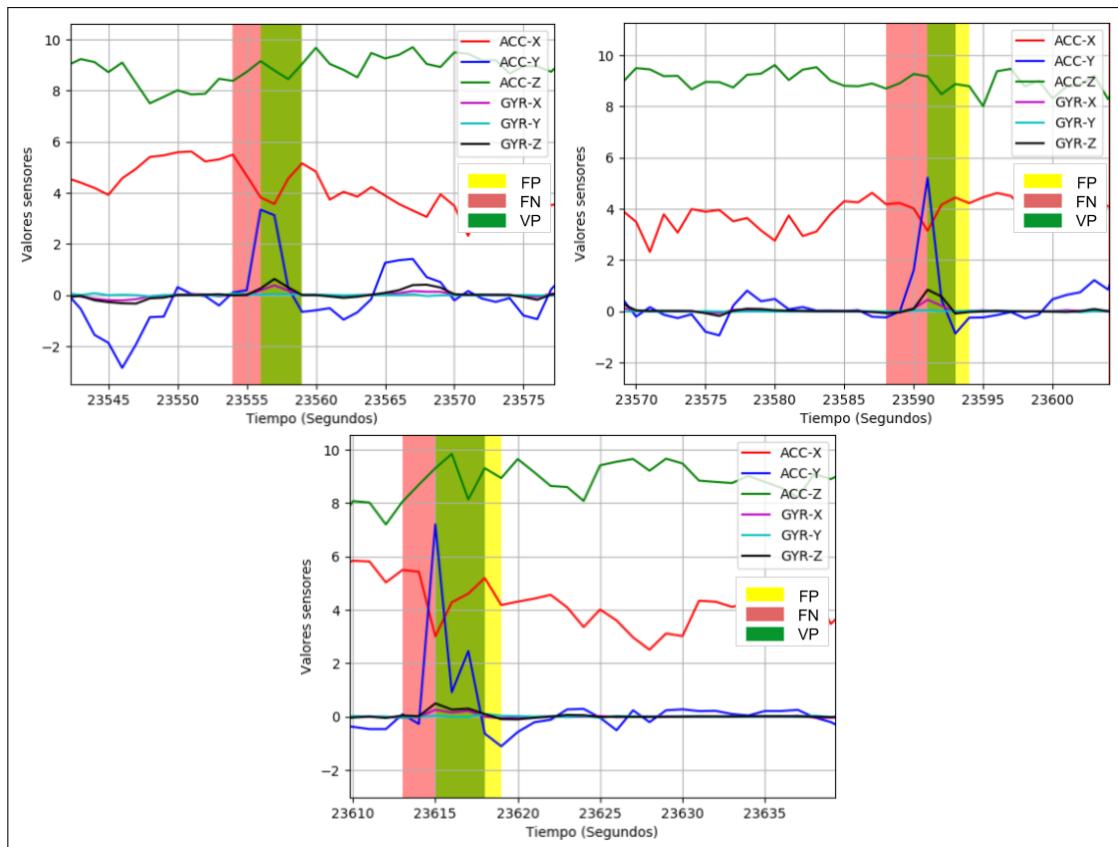


FIGURE 6.2: Resultados de la detección de anomalías del tipo giros a alta velocidad (Elaboración propia).

En el Cuadro 6.3 se presenta los resultados generales obtenidos para las anomalías del tipo Giros a alta velocidad, donde de 35 anomalías 23 fueron correctamente detectadas , es decir, el 65.71%.

Giros a alta velocidad		
VP	FN	Total
23	12	35

TABLE 6.3: Resultados del tipo Giros a alta Velocidad (Elaboración propia).

6.2.3 Detección de anomalías del tipo frenos en seco

Este tipo de anomalía suele ser uno de los valores atípicos más comunes que existen, ya que no sólo se presentan bajo los efectos del alcohol, drogas o fallas mecánicas, sino que también se presentan en contextos de distracción del conductor ya sea por el uso del celular u otro tipo de distracción, ante la aparición de un peatón o mascota que se presenta de manera repentina en la carril que conduce el agente, entre otros casos.

Los resultados de la detección de este tipo de anomalía se presentan en la Figura 6.3, donde al igual que el caso anterior este tipo de anomalía presenta una sección de falsos negativos, posteriormente un grupo de anomalías correctamente detectadas y finalmente falsos positivos; con lo cual es suficiente para generar alertas de manera oportuna y de esa manera poder evitar en lo posible algún accidente de tránsito.

