

# Similar Document Template Matching Algorithm

## Bounding Box Detection and Visualization

In this code, we demonstrate how to detect and visualize bounding boxes around text regions in an image using Python. We utilize the OpenCV library for image manipulation and pytesseract for text recognition.

### Code Explanation

- Import Required Libraries:** We import the necessary libraries, including cv2 for OpenCV, pytesseract for text extraction, and Output from pytesseract.
- Load the Image:** We load the image 'invoice\_0\_charspace\_1.jpg' into the 'img' variable using OpenCV's cv2.imread() function.
- Text Detection with pytesseract:** We use pytesseract's image\_to\_data function to detect text regions in the image. The result is stored in the dictionary 'd', containing information about each detected text region, such as its coordinates, size, and confidence score.
- Print Available Information:** We print the keys of the 'd' dictionary to see the available information about each detected text region.
- Bounding Box Visualization:** We iterate through the detected text regions and check if their confidence score is greater than 60. If the confidence is above the threshold, we draw a green bounding box around the text region using cv2.rectangle().
- Display the Image:** We create an OpenCV window named 'output', resize the image for better visualization, and display the image with the bounding boxes using cv2.imshow(). The cv2.waitKey(0) function waits for a key press before closing the window.
- Cleanup:** Finally, we destroy the OpenCV window to release resources.

Please make sure to adjust the image path if necessary and ensure that you have OpenCV and pytesseract properly installed in your environment before running this code.

```
In [ ]:
import cv2
import pytesseract
import Output

In [2]:
img = cv2.imread("/home/nee1/invoice_0_charspace_1.jpg")

In [3]:
d = pytesseract.image_to_data(img, output_type=Output.DICT)
# Print the keys of the resulting dictionary to see the available information
print(d.keys())

dict_keys(['level', 'page_num', 'block_num', 'par_num', 'line_num', 'word_num', 'left_t', 'top', 'width', 'height', 'conf', 'text'])

In [4]:
n_boxes = len(d['text'])
for i in range(n_boxes):
    if d['conf'][i] > 60:
        # Check if confidence score is greater than 60
        (x, y, w, h) = (d['left'][i], d['top'][i], d['width'][i], d['height'][i])
        img = cv2.rectangle(img, (x, y), (x + w, y + h), (0, 255, 0), 2)

import cv2
cv2.namedWindow("output", cv2.WINDOW_NORMAL)
img = cv2.imread('img.jpg')
ims = cv2.resize(img, (540, 540))
cv2.imshow("output", ims)
cv2.waitKey(0)
cv2.destroyAllWindows()

Out[4]:
225

In [ ]:
cv2.destroyAllWindows('img')
```

### Getting in Coordinates of Logo

- Importing Libraries:** The code starts by importing the necessary libraries.
- Spacy:** spacy is imported for natural language processing tasks, although it's not used in this code.
- cv2 (OpenCV):** cv2 is imported for image processing.
- numpy:** numpy is imported for numerical operations.
- Loading Pretrained YOLO Model:** The code loads a pretrained YOLO (You Only Look Once) model for object detection. YOLO is a deep learning model used for object detection tasks.
- Loading Class Names:** It reads a file named 'coco.names' that contains the class names used in the YOLO model. These class names include common objects and items that can be detected by the model.
- Loading and Preprocessing the Image:** The code loads an image named 'Original\_01\_page-0001.jpg' and preprocesses it for YOLO input. It resizes the image to 416x416 pixels, scales pixel values to a range between 0 and 1, and performs color channel swapping.
- Setting Input for the YOLO Model:** The preprocessed image is set as the input to the YOLO model.
- Getting Detections:** The code forward-connects the detected objects in the YOLO model to obtain object detections. The net.getUnconnectedOutLayersNames() method retrieves the output layer names of the model.
- Looping Over Detections:** The code then loops through the detected objects in the image.
  - It extracts scores, class IDs, and confidence values for each detected object.
  - If the confidence of a detected object is greater than 0.5, it assumes the object is a logo.
- Extracting Logo Coordinates:** For each detected object, the code extracts its coordinates and dimensions.
  - center\_x and center\_y are the coordinates of the center of the logo.
  - width and height represent the dimensions of the bounding box around the logo.
  - x1, y1, x2, and y2 calculate the coordinates of the top-left and bottom-right corners of the bounding box.
- Cropping the Logo Region:** It crops the region of the image containing the detected logo using the calculated coordinates.
- Checking Logo Region:** It checks if the cropped logo region contains non-black and non-white pixels, assuming an RGB image. This step helps verify if the region genuinely contains a logo.
- Printing Information:** If the region is considered a logo, the code prints the confidence score and the coordinates of the top-left corner of the logo bounding box.

