

EN.600.461 Computer Vision

Final Project

Recognizing and Translating Text from Images

Joon Hyuck (James) Choi
Senior Undergraduate
The Johns Hopkins University
3400 N Charles Street
Baltimore, MD 21218, USA
jchoi100@jhu.edu

Abstract

This work explores optical character recognition (OCR) in photos of printed and hand-written documents. It first explores basic preprocessing of photos of documents using OpenCV¹ functionalities for blurring, thresholding, and denoising. It next discusses the use of `tesseract-ocr`² to perform OCR. We then discuss the incorporation of the Google Cloud Translation API³ to translate the OCR results to different languages. It finally discusses template matching using the Johns Hopkins University (JHU) logo.

1 Introduction

This work achieved the following main goals stated in the original project proposal: 1) Given a photo of a document, convert it into a clean scanned version; 2) Take the scanned version and perform OCR. Optionally, this work took the OCR output and translated the text into different languages using the Google Cloud Translation API. This work also experimented with template matching using the JHU logo in order to determine whether the input document is an official JHU document or not.

2 Methods

In this section, we discuss the methods we took and external libraries used for each stage of our work.

2.1 Image Preprocessing

We implemented our program in Python using OpenCV. In order to feed the OCR algorithm clean input to achieve best performance, we used `cv2.medianBlur` to smoothen the input image with an aperture size of 5. Then, we passed the smoothened image through `cv2.adaptiveThreshold` with *adaptiveMethod*= `ADAPTIVE_THRESH_GAUSSIAN_C`, *threshold* = *binary*, *blockSize*=5x5, and *C*=2. The adaptive method we chose uses the weighted sum of the *blockSize* x *blockSize* neighborhood of pixel (*x*, *y*) - *C* as its threshold value. Parameter values were chosen empirically.

Finally, we performed denoising on the thresholded image to remove noise and make the output clean. We used `cv2.fastNlMeansDenoising` with *templateWindowSize*=7, *searchWindowSize* = 21, and *h*=7. Numbers here were also chosen empirically. The motivation behind denoising in our work was that document scans usually have many random specs here and there. Such noise can greatly reduce OCR accuracy because an OCR algorithm may confuse a random spec with punctuation marks or associate a spec with an actual character near by (e.g. confuse an *l* with an *i*).

2.2 OCR

In order to perform OCR on the cleaned inputs, we used the `tesseract-ocr` library through its python wrapper `pytesseract`. We used the raw output from `pytesseract`'s `image_to_string` method to pass into python's file writer and Google Cloud Translation API.

¹<http://opencv.org/>

²<https://github.com/tesseract-ocr>

³<https://cloud.google.com/translate/>

As a side experiment, we used the Keras ⁴ library to train the MNIST ⁵ dataset of 70,000 hand-written digits on three different convolutional neural networks (CNN). Results will be discussed in section 3.3. We referred to code online ⁶ to build the three CNNs.

2.3 Translation

We made use of Google Cloud's Translation API in order to translate the OCR processed document into several different languages based on user command line input. We referred to code posted online ⁷ and modified the details to fit our purpose.

2.4 Template Matching

Aside from OCR, which was the main focus of this project, we experimented with template matching in images. Specifically, we used the JHU logo to determine whether the input document was an official JHU document or not. We naively assumed that official JHU documents contained a JHU logo for our experiment's sake. We used `cv2.matchTemplate` with our input document as the source image, multiple JHU logo images (each with different sizes since `cv2.matchTemplate` is scale sensitive) as templates, and matching method `cv2.TM_CCOEFF_NORMED`. If the normalized `TM_CCOEFF` value came out to be greater than 0.51 for any of the input JHU logo templates of various sizes, we determined that the input document contained a JHU logo in it. The matching method and determinant value of 0.51 were chosen empirically.

3 Results

Here, we qualitatively discuss OCR results on typed and handwritten input documents.

3.1 Typed Documents

We ran our code on five document images with different font types and font sizes. The images were `cv_proj_description.png`,

⁴<https://keras.io/>

⁵<http://yann.lecun.com/exdb/mnist/>

⁶<http://machinelearningmastery.com/handwritten-digit-recognition-using-convolutional-neural-networks-python-keras/>

⁷<http://github.com/mouuff/mtranslate>

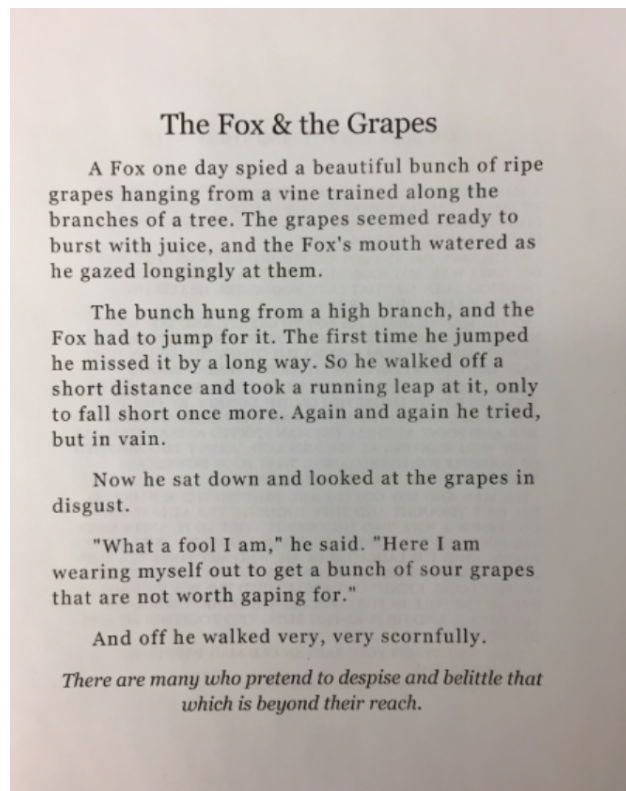


Figure 1: Sample typed document `aesop-fox.png`

`aesop-fox.png`, `aesop-miller.png`, `aesop-goose.png`, and `aesop-cat.png`.

