

Final project: Burial records classification and text extraction

Northeastern University

EAI6020: AI System Technologies

Daya Rudhramoorthi

July 3, 2021

Group 1

Xinyi Luo

Project Description

Introduction

Burial records are important evidences for tomb management. They often contain the name of the deceased and the detailed location and date of the burial. Because of the large time span since the burial service started, many dated records have not been electronicized, and there are handwritten records. This project aims to digitize these records. Digitalization of these records helps maintain them as well as analysis on historical populations.

Data description:

The data set consists of 120 scanned images in type ‘.jpg’. Each image is one burial record. There are two kinds of formats distinct from each other:

1. Organized form with detailed information (*Type I*)

DECEASED					INTERMENT NO.	
DABA, LEMANEE						
DATE OF INTERMENT	SECTION	LOT - TER	GR	AGE	VETERAN	
6/24/96	33	233ROW1	6	3	NO	
UNDERTAKER						
GATLING'S CHAPEL						
VAULT						
GRAVE BOX #30						
LOT OWNER			RELATION			
NEXT OF KIN			RELATION			
TADESSE DABA, FATHER						
7217 S. MAY, CHGO., IL 60621						
REMARKS						

Essential variables:

“Deceased” -Decedent name

“Lot” -Lot of cemetery

“Date of Interment” -Date of burial

“GR” -Grave spot

“Section” -Section of cemetery

Other variables: “Funeral Home”, “Vault”, “Next of Kin”, “Address”, “Age”..

2. Former form with only essential information (*Type 2*)

Name	Dafos, _____	Leroy _____
Lot		Sec.
Grave	155	Ter. G
Date of Burial	March 15 1916	
Age	4y 6m	
Date of Removal		To
MOUNT HOPE CEMETERY LIBRARY-BUSINESS DATA INDEX OF INTERMENTS		

Essential variables:

“Name” -Decedent name

“Lot” -Lot of cemetery

“Date of Burial” -Date of burial,

“GR” -Grave spot of cemetery

“Sec.” -Section of cemetery

Other variables: “Ter.”, “Date of Removal”, “Age”..

One thing to notice for type 2 forms is that there exist several forms with handwritten contents, which could heavily interfere with the machine's cognition ability as shown below.

Name	Daft. Edgar H	
Lot	248	Sec 5-
Grave	248	Ter 5
Date of Burial	Apr 2, 1897	
Age	33-5-28	
Date of Removal	Apr. 24, 1897	Form 134 / 135
		See 6

Methodology

This digitization project could be divided to two parts:

1. Image preprocess: train a classifier to recognize different types of forms and process the images separately through image preprocessing tools.

Techniques applied:

- CNN model for form classifier.
 - OpenCV for reading and improving image quality.
2. OCR progress: extract text content from treated images and record values into a dataframe file for remaining.

Techniques applied:

- Py-tesseract-OCR for text extraction.
- Azure Form Recognizer API for forms recognition.
- Power Automate, platform to generate workflow for FormRecognizer.
- Google Vision API

Form Classifier Building

Convolutional Neural Network(CNN) model:

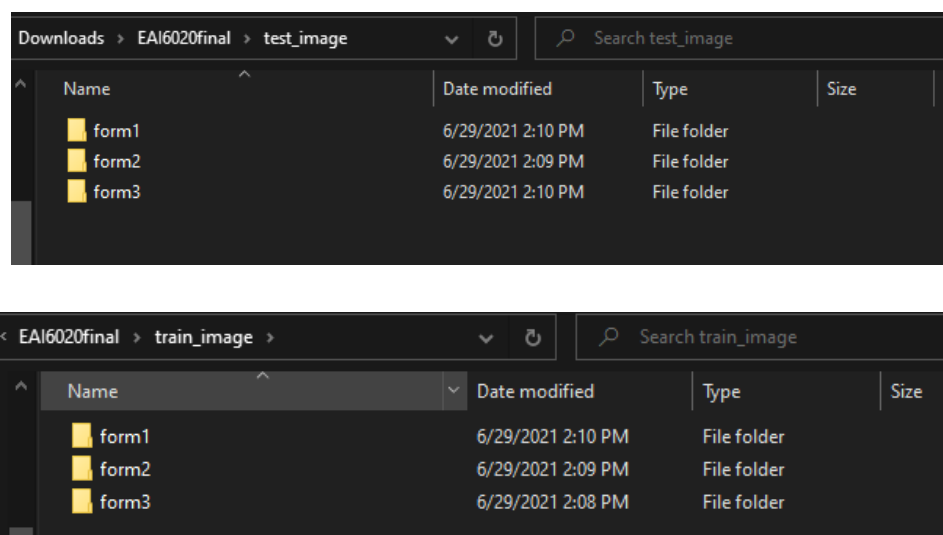
For dealing with images, CNN is known as one of the most powerful models for training with images. Therefore, we are going to build the image classifier based on CNN.