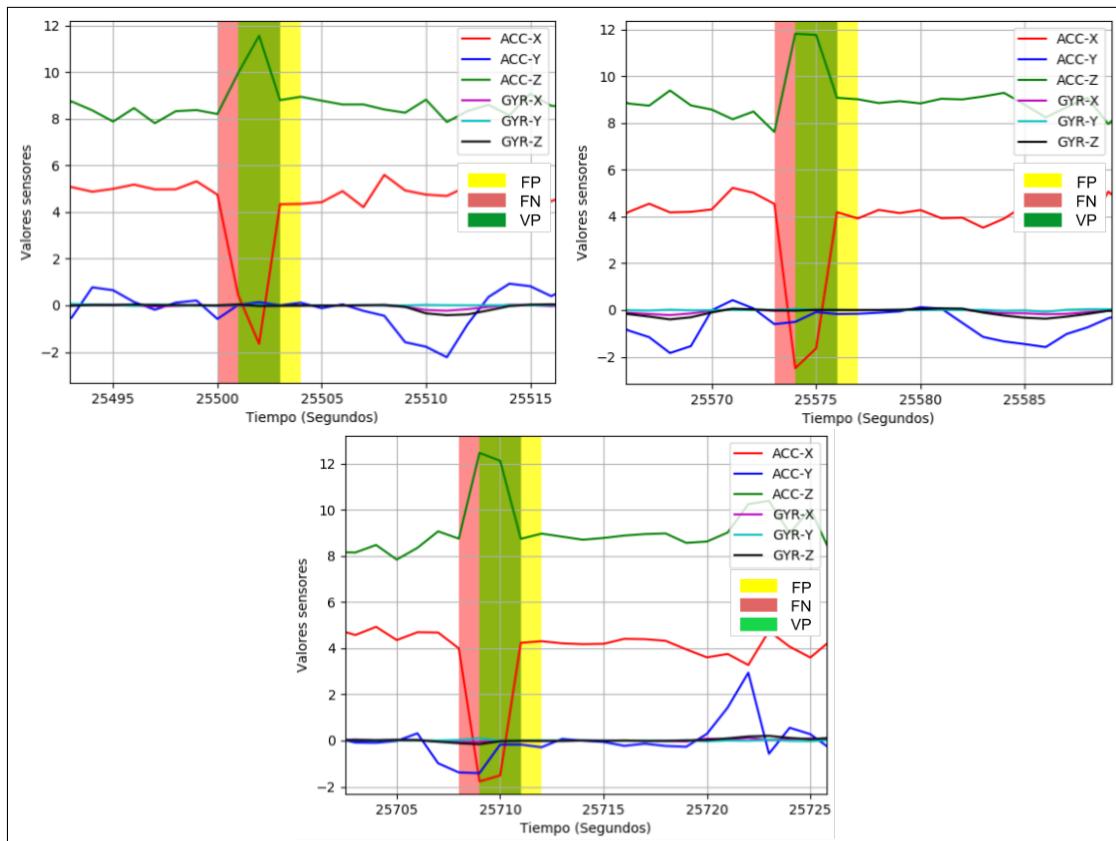


FIGURE 6.3: Resultados de la detección de anomalías del tipo frenos en seco
(Elaboración propia).

A continuación en el Cuadro 6.4 se puede ver que 19 anomalías de 24 fueron correctamente detectadas (79.17%), siendo así el tipo de anomalía que tiene el porcentaje de detección más elevado.

Frenos en seco		
VP	FN	Total
19	5	24

TABLE 6.4: Resultados del tipo Frenos en seco (Elaboración propia).

6.2.4 Detección de falsos positivos

Así como se detectó una gran cantidad de anomalías mediante este mecanismo, también se detectó una proporción considerable de falsos positivos, es decir, valores normales que fueron detectados erróneamente como valores atípicos.

De la misma forma que es importante conocer como este método detecta anomalías, también es importante saber en que casos el modelo propuesto falla; en la Figura 6.4 se presenta algunos casos donde el modelo falla, es decir, esta figura presenta algunos ejemplos de falsos positivos. La figura 6.4 ilustra claramente que estos falsos positivos se presentan generalmente de forma aislada, es decir, uno o dos valores detectados erróneamente como anomalías de forma continua, lo cual es un comportamiento diferente al de los verdaderos valores atípicos, ya que estos presentan una detección de tres valores atípicos de forma continua mínimamente.

