# PureOCR

github.com/michael-s-teodor/pureocr

Mihnea S. Teodorescu & Moe Assaf
michael.s.teodor@gmail.com & mormoassaf@gmail.com

December 14, 2019

## 1 Problem description

A program is required that can learn to scan and process handwritten (or typed) text, recognising its characters and outputting them as plain text. Every ASCII character from 0x21 to 0x7E should be identifiable. This includes digits, lower and upper case letters and most punctuation marks, adding up to a total of 94 characters.

## 2 Problem analysis

We interpret the problem as follows: A character exists in a given area of an image if and only if that area contains a cluster of non-white pixels. These areas can be easily identified recursively by an algorithm which maps each candidate cluster and sends it to the recognition stage. This interpretation is motivated by the following statements:

- Identifying candidate areas first provides a significant improvement in time complexity. This way, text recognition will not have to be conducted across the whole image, but rather in certain areas only.

- The problem can now be further simplified by applying a noise-removing filter function first. This reduces the number of required recursive steps exponentially.

Having mapped each cluster of pixels, it is now time to recognise each character. For this, we have set up a neural network of 10,000 inputs (for segments of $100 \times 100$ pixels) and 94 outputs – one per character. The network will be trained in accordance with each font to be recognised.

## 3 Design

In order to achieve the desired results, we have structured our algorithm into several stages:

- Pixel-continuity-oriented preprocessing approach for identifying letters and characters: our program goes through the image line by line and when we encounter a pixel we call a recursive functions which maps the character by going all directions and comparing the threshold value. Then we set the anchors of the letter and store them.

- Line recognition: After mapping all the characters we use a sorting algorithm to sort them by their y-axis coordinates. After that, we compare their coordinates consequently being able to recognise identify lines then we sort them by their x-coordinates and now we have achieved the right order of all letters.

- Merging neighbouring characters: There are occasions where we have letters like the letter i which has a dot and a line. This is one of the issues that we had because the dot or the line should not be recognised as one letter. This applies for question marks, semicolons, etc. We managed to fix this after implementing the line recognition because at allows us to compare their coordinates and merge the together.

- Neural network: we have the written a code for the neural network and it is able to feed forward the data and produce an output. However, we have not implemented the training algorithm yet.

- Postprocessing: Ultimately, the user will be provided with the option to apply a dictionary-based postprocessing algorithm. This aims to correct misspelled words, as well as incorrectly-recognised characters – should there be any.

# 4 Program code

main.py

```
1  from preprocessor import PreProcessor
2  from recogniser import Recogniser
3  from network import Network
4
5
6  img = PreProcessor("samples/5.jpg")
7  recogniser = Recogniser()
8  #===================================#
9  img.processImage()
10 img.convertToBinary()
11 img.previewImg()
12 #===================================#
```

