# Business Data Extraction Using a Programming Language

**Fabio Bragato Do Carmo, Vinícius Medeiros Magnani, Rafael Confetti Gatsios, Fabiano Guasti Lima**

School of Economics, Business Administration and Accounting at Ribeirão Preto, University of São Paulo, São Paulo, Brazil
Email: fabio_bragato@hotmail.com, viniciusmagnani@usp.br, rafaelgatsios@fearp.usp.br, fgl@usp.br

## Abstract

In the era of great informational quantity, the presence of technologies that assist in the extraction, transformation, and loading of data has become increasingly necessary. The term Big Data, usually used to describe this volume of information, requires the user to have knowledge of multiple tools such as Excel, VBA, SQL, Tableau, Python, Spark, AWS, and so on. In this context, the present work aims to study data extraction techniques using different methodologies. At the end of the work, a library of functions in the Python language will be made available that will deliver a compilation of stock price information available on the Yahoo Finance website as well as balance sheets from financial institutions released by Bacen. The main resource used will be Web Scraping, which is a method that aims to automate data collection via the web. Once the collection of functions has been structured, it will be made available for public enjoyment through the GitHub platform.

## Keywords

Data Extraction, Big Data, Python, Web Scraping

## 1. Introduction

In recent years, the impact caused by the advancement of technology in the financial market has become an essential phenomenon for us to understand the new behaviors and socioeconomic dynamics that accompany the business world. In this context, companies process a considerable volume of information that requires new data management tools to capture, store, and analyze (Laudon & Laudon, 2014).

The large set of stored data has been entitled "Big Data" by academics and

market professionals. According to Chen, Mao, and Liu (2014), Big Data refers to a set of data that cannot be read, managed, and processed by traditional IT tools and software/hardware within a tolerable time.

In consensus in the literature, Big Data is defined by five "Vs": velocity, volume, variety, veracity, and value (Naeem et al., 2022; Ferrara et al., 2014; Laender et al., 2002). Speed because companies need information to flow quickly, as close to real time as possible. Volume because of the large amount of data generated that needs to be stored. Variety due to the gathering of information from different sources. Veracity is related to the quality aspect of the information since it needs to be reliable. And finally, value, since it supports the decision making process of businesses, leading to the creation of new oportunities for companies (Chen, Mao, & Liu, 2014; Laudon & Laudon, 2014).

As a result of this technological expansion and the amount of interpreted information, then, programming languages, database managers, and cloud computing services have emerged as a foundation for the processing of incisive elements for decision making. One of the most important computational languages in the study of data is Python, created in 1991 by mathematician Guido Van Rossum. This language is object-oriented, interpreted, and written with the intention of easy learning and has several convenient modules to increase the productivity of the development of applications to users (Buriol, Marco, & Argenta, 2009).

This programming language includes several high-level structures such as lists, dictionaries, date/time, complexes, and others, in addition to a vast collection of ready-to-use modules and third-party frameworks that can be added. These circumstances provide sufficient resources to include it as one of the main components for analysis and decision making in data generation (Borges, 2014).

The Database Management System (DBMS), in turn, is responsible for the interface of the manipulation and organization of the collected information. Among such systems, a few common ones are MySQL, SQL Server, and Oracle. To use them, it is necessary to have knowledge of Structured Query Language (SQL). According to Miyagusuku (2008), the SQL language involves commands not only of data query but also of manipulation and definition of rules and operators that preserve the integrity and consistency of data in addition to allowing the implementation of procedures (stores procedures), functions (functions) and triggers (triggers), and applications developed in other languages. DBMSs are excellent informational processors and structures, but when it comes to massive amounts of data, activities need to be migrated to cloud computing; AWS (Amazon Web Services) is one example of a company that includes this service.

The term *Big Data*, commonly used to describe a vast mass of processed information, has also supported the emergence of new programming languages. Spark, for example, was built with a focus on the dimensionality of use, allowing the user, in the same application, to integrate different languages such as SQL. Because of these factors, today, there has been a growth in the number of profes-

sionals using Python, DBMSs, and cloud computing services to analyze, for example, the earnings of companies listed on the São Paulo Stock Exchange (B3).

Investors use many strategies to study the performance of a company, but it is essential that the first step before starting assessments is to build a database with all the information of the organization. Thus, the main objective of the present work is to structure the use of these new technologies for obtaining bulk data. Three distinct methods will be used to demonstrate the daily routine of a data scientist in the job market, the first being the search for prices of stock market negotiations and data analysis, the second the search for ".zip" format files extracted directly from the HyperText Markup Language (HTML) code of the site server page, and the third, the creation of "Comma-Separated Values" (CSV) format files with the objective of manipulating and organizing the data for further analysis. The codes are applied to information from DBMSs in Brazil; however the method and coding can be adapted to other DBMSs in other countries. Therefore, the novelty of this work is to highlight and disseminate routines for collecting and building a database applied to internet data in html, data in .zip format and in csv format.

## 2. Research Methodology

The work has a descriptive approach, with the main objective of providing evidence for the methodology used by a data scientist, from the creation of an account in a cloud platform for data storage to the construction and application of the code for data interpretation.

The objects of study are the historical prices of financial institutions traded on B3, as well as their income statements, during the periods of 2017, 2018, and 2019. The price information will be extracted directly from the *Application Programming Interface* (API) provided by the *Yahoo Finance* website, and the income statements will be collected directly from the Brazilian Central Bank website. To structure the information source, Amazon S3 (*Simple Storage Service*) will be used, which will store the files of the data that will be architected in CSV format. The programming language used will be *Python* in the *Jupyter* development environment, since there are function libraries that facilitate the extraction of the information to be collected.

Among the libraries available in the *Jupyter* development environment, the following will be used:

- *Pandas_datareader*: This includes functions that allow the user to extract historical stock price information from *Yahoo Finance*, one of the largest financial news sites in the world. A plausible alternative to obtaining the monetary reports of traded stocks is the *Quandl* library, a leading source of financial, economic, and alternative data sets for investment professionals.
- *NumPy*: This is mainly used to perform calculations on multidimensional arrangements. In addition, it provides a large set of functions and operations that help programmers easily perform numerical calculations.

- *Pandas*: This is the main tool for organizing data. It provides the simplicity of being able to clean, delete, or filter information according to a criterion as well as store tables and graphs in different file formats.
- *Requests*: This allows HTTP (*Hypertext Transfer Protocol*) *requests* to be made to the specified server of the webpage. It is the bridge for performing the extraction of data.
- *Zipfile*: This provides tools to create, read, write, add, and list files in the ".zip" format.
- *Io*: This is commonly used to speed up queries made via the web.

In the development of the methodology, all the cited collections will be used with the intention of contemplating the largest amount of possible resources allowed by Python in addition to having enough subisidies for the reader to learn and apply the obtained knowledge on his/her own having as focus the supply of enough subsidies to the reader to learn and apply by himself the obtained knowledge.

## 3. Development of Working Model

### 3.1. Creating an AWS Account

The AWS platform is the world's most widely adopted and comprehensive cloud platform, offering more than 175 comprehensive data center services worldwide. There are different uses for AWS. One of its main services is EMR, which facilitates the execution of Big Data frameworks such as Hadoop and Apache Spark. In the third data extraction method, only the function of S3 for storing CSV files in the cloud will be used. Therefore, to execute the codes, it is necessary to create an account on AWS.