Classifications determination:

Although the dataset has two main types of forms, the handwritten forms seem to cause high error because of their high variance. As a result, our classifier is going to be a multiple classifier for 3 form types:

- Form 1: Organized form
- Form 2: Older form with printed writing and lines
- Form 3: Older form with handwritten content

To prepare for model training, 120 images are divided into 'train_images' and 'test_images' with a ratio of '7/3', which to say the first 36 images are generated to the test dataset and 84 images are put into the train dataset. Both train and test images are divided into 3 folders with labels 'form1', 'form2', 'form3', as shown below.



Import image data through ImageGenerator and check the condition of both datasets. And the output should be a matrix with one-hotkey indicating which form the image is.

```
train_set = train_gen.flow_from_directory('train_image',
                                         target_size = (64, 64),
                                         batch_size = 4,
                                         class_mode = 'binary')

test_set = test_gen.flow_from_directory('test_image',
                                       target_size = (64, 64),
                                       batch_size = 4,
                                       class_mode = 'binary')
```

```
Found 84 images belonging to 3 classes.
Found 36 images belonging to 3 classes.
```

```
from keras.preprocessing import image
train_set.class_indices

{'form1': 0, 'form2': 1, 'form3': 2}
```

Model initialization:

For the CNN mode, 'relu' is going to be applied as the activation function.

Structure of the model:

- Two convolutional layers with activation function and followed by two pooling layers to feed forward the image data.
- One flatten() layer to connect all nodes.
- One dense layer with 'relu' function to generate the result of parameter weights.
- One Softmax() layer with 'softmax' function to sparse the result into one of the 3 classifications set before.
- Compile with 'RmsProp' optimization function.

```
# Initialising the CNN
cnnmodel = Sequential()

cnnmodel.add(Conv2D(32, (3, 3), input_shape = (64, 64, 3), activation = 'relu'))
cnnmodel.add(MaxPooling2D(pool_size = (2, 2)))

cnnmodel.add(Conv2D(32, (3, 3), activation = 'relu'))
cnnmodel.add(MaxPooling2D(pool_size = (2, 2)))

cnnmodel.add(Flatten())

cnnmodel.add(Dense(units = 128, activation = 'relu'))
cnnmodel.add(Dense(units = 3, activation = 'softmax'))

cnnmodel.compile(optimizer = 'rmsprop', loss = 'sparse_categorical_crossentropy', metrics = ['accuracy'])
```

Classifier Performance:

Training parameters: {train_set, steps_per_epoch = 21, epochs = 50, validation_data = test_set, validation_steps = 3}

With a validation dataset with 3*batch_size = 12 images, the model reaches 1.0 accuracy on the test set. The performance seems satisfying. Further image improvement will be based on the classification result.

```
Epoch 48/50
21/21 [=====] - 0s 23ms/step - loss: 0.1168 - accuracy: 0.9869 - val
y: 1.0000
Epoch 49/50
21/21 [=====] - 0s 23ms/step - loss: 3.1752e-04 - accuracy: 1.0000
racy: 1.0000
Epoch 50/50
21/21 [=====] - 0s 23ms/step - loss: 3.8014e-05 - accuracy: 1.0000
racy: 1.0000
WARNING:tensorflow:Your input ran out of data; interrupting training. Make sure that your da
t least `steps_per_epoch * epochs` batches (in this case, 100 batches). You may need to use
ng your dataset.
Accuracy of the model: [0.00018936669221147895, 1.0]
```

Image Preprocess

Image improvement with OpenCV:

Resize all images to a shape (900, 500). Crop all images with only 'essential variables' area.

