



UNIVERSIDADE FEDERAL DE SERGIPE
CENTRO DE CIÊNCIAS EXATAS E TECNOLOGIA
DEPARTAMENTO DE COMPUTAÇÃO

PreOCR - Trabalho de Processamento de Imagens T01 2023.2

Professora: Beatriz Trinchão Andrade

Grupo 13 - Everton Santos de Andrade Júnior



São Cristóvão – Sergipe

2024

Sumário

1	Introdução e Resultados	2
1.1	Resultados	2
2	Estrutura do Projeto	11
3	Requisitos	12
4	Métodos e Implementações	13
5	Resultados	14
6	Conclusão	15
	Referências	16

1

Introdução e Resultados

Neste trabalho, desenvolvemos um programa capaz de processar imagens binárias no formato PBM ASCII (PGM tipo P1), contendo texto dos tipos e tamanhos de fontes, variando de 10 à 40. Conseguimos determinar o básico do trabalho que são o número de linhas e palavras no texto e os retângulos dessas palavras, gerando uma imagem de saída com esses retângulos. Além disso, nosso grupo implementou mais funcionalidades, como a detecção de colunas e blocos do texto, utilizando o conceito de distância alinhada. Incluindo a *bouding box* (bbox) desses blocos e linhas que indicam as colunas.

Adicionalmente, para ilustrar nossos resultados de forma mais interativa, geramos uma série de imagens intermediárias que mostram o processo de detecção de blocos, colunas, palavras, e linhas destacando as regiões por um retângulo de diferentes cores. O formato dessas imagens é P3 (ascii RGB). Essas imagens foram usadas em sequência de frames para gerar um vídeo de saída usando o programa de linha de comando “ffmpeg”.

A nossa abordagem é dinâmica, baseado na altura das palavras mudamos os parâmetros enquanto o programa roda. Desse modo é possível identificar estruturas de texto em diferentes alinhamentos, como justificado, esquerda, centro e direita, além de lidar com diferentes tamanhos e estilos de fonte, incluindo Comic Sans, Impact, Cascadia Code, Arial e Times New Roman. Um exemplo de vídeo gerado pode ser encontrado em <https://youtu.be/uA45GeodGss?si=ZOgsA2av7EKZeG69>.

Na [Figura 3](#)

1.1 Resultados

Inicialmente, o algoritmo lê a imagem de entrada e a pré-processa aplicando filtro de medianos (nome da função é no código `median_blur`) e invertendo as cores. Essas etapas visam retirar possíveis ruídos do tipo sal de pimenta. O tamanho do filtro mediano é 3 de altura e largura

Figura 1 – Times New Roman em itálico, com 4 colunas, 6 blocos de texto.



Fonte: Autor.

já que foi confirmado que os ruídos não passam do tamanho de 1 pixel. A inversão da imagem é feita para aplicar algoritmos morfológicos, como a dilatação, para melhorar a conectividade dos componentes de texto e facilitar o processamento subsequente. Dependendo dos parâmetros especificados, o algoritmo pode aplicar operações adicionais como fechamento (closing) ou abertura (opening) para refinar ainda mais as regiões de texto. A ideia é utilizar um elemento estruturante, conhecido como SE (do inglês, *Structuring Element*), que connecte as letras; isso pode ser feito com um SE que seja horizontal. Foi testado com diferentes SEs, como fomarto circular, de cruz, vertical, mas o melhor é o horizontal. E faz sentido pois as letras vem logo à direta da anterior, então para conectalas, precisamos esticar na horizontal para "grudar" uma na

outra. Veja como é da uma "esticada" para horzinhotal ao aplicar um SE [000] [111] [000].

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

Na [Figura 3](#)

Aplicação de operações de dilatação horizontal para melhorar a conectividade entre caracteres adjacentes. Após as operações morfológicas, indetificamos as palavras individuais através componentes conectados dentro da imagem. Cada palavra é contida em uma bbox, que é desenhada retângulos ao redor delas na imagem de saída. Componentes conectados é um conceito de grafos que pode ser aplicado numa imagem binarias, pixels brancos são nós. Existem tipos de conectividade. Se consideramos os vizinhos de um pixel como apenas esquerda, direita, cima e baixo, então temos 4-conectividade ([GONZALEZ; WOODS, 2008, 2.5.2 Adjacency, Connectivity, Regions, and Boundaries](#)). Se adicionarmos as diagonais, agora temos 8-conectividade. O Algortimo para encontrar esses componentes conectados se basea numa busca profunda, marcado os visitados, toda vez que na imagem for não zero fazemos a busca até não encontrar mais caminhos, considerando as digonais. Se não encontramos mais caminhos quer dizer que esse é um compoenente. Repetimos até encontrar todos atualizando o ponto mais a em esquerda inferior e o maiis à direita superior do componente para encontrar o bbox da palavra. O código ilustra esse processo.