Figure 1 shows the AWS platform account creation table.



**Figure 1.** AWS account creation screen. Source: Own authorship.

## 3.2. Jupyter and Python Installation

Jupyter Notebook is an Integrated Development Environment (IDE). Briefly, it is the development environment for the applications made by a programmer. It is also one of the most used tools in the elaboration of Python files because it allows users to unite code and text in a single document. The installer provided by the Anaconda project on its official website (Anaconda, 2020) should be used to install the software. Figure 2 shows the images demonstrating the steps to follow after downloading the file:
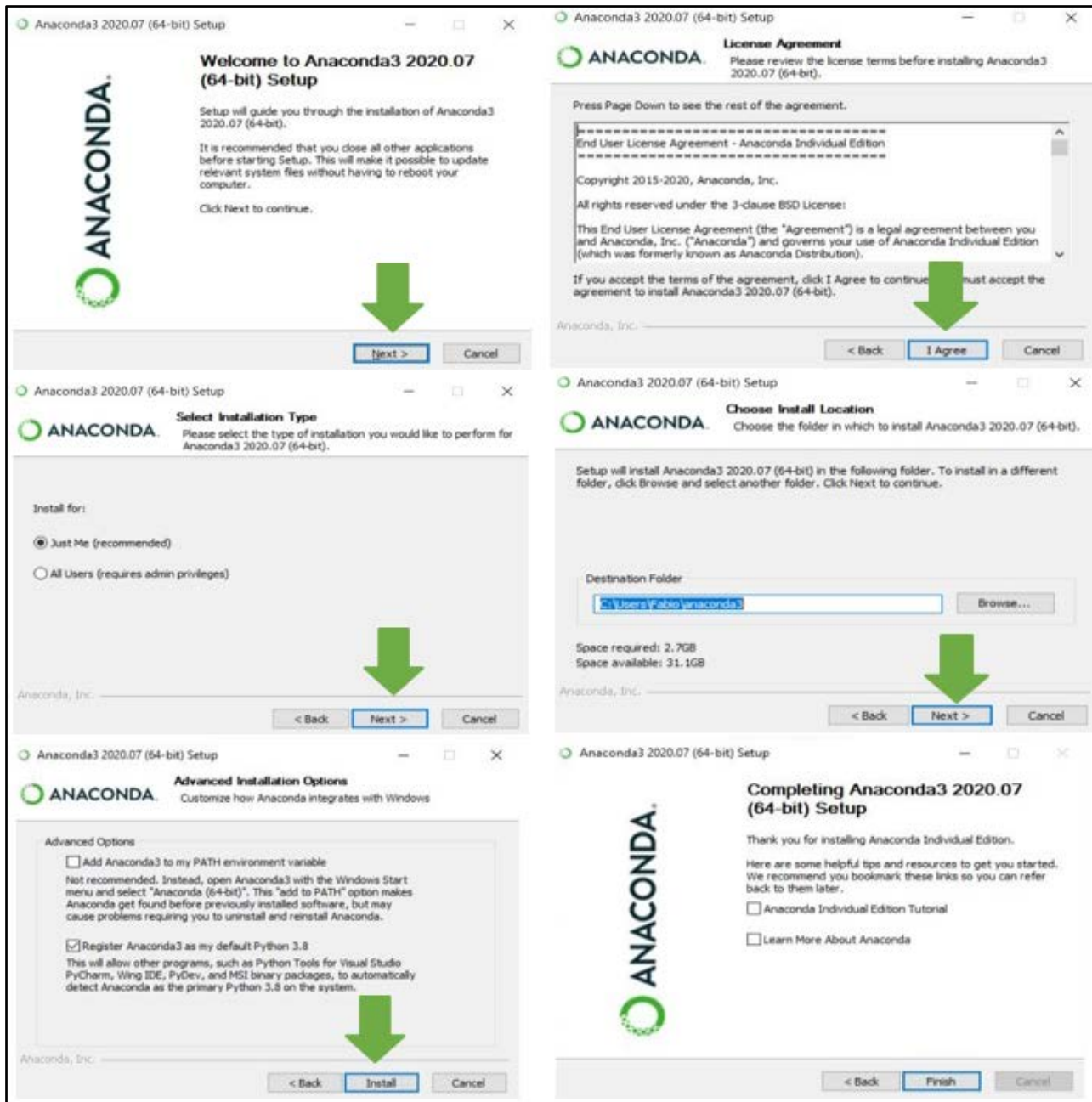


**Figure 2.** Steps to download file. Source: Own authorship.

### 3.3. Creating the First Python File

After installation, Jupyter Notebook should be started. A local server will run (by default, local output 8888), and, once loaded, a new tab in the browser will open, making it possible to see the computer Dashboard, Figure 3.

A list of files and subdirectories present on the computer is displayed. The user should create or select a folder to receive the new files. Afterward, the user should select the "New" option and choose "Python 3, Figure 4".

A new screen will open. This is where all the applications will be developed, Figure 5.

Figure 5 shows that the Jupyter Notebook is composed of the following elements:

1) File name: name of the file, which will also be displayed at the top of the screen.
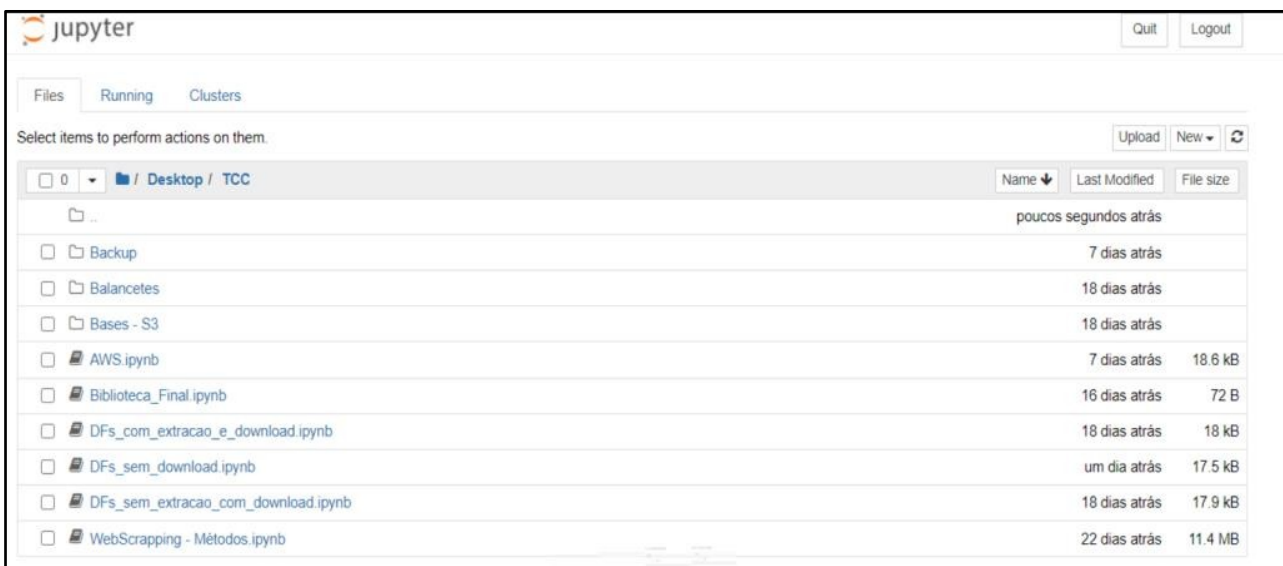


