

Universities of Burgos, León and  
Valladolid

Master's degree

# Business Intelligence and Big Data in Cyber-Secure Environments



Thesis of the Master's degree in  
Business Intelligence and Big Data in  
Cyber-Secure Environments

título del TFM

Presented by Adrián Riesco Valbuena  
in University of Burgos — February 9, 2022  
Tutor: Alvar Arnáiz González





# Universities of Burgos, León and Valladolid



## Master's degree in Business Intelligence and Big Data in Cyber-Secure Environments

Mr. Alvar Arnáiz González, professor of the department named Computer Engineering, area named Computer Languages and Systems.

Exposes:

That the student Mr. Adrián Riesco Valbuena, with DNI 71462231N, has completed the Thesis of the Master in Business Intelligence and Big Data in Cyber-Secure Environments titled NOMBRE TFM.

And that thesis has been carried out by the student under the direction of the undersigned, by virtue of which its presentation and defense is authorized.

In Burgos, February 9, 2022

Approval of the Tutor:

Mr. Alvar Arnáiz González





## Resumen

En este primer apartado se hace una **breve** presentación del tema que se aborda en el proyecto.

## Descriptores

Palabras separadas por comas que identifiquen el contenido del proyecto Ej: servidor web, buscador de vuelos, android ...

## **Abstract**

A **brief** presentation of the topic addressed in the project.

## **Keywords**

keywords separated by commas.



---

# Contents

---

<b>Contents</b>	<b>iii</b>
<b>List of Figures</b>	<b>v</b>
<b>List of Tables</b>	<b>vi</b>
 <b>Memory</b>	 <b>1</b>
<b>1. Introduction</b>	<b>3</b>
<b>2. Project objectives</b>	<b>5</b>
<b>3. Theoretical concepts</b>	<b>7</b>
3.1 API . . . . .	7
3.2 Orchestrator . . . . .	7
3.3 NoSQL Databases . . . . .	7
3.4 Containers . . . . .	7
3.5 Continuous Integration / Continuous Delivery . . . . .	7
3.6 Template engines . . . . .	8
3.7 Web Server Gateway Interface . . . . .	8
3.8 Tables . . . . .	8
<b>4. Techniques and tools</b>	<b>9</b>
4.1 GitHub . . . . .	9
4.2 APIs . . . . .	9
4.3 Postman . . . . .	9
4.4 Apache Airflow . . . . .	9

4.5	Apache Spark . . . . .	9
4.6	Cassandra . . . . .	10
4.7	Flask . . . . .	10
4.8	Bootstrap . . . . .	10
4.9	Docker . . . . .	10
<b>5.</b>	<b>Relevant aspects of the project</b>	<b>11</b>
<b>6.</b>	<b>Related works</b>	<b>13</b>
<b>7.</b>	<b>Conclusions and future work lines</b>	<b>15</b>
	<b>Appendix</b>	<b>16</b>
	<b>Appendix A Project Plan</b>	<b>19</b>
A.1	Introduction . . . . .	19
A.2	Temporary planning . . . . .	20
A.3	Feasibility study . . . . .	21
	<b>Appendix B Requirements</b>	<b>23</b>
B.1	Introduction . . . . .	23
B.2	General objectives . . . . .	23
B.3	Catalog of requirements . . . . .	23
B.4	Requirements specification . . . . .	23
	<b>Appendix C Design specification</b>	<b>25</b>
C.1	Introduction . . . . .	25
C.2	Data design . . . . .	25
C.3	Procedural design . . . . .	25
C.4	Architectural design . . . . .	25
	<b>Appendix D Programming technical documentation</b>	<b>27</b>
D.1	Introduction . . . . .	27
D.2	Directory structure . . . . .	27
D.3	Programmer's guide . . . . .	27
D.4	Compilation, installation and execution of the project . . . . .	27
D.5	System tests . . . . .	27
	<b>Appendix E User documentation</b>	<b>29</b>
E.1	Introduction . . . . .	29
E.2	User requirements . . . . .	29

E.3 Installation . . . . . 29

E.4 User’s manual . . . . . 29

---

## List of Figures

---

---

# List of Tables

---

3.1 Tools and technologies used . . . . .	8
---	---



# Memory





---

# Introduction

---

Description of the work, the structure of the memory and the rest of the material delivered.



---

## Project objectives

---

This section explains precisely and concisely what are the objectives pursued with the completion of the project. It is possible to distinguish between the objectives set by the requirements of the software to be built and the technical objectives that it poses when putting the project into practice.



---

# Theoretical concepts

---

In this section are covered the theoretical concepts in which the project has been based. All concepts are described in a detailed an simple way since this master's degree can be aimed at technical and non-technical students.

## 3.1 API

Section explaining API concepts.

## 3.2 Orchestrator

Section explaining Flow Orchestrator -> Airflow.

## 3.3 NoSQL Databases

Section explaining NoSQL Databases.

## 3.4 Containers

Section explaining Containers.

## 3.5 Continuous Integration / Continuous Delivery

Section explaining CI/CD.

Tools	App	AngularJS	API REST	BD	Memoria
HTML5		X			
CSS3		X			
BOOTSTRAP		X			
T <sub>E</sub> XMaker					X
Astah					X

Table 3.1: Tools and technologies used

### 3.6 Template engines

Section explaining Template engines -> Jinja.

### 3.7 Web Server Gateway Interface

Section explaining Web Server Gateway Interface (WSGI).

### 3.8 Tables

TablaSmall.

---

## Techniques and tools

---

In this section are presented the methodological techniques and development tools used to carry out the project.

### 4.1 GitHub

GitHub is the repository where the project was uploaded and its evolution was tracked.