Overall, this code snippet uses a pretrained YOLO model to detect logos in an image, extracts their coordinates, and checks if the region truly contains a logo based on pixel color. Detected logos meeting the confidence threshold are printed with their coordinates.

```
In [3]:
import spacy
# Load the english pre-trained model with NER
nlp = spacy.load("en_core_web_sm")

/home/nee1/anaconda3/lib/python3.8/site-packages/torch/cuda/_init_.py:546: UserWarning:
Can't initialize NVML
warnings.warn("Can't initialize NVML")

In [1]:
!python3 -m spacy download en_core_web_sm

/home/nee1/anaconda3/lib/python3.8/site-packages/torch/cuda/_init_.py:546: UserWarning:
Can't initialize NVML
warnings.warn("Can't initialize NVML")
Collecting en-core-web-sm==3.6.0
  Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.6.0/en_core_web_sm-3.6.0-py3-none-any.whl (12.8 MB)
12.8/12.0 MB 4.6 MB/s eta 0:00:00 eta 0:00:01
Requirement already satisfied: spacy<3.7.0,=>3.6.0 in ./anaconda3/lib/python3.8/site-packages (from en-core-web-sm==3.6.0) (3.6.1)
Requirement already satisfied: spacy-legacy<3.0,=>3.0.11 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0,=>1.0.0 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (1.0.0)
Requirement already satisfied: murmurhash<1.0,=>2.0.8 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2.0.7)
Requirement already satisfied: preshed<3.1.0,=>3.0.2 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (3.0.0)
Requirement already satisfied: catalogue<2.0,=>1.0.6 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (1.0.12)
Requirement already satisfied: srsly<3.0,=>2.4.3 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2.4.7)
Requirement already satisfied: catalogue<2.0,=>1.0.6 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (1.0.12)
Requirement already satisfied: typing-extensions<4.6, in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (4.0.0)
Requirement already satisfied: pathy<0.10.0,=>0.10.0 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (0.10.2)
Requirement already satisfied: smart-open<7.0,=>5.2.1 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (5.2.1)
Requirement already satisfied: tqdm<5.0,=>4.38.0 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (4.59.0)
Requirement already satisfied: numpy<1.15.0,=>1.15.0 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (1.20.1)
Requirement already satisfied: requests<3.0,=>2.13.0 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2.25.1)
Requirement already satisfied: jinja2<3.0,=>2.11.1 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2.11.3)
Requirement already satisfied: setuptools in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (68.2.1)
Requirement already satisfied: packaging<20.0,=>20.9 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (23.1)
Requirement already satisfied: langcodes<4.0,=>3.2.0 in ./anaconda3/lib/python3.8/site-packages (from spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (3.2.0)
Requirement already satisfied: annotated-types<0.4.0, in ./anaconda3/lib/python3.8/site-packages (from pydantic<1.8,=>1.8.1,=>3.0,=>1.7.4-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (0.6.0)