Também fazemos o agrupamento dessas palavras em blocos com base em sua proximidade espacial, possibilitando uma análise melhor da organização do texto podemos dizendo que cada bloco é um paragrafo.

Opcionalmente, aplicação de operações de fechamento e abertura para refinar as regiões de texto, dependendo dos parâmetros definidos.

Há uma dilatação adicional da imagem, se necessário, com base na altura média dos componentes de texto já filtrados.

Ao longo do processo, o algoritmo incorpora várias otimizações e ajustes de parâmetros para se adaptar a diferentes tipos de imagens de entrada e layouts de texto. Por exemplo, parâmetros como tamanho do kernel e contagens de iteração são ajustados dinamicamente com base nas características da imagem de entrada, como tamanho do texto e níveis de ruído.

Filtragem dos componentes de texto para eliminar bbox de pontuação e manter apenas bboxes de palavras. Agrupamento de palavras em blocos com base em sua proximidade espacial. Desenho de caixas delimitadoras ao redor das palavras e blocos identificados na imagem original. Geração de um vídeo para visualizar interativamente as etapas do algoritmo. Opcionalmente, escrita de imagens separadas para cada palavra identificada, para uso posterior em tentativas de reconhecimento óptico de caracteres (OCR). Escrita da imagem final com caixas delimitadoras de palavras e blocos para análise adicional.

Além disso, o algoritmo fornece opções para gerar imagens intermediárias e vídeos para visualizar as etapas de processamento, auxiliando na depuração e compreensão do comportamento do algoritmo. Essas visualizações melhoram a transparência e a interpretabilidade da operação do algoritmo.

Em resumo, o algoritmo implementado demonstra uma abordagem abrangente para o reconhecimento de texto em imagens binárias, utilizando uma combinação de técnicas de processamento de imagem, regras heurísticas e parâmetros adaptativos. Ao refinar iterativamente regiões de texto e analisar seus relacionamentos espaciais, o algoritmo alcança uma detecção precisa de palavras, linhas, colunas e blocos, estabelecendo as bases para tarefas de reconhecimento de texto mais avançadas, como reconhecimento óptico de caracteres (OCR).

Ela é especialmente útil para preencher pequenos espaços ou quebrar conexões entre objetos. No [Código 1](#) a função ‘understandable_dilate’ implementa uma dilatação morfológica de forma compreensível, percorrendo cada pixel da imagem. A variável ‘kernel’ representa o elemento estruturante, que define a forma e o tamanho da dilatação. A imagem é percorrida pixel a pixel, e para cada pixel, verifica-se se pelo menos um dos pixels na vizinhança definida pelo kernel é branco (valor 255). Se pelo menos um dos pixels na vizinhança é branco, o pixel central é definido como branco (255) na imagem resultante. Caso contrário, o pixel central é definido como preto (0) na imagem resultante. O parâmetro ‘iterations’ indica quantas vezes a dilatação deve ser aplicada à imagem. Quanto maior o número de iterações, mais ampliado será o objeto na imagem.

Já no [Código 2](#) ‘fast_dilate2’:** função implementa uma dilatação morfológica mais rápida, aproveitando as operações de matriz do NumPy. A imagem é expandida com um preenchimento adequado para evitar problemas de borda durante a aplicação do kernel. O kernel é então aplicado à imagem usando a função ‘np.maximum’, que calcula o máximo elemento a elemento entre a imagem e o subconjunto da imagem definido pelo kernel. Isso efetivamente realiza uma dilatação, onde o pixel central de uma vizinhança é definido como o valor máximo dessa vizinhança. Ambas as implementações alcançam o mesmo resultado de dilatação morfológica, mas a segunda função é mais eficiente computacionalmente devido ao uso de operações vetorizadas do NumPy. Isso resulta em uma execução mais rápida, especialmente em imagens grandes.