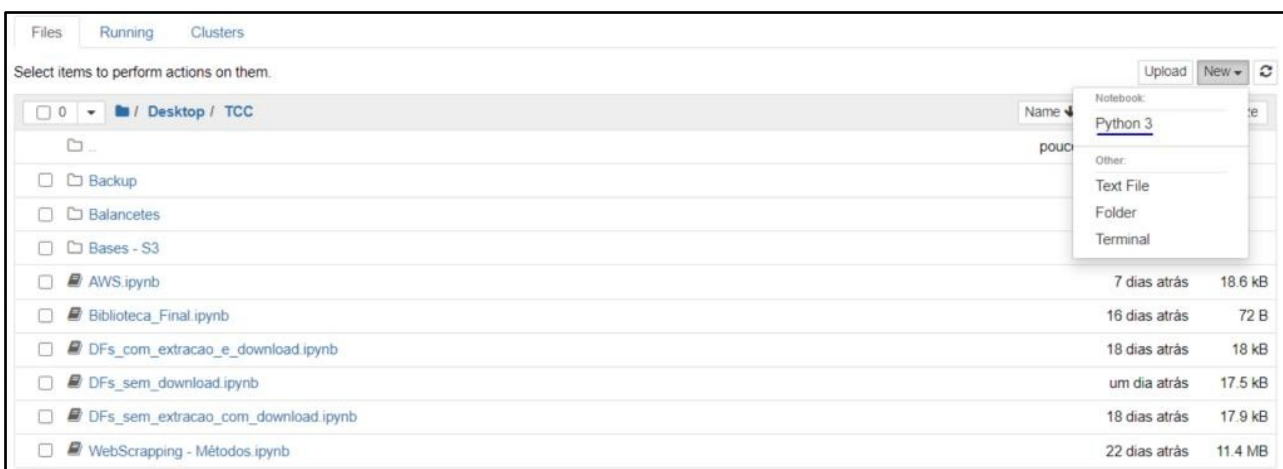Figure 3. Jupyter dashboard. Source: Own authorship.



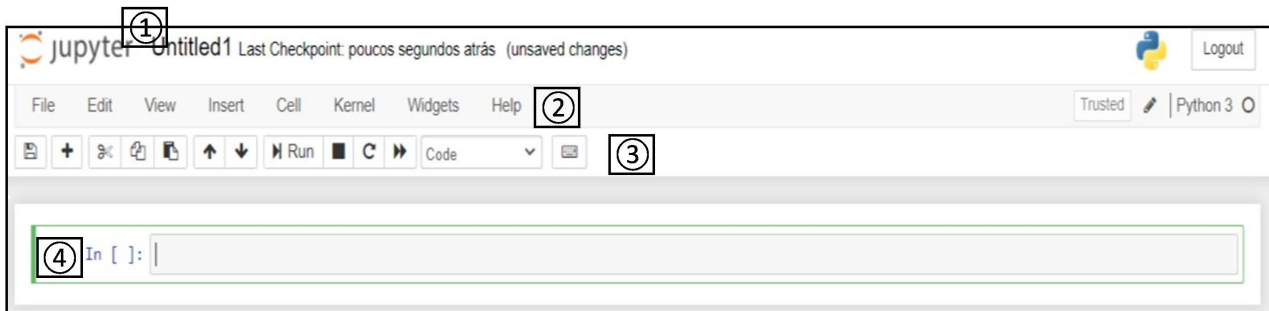Figure 4. Jupyter directories. Source: Own authorship.

**Figure 5.** Application development screen. Source: Own authorship.

2) Menu bar: different options that can be used to manipulate the operation of the notebook.

3) Toolbar: quick way to perform the most frequently used operations.

4) Code cell: allows the user you to edit or write new code.

## 3.4. Installing the Libraries

A library is a collection of modules accessible in a Python program that allows the user to simplify the programming process and remove the need to rewrite an already used function such as a calculation that extracts the square root of a number. We can call that function continuously just by modifying (or not modifying) the pre-existing parameters. Because of this, it is necessary to install open source libraries (that anyone has the right to study, modify, or distribute) that will help to create the data search methods. The execution is as follows, Figure 6.

## 3.5. Methods

### 3.5.1. Yahoo Finance

The first method used will be with the support of the library "pandas_datareader" to search the trading prices of shares on stock exchanges. The consulted data are extracted from the Yahoo Finance website and has a fifteen-minute lag time. The function created for the library called "data_actions_yahoo" was defined based on four parameters:

1st: Tickers, in which it is necessary to define the stock codes to be searched. By convention, all assets traded on B3 have the ending ".SA" (Joint Stock Company), as in "MGLU3.SA", "AZUL4.SA", "HBOR3.SA", and so on.

2nd: Start date, which will serve to reference the beginning of the consultation period. Written in a "year-month-day" pattern ("2020-01-01").

3rd: End date, which will serve to reference the end of the consultation period. Written in a "year-month-day" pattern ("2020-06-01").

4th: Type, which will reference the desired information among the possible types. "G"—All available information; "A"—Adjusted closing prices; "C"—Closing prices; "H"—Maximum prices; "L"—Minimum prices; "O"—Opening prices; "V"—Trading volume. Figure 7 shows the treatment of the function "dados_acoes_yahoo, Figure 7".

```
In [ ]:  pip install pandas_datareader

In [ ]:  pip install pandas

In [ ]:  pip install requests

In [ ]:  pip install zipfile

In [ ]:  pip install io
```

**Figure 6.** Installing the libraries. Source: Own authorship.

```
1 ▼ def dados_acoes_yahoo(tickers, data_inicial, data_final, tipo):
2
3 ▼     if tipo.upper() == 'G':
4           dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)
5           dados_finais = dados.stack()
6 ▼     elif tipo.upper() == 'A':
7           dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['Adj Close']
8           dados_finais = dados.stack()
9 ▼     elif tipo.upper() == 'C':
10          dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['Close']
11          dados_finais = dados.stack()
12 ▼    elif tipo.upper() == 'H':
13          dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['High']
14          dados_finais = dados.stack()
15 ▼    elif tipo.upper() == 'L':
16          dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['Low']
17          dados_finais = dados.stack()
18 ▼    elif tipo.upper() == 'O':
19          dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['Open']
20          dados_finais = dados.stack()
21 ▼    elif tipo.upper() == 'V':
22          dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['Volume']
23          dados_finais = dados.stack()
24      else:
25          return print('Tipo não encontrado. Tente novamente com um dos possíveis parâmetros ("G", "F", "C", "H", "L", "O" ou "V")')
26
27      return dados_finais
```

**Figure 7.** Handling of the "data_actions_yahoo" function. Source: Own authorship.

The second function, "plot_graphs_closing," aims to create a comparative chart of the stock prices evolution based on the parameters "tickers", "start_date", and "end_date". It is shown below **Figure 8**, demonstrating the function of graph plotting.

The third function, "simple_returns" serves to create a table with the simple stock returns in the specified period from the calculation and is evidenced in **Figure 9**. In Equation (1), the stock return calculation can be observed.

$$\left( \frac{\text{Share price in Tn}}{\text{Share Price in Tn} - 1} \right) - 1 \tag{1}$$

The parameters are "tickers", "start_date", and "end_date", **Figure 9**.