Requirement already satisfied: pydantic-core==2.6.3 in ./anaconda3/lib/python3.8/site-packages (from pydantic<1.8,=>1.8.1,=>3.0,=>1.7.4-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2.6.3)
Requirement already satisfied: typing-extensions<4.6, in ./anaconda3/lib/python3.8/site-packages (from pydantic<1.8,=>1.8.1,=>3.0,=>1.7.4-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (4.1.1)
Requirement already satisfied: charset-normalizer<3.0,=>3.0.2 in ./anaconda3/lib/python3.8/site-packages (from requests<3.0,=>2.13.0-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (3.0.0)
Requirement already satisfied: idna<3.0,=>2.5 in ./anaconda3/lib/python3.8/site-packages (from requests<3.0,=>2.13.0-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2.10)
Requirement already satisfied: urllib3<1.27,=>1.21.1 in ./anaconda3/lib/python3.8/site-packages (from requests<3.0,=>2.13.0-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (1.26.4)
Requirement already satisfied: certifi<=2017.4.17 in ./anaconda3/lib/python3.8/site-packages (from requests<3.0,=>2.13.0-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (2020.6.1)
Requirement already satisfied: click<9.0,=>7.1.1 in ./anaconda3/lib/python3.8/site-packages (from typer<0.10.0,=>0.3.0-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (7.1.2)
Requirement already satisfied: MarkupSafe<=0.22 in ./anaconda3/lib/python3.8/site-packages (from jinja2<3.0,=>2.11.1-spacy<3.7.0,=>3.6.0-en-core-web-sm==3.6.0) (1.1.1)
DEPRECATION: pydantic 4.0.0-unsupported has a non-standard version number. pip 23.1 will enforce this behaviour change. A possible replacement is to upgrade to a newer version of pip: run 'python -m pip install --upgrade pip'.
WARNING: You are using pip version 20.3.0, however you are currently using the default python interpreter (python3.8).
Discussion can be found at https://github.com/pypa/pip/issues/1205
Installing collected packages: en-core-web-sm
Attempting uninstall: en-core-web-sm
Found existing installation: en-core-web-sm 2.2.0
Uninstalling en-core-web-sm: 2.2.0
Successfully uninstalled en-core-web-sm-2.2.0
Successfully installed en-core-web-sm-3.6.0
# download and installation successful
You can now load the package via spacy.load('en_core_web_sm')
```

```
In [1]:
import cv2

/home/nee1/anaconda3/lib/python3.8/site-packages/torch/cuda/_init_.py:546: UserWarning:
Can't initialize NVML
warnings.warn("Can't initialize NVML")

In [2]:
nlp = spacy.load("en_core_web_sm")

In [3]:
with open('Invoice2.pdf', 'r', encoding='utf-8', errors='ignore') as f:
    text = f.read()
# Apply the NER model to the invoice text
doc = nlp(text)

In [4]:
print(doc)

In [5]:
# FitZ
# Title (document)
# Creator (WhiteToPDF 0.12.6.1)
# Producer (Qt 4.8.7)
# CreationDate (D:20230911190846Z)
3 0 obj
<<
/Type /Catalog
/Names 1 0 R
endobj
stream
xref
1 0 obj
<<
/Length 19 0 R
endobj
startxref
1235
endobj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SMask /None>
endobj
4 0 obj
<<
/Type /Page
/Parent 2 0 R
/Contents 9 0 R
/Resources 11 0 R
/Annotations 12 0 R
/MediaBox [ 0 0 612 792 ]
endobj
11 0 obj
<<
/ColorSpace <<
/PCS 4 0 R
/SP /DeviceRGB
/CS /DeviceGray
endobj
/SA false
/SM 0.02
/CA 1.0
/CA 1.0
/AS false
/SM
```