preprocessor.py

```
1
2  import cv2 #import the image processing
3
4  class PreProcessor():
5
6      def __init__(self,file):
7      #import a sample image also 0 is for grey scale and store its properties
8          self.file = file
9          self.original_img = cv2.imread(file,0)
10         self.processed_img = cv2.imread(file,0)
11         self.img = cv2.imread(file,0)
12
13         self.height, self.width = self.img.shape[0], self.img.shape[1]
14         self.threshold = -1
15
16         self.charAnchors = [] # [x_min,y_min,x_max,y_max]
17         self.charAnchorsRows = [] # [x_min,y_min,x_max,y_max]
18         self.data = [] #contains copies of charachters row/char
19
20     def processImage(self):
21         self.threshold = getThreshold(self.height,self.width,self.img)
22         self.charAnchors = findLetters(self.height,self.width,self.threshold,self.img)
```

```
23          self.charAnchors = bubbleSort(self.charAnchors,1)
24          self.charAnchorsRows = recogniseLines(self.charAnchors)
25          self.charAnchorsRows = sortLettersInOrder(self.charAnchorsRows)
26          self.charAnchorsRows = mergeNeighbours(self.charAnchorsRows)
27
28      def previewImg(self):
29          drawBoxesLines(self.charAnchorsRows,self.processed_img)
30          while True:
31              cv2.imshow("Preview (Press c to close)", self.processed_img)
32              key = cv2.waitKey(0) & 0xFF
33              # if the q key was pressed, break from the loop
34              if key == ord("c"):
35                  break
36              cv2.destroyAllWindows()
37
38      def convertToBinary(self):
39          if(self.threshold == -1):
40              self.threshold = getThreshold(self.height,self.width,self.img)
41          self.processed_img = convertToBinaryImg(self.height,self.width,self.threshold,self
                .processed_img)
42          return self.processed_img
43
44  #=================================================================================================
45
46  def getThreshold(height,width,img):
47      avg = 0
48      for i in range(height):
49          for j in range(width):
50              avg += img[i][j]
51      avg /= height*width
52      return 0.75*avg
53
54  def convertToBinaryImg(height,width,threshold,img):
55      for i in range(height):
56          for j in range(width):
57              if( img[i][j] > threshold):
58                  img[i][j] = 255
59              else:
60                  img[i][j] = 0
61      return img
62
63  def findLetters(height,width,threshold,img):
64      anchors = []
65      for i in range(height):
66          for j in range(width):
67              if( img[i][j] < threshold ):
68                  #map letter
69                  tempArr = [j,i,j,i]
70                  tempArr = mapLetter(j,i,tempArr,threshold,height,width,img)
71                  #store anchors
72                  anchors.append(tempArr)
73      return anchors
74
75  def mapLetter(x,y,tempArr,threshold,height,width,img):
76      if ( x < 0 or x >= width ) or ( y < 0 or y >= height ):
77          return tempArr
78
79      if( img[y][x] < threshold ):
```

```
80
81          #delete so that we dont encounter it again
82          img[y][x] = -1
83
84          #find min x and y and max x and y
85          if( x < tempArr[0] ):
86              tempArr[0] = x
87          if( y < tempArr[1] ):
88              tempArr[1] = y
89          if( x > tempArr[2] ):
90              tempArr[2] = x
91          if( y > tempArr[3] ):
92              tempArr[3] = y
93
94          for n in range(-1,2):
95              for m in range(-1,2):
96                  if not( n == 0 and m == 0 ):
97                      tempArr = mapLetter(x+n,y+m,tempArr,threshold,height,width,img)
98
99      return tempArr
100
101 def drawBoxesLines(anchorsRows,img):
102     #draw boxes around letters in rows with the same colour and alternate
103     for row in range(len(anchorsRows)):
104         img = drawBoxes(anchorsRows[row], row%2*100,img)
105
106 def drawBoxes(arr,colour,img):
107     for letter in range(len(arr)):
108         minX = arr[letter][0]
109         minY = arr[letter][1]
110         maxX = arr[letter][2]
111         maxY = arr[letter][3]
112         for i in range(minX,maxX+1):
113             img[minY][i] = colour
114             img[maxY][i] = colour
115         for j in range(minY,maxY+1):
116             img[j][minX] = colour
117             img[j][maxX] = colour
118     return img
119
120 def recogniseLines(charAnchors):
121     charAnchorsRows = []
122     row = 0
123     total = (charAnchors[0][1]+charAnchors[0][3])/2
124     n = 1
125     charAnchorsRows.append([])
126
127     deviation = 20 #NOTE: needs to be adjusted so that its adaptive
128
129     for i in range(len(charAnchors)):
130         avg = total/n
131
132         #calculate middle y
133         y = (charAnchors[i][1] + charAnchors[i][3])/2
134
135         if( y < avg + deviation and y > avg - deviation):
136             total += y
137             charAnchorsRows[row].append(charAnchors[i])
```

```
138            else:
139                row += 1
140                n = 0
141                total = y
142                charAnchorsRows.append([])
143                charAnchorsRows[row].append(charAnchors[i])
144
145            n+=1
146        return charAnchorsRows
147
148    def bubbleSort(arr,by):
149        length = len(arr)
150        for i in range(length-1):
151            for j in range(length-1):
152                if( arr[j][by] > arr[j+1][by] ):
153                    temp = arr[j+1]
154                    arr[j+1] = arr[j]
155                    arr[j] = temp
156        return arr
157
158    def sortLettersInOrder(charAnchorsRows):
159        for row in range(len(charAnchorsRows)):
160            #Sort by x
161            charAnchorsRows[row] = bubbleSort(charAnchorsRows[row], 0)
162        return charAnchorsRows
163
164    def mergeNeighbours(charAnchorsRows):
165
166        for row in range(len(charAnchorsRows)):
167            length = len(charAnchorsRows[row])-1
168            for char in range(length):
169                if(char + 1 < length ):
170                    c1 = charAnchorsRows[row][char]
171                    c2 = charAnchorsRows[row][char+1]
172                    #check if theyre on top of each other
173                    if( c1[1] >= c2[3] or c1[3] <= c2[1] ):
174                        #check if their x's are intersecting
175                        if( c2[0] <= c1[2] ):
176
177                            charAnchorsRows[row][char][0] = min(c1[0],c2[0])
178                            charAnchorsRows[row][char][1] = min(c1[1],c2[1])
179                            charAnchorsRows[row][char][2] = max(c1[2],c2[2])
180                            charAnchorsRows[row][char][3] = max(c1[3],c2[3])
181
182                            charAnchorsRows[row].remove(c2)
183                            length -=1
184
185        return charAnchorsRows
```
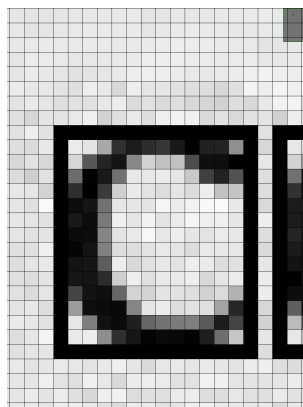
network.py

```
1   #### Libraries
2   # Third-party library
3   import numpy as np
4
5   class Network(object):
6
7       def __init__(self, sizes):
8           # sizes = number of neurons per layer [input, hidden, ..., hidden, output]
9           self.numLayers = len(sizes)
10          self.sizes = sizes
11
12          # Each neuron has one bias
13          # layer -> neuron -> random bias
14          self.biases = [np.random.randn(y, 1) for y in sizes[1:]]
15
16          # Each neuron from layer n has sizes[n-1] weights
17          # layer -> neuron -> random weights
18          self.weights = [np.random.randn(y, x) for x, y in zip(sizes[:-1], sizes[1:])]
19
20      def feedForward(self, a):
21          # This function returns the output of the network when a is applied
22          for b, w in zip(self.biases, self.weights):
23              a = sigmoid(np.dot(w,a) + b)
24          return a
25
26  #### Helper functions
27  def sigmoid(z):
28      return 1.0/(1.0 + np.exp(-z))
29
30  def sigmoid_derivative(z):
31      return np.exp(z) / (np.exp(z) + 1)**2
```
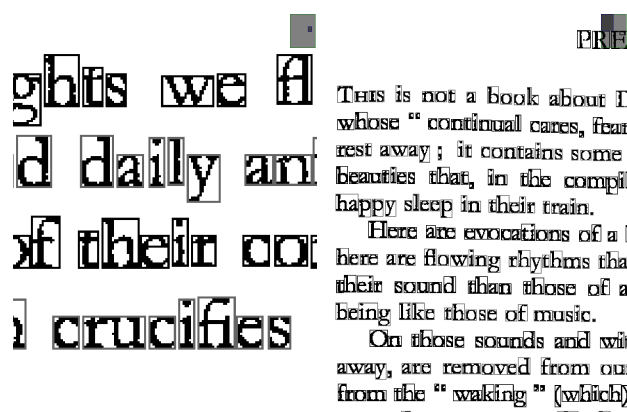
```
1   from network import Network
2
3   net = Network([10000,10,10,94])
4
5   class Recogniser():
6
7       def __init__(self):
8           pass
9
10      def retrieveChar(self,img):
11          output = net.feedForward(img)
12
13          max = 0
14          char = 0
15          for i in range(94):
16              if (max < output[i]):
17                  max = output[i]
18                  char = i
19          return chr(char + 33 )
20
21      def returnTexts(self,data):
22          texts = []
23          for img in range(len(data)):
24              texts.append(self.retrieveChar(img))
25          return texts
26
27
28      # Compute the output of the neural network
29
30      # Determine the most likely outcome and return its ASCII value
```