The last function, "retornos_log", has as objective of calculating the logarithmic returns of stocks in the specified period from the calculation:

$$\frac{\ln\left(\text{Share Price}_t\right)}{\ln\left(\text{Share Price}_{t-1}\right)} - 1 \tag{2}$$

The parameters are "tickers", "start_date", and "end_date", **Figure 10**.

Finally, many other functions can be created based on each user's need. However, the main goal of programming is to avoid repeating work and make it easier to write code.

```
1  def plotar_grafico_fechamentos(tickers, data_inicial, data_final):
2      dados = wb.DataReader(tickers, data_source='yahoo', start = data_inicial, end = data_final)['Adj Close']
3      dados_finais = dados.div(dados.iloc[0]/100.0)
4      dados_finais.plot(figsize=(16,8))
```

**Figure 8.** Handling of the "plot_graphs_closure" function. Source: Own authorship.

```
1  def retornos_simples(tickers, data_inicial, data_final):
2      returns = pd.DataFrame()
3      for ticker in tickers:
4          dados = wb.DataReader(ticker, data_source='yahoo', start = data_inicial, end = data_final)
5          returns[ticker] = (dados['Adj Close']/dados['Adj Close'].shift(1)) - 1
6
7      return returns
```

**Figure 9.** Treatment of the "simple_returns" function. Source: Own authorship.

```
1  def retornos_log(tickers, data_inicial, data_final):
2      returns = pd.DataFrame()
3      for ticker in tickers:
4          dados = wb.DataReader(ticker, data_source='yahoo', start = data_inicial, end = data_final)
5          returns[ticker] = np.log(dados['Adj Close'])/np.log(dados['Adj Close'].shift(1))
6
7      return returns
```

**Figure 10.** Treatment of the function "retornos_log." Source: Own authorship.

### 3.5.2. ZIP Files

The method of using ".zip" files is one of the most frequently utilized in data search. The purpose is to search the information of a site without installing any files. The information is extracted from the HTML code of the page itself or from the HTTP connection with the site server. The following example consolidates information from the balance sheets of financial institutions published by Bacen.

First of all, it is necessary to check the repository format, which is the ".zip" extension with CSV files. In addition, when you open the file in Excel, you can check the following structure, **Figure 11**.

As shown above, the data starts on the fourth line. This information will be important in view of the need to structure the data when writing the function. The code is as follows, **Figure 12**.

Analyzing the code line by line, we have the following steps:

- Line 1: Definition of the function that accepts two parameters, "start_harvest" and "end_harvest", both in the "year-month" format, for example, "202001".
- Line 3: Definition of the function to find the number of months between the initial and final harvests.
- Line 4: Creation of a list containing the crops to be searched.
- Line 5: Definition of a variable "i" with a value of zero that will serve as an auxiliary.
- Lines 7 to 17: A loop to add the crops chosen by the user.
- Line 19: Creation of a list that will contain the data from the CSV files to be extracted.

**Figure 11.** Excel structure. Source: Own authorship.

```python
def buscar_balancetes(safra_inicial, safra_final):

    qtd_meses = ((int(safra_final[:4]) - int(safra_inicial[:4])) * 12) + (int(safra_final[-2:]) - int(safra_inicial[-2:])) + 1
    lista_datas = []
    i = 0

    for mes in range(qtd_meses):
        if i == 0 and int(safra_inicial[-2:]) == 12:
            data = str(((int(safra_inicial[:4]) + 1) * 100 + 1))
        elif i == 0 and int(safra_inicial[-2:]) != 12:
            data = safra_inicial
        elif i > 0 and int(data[-2:]) != 12:
            data = str(int(data) + 1)
        else:
            data = str(((int(data[:4]) + 1) * 100 + 1))
        lista_datas.append(data)
        i = 1

    appended_data = []

    for dt in lista_datas:
        url = 'http://www4.bcb.gov.br/fis/cosif/cont/balan/bancos/'+ dt + 'BANCOS.ZIP'
        r = requests.get(url)
        z = zipfile.ZipFile(io.BytesIO(r.content))
        unzipped_file = z
        df = pd.read_csv(unzipped_file.open(zipfile.ZipFile.namelist(unzipped_file)[0]), sep=';', engine='python', skiprows=3)
        appended_data.append(df)

    appended_data = pd.concat(appended_data)
    return appended_data
```

**Figure 12.** Handling the "fetch_balancets" function. Source: Own authorship.

- Lines 21 to 27: The loop to capture the file information is built.
- Line 22: Definition of the URL of the site where we will collect the information from the balance sheets.
- Line 23: Use of the method "requests.get" to make the connection with the site.
- Line 24: Use of a compilation of methods. Analyzing from the inner layer outward, we have the use of "r.content", which has the function of collecting all the information contained in the connection made in line 23. It also uses the "io.BytesIo" to increase the search speed and, lastly, the "zipfile.Zipfile" to handle the ".zip" file.
- Line 25: A variable is set to capture the unzipped ".zip" file.
- Line 26: A "DataFrame" is created and the method "pd.read_csv" is used to read the CSV file contained in the decompressed archive. It is important to observe the parameters used in the function.
- "sep = ';'": CSV file with separator ";".

- "skiprows = 3": Skips three rows as indicated in the Excel figure above to start reading the information from the row that has the data pertinent to the trial balances.
- "engine = 'python'": An optional parameter that defines the parsing engine. The C engine is faster, while the Python engine is currently more resourceful.
- Line 27: Addition of the extracted data to the list of data created.
- Line 29: Concatenation of all extracted data into a single table.
- Line 30: Return of the extracted data.

### 3.5.3. CSV in AWS S3

The CSV file method in AWS S3 is an alternative way to manipulate and organize data. As an example, we will use the financial statement data provided by the Central Bank on the IFDATA website (https://www3.bcb.gov.br/ifdata/). First, the user must download the desired files. In the institution type, the user must select the option "Financial Conglomerates and Independent Institutions" in addition to choosing one of the three possible reports (assets, liabilities, or income statement), Figure 13.

Once the files have been obtained, the user logs in to the AWS website and selects the "AWS Management Console" option to log into the account, Figure 14.

Upon logging in, the user must click on services in the top menu and select "S3" in the storage group, Figure 15.

This is the screen where the user can create several directories ("Buckets") that must have unique names. To do so, press "Create Bucket" and configure it the way you prefer, Figure 16.



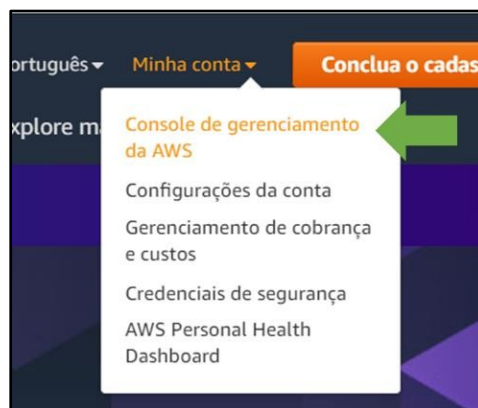**Figure 13.** IFDATA data selection screen. Source: Own authorship.