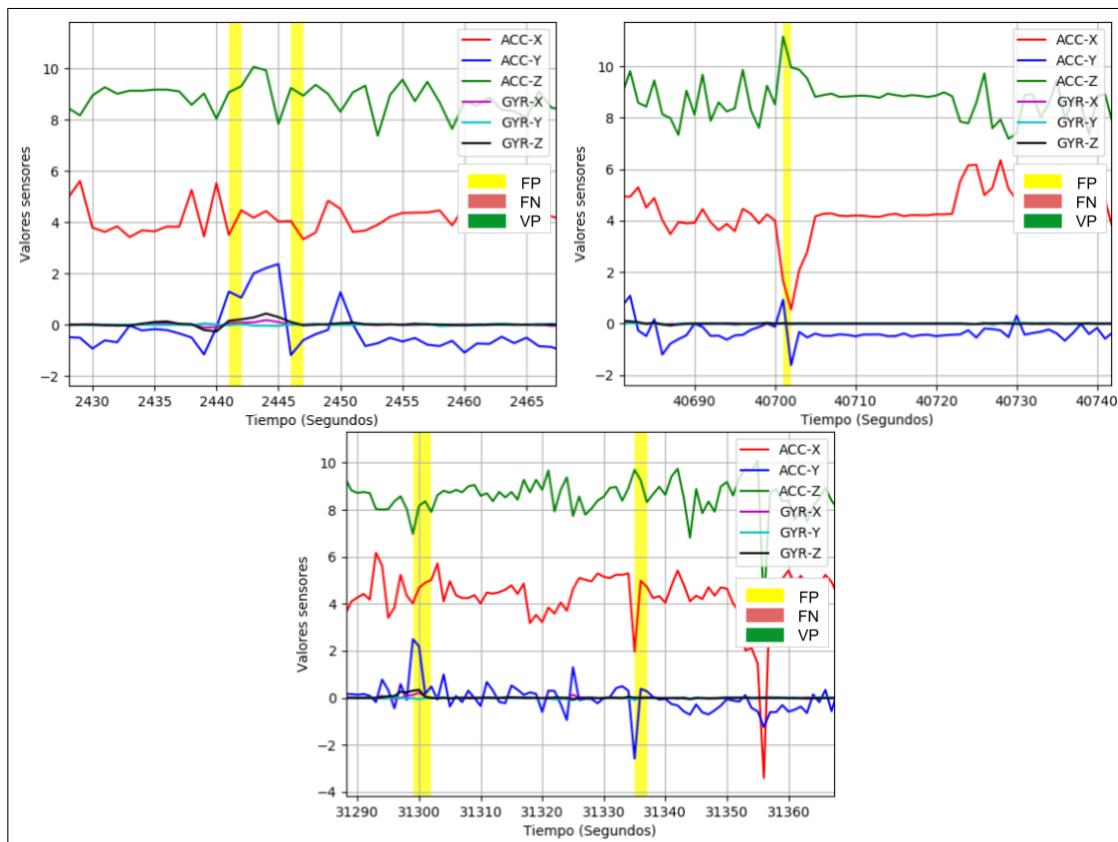


FIGURE 6.4: Resultados de la detección de falsos positivos (Elaboración propia).

Chapter 7

CONCLUSIONES Y TRABAJOS FUTUROS

Después de haber realizado el procedimiento descrito en los anteriores capítulos, con el objetivo de comprobar la hipótesis establecida en la presente investigación, se generó un mecanismo (modelo) capaz de detectar anomalías de conducción. De esta forma se puede decir que se ha cumplido a cabalidad con los objetivos propuestos en esta investigación. A continuación se presentará las conclusiones a las que se llegó, así como también aquellas nuevas ideas e inquietudes que surgieron durante el proceso de desarrollo, las cuales podrían mejorar los resultados obtenidos por el presente trabajo.

7.1 Conclusiones

El objetivo fundamental de este trabajo de investigación fue desarrollar un mecanismo capaz de detectar anomalías de conducción, tal que, se aporte con una solución para alertar de forma oportuna el hallazgo de patrones anómalos en la conducción de agentes, ya sean humanos o autónomos, independizando cada modelo según la experiencia y el ambiente por el que recorre cada agente.

Así pues, el principal aporte de este estudio consiste en la implementación de un mecanismo capaz de identificar anomalías a partir de los datos de conducción normal de cada agente, sin intervención humana, es decir, el modelo detector no requiere que un humano intervenga para generarlo, sin embargo, este puede ser optimizado por medio del ajuste del hiperparámetro *Contaminación* con el fin de definir cuán sensible a las anomalías será dicho detector. Por otra parte, se puede decir que el mecanismo de detección de este trabajo de investigación, además de ser novedoso, es uno de los pocos trabajos que se realizaron con un enfoque “semi-supervisado”, ya que la mayoría de los trabajos realizados a la fecha fueron realizados mediante un enfoque supervisado.

Las conclusiones que se derivan de este trabajo de investigación se hicieron en base a los diferentes experimentos realizados, dichas conclusiones se exponen a continuación.

- Se comprueba, a partir del análisis de resultados de este estudio, la capacidad con la que cuentan los sensores iniciales de un dispositivo móvil para representar correctamente el movimiento de un automóvil y de esa manera ser capaz de alimentar, con un previo pre-procesamiento, un mecanismo de detección de anomalías.
- En este trabajo se compararon diferentes arquitecturas de redes neuronales para generar el modelo del comportamiento normal, donde la red más simple logró los mejores resultados tanto en precisión como en el tiempo empleado durante el proceso de predicción; demostrando así, que no siempre las redes más complejas interpretan mejor los conjuntos de datos.
- Por otra parte, se comparó diferentes técnicas para definir un método de detección de anomalías adecuado al contexto de la presente investigación, donde por la simplicidad de su entrenamiento y por su robusto resultado se optó por la elección de la técnica de bosques de aislamiento, con un valor de 0.0075 para el hiperparámetro *Contaminación*.
- Integrando el modelo del comportamiento normal y el método de detección de anomalías, los cuales sólo fueron entrenados con el conjunto de datos "normal", se logra la creación de un mecanismo capaz de identificar valores atípicos de la conducción de cada agente.
- Finalmente se evaluó la capacidad del mecanismo de detección, mediante el conjunto de evaluación el cuál presenta muestras anómalas, dando como resultado la correcta detección del 67.68% de las muestras, que presentan anomalías en el conjunto de evaluación, así como también presenta una tasa de tan sólo 0.80% de muestras normales detectadas como anomalías. Siendo un gran avance en el ámbito de la detección de anomalías con un enfoque semi-supervisado.

