

MTP2

K Laxman 2018CS50408

COD892 Project Proposal

Advisors: Prof.M.Balakrishnan, Prof.Volker Sorge, Abhishek Baghel

Supervisor: Prof. Vireshwar Kumar

January 04, 2023



Extracting table from pdf using layout parser and GCV(Google Cloud vision)

Automated tool for converting documents containing the pdf files to extract text out of the image per page

- ❶ **Problem** : extract structured data from a PDF file containing tables and text .
- ❷ **Algorithm** :
- ❸ Import required files and Load PDF file using pdf2image.
- ❹ For each page in the PDF file:
 - ❺ a. Perform layout analysis using Detectron2 LayoutModel from layout parser.
 - ❻ b. Draw boxes around detected elements using lp.drawbox function.
 - ❼ c. Extract tables from the layout analysis using lp.Layout() function.
 - ❽ d. Perform text recognition on the page image using the GCV API through the layoutparser library.

Extracting table from pdf using layout parser and GCV(Google Cloud vision)

Automated tool for converting documents containing the pdf files to extract text out of the image per page

- ① e. Draw text boxes around detected text using the `lp.drawtext()` function.
- ② f. Extract text content and location coordinates using OCR agent from layoutparser.
- ③ Group text content into columns based on the x-coordinates
- ④ Convert the extracted table data into a pandas DataFrame and save as a CSV file.
- ⑤ end of loop and end of program

Extracting table from pdf using layout parser and GCV(Google Cloud vision)

Group text content into columns based on the x-coordinates .How does it work ?

- ❶ Initialize an empty dictionary to hold the columns
- ❷ Loop through each text block:.
- ❸ a. Determine the x-coordinate of the text block by finding the average of the minimum and maximum x-values of the block's bounding box
- ❹ b. Check if there is already a column for that x-coordinate:
- ❺ i. If there is, append the text content to the existing column
- ❻ ii. If there is not, create a new column and add the text content to it
- ❼ Determine the maximum length of any column
- ❽ Remove any columns with fewer elements than the maximum length minus 3
- ❾ Fill any missing values in each column with "None"
- ❿ Create a pandas DataFrame with the resulting column-wise text content and save as CSV