

IBM Python Project for Data Engineering

Role: Data Engineer

Final Project: Acquiring and Processing Information on the World's Largest Banks

In this project, you will put all the skills acquired throughout the course and your knowledge of basic Python to test. You will work on real-world data and perform the operations of Extraction, Transformation, and Loading (ETL) as required.

Disclaimer:

Cloud IDE is not a persistent platform, and you will lose your progress every time you restart this lab. We recommend saving a copy of your file on your local machine as a protective measure against data loss.

Project Scenario:

You have been hired as a data engineer by research organization. Your boss has asked you to create a code that can be used to compile the list of the top 10 largest banks in the world ranked by market capitalization in billion USD. Further, the data needs to be transformed and stored in GBP, EUR and INR as well, in accordance with the exchange rate information that has been made available to you as a CSV file. The processed information table is to be saved locally in a CSV format and as a database table.

Your job is to create an automated system to generate this information so that the same can be executed in every financial quarter to prepare the report.

Particulars of the code to be made have been shared below.

Parameter	Value
Code name	<code>banks_project.py</code>
Data URL	<code>https://web.archive.org/web/20230908091635/https://en.wikipedia.org/wiki/List_of_largest_banks</code>
Exchange rate CSV path	<code>https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMSkillsNet-work-PY0221EN-Coursera/labs/v2/exchange_rate.csv</code>
Table Attributes (upon Extraction only)	<code>Name, MC_USD_Billion</code>
Table Attributes (final)	<code>Name, MC_USD_Billion, MC_GBP_Billion, MC_EUR_Billion, MC_INR_Billion</code>
Output CSV Path	<code>./Largest_banks_data.csv</code>
Database name	<code>Banks.db</code>

Table name	Largest_banks
Log file	code_log.txt

Project tasks

Task 1:

Write a function `log_progress()` to log the progress of the code at different stages in a file `code_log.txt`. Use the list of log points provided to create log entries as every stage of the code.

Task 2:

Extract the tabular information from the given URL under the heading 'By market capitalization' and save it to a dataframe.

- Inspect the webpage and identify the position and pattern of the tabular information in the HTML code
- Write the code for a function `extract()` to perform the required data extraction.
- Execute a function call to `extract()` to verify the output.

Task 3:

Transform the dataframe by adding columns for Market Capitalization in GBP, EUR and INR, rounded to 2 decimal places, based on the exchange rate information shared as a CSV file.

- Write the code for a function `transform()` to perform the said task.
- Execute a function call to `transform()` and verify the output.

Task 4:

Load the transformed dataframe to an output CSV file. Write a function `load_to_csv()`, execute a function call and verify the output.

Task 5:

Load the transformed dataframe to an SQL database server as a table. Write a function `load_to_db()`, execute a function call and verify the output.

Task 6:

Run queries on the database table. Write a function `load_to_db()`, execute a given set of queries and verify the output.

Task 7:

Verify that the log entries have been completed at all stages by checking the contents of the file `code_log.txt`.

Preliminaries: Installing libraries and downloading data

Before building the code, you need to install the required libraries.

The libraries needed for the code are:

`requests` - The library used for accessing the information from the URL.

`bs4` - The library containing the BeautifulSoup function used for webscraping.

`pandas` - The library used for processing the extracted data, storing it in required formats, and communicating with the databases.

`sqlite3` - The library required to create a database server connection.

`numpy` - The library required for the mathematical rounding operations.

`datetime` - The library containing the function datetime used for extracting the timestamp for logging purposes.

Install the required libraries from the terminal window. The command syntax is:

```
python3.11 -m pip install <library_name>
```

Also, download the required exchange rate file using the terminal command:

wget

https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMSkillsNetwork-PY0221EN-Coursera/labs/v2/exchange_rate.csv

Code Structure

Create the file `banks_project.py` in the path `\home\project\`

Task 1: Logging function

Write the function to log the progress of the code, `log_progress()`. The function accepts the message to be logged and enters it to a text file `code_log.txt`.

The format to be used for logging must have the syntax:

`<time_stamp> : <message>`

Each log entry must happen in the next line in the text file.

You must associate the correct log entries with each of the executed function calls. Use the following table to note the logging message at the end of each function call that follows.

Task	Log message on completion
Declaring known values	Preliminaries complete. Initiating ETL process
Call extract() function	Data extraction complete. Initiating Transformation process

Call transform() function	Data transformation complete. Initiating Loading process
Call load_to_csv()	Data saved to CSV file
Initiate SQLite3 connection	SQL Connection initiated
Call load_to_db()	Data loaded to Database as a table, Executing queries
Call run_query()	Process Complete
Close SQLite3 connection	Server Connection closed

At this stage, you should now make the first log entry from the table above.

Peer graded assignment prompt:

Take a screenshot of the code, as created for the `log_progress()` function and save it to your local machine as `Task_1_log_function.png`

Task 2 : Extraction of data

Analyze the webpage on the given URL:

https://web.archive.org/web/20230908091635/https://en.wikipedia.org/wiki/List_of_largest_banks

Identify the position of the required table under the heading **By market capitalization**. Write the function `extract()` to retrieve the information of the table to a Pandas data frame.

Note: Remember to remove the last character from the **Market Cap** column contents, like, '\n', and typecast the value to float format.

Write a function call for `extract()` and print the returning data frame.

Make the relevant log entry.