Este trabajo de investigación antes que presentar una solución final sienta las bases para el desarrollo de sistemas de detección de anomalías de la conducción de los agentes, mediante el uso de técnicas de Inteligencia Artificial, resaltando la capacidad y alcance que conlleva este estudio, ya que no sólo se enfoca en la conducción de agentes humanos, sino que es igual de capaz de ser aplicado en un enfoque de conducción autónomo.

7.2 Trabajos futuros

Una vez concluido el trabajo de investigación, se considera interesante investigar sobre diferentes aspectos de la detección de anomalías y se propone:

- Agregar la velocidad del vehículo como un nuevo parámetro del conjunto de datos, debido a que esto podría brindar un mejor entendimiento del comportamiento normal de conducción, así como también de las anomalías.
- En lugar de trabajar con los datos en crudo, usar la diferencia entre un dato capturado en el tiempo t y un dato capturado en $t - 1$ ($diff_t = dato_t - dato_{t-1}$), la aplicación de éste pre-procesamiento de datos podría maximizar la detección de aquellas anomalías que presentan elevadas diferencias entre los datos consecutivos.

- Validar el modelo con nuevos tipos de anomalías como por ejemplo: derrapes, choques, giros en U a alta velocidad, entre otros. Esto debido a que el estudio se limitó al reconocimiento de sólo tres tipos de anomalías por la dificultad y peligro que conlleva su captura.
- Extender el modelo para que sea capaz de determinar no sólo una anomalía sino también el tipo al que dicha anomalía pertenece.
- Migrar el modelo del comportamiento normal de keras a Tensorflow 2.0.
- Implementar un sistema de información para monitorear las anomalías mediante el mecanismo de detección propuesto en el presente trabajo.

BIBLIOGRAPHY

References

- 0.22, S. (n.d.). *Novelty and outlier detection*. https://scikit-learn.org/stable/modules/outlier_detection.html. (Último acceso en 19 de diciembre de 2019)
- Alashwal, H., Bin D., S., & Othman, R. (2006, 01). One-class support vector machines for protein protein interactions prediction. *Int J Biomed Sci*, 1.
- Araujo, R., Igreja, A., R., D. C., & Araujo, R. (2012, 6). Driving coach: A smartphone application to evaluate driving efficient patterns. *Proceedings of the 2012 IEEE Intelligent Vehicles Symposium (IV); Alcala de Henares, España.*(3–7), 1005–1010. Retrieved from https://www.researchgate.net/publication/261309792_Driving_coach_A_smartphone_application_to_evaluate_driving_efficient_patterns
- Bellman, R. E. (2003). *Dynamic programming*. Courier Dover Publications ISBN.
- Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Netw*, 5, 157–166. Retrieved from <https://ieeexplore.ieee.org/document/279181/authors#authors>
- Bhoyar, V., Lata, P., Katkar, J., Patil, A., & Javale, D. (2013, 3–4). Symbian based rash driving detection system. *International Journal of Emerging Trends & Technology in Computer Science (IJETTCS)*, 2, 124–126. Retrieved from <https://www.ijettcs.org/Volume2Issue2/IJETTCS-2013-03-28-046.pdf>
- Bishop, M. C. (2006). *Pattern recognition and machine learning*. Springer Science+Business Media, LLC.
- Boonmee, S., & Tangamchit, P. (2009, 5). Portable reckless driving detection system. *Proceedings of the 6th IEEE International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology; Pattaya, Chonburi, Tailandia.*(6–9), 412–415. Retrieved from <https://ieeexplore.ieee.org/abstract/document/5137037>
- Chen, Z., Yu, J., Zhu, Y., Chen, Y., & Li, M. (2015, 6). D3: Abnormal driving behaviors detection and identification using smartphone sensors. *Proceedings of the 12th Annual IEEE International Conference on Sensing, Communication, and Networking (SECON); Seattle, WA, USA.*, 20–25. Retrieved from <http://www.winlab.rutgers.edu/~yychen/papers/D3-Abnormal%20Driving%20Behaviors%20Detection%20and%20Identification%20Using%20Smartphone%20Sensors.pdf>
- Cho, K., Merriënboer, B. V., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014a). Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv*, 1406.1078. Retrieved from <https://www.aclweb.org/anthology/D14-1179.pdf>