We print the PDF files that belong to each cluster, helping us understand how the documents are grouped.

Plotting Results

Finally, we use Principal Component Analysis (PCA) to reduce the dimensions of the data for visualization purposes. We then create a scatter plot to visualize the clustering results in a 2D space, coloring the points based on their cluster assignments.

The code provides a way to analyze and group PDF documents based on both text and image content, which can be useful in various document categorization and clustering tasks.

```
In [ ]:
import fitz
import cv2
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import make_pipeline
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA

# Function to extract text from a PDF file
def extract_text_from_pdf(pdf_file):
    doc = fitz.open(pdf_file)
    text = ""
    for page_num in range(len(doc)):
        page = doc[page_num]
        text += page.get_text()
    return text

# Function to extract images from a PDF file
def extract_images_from_pdf(pdf_file):
    doc = fitz.open(pdf_file)
    images = []
    for page_num in range(len(doc)):
        page = doc[page_num]
        xref_list = page.get_images(full=True)
        for xref in xref_list:
            base_image = doc.extract_image(xref[0])
            image_data = base_image["image"]
            images.append(image_data)
    return images

# Function to compute TF-IDF similarity between two texts
def compute_similarity(text1, text2):
    tfidf_vectorizer = TfidfVectorizer()
    tfidf_matrix = tfidf_vectorizer.fit_transform([text1, text2])
    similarity_matrix = cosine_similarity(tfidf_matrix, tfidf_matrix)
    return similarity_matrix[0][1]

# Function to compare images using structural similarity index (SSIM)
def compare_images(image1, image2):
    image1 = cv2.imdecode(np.frombuffer(image1, np.uint8), -1)
    image2 = cv2.imdecode(np.frombuffer(image2, np.uint8), -1)
    if image1.shape != image2.shape:
        return 0.0
    return cv2.matchTemplate(image1, image2, cv2.TM_CCOEFF_NORMED)[0][0]

# Load your dataset of PDF files here
# The base path for the PDF files
base_path = 'Samples of electronic invoices/Dataset with valid information/invoice_'

# List to store the file paths for 500 images
pdf_files = []

# Loop to generate file paths for 500 images
for i in range(50):
    pdf_file_path = f'{base_path}{i}.pdf'
    pdf_files.append(pdf_file_path)

# Create feature vectors for text and image similarity
text_similarity_matrix = np.zeros((len(pdf_files), len(pdf_files)))
image_similarity_matrix = np.zeros((len(pdf_files), len(pdf_files)))

for i, pdf_file1 in enumerate(pdf_files):
    text1 = extract_text_from_pdf(pdf_file1)
    images1 = extract_images_from_pdf(pdf_file1)

    for j, pdf_file2 in enumerate(pdf_files):
        text2 = extract_text_from_pdf(pdf_file2)
        images2 = extract_images_from_pdf(pdf_file2)

        text_similarity_matrix[i][j] = compute_similarity(text1, text2)
        if (len(images1) == 0 | len(images2) == 0):
            image_similarity_matrix[i][j] = 0;
        else:
            image_similarity_scores = []
            for image1 in images1:
                for image2 in images2:
                    similarity_score = compare_images(image1, image2)
                    image_similarity_scores.append(similarity_score)
            image_similarity_matrix[i][j] = max(image_similarity_scores)

# Combine text and image similarity scores
combined_similarity_matrix = text_similarity_matrix + image_similarity_matrix

# Normalize the combined similarity matrix
scaler = StandardScaler()
normalized_combined_similarity_matrix = scaler.fit_transform(combined_similarity_matrix)

# Apply K-Means clustering
num_clusters = 5 # You can adjust this based on your dataset and requirements
kmeans = KMeans(n_clusters=num_clusters)
clusters = kmeans.fit_predict(normalized_combined_similarity_matrix)

# Print clusters
for i in range(num_clusters):
    print(f'Cluster {i + 1}:')
    for j, pdf_file in enumerate(pdf_files):
        if clusters[j] == i:
            print(pdf_file)

# Plot the results (for visualization purposes, using the first two dimensions)
pca = PCA(n_components=2)
reduced_features = pca.fit_transform(normalized_combined_similarity_matrix)

plt.scatter(reduced_features[:, 0], reduced_features[:, 1], c=clusters, cmap='viridis')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.title('K-Means Clustering of Documents')
plt.show()
```



In [ ]: