

README

for Code to Extract Data from PDFs, Scrape Web Data, and Store in Database

Overview

The provided code performs the following tasks:

1. Extracts well data and stimulation data from a list of PDF files.
2. Uses web scraping to gather additional information about each well based on its well name.
3. Preprocesses the extracted and scraped data.
4. Stores the final data in a MySQL database.

Prerequisites

Before you can run the code, you need to have the following software and Python libraries installed:

- Python (3.6 or later)
- ChromeDriver (compatible with your Chrome version)
- PyPDF2
- pdf2image
- pytesseract
- selenium
- webdriver_manager
- pandas
- numpy
- BeautifulSoup
- glob, os, re, time
- mysql.connector

You can install the necessary libraries using pip:

```
1 | pip install PyPDF2 pdf2image pytesseract selenium webdriver_manager pandas numpy  
   | beautifulsoup4 mysql-connector-python
```

Also, you need to have MySQL set up and running on your system.

Running the Code

1. Place all the PDF files from which you want to extract data in a folder named `DSCI560_Lab5`.
2. Open a terminal or command prompt.
3. Navigate to the directory containing `Lab5_Part1.py`.
4. Run the script:

```
1 python Lab5_Part1.py
```

Outputs

1. `Task_PDF_original.csv`: Contains the data extracted from the PDF files.
2. `Task1_PDF_preprocessed.csv`: Contains the preprocessed version of the extracted data.
3. Data stored in the MySQL database in the table named `wells_data`.

```
oil_wells_analysis.ipynb > !pip install PyPDF2
+ Code + Markdown | ▶ Run All | ✕ Clear All Outputs | ☰ Outline | ... | Select Kernel

... Mounted at /content/drive
/content/drive/Shared drives/560_GROUP/DSCI560_Lab5/W11920.pdf
Extract well data from pdfs: {'Well_name': 'MAGNUM 3-36-25H', 'API#': '33-053-04069'}
Extract well data from pdfs: {'Well_name': 'MAGNUM 2-36-25H', 'API#': '33-053-03944'}
Extract well data from pdfs: {'Well_name': 'Colville 5301 44-12T', 'API#': '33-053-04981'}
Extract well data from pdfs: {'Well_name': 'DAHL FEDERAL 2-15H', 'API#': '33-002-58854'}
Extract well data from pdfs: {'Well_name': 'Magnum 1-36-25H', 'API#': '33-053-03943'}
Extract well data from pdfs: {'Well_name': '(see details', 'API#': '33-053-06010'}
Extract well data from pdfs: {'Well_name': 'DAHL 15-11H', 'API#': '33-002-58854'}
Extract well data from pdfs: {'Well_name': 'BRAY 5301 43-12H', 'API#': '33-053-03911'}
Extract well data from pdfs: {'Well_name': 'FOLEY FEDERAL 5301 43-12H', 'API#': '33-000-10000'}
Extract well data from pdfs: {'Well_name': 'Chalmers Wade Federal 5300 44-24 12TXR', 'API#': '33-053-06012'}
Extract well data from pdfs: {'Well_name': 'Cc lumbus Federal 1-16H', 'API#': '33-053-04852'}
Extract well data from pdfs: {'Well_name': 'Gamma Federal 5300 41-3113T2', 'API#': '33-053-06232'}
/content/drive/Shared drives/560_GROUP/DSCI560_Lab5/W15358.pdf
Extract well data from pdfs: {'Well_name': 'Kline Federal 5300 31-18 6B', 'API#': '33-053-06057'}
Extract well data from pdfs: {'Well_name': '(see details', 'API#': '33-053-06011'}
Extract well data from pdfs: {'Well_name': 'Atlanta 4-6H', 'API#': '33-105-02729'}
Extract well data from pdfs: {'Well_name': 'Columbus Federal 3-16H', 'API#': '33-053-04856'}
Extract well data from pdfs: {'Well_name': 'Columbus Federal 2-16H', 'API#': '33-053-04855'}
Extract well data from pdfs: {'Well_name': '(see details', 'API#': '33-053-05924'}
Extract well data from pdfs: {'Well_name': 'Innoko 5301 43-12T', 'API#': '33-053-03911'}
Extract well data from pdfs: {'Well_name': 'Atlanta 2-6H I Approximate Start Date D Drilling Prognosis', 'API#': '33-105-02731'}
Extract well data from pdfs: {'Well_name': 'Atlanta Federal 8-6H', 'API#': '33-105-02725'}
Extract well data from pdfs: {'Well_name': 'Atlanta Federal 10-6H', 'API#': '33-105-02723'}
Extract well data from pdfs: {'Well_name': 'Atlanta #1 SWD I Approximate Start Date D Drilling Prognosis', 'API#': '33-105-90258'}
Extract well data from pdfs: {'Well_name': 'Carson SWD 5301 12-24', 'API#': '33-053-90329'}
Extract well data from pdfs: {'Well_name': 'Gamma Federal 5300 41-31 12B', 'API#': '33-053-06231'}
/usr/local/lib/python3.10/dist-packages/PyPDF2/cmap.py:142: PdfReadWarning: Advanced encoding /90ms-RKSJ-H not implemented yet
warnings.warn(
/usr/local/lib/python3.10/dist-packages/PyPDF2/cmap.py:142: PdfReadWarning: Advanced encoding /90ms-RKSJ-V not implemented yet
warnings.warn(
Extract well data from pdfs: {'Well_name': 'Atlanta 13-6H', 'API#': '33-105-02720'}
Extract well data from pdfs: {'Well_name': 'AUmta 12-6H Sec 5, 6, 7, & 8 T153N R101W', 'API#': '33-105-02721'}
Extract well data from pdfs: {'Well_name': 'Atlanta Federal 7-6H', 'API#': '33-105-02726'}
Extract well data from pdfs: {'Well_name': 'Atlanta Federal 9-6H', 'API#': '33-105-02724'}
Extract well data from pdfs: {'Well_name': 'Tallahassee 3-16H', 'API#': '33-053-04853'}
Extract well data from pdfs: {'Well_name': 'BUCK SHOT SWD 5300 31-31', 'API#': '33-053-90244'}
Extract well data from pdfs: {'Well_name': 'Atlanta Federal 5-6H', 'API#': '33-105-02728'}
Extract well data from pdfs: {'Well_name': 'Atlanta 1-6H', 'API#': '33-105-02732'}
Extract well data from pdfs: {'Well_name': 'Wade Federal 5300 21-30 12T', 'API#': '33-053-06129'}
```