**Figure 14.** AWS login screen. Source: Own authorship.

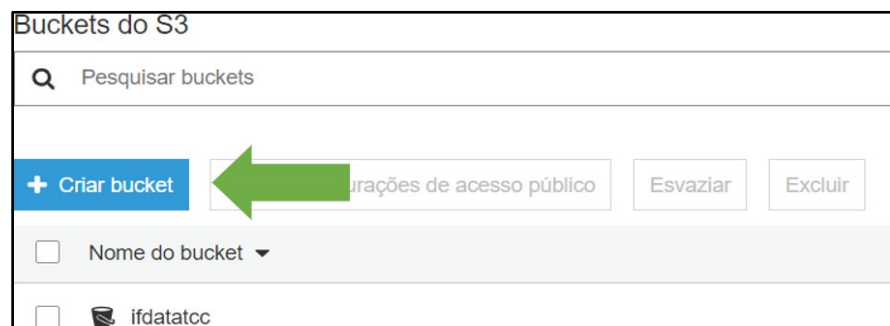**Figure 15.** AWS selection menu. Source: Own authorship.



**Figure 16.** Creating the bucket. Source: Own authorship.

After creating the Bucket, creating a folder for each report in order to maintain the organization of the files is suggested, **Figure 17**.

The user must upload the files that he/she downloaded in the created folder, and it is essential that, for the functioning of the library, the file is named with the standard "year-quarter", for example, "202001" referring to the first quarter of 2020, **Figure 18**.

Finally, the user must select all the files that were uploaded to the folder and in the "Actions" option, select "Make Public", **Figure 19**.

After finishing the S3 step, the library is ready for operation. The image below shows its construction part by part and contains more than one function for use. In order to facilitate understanding, the code will be separated into four parts:

Function "buscar_relatorio": This function's aim is to structure a report of a specific period in a grid, **Figure 20**.

**Part 1:** construction of the asset report.

- Line 1: Define the function "buscar_relatorio", which has the parameters name_bucket (ex: "ifdatatcc"), account (ex: "assets", "liabilities", or "DRE"), year (ex: 2020) and quarter (ex: 1, 2, 3, or 4).
- Line 2: The "try" method is used for error handling, as will be shown at the end of the code.
- Line 4: A conditional is created to check if one of the reports was selected.
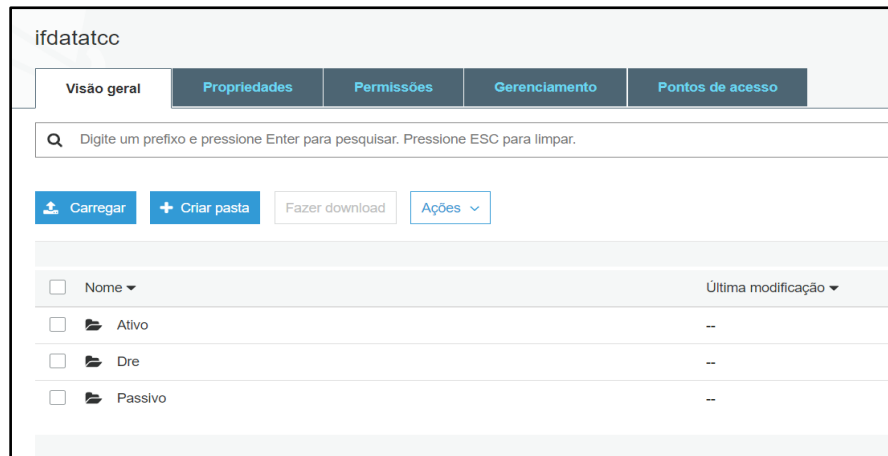- Line 7: This shows a link where the CSV file is stored in S3.

**Figure 17.** Bucket structure. Source: Own authorship.



**Figure 18.** Structured files in AWS. Source: Own authorship.
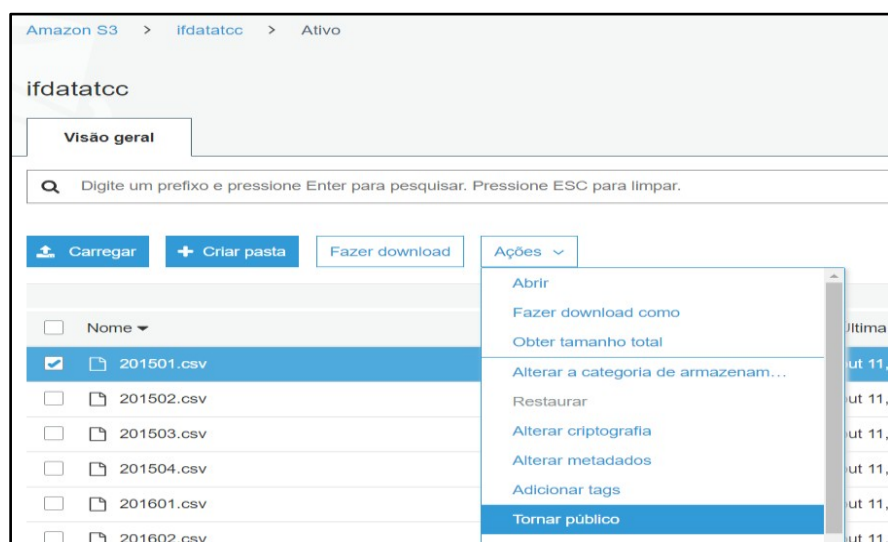


**Figure 19.** Making the archives public. Source: Own authorship.

```
1   def buscar_relatorio(nome_bucket, conta, ano, trimestre):
2       try:
3
4           if conta.upper() == 'ATIVO' or conta.upper() == 'PASSIVO' or conta.upper() == 'DRE':
5
6               # Link do CSV
7               url = 'https://{}.s3-sa-east-1.amazonaws.com/{}/{}0{}.csv'.format(nome_bucket, conta.title(), ano, trimestre)
8
9               # Função para ler o CSV
10              df = pd.read_csv(url, sep=";")
11
12
13              if conta.upper() == 'ATIVO':
14
15                  # Renomeação das colunas
16                  df = df.rename(
17                      columns={
18                          'Disponibilidades (a)': 'Disponibilidades',
19                          'Aplicações Interfinanceiras de Liquidez (b)': 'Aplicações Interfinanceiras de Liquidez',
20                          'TVM e Instrumentos Financeiros Derivativos (c)': 'TVM e Instrumentos Financeiros Derivativos',
21                          'Operações de Crédito (d1)': 'Operações de Crédito',
22                          'Unnamed: 13': 'Provisão sobre Operações de Crédito',
23                          'Unnamed: 14': 'Operações de Crédito Líquidas de Provisão',
24                          'Arrendamento Mercantil': 'Arrendamento Mercantil a Receber',
25                          'Unnamed: 16': 'Imobilizado de Arrendamento',
26                          'Unnamed: 17': 'Credores por Antecipação de Valor Residual',
27                          'Unnamed: 18': 'Provisão sobre Arrendamento Mercantil',
28                          'Provisão sobre Arrendamento Mercantil (e4)': 'Provisão sobre Arrendamento Mercantil',
29                          'Unnamed: 19': 'Arrendamento Mercantil Líquido de Provisão',
30                          'Outros Créditos - Líquido de Provisão (f)': 'Outros Créditos - Líquido de Provisão',
31                          'Outros Ativos Realizáveis (g)': 'Outros Ativos Realizáveis',
32                          'Permanente Ajustado (h)': 'Permanente Ajustado',
33                          'Ativo Total Ajustado (i) = (a) + (b) + (c) + (d) + (e) + (f) + (g) + (h)': 'Ativo Total Ajustado',
34                          'Credores por Antecipação de Valor Residual (j)': 'Credores por Antecipação de Valor Residual',
35                          'Ativo Total (k) = (i) - (j)': 'Ativo Total'
36                      })
37
38                  # Exclusão da coluna vazia
39                  df = df.drop(columns='Unnamed: 26')
40
41                  # Excluí os dados da primeira linha
42                  df = df.drop([0])
43
44                  # Encontra o índice do fim dos dados das entidades financeiras
45                  for i in range(len(df)):
46                      if df.iloc[i][0] == 'TCB - Tipo de Consolidado Bancário':
47                          marcador = i
48
49                  # Excluí os dados
50                  df = df.drop(df.index[marcador:len(df)])
51
52                  # Reseta os índices para posiciona-los de maneira correta
53                  df = df.reset_index(drop=True)
54
55                  # Retorna o dataframe como resultado da função
56                  return df
57
```

