

基於OpenCV的表格文本內容提取

小白 AI算法與圖像處理 2022-01-25 17:00

點擊下方“ AI算法與圖像處理”，一起進步！

重磅乾貨，第一時間送達



AI算法與圖像處理

考研逆襲985，非科班跨行AI，目前從事計算機視覺的工業和商業相關應用的工作。分...
248篇原創內容

公眾號

小伙伴們可能會覺得從圖像中提取文本是一件很麻煩的事情，尤其是需要提取大量文本時。PyTesseract是一種光學字符識別（OCR），該庫提供了文本圖像。

PyTesseract確實有一定的效果，用PyTesseract來檢測短文本時，結果相當不錯。但是，當我們用它來檢測表格中的文本時，算法執行失敗。

LAPORAN HARIAN OTG, ODP, PDP, Dan KONFIRMASI																			
PROVINSI SULAWESI TENGAH																			
TAHUN 2020																			
Tanggal, 18 Juni 2020																			
No	Kabupaten/Kota	OTG					ODP					PDP					Konf (+)	Sembuh	Meninggal
		Kasus Baru	Kasus Lama	Selesai Masa	(KL + KB) - Selesai Masa	Bukan Covid-19	Kasus Baru	Kasus Lama	Selesai Masa	(KL + KB) - Selesai Masa	Bukan Covid-19	Kasus Baru	Kasus Lama	Selesai Masa PDP	(KL + KB) - Selesai Masa	Bukan Covid-19			
1	Banggai	0	0	0	0	0	0	4	0	4	49	0	0	0	0	10	10	4	0
2	Banggai Kepulauan	0	1	0	0	28	0	0	0	0	9	0	0	0	0	0	1	1	0
3	Banggai Laut	0	13	0	13	16	0	0	0	0	11	0	0	0	0	0	1	0	0
4	Buol	7	9	0	16	0	0	33	9	24	70	0	0	0	0	6	57	50	0
5	Donggala	0	70	0	70	38	0	0	0	0	32	1	0	0	1	7	1	0	0
6	Morowali	0	84	0	84	249	0	0	0	0	33	1	4	0	5	16	12	10	0
7	Morowali Utara	0	21	3	18	249	0	1	1	0	14	0	5	0	5	20	14	13	1
8	Parigi Moutong	0	87	0	87	5	0	7	1	6	106	0	0	0	0	6	1	0	0
9	Poso	1	49	11	39	234	0	2	0	2	45	0	3	0	3	19	13	12	0
10	Sigi	0	15	1	14	9	0	1	0	1	11	0	1	1	0	11	5	3	0
11	Tojo Una-una	0	7	0	7	18	0	2	0	2	31	0	1	0	1	8	0	0	0
12	Toli-toli	20	23	0	43	27	2	202	43	161	77	0	7	5	2	15	18	12	0
13	Kota Palu	0	247	203	44	431	5	32	0	37	167	0	1	1	0	54	39	25	3
	Provinsi	28	577	218	387	1304	7	284	54	237	655	2	22	7	17	172	172	130	4

圖1. 直接使用PyTesseract檢測表中的文本

圖1描繪了文本檢測結果，綠色框包圍了檢測到的單詞。可以看出算法對於大部分文本都無法檢測，尤其是數字。而這些數字卻是展示了每日COVID-19病例的相關信息。那麼，如何提取這些信息？

簡介

在编写算法时，我们通常应该以我们人类理解问题的方式来编写算法。这样，我们可以轻松地将想法转化为算法。

当我们阅读表格时，首先注意到的就是单元格。一个单元格使用边框（线）与另一个单元格分开，边框可以是垂直的也可以是水平的。识别单元格后，我们继续阅读其中的信息。将其转换为算法，您可以将过程分为三个过程，即**单元格检测**、**区域（ROI）选择**和**文本提取**。

在执行每个任务之前，让我们先导入必要内容

```
1 import cv2 as cv
2 import numpy as np
3 filename = 'filename.png'
4 img = cv.imread(cv.samples.findFile(filename))
5 cImage = np.copy(img) #image to draw lines
6 cv.imshow("image", img) #name the window as "image"
7 cv.waitKey(0)
8 cv.destroyAllWindows() #close the window
```

单元格检测

查找表格中的水平线和垂直线可能是最容易开始的。有多种检测线的方法，这里我们采用OpenCV库中的Hough Line Transform。

在应用霍夫线变换之前，需要进行一些预处理。第一是将存在的RGB图像转换为灰度图像。因为灰度图像对于Canny边缘检测而言非常重要。

```
1 gray = cv.cvtColor(img, cv.COLOR_BGR2GRAY)
2 cv.imshow("gray", gray)
3 cv.waitKey(0)
4 cv.destroyAllWindows()
5 canny = cv.Canny(gray, 50, 150)
6 cv.imshow("canny", canny)
```

```
7 cv.waitKey(0)
8 cv.destroyAllWindows("canny")
```

