

# Project Report: Invoice Data Extraction and Accuracy Evaluation

## Introduction

This project automates the extraction of structured information from both regular and scanned PDF invoices. The primary goal is to reduce manual entry errors and improve efficiency by using advanced text extraction techniques and Optical Character Recognition (OCR).

## Components Overview

The project consists of several key modules, each tailored to specific aspects of the extraction and evaluation process:

1. **Data Extraction: (`data_extraction.py`)** Handles the reading and extraction of text from PDF files using both direct methods and OCR.
  - a. To retrieve text from PDF files.
  - b.
    - i. **Regular PDFs:** Utilizes `PyMuPDF (fitz)` to directly extract text.
    - ii. **Scanned PDFs:** Converts PDF pages to images using `pdf2image`, then applies `pytesseract` for OCR to extract text. Image quality is enhanced using `OpenCV` for better OCR accuracy.
  - c. Differentiating between regular and scanned PDFs is essential as text extraction methods vary significantly in approach and complexity, affecting the accuracy and efficiency of data retrieval.
  - d. Module Description:

**`extract_text_from_pdf`:** This function first attempts to extract text directly from the PDF using `fitz` (PyMuPDF). If it finds that the PDF is mostly images or if the text extraction does not yield usable results, it falls back to OCR.

**`extract_text_from_scanned_pdf`:** This function handles cases where the PDF is composed of images, which is common with scanned documents. It converts each page to an image using `pdf2image`, processes it for OCR readiness by converting to grayscale and applying thresholding, and then uses `pytesseract` to perform OCR and extract text.

This setup ensures comprehensive handling of different types of PDFs, maximizing the chances of successful data extraction from both digitally created and scanned documents. The fallback mechanism from direct text extraction to OCR ensures that the module can handle a wide variety of document formats efficiently.

2. **Text Processing (`text_processing.py`)** Uses regular expressions to parse and extract specific data points from the raw text.
- a. To parse the raw text extracted from PDFs and identify structured information such as invoice numbers, dates, amounts, etc.
  - b. Employs regular expressions to search and capture relevant data based on patterns that define the typical formatting of invoice data.
  - c. Parsing transforms unstructured text into structured data, which is essential for subsequent processing and analysis, ensuring only relevant information is extracted and used.
  - d. Description:

**`parse_invoice_data`**: This function applies regular expressions to raw text to extract structured information. Each field (like invoice number, dates, customer details, etc.) is identified with a unique regex pattern. If the pattern finds a match in the text, it's added to the resulting dictionary; otherwise, the field is set to `None`.

**`extract_fields_from_texts`**: This function processes a list of text blocks from multiple invoices, using `parse_invoice_data` to parse each one. It accumulates the results in a list, allowing for batch processing of multiple invoice texts.

This module simplifies the task of converting unstructured text from invoices into a structured format, making it easier to handle and prepare for further processing, such as data cleaning and analysis.

3. **Data Processing (`data_processing.py`)** Manages data loading, cleaning, and merging to prepare for accuracy evaluation.
- a. To clean and merge the extracted data with a ground truth dataset for comparison.
  - b. Standardizes the formatting of text fields (e.g., trimming whitespace, normalizing case) and merges datasets on common keys (e.g., Invoice Number) to align the extracted data with its corresponding ground truth.
  - c. Cleaning reduces discrepancies caused by formatting differences, while merging aligns the extracted data with its validation counterpart, essential for accurate evaluation.
  - d. Description:

**`load_data`**: This function reads data from a CSV file into a DataFrame. It includes basic error handling to manage file not found errors and other I/O issues.

**merge\_data**: This function merges two datasets based on a specified key, such as 'Invoice Number'. It allows for inner join operations, ensuring that only matching records are kept.

**clean\_data**: This function standardizes the formatting of specified columns in the DataFrame to ensure uniformity, which is crucial for accurate data comparison. It handles tasks such as stripping excess whitespace, converting text to lowercase, and replacing multiple spaces with a single space.

These functions collectively prepare the dataset for accurate and efficient data analysis, handling fundamental data management tasks required before the data can be evaluated for extraction accuracy.

4. **Accuracy Evaluation :([accuracy\\_evaluation.py](#))** Compares extracted data against a ground truth dataset to measure extraction accuracy.
  - a. To compute and report the accuracy of the extracted data by comparing it against the ground truth.
  - b. Uses both direct comparison and fuzzy matching (for fields with expected textual variations) to evaluate the accuracy of each field.
  - c. Provides quantitative metrics on the performance of the extraction process, highlighting potential areas for improvement and confirming the reliability of the system.
  - d. Description:

**clean\_text**: This function normalizes text to ensure that comparisons are not affected by differences in formatting or case.

**evaluate\_accuracy**: This function compares extracted and ground truth data for each column where exact matching is expected, calculating the percentage of correct matches.

**evaluate\_fuzzy\_matching**: This function uses fuzzy logic to evaluate fields where variations in text representation might occur, useful for fields like addresses or detailed descriptions.

This setup ensures that the evaluation of extracted data can be done robustly, handling both straightforward and complex cases, thus providing a comprehensive assessment of the extraction system's performance.

## Accuracy Evaluation Metrics

The system's accuracy in extracting various fields from invoices is as follows:

- **Accuracy for Invoice Date:** 1.00 (100%)
- **Accuracy for Due Date:** 1.00 (100%)

- **Accuracy for Customer Details:** 1.00 (100%)
- **Accuracy for Place of Supply:** 1.00 (100%)
- **Accuracy for Taxable Amount:** 1.00 (100%)
- **Accuracy for Total:** 1.00 (100%)
- **Accuracy for Total Discount:** 1.00 (100%)
- **Accuracy for Bank Details:** 1.00 (100%)

These metrics indicate that the project achieves perfect accuracy across all evaluated fields.

## Methodology

1. **Direct Text Extraction from Text-Based PDFs:**Text is directly extracted from PDFs where text layers are accessible, using PyMuPDF (fitz). This method is preferred for its efficiency and accuracy in handling digital PDFs.
2. **OCR for Image-Based PDFs:**For PDFs that are essentially scanned images, the project employs `pdf2image` to convert pages to images, which are then processed using `pytesseract` after image enhancement techniques (grayscale conversion and thresholding) to extract text.
3. **Text Parsing:**Extracted text is parsed using predefined regular expressions that identify specific invoice-related information, ensuring that only relevant data is captured.
4. **Data Cleaning and Merging:**Extracted data is cleaned for consistency (e.g., removing extra spaces, standardizing text formats) and merged with ground truth data for comparison.
5. **Accuracy Evaluation:**The accuracy of extracted fields is assessed by comparing them against the ground truth, with perfect scores indicating a highly reliable extraction process.

## Conclusion

The Invoice Data Extraction and Accuracy Evaluation project demonstrates exceptional performance with 100% accuracy across all fields, underscoring the effectiveness of the methodologies used. This automation not only enhances data processing efficiency but also ensures high reliability, crucial for applications that require precise data extraction from various document formats.