Execute the code using the command:

```
python3.11 banks_project.py
```

Quiz question prompt:

While inspecting the webpage, note the attributes of the data in the first row. There will be a quiz questions based on these attributes.

Peer graded assignment prompt:

Take a screenshot of the html code of the table, as obtained by inspecting the webpage. Make sure that the contents of at least the first row of the table, as entered in the HTML code, are completely visible. Save this screenshot to your local machine as `Task_2a_extract.png`

Take a screenshot of the code, as created for the `extract()` function and save it to your local machine as `Task_2b_extract.png`.

Take a screenshot of the output, as obtained upon execution in the terminal, and save it to your local machine as `Task_2c_extract.png`.

Task 3 : Transformation of data

The Transform function needs to perform the following tasks:

1. Read the exchange rate CSV file and convert the contents to a dictionary so that the contents of the first columns are the keys to the dictionary and the contents of the second column are the corresponding values.

You can use the below-mentioned syntax to achieve this. Remember to modify the statement as per your code.

```
dict = dataframe.set_index('Col_1_header').to_dict()['Col_2_header']
```

2. Add 3 different columns to the dataframe, viz. `MC_GBP_Billion`, `MC_EUR_Billion` and `MC_INR_Billion`, each containing the content of `MC_USD_Billion` scaled by the corresponding exchange rate factor. Remember to round the resulting data to 2 decimal places.

A sample statement is being provided for adding the `MC_GBP_Billion` column. You can use this to add the other two statements on your own.

```
df['MC_GBP_Billion'] = [np.round(x*exchange_rate['GBP'],2) for x in df['MC_USD_Billion']]
```

Write the function call for `transform()` and print the contents of the returning data frame. Comment out all previous print statements.

Make the relevant log entry and execute the code.

Quiz question prompt:

1. Experiment with the statement provided for adding the transformed columns to the dataframe. There will be a question on this in the quiz.
2. Print the contents of `df['MC_EUR_Billion'][4]`, which is the market capitalization of the 5th largest bank in billion EUR. Note this value, as it will be the answer to a question in the final quiz.

Peer graded assignment prompt:

Take a screenshot of the code, as created for the `transform()` function, and save it to your local machine as `Task_3a_transform.png`.

Take a snapshot of the output and save it as `Task_3b_tranform.png`.

Task 4: Loading to CSV

Write the function to load the transformed data frame to a CSV file, like `load_to_csv()`, in the path mentioned in the project scenario.

Make the relevant log entry.

Peer graded assignment prompt:

Double-click the created CSV file in the `Explorer` tab on the left ribbon of the programming pane in Cloud IDE. Note that its contents are displayed on the editor screen. Take a snapshot of this screen and save it as `Task_4_CSV.png`.

Task 5: Loading to Database

Write the function to load the transformed data frame to an SQL database, like, `load_to_db()`. Use the database and table names as mentioned in the project scenario.

Before calling this function, initiate the connection to the SQLite3 database server with the name `Banks.db`. Pass this connection object, along with the required table name `Largest_banks` and the transformed data frame, to the `load_to_db()` function in the function call.

Make the relevant log entry.

Upon successful function call, you will have loaded the contents of the table with the required data and the file `Banks.db` will be visible in the `Explorer` tab of the IDE under the `project` folder.

Peer graded assignment prompt:

Take a single screenshot of the code, as created for `load_to_csv()` and `load_to_db()` functions, and save it to your local machine as `Task_4_5_save_file.png`.

Task 6: Function to Run queries on Database

Write the function `run_queries()` that accepts the query statement, and the SQLite3 Connection object, and generates the output of the query. The query statement should be printed along with the query output.

Execute 3 function calls using the queries as mentioned below.

1. Print the contents of the entire table

Query statement:

```
SELECT * FROM Largest_banks
```

2. Print the average market capitalization of all the banks in Billion USD.

Query statement:

```
SELECT AVG(MC_GBP_Billion) FROM Largest_banks
```

3. Print only the names of the top 5 banks

Query statement:

```
SELECT Name from Largest_banks LIMIT 5
```

Make the relevant log entry.

Peer graded assignment prompt:

Take the snapshot of the output and save it as `Task_6_SQL.png`. Please adjust the size of the terminal prompt in order to take a single screenshot that captures all three outputs together.

Quiz question prompt:

There will be a quiz question on the output of these queries.

Task 7: Verify log entries

After updating all the `log_progress()` function calls, you have to run the code for a final execution. However, you will first have to remove the `code_log.txt` file, that would have been created and updated throughout the multiple executions of the code in this lab. You may remove the file using the following command on a terminal.

```
rm code_log.txt
```

Once the existing file is removed, now run the final execution. Upon successful completion of execution, open the `code_log.txt` file by clicking on it in the `Explorer` tab of the toolbar on left side of the programming pane of the IDE, under the `project` folder. You should see all the relevant entries made in the text file in relation to the stages of code execution.

Peer graded assignment prompt:

Take the snapshot of the file contents and save it as `Task_7_log_content.png`.

Conclusion

Congratulations on completing this project!

With this, you are now trained to perform ETL operations on real-world data and make the processed information available for further use in different formats.

You should now be able to:

- Use Webscraping techniques to extract information from any website as per requirement.
- Use Pandas data frames and dictionaries to transform data as per requirement.

- Load the processed information to CSV files and as Database tables
- Query the database tables using SQLite3 and pandas libraries
- Log the progress of the code properly