With the exception of `aesop-goose.png`, performance on all documents were near-perfect. A sample image is shown in Figure 1. The reason that `aesop-goose.png` did not produce satisfactory results was that the font in the document was extremely thick and bold. Therefore, after preprocessing, the alphabet characters in the thresholded document were hollow with just the edges remaining. Therefore, the OCR algorithm could not recognize any of the characters in this text.

We ran our code on additional two document images that were taken under sub-optimal conditions: irregular lighting, partial skew, and/or glossy surfaces. Moreover, these documents had mixed font types and sizes, and the structure of the document was more complicated than the previous inputs. These images were `cv-jchoi-proposal.png` and `lav-notes1.png`. Despite the less optimal settings, we correctly denoised and thresholded the images and got similar results as before.

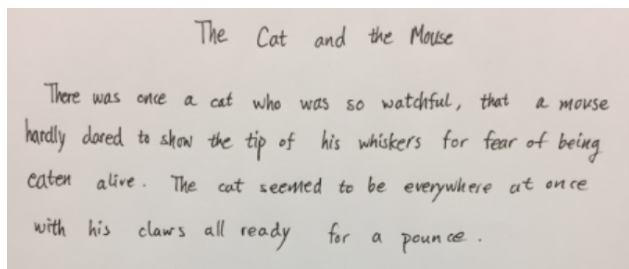


Figure 2: Reasonable handwritten doc: hand_cat.png

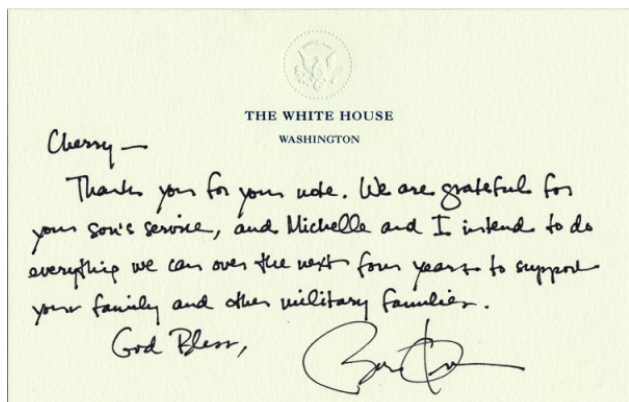


Figure 3: Sample handwritten document: obama1.png

3.2 Handwritten Documents

We applied our code on photos of handwritten documents. Performance on very neatly written documents (hand_fox.png and hand_goose.png) was nowhere close to what we saw for the printed dataset. Moreover, any reasonably neat handwritten documents (hand_miller.png and hand_cat.png) showed even poorer performance. Figure 2 shows one of our reasonably neat images hand_cat.png. Furthermore, our code was not able to detect a single word from document images with much less neat handwriting. Such images were obama1.png, obama2.png, obama3.png, bush1.png, bush2.png, clinton1.png, and clinton2.png. Figure 3 shows obama1.png.

4 Conclusion

We explored optical character recognition on photographed documents in various settings and formats. We tested with optimal and suboptimal photography settings, with documents using various fonts, and with printed and handwritten documents.

For printed documents, we saw highly accurate results for both optimal and suboptimal settings, even if they contained various font sizes and font types. However, we were not able to achieve accurate results with handwritten document photos, disregarding the neatness of the handwriting. For sufficiently small portions of handwritten documents (e.g. a single word or a short sentence), the OCR algorithm was able to produce correct output occasionally. However, even this was not often enough to be considered reliable.

Nonetheless, on typed documents, we were able to effectively preprocess the document image and perform OCR on the cleaned data. Moreover, we leveraged the Google Cloud Translation API to translate the document into several different languages. Our short experiment with template matching also deserves future development into a project of its own.

How to Run the Code

Please extract FINAL_PROJECT-jchoi100.zip to view and run the code.

Assuming that the user has Python 2.5+, OpenCV, pytesseract, and tesseract installed, run the following command on the command line in the same directory as driver.py.

```
$python driver.py {path-to-image-file} {list of target languages separated by single space}
```

Sample usage: \$python driver.py data/doc7.png es fr

performs OCR on doc7.png, outputs the thresholded image, the text in the original input language, and the text translated in Spanish and French.

Saved Output from Experiments

The output that we created by running our code on all the images in our dataset are saved under FINAL_PROJECT-jchoi100/code/data.

Acknowledgments

We thank the image processing and OCR libraries, the code we referred to for constructing CNNs using Keras, and MNIST for the digit dataset.

And thanks to Professor Reiter and the TAs!