# BIG DATA PROJECT REPORT

## IMPLEMENTATION OF A DATALAKE ARCHITECTURE FOR THE CONSOLIDATION AND ANALYSIS OF PUBLIC DATA

*Prepared by : Aaraba Anass and Omar Farouk Lafdil*

*supervised by : Mr. Benelallam Imade*

# ABSTRACT :

In this project, we implemented a Big Data architecture in the form of a data warehouse lake using a selection of technologies, including Hadoop, Hbase, Spark, web scraping, and AI models. The goal of this project was to consolidate and analyze public data from the website "marchespublics.gov.ma". To achieve this, we used web scraping techniques to extract various document formats, including .pdf, .png, .jpg, .jpeg, .doc, .docx, .ppt, and .zip files. We then applied data cleansing and transformation techniques to convert these documents into .txt files and used AI models and natural language processing tools to extract relevant data. The resulting data was stored in .csv files and loaded into the Hadoop ecosystem using HDFS commands. We also configured Anaconda and a virtual environment to run a Jupyter notebook, and used PySpark to conduct further analysis. This project demonstrates the effectiveness of using a diverse set of tools and techniques to build a robust Big Data architecture for the consolidation and analysis of public data.

# GENERAL PLAN :

1 / INTRODUCTION

2 / THE DATA

3 / BIG DATA ANALYTICS ARCHITECTURE

4 / PROJECT DEMONSTRATION

5 / CONCLUSION

# INTRODUCTION :

The goal of this project was to build a Big Data architecture in the form of a data warehouse lake using a selection of technologies, including Hadoop, Hbase, Spark, web scraping, and AI models. The objective was to consolidate and analyze public data from the website "marchespublics.gov.ma" in order to extract all types of data, such as the subject, the project owner, the contract winner, and some relevant data. To achieve this, we applied web scraping techniques to extract various document formats from the website, including .pdf, .png, .word, and .zip files. We then used data cleansing and transformation techniques to convert these documents into .txt files and applied AI models and natural language processing tools to extract relevant data. The resulting data was stored in .csv files and loaded into the Hadoop ecosystem using HDFS commands. We also configured Anaconda and a virtual environment to run a Jupyter notebook and used PySpark to conduct further analysis. In this report, we will describe the data, methods, and results of this project, as well as the implications and recommendations for future work.

# THE DATA:

## 1 / DATA SOURCES

The data for this analysis was collected from the website "marchespublics.gov.ma", which is a public portal for procurement and tendering information in Morocco. The website contains a range of documents and information related to various procurement processes, including .pdf, .png, .word, and .zip files. We used web scraping techniques to extract these documents from the website and used data transformation techniques to convert them into .txt files for further analysis. Unfortunately, we had a capacity issue, so we limited the period to 3 months, from January to March 2019.

# THE DATA:

## 2 / DATA EXTRACTION

To extract the data from "marchespublics.gov.ma", we used web scraping techniques using the BeautifulSoup library and web scraping using Selenium webdriver. We wrote a script to scrape the website and extract the relevant documents and information. This process involved navigating to the "https://www.marchespublics.gov.ma/index.php?page=entreprise.EntrepriseAdvancedSearch&AvisExtraitPV" link, configuring the download directory, and removing duplicate files. We encountered some challenges during the data extraction process, such as dealing with CAPTCHAs. To address these challenges, we used the undetected_webdriver to ensure that CAPTCHAs don't recognize our bot. Overall, the data extraction process required more than 7 hours and a storage capacity of 3GB.

| Nom | Modifié le | Type | Taille |
|---|---|---|---|
| _AO 2-2020 Extrait-PV &amp; Réultat Dé... | 14/12/2022 01:50 | Adobe Acrobat D... | 684 Ko |
| ~$trait du procés verbal de l'A.O.O n°47.... | 15/12/2022 02:04 | Document Micros... | 1 Ko |
| 001 | 13/12/2022 23:36 | Fichier JPG | 529 Ko |
| 01-2021 | 13/12/2022 23:48 | Adobe Acrobat D... | 857 Ko |
| 01-2021-DA | 13/12/2022 23:32 | Adobe Acrobat D... | 269 Ko |
| 01-e-2021 extrait PV | 13/12/2022 23:34 | Document Micros... | 79 Ko |
| 1 | 15/12/2022 04:17 | Fichier JPG | 567 Ko |
| 1tapis | 13/12/2022 16:53 | Adobe Acrobat D... | 587 Ko |
| 02 2021 | 13/12/2022 23:58 | Fichier JPG | 5 024 Ko |
| 02-2020 | 14/12/2022 01:52 | Adobe Acrobat D... | 328 Ko |
| 02-c-2021 extrait PV | 13/12/2022 23:33 | Document Micros... | 81 Ko |
| 2 | 15/12/2022 04:17 | Fichier JPG | 631 Ko |
| 2 | 14/12/2022 01:59 | WinRAR ZIP archive | 1 565 Ko |
| 003 | 13/12/2022 17:52 | Adobe Acrobat D... | 188 Ko |
| 03 | 14/12/2022 01:43 | Fichier JPG | 354 Ko |

# THE DATA:

## 3 / DATA PREPROCESSING AND CLEANSING

This part includes multiple data processing steps in order to convert every file to a .txt file :

### 1- UNZIP THE .ZIP FILES

- The first step was to unzip the .zip files in order to maximize our data lake. In this step, we used the zipfile and glob libraries

### 2- DEALING WITH .PDF FILES

- The next step was dealing with .pdf files by converting every .pdf file to an image using the pdf2image library.
- After that, we stored every converted image to a directory, then they were deleted.

### 3- DEALING WITH IMAGES

- The next step was converting every image stored in the images directory and other .png, .jpg, .jpeg... files to .txt files.
- This process is done by using a cvttext.py script that use PIL and pytesseract libraries.
- This process was done using a cvttext.py script that used the PIL and pytesseract libraries.

### 4- DEALING WITH WORD FILES

- The next step was dealing with every .doc or .docx file by converting them to image files using the aspose.words API, which handles every version of Word and converts it to a .png image.
- Then we dealt with every image to get a .txt file for data mining purposes.

# THE DATA:

## 4 / DATA MINING

This part includes multiple data mining steps in order to insert all valuable data into a .csv file.

### 1- EXTRACTING VALUABLE DATA

- The first step was to extract the data we needed for our analysis, including the subject, the project owner, the winner of the contract, the price of the contract, the competitors, and their offers.
- The output was stored in a .txt file in a directory for text mining.
- For this step, we used the OpenAI API to extract all useful data.

### 2- DATA EXTRACTION

- This step was about extracting the data from the output of every OpenAI API request and finally generating a data frame containing the specific data we needed
- In this step, we used various text mining libraries.

### 3- CREATION OF A .CSV FILE

- The next step was taking the data frame and converting it to a .csv file.
- For some financial reasons, we were not able to pursue requests for every file at the same time. There is a rate limit in the OpenAI API on the number of requests per minute.

### 4- MERGING DATA IN ONE .CSV FILE

- After completing the process of extraction for every file in the data lake, we obtained multiple .csv files
- Finally, we combined the .csv files into one .csv file.

# THE DATA:

## 5 / FINAL DATASET CHARACTERISTICS

Our final data set is a .csv file containing the valuable data we sought to use in our analysis. Our data set includes:

### 1- OBJET AND MAITRE_OUVRAGE COLUMNS

- The first column is about the subject of every minutes of the open call for tenders document.

- The second column is about the project owner

### 2- GAGNANT AND OFFRE_GAGNANT

- These two columns are about the winner of the project and it's offer in moroccan dirham currency.

### 3- CONCURRENTS COLUMN

- This column is about the competitors of the winner.
- This column contains a list of the competitors.

### 4- OFFRES_CONCURRENTS COLUMN

- Finally, this column contains the competitors' offers in Moroccan dirham currency.
- This column also contains a list of the offers.