- Dang-Nhac, L., Duc-Nhan, N., Thi-Hau, N., & Ha-Nam, N. (2018, 4). Vehicle mode and driving activity detection based on analyzing sensor data of smartphones. *Sensors*, 18(4), 1036. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5948751/>
- Dauphin, Y. N., Fan, A., Auli, M., & Grangiera, D. (2017, 9). Language modeling with gated convolutional networks. *arXiv*, 1612.08083v3(8). Retrieved from <https://arxiv.org/pdf/1612.08083.pdf>
- de Salud (OMS), O. M. (n.d.). *Informe de la situación mundial de la seguridad vial 2015*. https://www.who.int/violence_injury_prevention/road_safety_status/2015/Summary_GSRRS2015_SPA.pdf?ua=1. (Último acceso en 19 de diciembre de 2019)
- Elman, J. (1990). Finding structure in time. *Cognitive Science*, 14(2), 179–211. Retrieved from <https://crl.ucsd.edu/~elman/Papers/fsit.pdf>
- Eren, H., Makinist, S., Akin, E., & Yilmaz, A. (2012, 6). Estimating driving behavior by a smartphone. *Proceedings of the Intelligent Vehicles Symposium. Alcalá de Henares, España.*(3–7), 234—239.
- Ferreira, J., Carvalho, E., Ferreira, B., De Souza, C., Suhara, Y., Pentland, A., & Pessin, G. (2017). Driver behavior profiling: An investigation with different smartphone sensors and machine learning. *PLoS One*, 12(4), e0174959. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5386255/>
- Galarneck, M. (n.d.). *Explaining the 689599.7 rule for a normal distribution*. <https://towardsdatascience.com/understanding-the-68-95-99-7-rule-for-a-normal-distribution-b7b7cbf760c2>. (Último acceso en 19 de diciembre de 2019)
- Guo, Y., Liao, W., Wang, Q., Yu, L., Ji, T., & Li, P. (2018). Multidimensional time series anomaly detection: A gru-based gaussian mixture variational autoencoder approach. *Proceedings of Machine Learning Research*, 95, 97–112. Retrieved from <http://proceedings.mlr.press/v95/guo18a/guo18a.pdf>
- Hawkins, S., He, H., Williams, G., & Baxter, R. (2002, 9). Outlier detection using replicator neural networks. *International Conference on Data Warehousing and Knowledge Discovery*, 170—180. Retrieved from <https://togaware.com/papers/dawak02.pdf>
- Hochreiter, S., & Schmidhuber, J. (1997, 11). Long short-term memory. *Neural Comput*, 9(8)(15), 1735—1780. Retrieved from <https://www.mitpressjournals.org/doi/abs/10.1162/neco.1997.9.8.1735>
- Hsu, A., & Griffiths, T. (2010). Effects of generative and discriminative learning on use of category variability. Retrieved from <https://cocosci.princeton.edu/tom/papers/discgencat.pdf>
- Jayesh, B. (n.d.). *The artificial neural networks handbook: Part 4*. <https://medium.com/@jayeshbahire/the-artificial-neural-networks-handbook-part-4-d2087d1f583e>. (Último acceso en 19 de diciembre de 2019)
- Jing, Y., & Guanci, Y. (2018, 03). Modified convolutional neural network based on dropout and the stochastic gradient descent optimizer. *Algorithms*, 11, 28.
- Johnson, D., & Trivedi, M. (2011, 10). Driving style recognition using a smartphone as a sensor platform. *14th International IEEE Conference en Intelligent Transportation Systems (ITSC)*, 1609—1615. Retrieved from http://cvrr.ucsd.edu/publications/2011/Johnson_ITSC2011.pdf
- Klos, M., & Waszczyszyn, Z. (2011). Modal analysis and modified cascade neural networks in identification of geometrical parameters of circular arches. *Computers & Structures*, 89,

- 581–589. Retrieved from <http://cames.ippt.gov.pl/index.php/cames/article/view/110>
- Koh, D. W., & Kang, H. (2015, 7). Smartphone-based modeling and detection of aggressiveness reactions in senior drivers. *Proceedings of the IEEE Intelligent Vehicles Symposium; Seoul, Korea.*(1), 12–17. Retrieved from <https://ieeexplore.ieee.org/abstract/document/7225655>
- Kridalukmana, R., Yan-Lu, H., & Naderpour, M. (2017). An object oriented bayesian network approach for unsafe driving maneuvers prevention system. *12th International IEEE Conference*. Retrieved from <https://opus.lib.uts.edu.au/bitstream/10453/122196/4/633634.pdf>
- Lecun, Y., Bengio, Y., & Hinton, G. (2015, 5). Deep learning. *Nature*, 521(7553)(27), 436—444. Retrieved from <https://doi.org/10.1038/nature14539>
- Lecun, Y., Jackel, L., Boser, B., Denker, J., Graf, H., Guyon, I., ... Hubbard, W. (1998). Handwritten digit recognition : Applications of neural networks chips and automatic learning. *Proceedings of the IEEE*, 86(11), 2278–2324. Retrieved from <http://yann.lecun.com/exdb/publis/pdf/lecun-89c.pdf>
- Mass, A., Hannun, A., & Ng, A. (2013). Rectifier nonlinearities improve neural network acoustic models. *International Conference on Machine Learning (icml)*. Retrieved from https://ai.stanford.edu/~amaas/papers/relu_hybrid_icml2013_final.pdf
- Michael, A. (2015). *Neural networks and deep learning*. Determination Press. Retrieved from <http://neuralnetworksanddeeplearning.com/index.html>
- Moindrot, O., & Genthial, G. (2018, 1). *Splitting into train, dev and test sets*. <https://cs230-stanford.github.io/train-dev-test-split.html>. (Último acceso en 16 de octubre de 2019)
- Muhammad, R. (n.d.). *Convolutional neural network. in a nut shell*. <https://engmrk.com/convolutional-neural-network-3/>. (Último acceso en 19 de diciembre de 2019)
- Olah, C. (n.d.). *Understanding lstm networks*. <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>. (Último acceso en 11 de octubre de 2019)
- Özler, H. (n.d.). *Accuracy trap! pay attention to recall, precision, f-score, auc*. <https://medium.com/datadriveninvestor/accuracy-trap-pay-attention-to-recall-precision-f-score-auc-d02f28d3299c>. (Último acceso en 16 de octubre de 2019)
- Pascanu, R., Mikolov, T., & Bengio, Y. (2013, 6). On the difficulty of training recurrent neural networks. In *Proceedings of the International Conference on Machine Learning; Atlanta, GA, USA*(16–21), 1310–1318. Retrieved from <http://proceedings.mlr.press/v28/pascanu13.pdf>
- Pyle, D. (1999). *Data preparation for data mining*. Morgan Kaufmann Publishers, Inc. Retrieved from <https://pdfs.semanticscholar.org/470a/828d5e3962f2917a0092cc6ba46ccfe41a2a.pdf>
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by backpropagating errors. *Nature*, 323(6088), 533–536. Retrieved from <https://psycnet.apa.org/record/1987-33645-001>
- Russakovsky, O. e. a. (2014). Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*. Retrieved from <http://link.springer.com/article/10.1007/s11263-015-0816-y#>
- Schölkopf, B., & Smola, A. J. (2002). *Support vector machines, regularization, optimization, and beyond*. MIT Press.