### 4.2 APIs

During this project there were used APIs from two different providers to gather the information: Twitter API and Spotify API.

### 4.3 Postman

Postman is a tool that allows the user to send HTTPS request in a simple way.

### 4.4 Apache Airflow

Apache Airflow is a flow orchestrator that allows the user to...

### 4.5 Apache Spark

Apache Spark is...

## **4.6 Cassandra**

Cassandra is a NoSQL database that...

## **4.7 Flask**

Flask is...

## **4.8 Bootstrap**

Bootstrap is...

## **4.9 Docker**

Docker is...



---

## Relevant aspects of the project

---

The first step of the project was the feasibility and viability analysis of the concept devised. The author was looking to use two data sources with:

- Real and updated data, preferable related to the social interest.
- The possibility of getting a stream data flow.
- The potential to combine both to get an added value.

Considering the previous points, the author found an interesting option on Twitter and Spotify providers. Both of them provides solid APIs for a fluid development and have the characteristics needed to combine the data collected. Consequently, the author designed the following use case:

1. The Twitter API is consulted to gather the *tweets* with the hashtag *#NowPlaying*.
2. The tweet is cleaned, removing the stopwords and the other hashtags and getting the song name and artist as isolated as possible.
3. The Spotify API is consulted to gather the information of the song identified.
4. The vector values of the cleaned Twitter data and the name of the song returned by Spotify are compared to ensure they are the same.
5. The data is moved to the database, ready to be stored and visualized.

The design phase...

The project development was undertaken following an Agile methodology.

---

## Related works

---

This section would be similar to a state of the art of a thesis or dissertation. In a final master's thesis, its presence does not seem so obligatory, although it can be left to the tutor's judgment to include a small commented summary of the works and projects already carried out in the field of the current project.



---

## Conclusions and future work lines

---

Every project must include the conclusions derived from its development. These can be of a different nature, depending on the type of project, but normally there will be a set of conclusions related to the results of the project and a set of technical conclusions. In addition, it is very useful to make a critical report indicating how the project can be improved, or how work can continue along the lines of the completed project.



# Appendix





## Appendix A

---

# Project Plan

---

### A.1 Introduction

The project planning was decided in an initial meeting between the author and its tutor. It was based in an Agile methodology, with two-weeks *sprints* and meetings between the author and his tutor conditioned to their availability.

The project repository was stored in GitHub under the url <https://github.com/AdrianRiesco/Data-Engineer-project>. Each *sprint* was created as an *milestone*, with the *issues* contained there being the tasks assigned. The *issues* were created to reflect tasks at most eight hours, allowing the author segregate his work and manage each *sprint* better. The author closed an *issue* when the task was finished and a *milestone* when the *sprint* was over, regardless of its state. If a task remained in an open state when a *sprint* reached its planned end date, the *issue* was transfered to the next *milestone*.

A meeting was held by the author and his tutor at the end of each sprint. During these meetings, both of them reviewed the state and development of the tasks of the corresponding sprint and planned the tasks of the next sprint. All the *milestones* and *issues* can be consulted in the project repository.

## A.2 Temporary planning

The sprints carried out for the development of the project are described below with they correspondant dates:

**Initial meeting.** Held on Monday January 31st, it was the start point for the first sprint. During this meeting, the objective of the project, the data source and the tools to be used were validated by both the author and his tutor. The author previously made a research and came with an idea and the tutor exposed his point of view to create the final goal.

**Sprint 1** . Weeks of January 31st and February 7th. This Sprint had the following tasks assigned:

- Configure the work environment.
- Configure the project memory template.
- Write a draft of the objectives and main goals.
- Write a brief description of the tools selected.
- Write a brief explanation of the selected tools and the work methodology.
- Inspect Twitter and Spotify APIs.

The end-of-sprint meeting was held on M— February –th.

**Sprint 2** . Weeks of February 14th and February 21st. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— February –th.

**Sprint 3** . Weeks of February 28th and March 7th. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— March –th.

**Sprint 4** . Weeks of March 14th and March 21st. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— March –th.

**Sprint 5** . Weeks of March 28th and April 4th. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— April –th.

**Sprint 6** . Weeks of April 11th and April 18th. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— April –th.

**Sprint 7** . Weeks of April 25th and May 2nd. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— May –th.

**Sprint 8** . Weeks of May 9th and May 16th. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— May –th.

**Sprint 9** . Weeks of May 23rd and May 30th. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— June –th.

**Sprint 10** . Weeks of June 6th and May 13th. This Sprint had the following tasks assigned:

- Task1.

The end-of-sprint meeting was held on M— June –th.

## A.3 Feasibility study

The architecture of the project and the use case were designed to ensure its feasibility.

**Economic feasibility**

The project is based on open-source platforms to ensure its economic and legal feasibility. The APIs where the information was gathered are free to use if the developer keeps his queries under specific limit rates.

**Legal feasibility**

The project is based on open-source platforms to ensure its economic and legal feasibility.

## *Appendix B*

---

# **Requirements**

---

- B.1 Introduction**
- B.2 General objectives**
- B.3 Catalog of requirements**
- B.4 Requirements specification**



## *Appendix C*

---

# **Design specification**

---

**C.1** Introduction

**C.2** Data design

**C.3** Procedural design

**C.4** Architectural design





## *Appendix D*

---

# **Programming technical documentation**

---

- D.1 Introduction
- D.2 Directory structure
- D.3 Programmer's guide
- D.4 Compilation, installation and  
execution of the project
- D.5 System tests



## *Appendix E*

---

# **User documentation**

---

**E.1 Introduction**

**E.2 User requirements**

**E.3 Installation**

**E.4 User's manual**