Process all images through grayscale, denoise and filtering of OpenCV.

Recheck the size of processed images: 'Form1': all sized (300, 1800)

'Form3': most sized (380,1620). However, 'Form2' have sizes distinct in a wide range.

```
[[0. 1. 0.]]
This is a form 2 image
the shape of processed image: (760, 3240)
Decedent D ahleen Myrtle
```

```
[[0. 1. 0.]]
This is a form 2 image
the shape of processed image: (1520, 6480)
Decedent atin Vietor
```

Therefore, to cut small fields of a single variable, we need to cut with ratio. (For example: How much percentage does the 'Decedent' block covering the processed Form 1 image)

```
# Cut needed blocks for different information by ratios\n",
def imacrop(type,attribute,image):
    height,width = image.shape[:2]

    if(type=='f1'):
        if(attribute=='decedent'):
            startrow,startcol = int(height*0.1923),int(width *0.01)
            endrow,endcol = int(height *0.67),int(width *0.6591)
            return image[startrow:endrow,startcol:endcol]
```

Text Extraction with Py-tesseract:

Install the tesseract in system and direct to the commander:

```
import pytesseract
pytesseract.pytesseract.tesseract_cmd = r'C:\\Program Files\\Tesseract-OCR\\tesseract.exe'
```

For a single variable, extract the string from the cropped block and tune slightly with re.sub().

```
image_placeholder = form1process(image_red,a)

f1_decident = pt.image_to_string(imacrop('f1','decident',image_placeholder), lang = 'eng')
f1_decident = re.sub('[^A-Za-z0-9]+',' ', f1_decident)
f1_decident = f1_decident.replace('DECEDEMENT','')
f1_decident = f1_decident.replace('DECEASED','')
```

Text Extraction Performance:

Surely, tesseract could extract text from the designed area. However, the results obviously need much more tuning and improvement. I was trying to include as much information as possible to adapt the most images. That causes there is much noises text involved and extracted with needed information, as shown below:


```

[[0. 0. 1.]]
This is a form 3 image
the shape of processed image: (380, 1620)
Decedent Dahike Reinhold
Date nt "July 13, 1945
Section
Lot NEt 190
Grave

[[0. 0. 1.]]
This is a form 3 image
the shape of processed image: (380, 1620)
Decedent Dahlman William 3
Date Mareh 10 1919
Section 18
Lot Vis 87
Grave

[[0. 1. 0.]]
This is a form 2 image
the shape of processed image: (760, 3240)
Decedent Jahlman William HB
Date farch 10 1919
Section
Lot B87
Grave

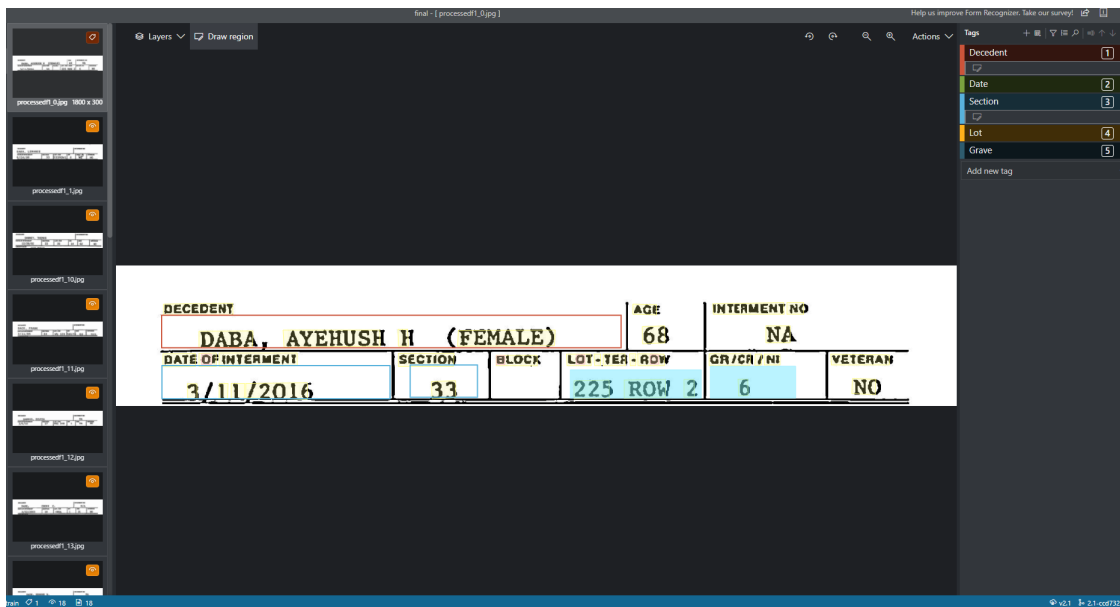
[[1. 0. 0.]]
This is a form 1 image
the shape of processed image: (300, 1800)
Decedent OEEceasid DAILY DOLORES M
Date Sr aw ee hadi1/2/2001
Section
Lot INW266
Grave