- Shai, S., & Shai, B. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge University Press. Retrieved from <https://www.cs.huji.ac.il/~shais/UnderstandingMachineLearning/understanding-machine-learning-theory-algorithms.pdf>
- Smits, P., Dellepiane, S., & Schowengerdt, R. (1999). Quality assessment of image classification algorithms for land-cover mapping: a review and a proposal for a cost-based approach. *International journal of remote sensing* 20, 8, 1461—1486.
- Suad, A., & Wesam, S. (2017). Review of data preprocessing techniques in data mining. *Review of Scientific Instruments*, 12(16), 4102–4107. Retrieved from <http://docsdrive.com/pdfs/medwelljournals/jeasci/2017/4102-4107.pdf>
- Varun, C., & Arindam, K., B.and Vipin. (2009). *Anomaly detection: A survey*. Universidad de Minnesota. Retrieved from https://www.researchgate.net/publication/220565847_Anomaly_Detection_A_Survey
- Werbos, P. J. (1988). Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 1(4), 339—356. Retrieved from <https://www.sciencedirect.com/science/article/abs/pii/089360808890007X>
- Who-Lee, K., Sik-Yoon, H., Min-Song, J., & Ryoung-Park, K. (2018, 4). Convolutional neural network-based classification of driver's emotion during aggressive and smooth driving using multi-modal camera sensor. *Sensors*, 18(4), 957. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5948584/>
- Wikipedia. (n.d.). *Neuron*. <https://simple.wikipedia.org/wiki/Neuron>. (Último acceso en 19 de diciembre de 2019)
- Williams, G., & Baxter, R. (2002, 12). A comparative study of rnn for outlier detection in data mining. *IEEE International Conference on Data Mining*, 1–16. Retrieved from <https://togaware.com/papers/tr02102.pdf>
- Wolpher, M. (n.d.). *Anomaly detection in unstructured time series data using an lstm autoencoder*. <http://www.diva-portal.org/smash/get/diva2:1225367/FULLTEXT01.pdf>. (Último acceso en 11 de octubre de 2019)
- Xue, Z., Shang, Y., & Feng, A. (2010, 5). Semi-supervised outlier detection based on fuzzy rough c-means clustering. *Mathematics and Computers in Simulation*, 80(9), 1911—1921. Retrieved from https://www.researchgate.net/publication/220348246_Semi-supervised_outlier_detection_based_on_fuzzy_rough_C-means_clustering
- Yan, S. (n.d.). *Understanding lstm and its diagrams*. <https://medium.com/mlreview/understanding-lstm-and-its-diagrams-37e2f46f1714>. (Último acceso en 11 de octubre de 2019)
- Zaldivar, J., Calafate, C., Cano, J., & Manzoni, P. (2011, 10). Providing accident detection in vehicular networks through obd-ii devices and android-based smartphones. *Proceedings of the 2011 IEEE 36th Conference on Local Computer Networks; Bonn, Alemania.(4–7)*, 813—819.
- Zeiler, M. D., Ranzato, M., Monga, R., Mao, M., Yang, K., Le, Q. V., & Hinton, G. E. (2013). On rectified linear units for speech processing. *International Conference on Acoustics, Speech and Signal Processing. IEEE*, 3517—3521. Retrieved from <https://static.googleusercontent.com/media/research.google.com/es//pubs/archive/40811.pdf>
- Zenon, W. (2011). Artificial neural networks in civil engineering: another five years of research in poland. *Computer Assisted Mechanics and Engineering Sciences*, 18, 131–146.