下面的两幅图分别显示了灰度图像和Canny图像。

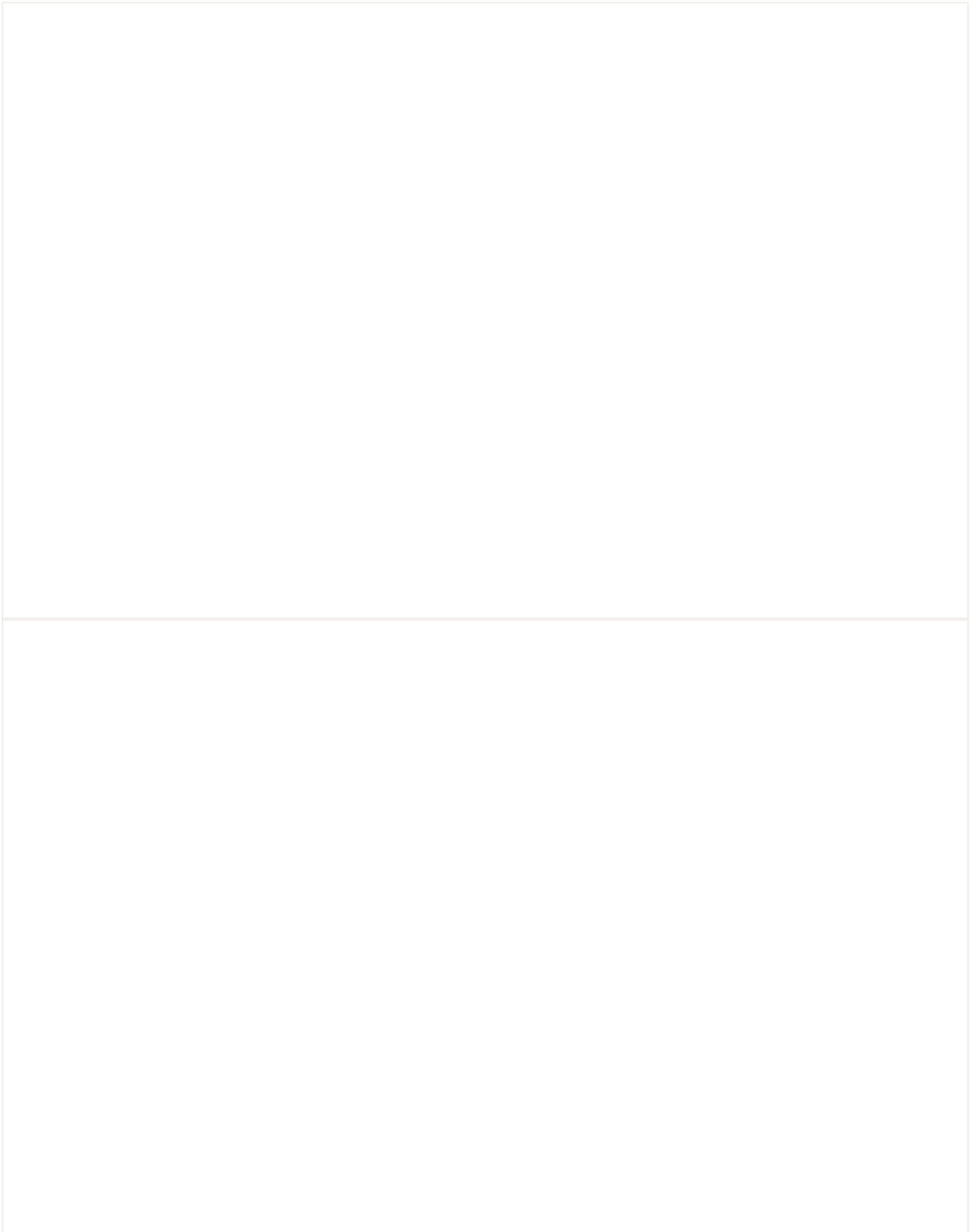


图2. 灰度和Canny图像

霍夫线变换

在OpenCV中，此算法有两种类型，即标准霍夫线变换和概率霍夫线变换。标准变换为我们提供直线方程，因此我们无法得知直线的起点和终点。概率变换将为我们提供线列表，即直线起点与终点的坐标值列表。我们优先选用的是概率变化。