**Figure 20.** Handling the assets of the "fetch_report" function. Source: Own authorship.

- Line 13: This creates a conditional for when the "Asset" report is selected.
- Lines 16 to 53: These lines contain data processing, since it is necessary to process the information in order to use it.
- Line 56: This line returns the processed data to the function user.

**Part 2:** construction of the liability report, **Figure 21**.

The same steps are performed as in the previous method, but the data are treated according to their peculiarities. The final return is given in line 104.

**Part 3:** construction of the DRE (Statement of Income for the Year) report, **Figure 22**.

This has the same logic, but is also treated according to its differences. The final return is on line 161.

**Part 4:** final procedure of the code for handling exceptions, **Figure 23**.

- Line 163: In case the name of the selected report is not found, a phrase will be returned to the user in order to correct the value passed as parameter.

```
58          elif conta.upper() == 'PASSIVO':
59
60              # Renomeação das colunas
61              df = df.rename(
62              columns={
63                  'Captações': 'Depósitos à Vista',
64                  'Unnamed: 10': 'Depósitos de Poupança',
65                  'Unnamed: 11': 'Depósitos Interfinanceiros',
66                  'Unnamed: 12': 'Depósitos a Prazo',
67                  'Unnamed: 13': 'Outros Depósitos',
68                  'Unnamed: 14': 'Depósito Total',
69                  'Unnamed: 15': 'Obrigações por Operações Compromissadas',
70                  'Unnamed: 16': 'Letras de Crédito Imobiliário',
71                  'Unnamed: 17': 'Letras de Crédito do Agronegócio',
72                  'Unnamed: 18': 'Letras Financeiras',
73                  'Unnamed: 19': 'Obrigações por Títulos e Valores Mobiliários no Exterior',
74                  'Unnamed: 20': 'Outros Recursos de Aceites e Emissão de Títulos',
75                  'Unnamed: 21': 'Recursos de Aceites e Emissão de Títulos',
76                  'Unnamed: 22': 'Obrigações por Empréstimos e Repasses',
77                  'Unnamed: 23': 'Captações',
78                  'Instrumentos Derivativos (f)': 'Instrumentos Derivativos',
79                  'Outras Obrigações (g)': 'Outras Obrigações',
80                  'Passivo Circulante e Exigível a Longo Prazo (h) = (e) + (f) + (g)': 'Passivo Circulante e Exigível a Longo Prazo',
81                  'Resultados de Exercícios Futuros (i)': 'Resultados de Exercícios Futuros',
82                  'Patrimônio Líquido (j)': 'Patrimônio Líquido',
83                  'Passivo Total (k) = (h) + (i) + (j)': 'Passivo Total'
84              })
85
86              # Exclusão da coluna vazia
87              df = df.drop(columns='Unnamed: 30')
88
89              # Excluí os dados da primeira linha e segunda linha
90              df = df.drop([0, 1])
91
92              # Encontra o índice do fim dos dados das entidades financeiras
93              for i in range(len(df)):
94                      if df.iloc[i][0] == 'TCB - Tipo de Consolidado Bancário':
95                          marcador = i
96
97              # Excluí os dados
98              df = df.drop(df.index[marcador:len(df)])
99
100             # Reseta os índices para posiciona-los de maneira correta
101             df = df.reset_index(drop=True)
102
103             # Retorna o dataframe como resultado da função
104             return df
```

**Figure 21.** Handling the liabilities of the "fetch_report" function. Source: Own authorship.

```
106         elif conta.upper() == 'DRE':
107
108             # Renomeação das colunas
109             df = df.rename(
110             columns={
111                 'Resultado de Intermediação Financeira': 'Rendas de Operações de Crédito',
112                 'Unnamed: 10': 'Rendas de Operações de Arrendamento Mercantil',
113                 'Unnamed: 11': 'Rendas de Operações com TVM',
114                 'Unnamed: 12': 'Rendas de Operações com Instrumentos Financeiros Derivativos',
115                 'Unnamed: 13': 'Resultado de Operações de Câmbio',
116                 'Unnamed: 14': 'Rendas de Aplicações Compulsórias',
117                 'Unnamed: 15': 'Receitas de Intermediação Financeira',
118                 'Unnamed: 16': 'Despesas de Captação',
119                 'Unnamed: 17': 'Despesas de Obrigações por Empréstimos e Repasses',
120                 'Unnamed: 18': 'Despesas de Operações de Arrendamento Mercantil',
121                 'Unnamed: 19': 'Resultado de Operações de Câmbio',
122                 'Unnamed: 20': 'Resultado de Provisão para Créditos de Difícil Liquidação',
123                 'Unnamed: 21': 'Despesas de Intermediação Financeira',
124                 'Unnamed: 22': 'Resultado de Intermediação Financeira',
125                 'Outras Receitas/Despesas Operacionais': 'Rendas de Prestação de Serviços',
126                 'Unnamed: 24': 'Rendas de Tarifas Bancárias',
127                 'Unnamed: 25': 'Despesas de Pessoal',
128                 'Unnamed: 26': 'Despesas Administrativas',
129                 'Unnamed: 27': 'Despesas Tributárias',
130                 'Unnamed: 28': 'Resultado de Participações',
131                 'Unnamed: 29': 'Outras Receitas Operacionais',
132                 'Unnamed: 30': 'Outras Despesas Operacionais',
133                 'Unnamed: 31': 'Outras Receitas/Despesas  Operacionais',
134                 'Resultado Operacional (e) = (c) + (d)': 'Resultado Operacional',
135                 'Resultado Não Operacional (f)': 'Resultado Não Operacional',
136                 'Resultado antes da Tributação, Lucro e Participação (g) = (e) + (f)': 'Resultado antes da Tributação, Lucro e Participação',
137                 'Imposto de Renda e Contribuição Social (h)': 'Imposto de Renda e Contribuição Social',
138                 'Participação nos Lucros (i)': 'Participação nos Lucros',
139                 'Lucro Líquido (j) = (g) + (h) + (i)': 'Lucro Líquido',
140                 'Juros Sobre Capital Social de Cooperativas (k)': 'Juros Sobre Capital Social de Cooperativas'
141             })
142
143             # Exclusão da coluna vazia
144             df = df.drop(columns='Unnamed: 39')
145
146             # Excluí os dados da primeira linha e segunda linha
147             df = df.drop([0, 1])
148
149             # Encontra o índice do fim dos dados das entidades financeiras
150             for i in range(len(df)):
151                     if df.iloc[i][0] == 'TCB - Tipo de Consolidado Bancário':
152                         marcador = i
153
154             # Excluí os dados
155             df = df.drop(df.index[marcador:len(df)])
156
157             # Reseta os índices para posiciona-los de maneira correta
158             df = df.reset_index(drop=True)
159
160             # Retorna o dataframe como resultado da função
161             return df
```

**Figure 22.** Handling the DRE of the "fetch_report" function. Source: Own authorship.

- Lines 165 to 170: Error handling related to searching for non-existent periods and when the required library for operation is not imported.