Retrieved from <http://cames.ippt.gov.pl/index.php/cames/article/view/110>
Zhang, A., Lipton, Z., Li, M., & Smola, A. (2019, 9). *Dive into deep learning*. <https://en.d2l.ai/d2l-en.pdf>. (Último acceso en 11 de octubre de 2019)

Appendix A

Experimentos de diferentes arquitecturas para los autoencoders

A.1 Redes densas

A.1.1 Redes densas para 3 componentes

Arquitecturas Densas				
	Tipo	Salida	Activacion	# Parametros
NN_33	Input	(3,3)		0
	Flatten	9		0
	Dense	8	elu	80
	Dense	5	elu	45
	Dense	8	elu	48
	Dense	9	tanh	81
	Reshape	(3,3)		0
NN_43	Input	(4,3)		0
	Flatten	12		0
	Dense	10	elu	130
	Dense	5	elu	55
	Dense	8	elu	48
	Dense	12	tanh	108
	Reshape	(4,3)		0
NN_53	Input	(5,3)		0
	Flatten	15		0
	Dense	10	elu	160
	Dense	6	elu	66
	Dense	11	elu	77
	Dense	15	tanh	180
	Reshape	(5,3)		0

TABLE A.1: Arquitectura densa para 3 componentes principales (Elaboración propia).

A.1.2 Evaluación redes densas

Red	val_loss	val_acc	val_f1score	loss	acc	f1score
NN_33	0.003956	0.899894	0.287709	0.003898	0.900735	174173.640890
NN_43	0.006572	0.869167	0.233547	0.006104	0.869548	87092.835609
NN_53	0.006400	0.846857	101587.687459	0.006226	0.850568	0.283059

TABLE A.2: Tabla de evaluación de redes densas (Elaboración propia).

A.2 Redes convolucionales

A.2.1 Redes convolucionales para 3 componentes

Arquitecturas Convolucionales					
		Tipo	Salida	Activacion	# Parametros
Redes Conv - 3 componentes principales	CNN_33	Input	(3,3)		0
		Conv1D	(3,2)	elu	20
		MaxPooling1D	(2,2)		0
		Conv1D	(2,4)	elu	28
		MaxPooling1D	(1,4)		0
		Conv1D	(1,6)	elu	54
		UpSampling1D	(3,6)		0
		Conv1D	(3,3)	tanh	57
Redes Conv - 3 componentes principales	CNN_43	Input	(4,3)		0
		Conv1D	(4,2)	elu	20
		MaxPooling1D	(2,2)		0
		Conv1D	(2,4)	elu	26
		MaxPooling1D	(1,4)		0
		Conv1D	(1,6)	elu	54
		UpSampling1D	(4,6)		0
		Conv1D	(4,3)	tanh	57
Redes Conv - 3 componentes principales	CNN_53	Input	(5,3)		0
		Conv1D	(5,2)	elu	20
		MaxPooling1D	(3,2)		0
		Conv1D	(3,4)	elu	28
		MaxPooling1D	(2,4)		0
		Conv1D	(2,5)	elu	45
		Reshape	(5,2)		0
		Conv1D	(5,3)	tanh	15

TABLE A.3: Arquitectura convolucional para 3 componentes principales (Elaboración propia).

A.2.2 Evaluación redes convolucionales

Red	val_loss	val_acc	loss	acc
CNN_33	0.0070	0.8453	0.0067	0.8434
CNN_43	0.0100	0.8136	0.0097	0.8152
CNN_53	0.0102	0.7951	0.0108	0.7907

TABLE A.4: Tabla de evaluación de redes convolucionales (Elaboración propia).

A.3 Redes recurrentes

A.3.1 Redes recurrentes para 3 componentes

Arquitecturas Recurrentes				
	Tipo	Salida	Activacion	# Parametros
RNN_33	Input	(3,3)		0
	LSTM	(3,9)	elu	468
	LSTM	6	elu	384
	Reshape	(3,2)		0
	LSTM	(3,3)	elu	72
	LSTM	(3,9)	elu	468
	TimeDistributed(Dense(3))	(3,3)	tanh	30
RNN_43	Input	(4,3)		468
	LSTM	(4,9)	elu	384
	LSTM	6	elu	0
	Reshape	(3,2)		216
	LSTM	(3,6)	elu	576
	LSTM	(3,9)	elu	40
	TimeDistributed(Dense(3))	(3,4)	tanh	0
	Reshape	(4,3)		0
RNN_53	Input	(5,3)		0
	LSTM	(5,9)	elu	468
	LSTM	6	elu	384
	Reshape	(3,2)		0
	LSTM	(3,3)	elu	72
	LSTM	(3,9)	elu	468
	TimeDistributed(Dense(3))	(3,5)		50
	Reshape	(5,3)	tanh	0

TABLE A.5: Arquitectura recurrente para 3 componentes principales (Elaboración propia).

A.3.2 Evaluación redes recurrentes

Red	val_loss	val_acc	loss	acc
RNN_33	0.0039	0.8855	0.0037	0.8900
RNN_43	0.0108	0.7871	0.0102	0.7884
RNN_53	0.0295	0.2986	0.0290	0.3370

TABLE A.6: Tabla de evaluación de redes recurrentes (Elaboración propia).