Description of Functions

1. `well_name_func(pdf)` : Determines the pattern to extract the well name based on the PDF filename.
2. `api_func(pdf)` : Determines the pattern to extract the API number based on the PDF filename.
3. `api_fix_func(pdf, API_pattern_match)` : Fixes the API number format based on the matched API pattern.
4. `extract_well_data_from_pdfs(folder_path)` : Extracts well-related data from the PDF files.
5. `extract_text_from_pdf(pdf_path, page_number)` : Extracts text from a specific page of a PDF using OCR.
6. `extract_row1_data(row1_string)`, `extract_row2_data(row2_string)` : Extract data from specific strings.
7. `extract_stimulation_data_from_pdfs(pdf_files)` : Extracts stimulation data from the list of PDF files.
8. `scrape_well_data(df)` : Gathers well-related information from a website using web scraping.
9. `merge_and_save_dataframes(files_dict, stimulation_dict, output_folder)` : Merges extracted data and saves it as a CSV.
10. `preprocess_and_save_dataframes(df, output_folder, task)` : Preprocesses the data and saves it as a CSV.
11. `store_db(df)` : Stores the final data in a MySQL database.

Notes

- The PDF extraction process uses patterns matched with regular expressions. If the structure of the PDFs changes, the patterns might need adjustments.
- Web scraping is done using the Selenium library and the Chrome browser in headless mode. If the structure of the target website changes, the scraping logic might need adjustments.
- Before running the code, ensure that the MySQL server is up and running and that you have provided the correct database credentials.

Always make backups of your data and run the code on a small set of files first to ensure it works as expected.