```

Apply Azure FormRecognizer API

FormRecognizer(FR) is an AI-powered text extractor deployed by Microsoft on their cloud platform, Azure. FormRecognizer could be trained to recognize the layout of certain forms and automatically drag paired information from forms with the same layout.

Step 1: Manually draw regions on the train set images and attach labels to them.



Step 2: Train the model

Train Result

Model ID: af13f605-ee93-4817-a46b-e08ebaa410c8

Tag	Estimated Accuracy
Date	85.70%
Decedent	57.10%
Grave	85.70%
Lot	85.70%
Section	85.70%

Train a new model

Model name:

[Train](#)

Training record

Model information

Model ID: af13f605-ee93-4817-a46b-e08ebaa410c8

Model Name: f1

Created date and time: 7/3/2021, 1:45:48 AM

Average accuracy: 80.00%

[Learn more about improving model accuracy.](#)

[Download JSON file](#)

Step 3: Check the performance by predicting a new image.

The screenshot displays the Azure AI Studio interface for a Form Recognizer model. On the left, a sample image of a cemetery record form is shown. The form contains the following text:

- DECEASED** (highlighted in red)
- DABA, LEMANEE** (highlighted in red)
- DATE OF INTERMENT** 6/24/96
- SECTION** 33
- LOT - TER** 233ROW1
- GR** 6
- AGE** 3 (with handwritten 'NA' below it)
- VETERAN** NO

On the right, the 'Predictions' section shows the model's output:

- modelId:** cdddcce-e238-4047-8629-37c1aeceb9
- docTypeConfidence:** 96.70%

Page #	Field name	Value	Confidence
1	Date	6/24/96	99.00%
1	Section	33	99.40%
1	Grave	6	99.00%
1	Decedent	DABA, LEMANEE	91.90%
1	Lot	233ROW1	94.40%

I used the cloud sample trainer to train my FormRecognizer model. And generate a workflow on Power Automate platform to run the pipeline through processed images to excel data.