# 5 Test results

- This image shows how the threshold value is responsible for determining whether a pixel is part of a character or not.



- Here the image shows how the outcome turned out to be. As it shows in the picture, we draw a box around each character and it seems to be working perfectly. The image also shows how we alternate between grey and black to show how lines are being recognised. It also shows how the i's are merged together and this also applies for question marks, semicolon, question marks, etc.



# 6 Evaluation

- Threshold value: the way we determine the threshold value in our code depends on the average intensity of all pixels and we take a percentage out of that. However, in certain circumstances this can be quite deterministic and ambiguous. Thus, we need to implement a new way of find the threshold value and possibly make it dynamic in some parts of the image where there might be some shades which could cause the program to fail.

- Merged letters: there will be occasions where the letters will be connected together. Consequently, the preprocessor program will mark it as one single character making it hard for the neural network to understand causing an error. When the training part is implemented, the neural network will be trained enough to recognise merged letters and apply a different approach to those specific merged letters.

- Noise: there happens to be random chunks of letters which are not connected to the actual letter itself causing them to be recognised as characters in the preprocessor. We need to find a way of recognising them as errors and simply return nothing.

- Shades: shades can really mess up the threshold value and even cause the recursive function to reach its maximum depth due to python having a limit. We need to find a good procedure of either removing those shades from the image completely or making the threshold dynamic. Otherwise, the program would recognise shades as letters or just crash due to the recursion depth limit.

- Orientation and dimensions: we have not implemented any sort of procedure to the change the dimensions of the image so that texts can be aligned properly for the preprocessor to analyse.