““

Código 1 – .

```
1 def understandable_dilate(image, kernel):
2     result = np.zeros(image.shape)
3     height = image.shape[0]
4     width = image.shape[1]
5     kernel_height = kernel.shape[1]
6     kernel_width = kernel.shape[0]
7     kernel_width_delta = kernel_width // 2
8     kernel_height_delta = kernel_height // 2
9     for y in range(height):
10         for x in range(width):
11             all_good = False
12             for j in range(kernel_height):
13                 for i in range(kernel_width):
14                     i_offset = i - kernel_width_delta
15                     j_offset = j - kernel_height_delta
16                     color = 0
17                     if (
18                         (x + i_offset) >= 0
19                         and (x + i_offset) < width
20                         and y + j_offset >= 0
21                         and y + j_offset < height
22                     ):
23                         color = image[y + j_offset, x + i_offset]
24                         kcolor = kernel[j, i]
25                         if int(kcolor) * int(color):
26                             all_good = True
27                             break
28             if all_good:
29                 break
30         if all_good:
31             result[y, x] = 255
32         else:
33             result[y, x] = 0
34     return result
```

Figura 2 – Cacadia Code em negrito com tamanho 10, 2 colunas e 5 blocos de texto. A entrada possui baixa resolução e ruídos.



Fonte: Autor.

Figura 3 – Times New Roman em itálico, com 4 colunas, 6 blocos de texto.

<p><i>>Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud</i></p> <p><i>Exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.</i></p>	<p><i>Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi</i></p> <p><i>Beatae vitae dicta sunt explicabo. Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia magni dolores eos qui ratione voluptatem sequi nesciunt. Neque porro</i></p>	<p><i>quisquam est, qui dolorem ipsum quia dolor sit amet, consectetur adipiscing velit, sed quia non numquam eius modi tempora incidunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim ad minima veniam, quis nostrum exercitatione m ullam corporis</i></p>	<p><i>Suscipit laboriosam, nisi ut aliquid ex ea commodi consequatur? Quis autem vel eum iure reprehenderit qui in ea voluptate velit esse quam nihil molestiae consequatur, vel illum qui dolorem eum fugiat quo voluptas nulla pariatur?</i></p>
---	--	---	--

Fonte: Autor.

Figura 4 – Cacadia Code em negrito com tamanho 10, 2 colunas e 5 blocos de texto. A entrada possui baixa resolução e ruídos.



Fonte: Autor.

Código 2 – .

```
1 def fast_dilate2(image, kernel):
2     global counter
3     height, width = image.shape
4     kernel_height, kernel_width = kernel.shape
5
6     kernel_width_delta = kernel_width // 2
7     kernel_height_delta = kernel_height // 2
8
9     # We pad by the kernel delta, top, bottom, left and right
10    padded_image = pad(
11        image,
12        kernel_height_delta,
13        kernel_height_delta,
14        kernel_width_delta,
15        kernel_width_delta,
16    )
17
18    dilated = np.zeros(image.shape, dtype=np.uint8)
19    for j in range(kernel_height):
20        for i in range(kernel_width):
21            if kernel[j, i] == 1:
22                shifted_sub_image = padded_image[j : j + height, i :
23                    i + width]
24                dilated = np.maximum(
25                    dilated, shifted_sub_image
26                )
27    return dilated
```

2

Estrutura do Projeto

3

Requisitos

```
python3 src/main.py assets/grupo_13_arial_esquerda_tamanho_16_colunas_2_blocos_4_7
```

4

Métodos e Implementações

Seguindo as formulas de dilatação em ([GONZALEZ; WOODS, 2008](#), capítulo 9)

5

Resultados

implementação <<https://www.youtube.com/watch?v=uA45GeodGss>> Descrição da implementação do programa. Destaque para soluções desenvolvidas para problemas específicos encontrados durante o desenvolvimento. Resultados:

Apresentação dos resultados obtidos. Inclusão de exemplos de imagens de entrada e saída. Discussão sobre a eficácia do programa e eventuais limitações.

6

Conclusão

Sumarização dos principais resultados e contribuições do trabalho. Reflexão sobre o aprendizado durante o desenvolvimento do programa. Sugestões para trabalhos futuros ou melhorias no programa.

Referências

GONZALEZ, R. C.; WOODS, R. E. *Digital Image Processing*. Upper Saddle River, NJ: Pearson Prentice Hall, 2008. Citado 2 vezes nas páginas 4 e 13.