Step 1: Upload data to Azure Blob storage

The screenshot shows the Power Automate workflow editor. The workflow is titled 'finalOCR'. The steps are as follows:

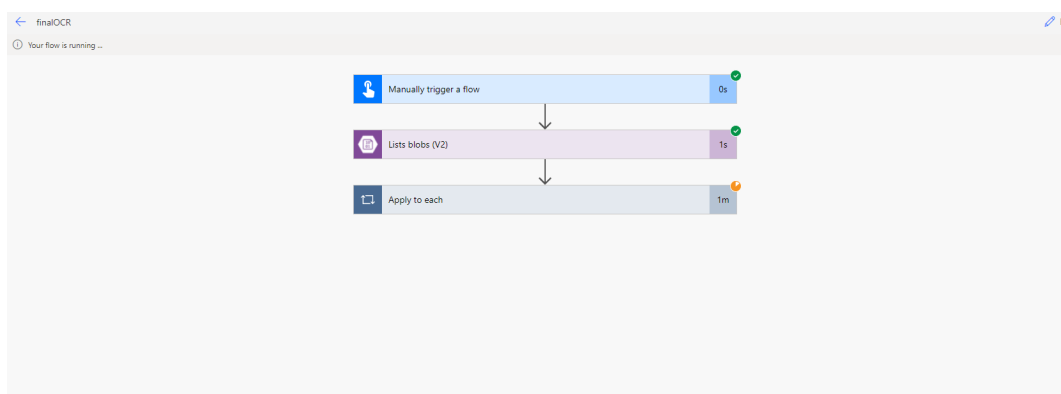
- Lists blobs (V2) (Preview)**
- Apply to each** (with a loop over 'value')
- Copy blob (V2) (Preview)**
- Get blob content (V2) (Preview)**
- Form Recognizer** (with the following configuration):
 - Connection name: cdddcce-e238-4047-8629-37c1aeceb9
 - Endpoint URL: https://final3.cognitiveservices.azure.com/
 - Account Key: [redacted]

Step 2: Conduct workflows with function nodes on Power Automate

The screenshot shows the continuation of the Power Automate workflow. The steps are as follows:

- Apply to each 2** (with a loop over 'pageResults')
- Create CSV table**

Step 3: Run the flow and export the result.



FormRecognizer Performance:

FR works better than the original tesseract. Results come in a well-organized format with clean format. That is because FR firstly recognizes text content clusters in the whole image, then starts to learn the layout through text clustering. That helps it avoid many noise pixels. And I'm using this dataset as my final submission.

A	B	C	D	E
Decedent	Date	Section	Lot	Grave
DABA AYEHUSH	3/11/2016	33	225 ROW 2	6
DABA LEMANE	6/24/96	33	233ROW1	6
DABA SHORRO	9/19/96	33	233ROW1	2
DABE PATRICIA	5/11/96			510
DABNEY ABRAH	12/8/97	28	16	6
DABNEY ERWIN	11/28/2015	30	30 ROW 2	5

Other text extraction alternatives:

Google Vision API:

For all visual inspection needs in manufacturing, we now have a dedicated solution. [Learn more](#)

Vision AI

Derive insights from your images in the cloud or at the edge with AutoML Vision or use pre-trained Vision API models to detect emotion, understand text, and more.

[Try Vision AI free](#)

- Use machine learning to understand your images with industry-leading prediction accuracy
- Train machine learning models that classify images by your custom labels using AutoML Vision
- Detect objects and faces, read handwriting, and build valuable image metadata with Vision API

VIDEO

Fortune 500 global power company AES drives green energy with AutoML Vision

BENEFITS

Detect objects automatically

Reduce purchase friction

With Vision API's [vision product search](#), retailers can create an

Gain intelligence at the edge

It is also a trendy tool for computer vision tasks.