# THE DATA:

## 6 / DATA ANALYSIS USING APACHE HADOOP

In order to achieve our analysis goals, we followed multiple steps

### 1- DATA TRANSFER TO HADOOP SYSTEM

- The first step was transferring our data set into Cloudera Hadoop using a Vmware virtual machine.

### 2- HDFS STORING

- The next step was connecting to the Hadoop system using the 'hdfs' command.
- Then, we imported the csv file into the Hadoop system using the '-put' command.
- Finally, we checked that the file had been imported into the Hadoop file system.

### 3- ANACONDA3 INSTALLATION

- Due to some issues with the Spark default installation, we were not able to interact with the pyspark shell.
- So, we installed the Anaconda3 version that is compatible with our Unix system.
- After installing Anaconda3 and setting up the environment, we were able to start a Jupyter-Notebook.

### 4- PYSPARK INSTALLATION

- After starting the Jupyter notebook, we installed Pyspark using the 'pip' command.
- After installing Pyspark, we were able to begin the data analysis.

# BIG DATA ANALYTICS ARCHITECTURE

## 1 / LOGICAL ARCHITECTURE

In our big data analysis, we implemented the following logical architecture:

# BIG DATA ANALYTICS ARCHITECTURE

## 2 / TECHNICAL ARCHITECTURE

In our big data analysis, we implemented the following technical architecture:

# PROJECT DEMONSTRATION

## 1 / DATA SCRAPING

The first task of our big data analytics was scraping data from the public website:

This Pycharm project contains all the work from data scraping to moving the data to the Hadoop ecosystem.
The file in charge of scraping the data is "data_scraping.ipynb".

The data lake resulting from the data scraping process includes 1.49 GB of data in multiple formats (.pdf, .png, .jpg, etc.).

# PROJECT DEMONSTRATION

## 2 / DATA PREPROCESSING AND CLEANSING

The next task of our big data analytics was processing and cleansing the data scraped in the first task.

This Pycharm project contains all the work from data scraping to moving the data to the Hadoop ecosystem. Multiple python scripts are responsible for handling this task, such as cvttext.py, script2.py, imgha.dler.py, and data_processing.py which is in charge of orchestrating this process.



The result of this process is a data lake that includes data in .txt format.

# PROJECT DEMONSTRATION

## 2 / DATA MINING

The next task of our big data analytics was about the data mining in order to extract the valuable data we need:

This Pycharm project contain all the works from the data scraping to moving the data to the Hadoop ecosystem: Multiple python scripts are in charge of handling this task such as datamining.py, data_extraction.py and data_merging.py



The result of this process is our data warehouse in a .csv format. this .csv file contains all the extracted data we need for our analytics.

# PROJECT DEMONSTRATION

## 2 / DATA ANALYSIS USING APACHE HADOOP

The next task of our big data analytics was about the analytics using Pyspark in Apache Hadoop:

First we prepare the Anaconda3 environement to performe the analytics with Pyspark using a Jupyter-notebook:

# PROJECT DEMONSTRATION

## 2 / DATA ANALYSIS USING APACHE HADOOP

Next we copy our .csv file using WinSCP to transfer the CSV file from our Windows machine to the Cloudera Hadoop VM. Then we store our data into the Hadoop Distributed File System :

# PROJECT DEMONSTRATION

## 2 / DATA ANALYSIS USING APACHE HADOOP



Finally we have proceeded to the data analytics using Pyspark in a Jupyter-notebook using the data we have stored in the hdfs :

# CONCLUSION

In this report, we analyzed a large dataset of procurement information from "marchespublics.gov.ma" using a range of big data tools and techniques, including web scraping, data transformation, AI modeling, and data analysis.

These findings have several practical and theoretical implications for understanding procurement processes in Morocco.

While our analysis provides valuable insights into the procurement processes on "marchespublics.gov.ma", it is important to note that there were several limitations and challenges that we faced such as the storage limite.

Overall, this report provides a detailed analysis of procurement processes on "marchespublics.gov.ma" using big data tools and techniques, and offers valuable insights into the procurement landscape in Morocco. Our findings contribute to the field of big data analysis and have practical and theoretical implications for understanding procurement processes in the region.