Function "buscar_relatorios": Its aim is to consolidate in a single table more than one period specified by the function "buscar_relatorio.", Figure 24.

To build this code, the function previously built was used to help in the consolidation of dates.

- Line 1: The function "buscar_relatorios" is defined; it the parameters "name_bucket" (ex: "ifdatatcc"), account (ex: "assets", "liabilities", or "DRE"), "initial_year" (ex: 2019), "initial_quarter" (ex: 1, 2, 3, or 4), "final_year" (ex: 2020), and "final_quarter" (ex: 1, 2, 3, or 4).
- Line 2: A calculation is created to find the number of reports to be extracted.
- Lines 3 to 6: A loop is needed to find the position of "start_date" in the date list.
- Lines 8 to 14: A loop runs that adds the number of reports calculated in line 2 from the start date. At the end of the procedure, the data are merged into a single table.
- Line 16: Returns the final information to the user.

Function "periodos_disponible": It has as objective to show the user which is the available report periods. It does not have parameters, Figure 25.

```
163        else:
164            print('Nome da conta inserido incorretamente. Tente novamente passando um dos possíveis valores: "ativo", "passivo" ou "DRE".')
165    except:
166        # Tratamento de erro
167        return print('Arquivo não encontrado. Possível motivo: \n'
168                     '- Período procurado inexistente. \n'
169                     '- Biblioteca "Pandas" não importada. \n'
170                     )
```

Figure 23. Error handling of the "fetch_report" function. Source: Own authorship.

```
1   def buscar_relatorios(nome_bucket, conta, ano_inicial, trimestre_inicial, ano_final, trimestre_final):
2       nr_relatorios = (ano_final - ano_inicial) * 4 + (trimestre_final - trimestre_inicial) + 1
3       for i in lista_datas:
4           if i == (ano_inicial * 100) + trimestre_inicial:
5               posicao = lista_datas.index(i)
6               break
7
8       for j in range(nr_relatorios):
9           if j == 0:
10              dado = buscar_relatorio(nome_bucket, conta, int(str(lista_datas[posicao + j])[:4]), int(str(lista_datas[posicao + j])[-1:]))
11              dados = dado
12          else:
13              dado = buscar_relatorio(nome_bucket, conta, int(str(lista_datas[posicao + j])[:4]), int(str(lista_datas[posicao + j])[-1:]))
14              dados = dados.append(dado)
15
16      return dados
```

Figure 24. Handling the "fetch_reports" function. Source: Own authorship.

```
1   def periodos_disponiveis():
2       print('Dados disponíveis: \n' +
3             str(min(lista_datas))[-1:] +
4             'º Trimestre/' +
5             str(min(lista_datas))[:4] +
6             ' à ' +
7             str(max(lista_datas))[-1:] +
8             'º Trimestre/' +
9             str(max(lista_datas))[:4] +
10            '.'
11            )
```

Figure 25. Handling the "available_periods" function. Source: Own authorship.

## 4. Results and Findings

After the construction of the three methods, the information was consolidated into a single file from all the functions created in order to establish a vast repertoire of utilities to the users of the elaborated library. The result of each function is shown below, Figure 26.

- "data_actions_yahoo".
- "plot_graphic_closures", Figure 27.
- "returns_simple", Figure 28.
- "returns_log", Figure 29.
- "fetch_balancets", Figure 30.
- "fetch_report", Figure 31.
- "fetch_reports", Figure 32.
- "periods_available", Figure 33.



**Figure 26.** Return of the function "dados_acoes_yahoo". Source: Own authorship.



**Figure 27.** Return of the function "plot_graphic_closures". Source: Own authorship.

**Figure 28.** Return from the "simple_returns" function. Source: Own authorship.



**Figure 29.** Return of the function "retornos_log". Source: Own authorship.



**Figure 30.** Return of the function "fetch_balancets". Source: Own authorship.

```
buscar_relatorio('ifdatatcc', 'ativo', 2020, 1)
```

| | Instituição financeira | Código | TCB | SR | TD | TC | Cidade | UF | Data | Disponibilidades | ... | Imobilizado de Arrendamento | Credores por Antecipação de Valor Residual | Provisão sobre Arrendamento Mercantil | Arrendame Merca Líquido Provis |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | ITAU | 10069.0 | b1 | S1 | C | 2.0 | SAO PAULO | SP | 03/2020 | 37.685.356 | ... | 567.954 | 376.783 | -259.186 | 8.607.5 |
| 1 | BB | 49906.0 | b1 | S1 | C | 1.0 | BRASILIA | DF | 03/2020 | 15.171.888 | ... | 228.259 | 75.303 | -3.864 | 334.1 |
| 2 | CAIXA ECONOMICA FEDERAL | 360305.0 | b1 | S1 | I | 1.0 | BRASILIA | DF | 03/2020 | 11.207.090 | ... | 0 | 0 | 0 | |
| 3 | BRADESCO | 10045.0 | b1 | S1 | C | 2.0 | OSASCO | SP | 03/2020 | 22.559.650 | ... | 4.321.972 | 1.415.538 | -148.768 | 5.638.5 |
| 4 | SANTANDER | 30379.0 | b1 | S1 | C | 3.0 | SAO PAULO | SP | 03/2020 | 13.889.118 | ... | 5.017.954 | 2.230.363 | -38.192 | 7.223.5 |

**Figure 31.** Return of the function "buscar_relatorio". Source: Own authorship.

```
buscar_relatorios('ifdatatcc', 'ativo', 2020, 1, 2020, 2)
```

| | Instituição financeira | Código | TCB | SR | TD | TC | Cidade | UF | Data | Disponibilidades | ... | Imobilizado de Arrendamento | Credores por Antecipação de Valor Residual | Provisão sobre Arrendamento Mercantil | Arrendame Merca Líquido Provis |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | ITAU | 10069.0 | b1 | S1 | C | 2.0 | SAO PAULO | SP | 03/2020 | 37.685.356 | ... | 567.954 | 376.783 | -259.186 | 8.607.5 |
| 1 | BB | 49906.0 | b1 | S1 | C | 1.0 | BRASILIA | DF | 03/2020 | 15.171.888 | ... | 228.259 | 75.303 | -3.864 | 334.1 |
| 2 | CAIXA ECONOMICA FEDERAL | 360305.0 | b1 | S1 | I | 1.0 | BRASILIA | DF | 03/2020 | 11.207.090 | ... | 0 | 0 | 0 | |
| 3 | BRADESCO | 10045.0 | b1 | S1 | C | 2.0 | OSASCO | SP | 03/2020 | 22.559.650 | ... | 4.321.972 | 1.415.538 | -148.768 | 5.638.5 |
| 4 | SANTANDER | 30379.0 | b1 | S1 | C | 3.0 | SAO PAULO | SP | 03/2020 | 13.889.118 | ... | 5.017.954 | 2.230.363 | -38.192 | 7.223.5 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 1299 | COOPERATIVA DE ECONOMIA E CRÉDITO MÚTUO DOS TR... | 7494300.0 | b3S | S5 | I | 2.0 | ERECHIM | RS | 06/2020 | NI | ... | NI | NI | NI | |
| 1300 | COOPERATIVA DE CRÉDITO RURAL DO AGRESTE CENTRA... | 8202793.0 | b3S | S5 | I | 2.0 | ARAPIRACA | AL | 06/2020 | NI | ... | NI | NI | NI | |