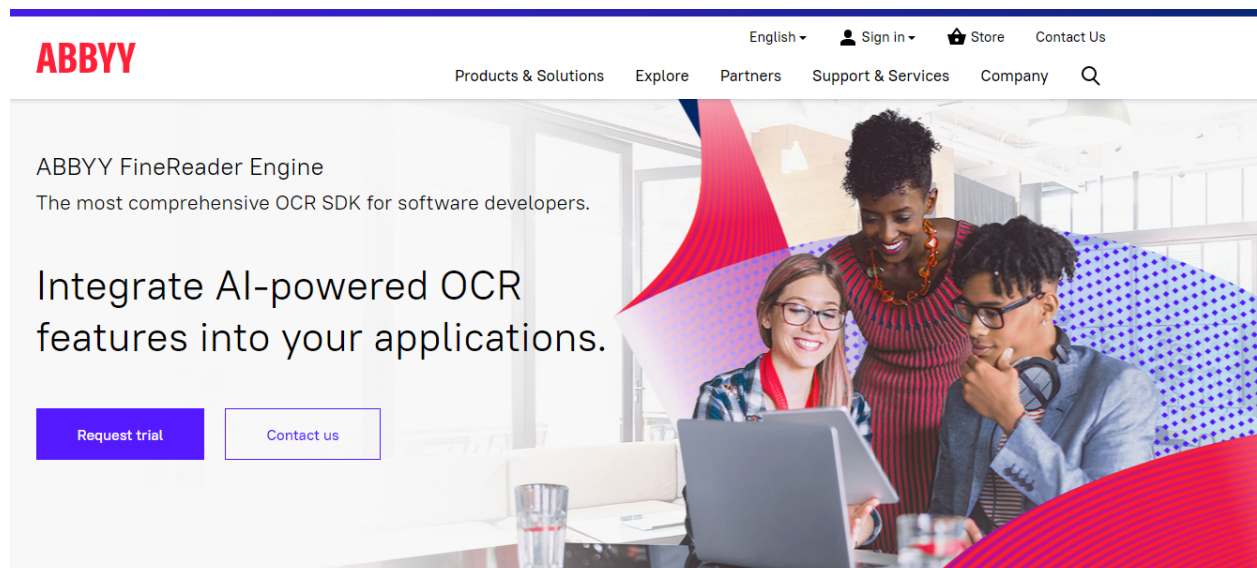
To apply the Google Vision, the most simple way is to deploy a script with access to Google Cloud Platform.

Vision API could be installed and imported as an extension.

With a GCP account, users could apply the API directly to the data and generate text from results.

ABBY Reader Engine:

ABBYY is a company who releases a series of OCR document reader applications. The engine provides an integrated end-to-end text extraction service. But it cost a little much to purchase an engine like that. It is more suitable for business.



The screenshot shows the ABBYY website homepage. At the top, there is a navigation bar with the ABBYY logo on the left and links for English, Sign in, Store, and Contact Us on the right. Below this, a secondary navigation bar includes links for Products & Solutions, Explore, Partners, Support & Services, and Company, along with a search icon. The main content area features a large hero image of three people (two women and one man) working together at a desk with a laptop. Overlaid on the left side of the hero image is the text: "ABBYY FineReader Engine", "The most comprehensive OCR SDK for software developers.", and "Integrate AI-powered OCR features into your applications." Below this text are two buttons: "Request trial" and "Contact us".

References

Microsoft Azure. (). *Form Recognizer – Automated Data Processing Systems*

<https://azure.microsoft.com/en-us/services/form-recognizer/>

Microsoft Docs. (2021). *IDs - Form Recognizer - Azure Applied AI Services*

<https://docs.microsoft.com/en-us/azure/cognitive-services/form-recognizer/concept-identification-cards>

Microsoft Power Platform. (2021). *Microsoft Power Automate Portal*

<https://us.flow.microsoft.com/en-us/>

Issagha BA. (2020). *How to automate form processing with Azure Form Recognizer* - Microsoft Industry Blogs - United Kingdom.

<https://cloudblogs.microsoft.com/industry-blog/en-gb/technetuk/2020/10/22/how-to-automate-forms-processing-with-azure-form-recognizer/>

Lee, M (2021). Pytesseract · PyPI. <https://pypi.org/project/pytesseract/>

Microsoft Azure. (2021). *Home - Form OCR Testing Tool (fott-2-1.azurewebsites.net)*

<https://fott-2-1.azurewebsites.net/>

Google. (2021). *Vision AI | Derive Image Insights via ML | Cloud Vision API (google.com)*

<https://cloud.google.com/vision/>

ABBYY. (2021). *The Digital Intelligence Company* <https://www.abbyy.com/>