Appendix B

Arquitectura del Sistema de Demostración

Si bien el trabajo presentado no promete un sistema que implemente el mecanismo propuesto, es necesario tener un prototipo para demostrar el correcto funcionamiento del detector de anomalías, por lo que este anexo se centra en presentar la arquitectura tanto física como lógica del prototipo.

B.1 Arquitectura Física

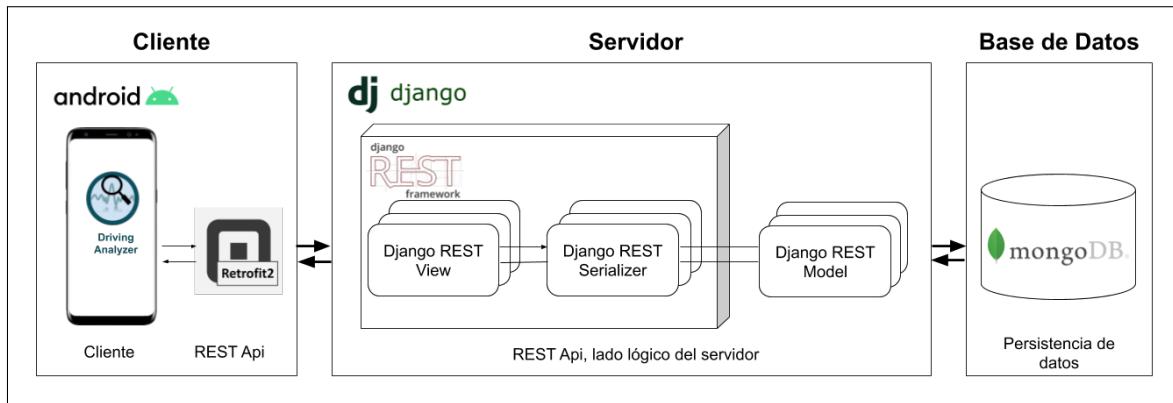


FIGURE B.1: Arquitectura física del Sistema de Demostración (Elaboración propia).

La arquitectura física del prototipo cuenta con tres partes principales: el cliente, el servidor y la base de datos, como se puede observar en la Figura B.1. En primer lugar el cliente está compuesto por una aplicación móvil para Android encargada de capturar los datos de los sensores iniciales, para posteriormente enviarlas al servidor usando Retrofit 2, posteriormente los datos ingresan al servidor realizado en Django, una vez el servidor recibe estos datos, se encarga de preprocesarlos para posteriormente predecir si los datos enviados son anomalías o no lo son, y por último los datos recibidos son almacenados en la Base de Datos de MongoDB con su respectiva predicción con el objetivo de hacer un monitoreo de la conducción del usuario.

B.2 Arquitectura Lógica

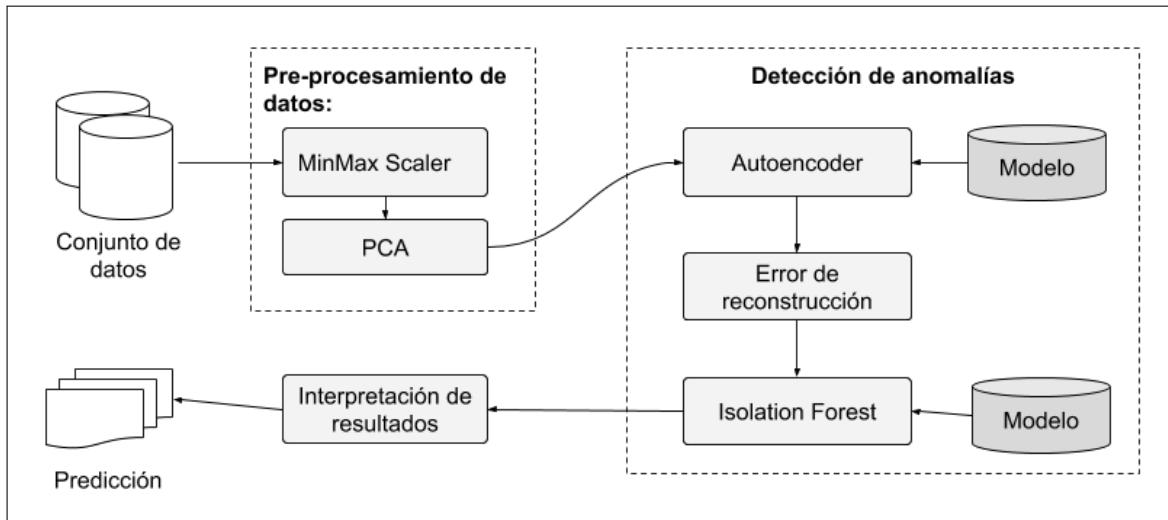


FIGURE B.2: Arquitectura Lógica del Sistema de Demostración (Elaboración propia).

En cuanto a la arquitectura lógica se observa las dos principales partes del mecanismo de detección propuesto: el pre-procesamiento y la detección de anomalías, como se puede observar en la Figura B.2.