**Figure 32.** Return of the function "buscar_relatorios". Source: Own authorship.

```
periodos_disponiveis()

Dados disponíveis:
1º Trimestre/2015 à 2º Trimestre/2020.
```

**Figure 33.** Return of the function "periods_available". Source: Own authorship.

## 5. Conclusion

In summary, it is possible to conclude that data extraction tools are currently of extreme importance, in view of the large informational volume processed and the different sources of information available to users. It is interesting to point out that the profession of data analyst has become increasingly essential for companies (Lima, 2018). There is applicability for the techniques of capturing information in several fields, as in the area of health, in which evidence can be collected for research proof (Bradley, Curry, & Devers, 2007), as in the Business and Competitive Intelligence, where a company can probe the Web to acquire

and analyze information about the activity of its competitors (Baumgartner et al., 2005; Chen et al., 2002), and as in crawling os Social Web platforms, to map consumer desires and understand, model and predict human behavior (Gkotsis et al., 2013; Catanese et al., 2011).

This work explored the *Web Scraping* technique, which is a main method in the context of the digital era. The extraction can be made from simple texts, as well as, with the use of more advanced techniques, from audio and even video material using the aid of artificial intelligence. Thus, the main contribution of the work was to structure the use of three distinct methods to structure a data collection process, the first being the search for stock trading prices on the stock exchange and graphically analyzing the data, the second the search for files in ".zip" format extracted directly from the HyperText Markup Language (HTML) code of the website's server page, and the third, the creation of files in the "Comma-Separated Values" (CSV) format in order to manipulate and organize data for further analysis. The codes are applied to information from SGBDs in Brazil; however, the method and coding can be adapted to other DBMSs in other countries. Therefore, the code of the work used is available for public enjoyment on the GitHub platform (Github, 2020). The results show that the construction of programs for data extraction can be a great ally in the efforts of information users. Once created, the time saving of such tools for the structuring of reports that bring value generation becomes obvious.

The main limitation of this study is related to the need to check whether downloading data from the site with scripts is allowed and does not violate local data protection laws. Finally, regarding future work for the evolution of the solutions developed in this project, we would like to mention the study of other techniques using resources such as the AWS EMR to create a virtual machine capable of processing large-scale files as well as the demonstration of the use of *machine learning* and *deep learning* techniques.

## Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

## References

Anaconda (2020) Individual Edition. https://www.anaconda.com/products/individual

Baumgartner, R., Fröhlich, O., Gottlob, G., Harz, P., Herzog, M., & Lehmann, P. (2005). *Web Data Extraction for Business Intelligence: The Lixto Approach*. Gesellschaft für Informatik eV.

Borges, L. E. (2014). *Python para desenvolvedores: Aborda Python 3.3* (p. 14). Novatec Editora. https://books.google.com.br/books?hl=pt-BR&lr=lang_pt&id=eZmtBAAAQBAJ&oi=fnd&pg=PA14&dq=python&ots=VDSosqIkmu&sig=TZ0MbKn058lnRJ9zgZrNLmoOFh4#v=onepage&q=python&f=false

Bradley, E. H., Curry, L. A., & Devers, K. J. (2007). Qualitative Data Analysis for Health Services Research: Developing Taxonomy, Themes, and Theory. *Health Services Research, 42,* 1758-1772. https://doi.org/10.1111/j.1475-6773.2006.00684.x

Buriol, T. M., Marco, B., & Argenta, M. A. (2009). *Acelerando o desenvolvimento eo processamento de análises numéricas computacionais utilizando python e cuda.* https://www.researchgate.net/profile/Marco_Argenta/publication/228683446_ACELER ANDO_O_DESENVOLVIMENTO_EO_PROCESSAMENTO_DE_ANALISES_NUME RICAS_COMPUTACIONAIS_UTILIZANDO_PYTHON_E_CUDA/links/5630d6c908 ae0530378cdf06.pdf

Catanese, S. A., De Meo, P., Ferrara, E., Fiumara, G., & Provetti, A. (2011, May). Crawling Facebook for Social Network Analysis Purposes. In *Proceedings of the International Conference on Web Intelligence, Mining and Semantics* (pp. 1-8). https://doi.org/10.1145/1988688.1988749

Chen, H., Chau, M., & Zeng, D. (2002). CI Spider: A Tool for Competitive Intelligence on the Web. *Decision Support Systems, 34,* 1-17. https://doi.org/10.1016/S0167-9236(02)00002-7

Chen, M., Mao, S., & Liu, Y. (2014). Big Data: A Survey. *Mobile Networks and Applications, 19,* 171-209. https://doi.org/10.1007/s11036-013-0489-0

Ferrara, E., De Meo, P., Fiumara, G., & Baumgartner, R. (2014). Web Data Extraction, Applications and Techniques: A Survey. *Knowledge-Based Systems, 70,* 301-323. https://doi.org/10.1016/j.knosys.2014.07.007

Github (2020) Fabiobragato/Finances. https://github.com/fabiobragato/finances.git

Gkotsis, G., Stepanyan, K., Cristea, A. I., & Joy, M. (2013, July). Self-Supervised Automated Wrapper Generation for Weblog Data Extraction. In *British National Conference on Databases* (pp. 292-302). Springer. https://doi.org/10.1007/978-3-642-39467-6_26

Laender, A. H., Ribeiro-Neto, B. A., Da Silva, A. S., & Teixeira, J. S. (2002). A Brief Survey of Web Data Extraction Tools. *ACM SIGMOD Record, 31,* 84-93. https://doi.org/10.1145/565117.565137

Laudon, K. C., & Laudon, J. P. (2014). *Management Information Systems: Managing the Digital Firm* (13th ed., p. 37). Prentice Hall. http://dinus.ac.id/repository/docs/ajar/Kenneth_C.Laudon,Jane_P_.Laudon_-_Manage ment_Information_Sysrem_13th_Edition_.pdf

Lima, F. G. (2018). Análise de Riscos. Ed. Atlas.

Miyagusku, R. (2008). Curso Prático de SQL. Guia de Referência Completo Para Usar a linguagem SQL nos Bancos de Dados: MS SQL Server, Oracle, PostgreSQL, MySQL.

Naeem, M., Jamal, T., Diaz-Martinez, J., Butt, S. A., Montesano, N., Tariq, M. I., De-La-Hoz-Valdiris, E. et al. (2022). Trends and Future Perspective Challenges in Big Data. In *Advances in Intelligent Data Analysis and Applications* (pp. 309-325). Springer. https://doi.org/10.1007/978-981-16-5036-9_30