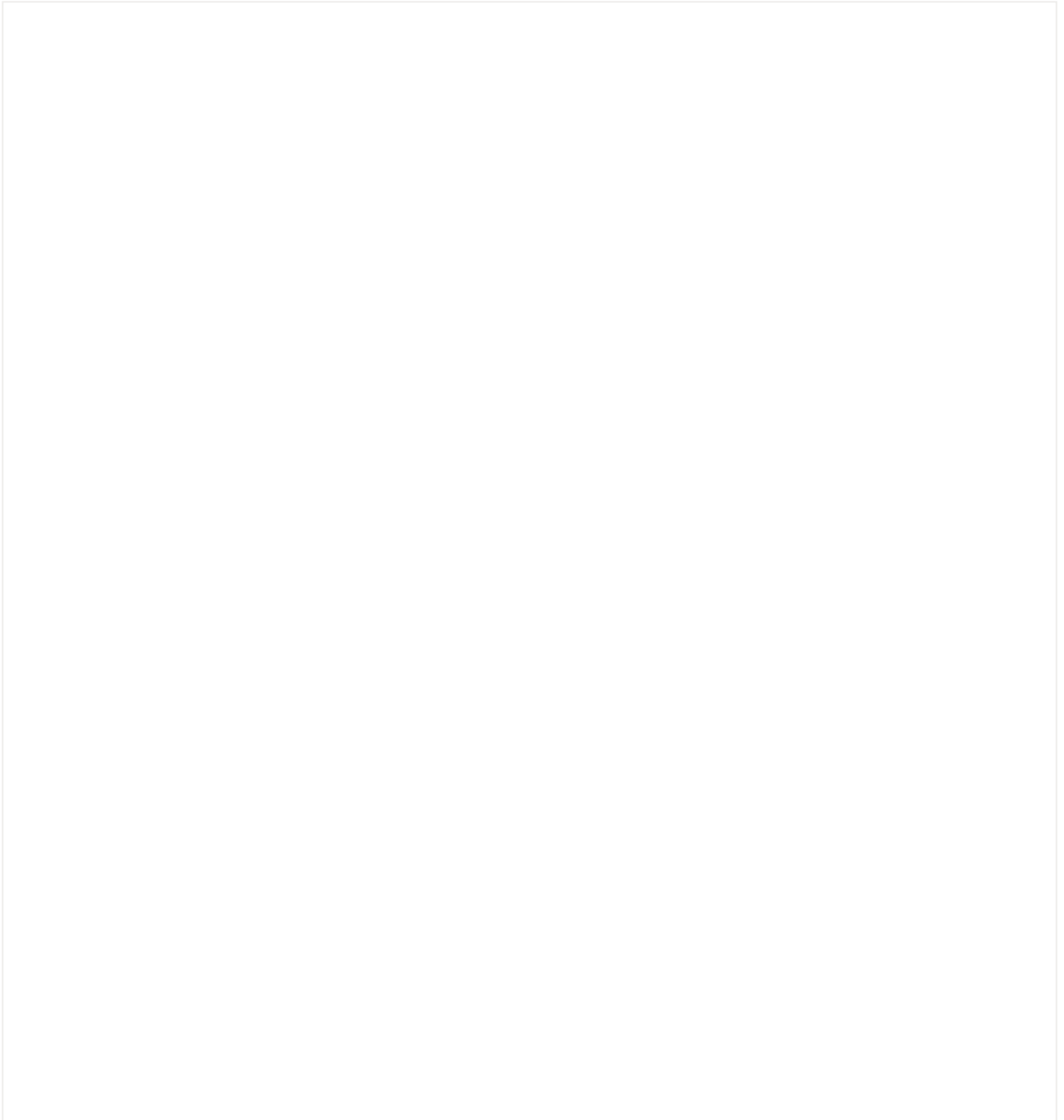


图3. 霍夫线变换结果示例（来源：OpenCV）

对于HoughLinesP函数，有如下几个输入参数：

- **image** -8位单通道二进制源图像。该图像可以通过该功能进行修改。
- **rho** —累加器的距离分辨率，以像素为单位。
- **theta** —弧度的累加器角度分辨率。
- **threshold**-累加器阈值参数。仅返回那些获得足够投票的行
- **line** — 线的输出向量。这里设置为无，该值保存到linesP

- **minLineLength** —最小行长。短于此的线段将被拒绝。
- **maxLineGap** —同一线上的点之间允许链接的最大间隙。

```

1 # cv.HoughLinesP(image, rho, theta, threshold[, lines[, minLineLength[, maxLine
2 rho = 1
3 theta = np.pi/180
4 threshold = 50
5 minLinLength = 350
6 maxLineGap = 6
7 linesP = cv.HoughLinesP(canny, rho , theta, threshold, None, minLinLength, maxL

```

为了区分水平线和垂直线，我们定义了一个函数并根据该函数的返回值添加列表。

```

1 def is_vertical(line):
2     return line[0]==line[2]
3 def is_horizontal(line):
4     return line[1]==line[3]
5 horizontal_lines = []
6 vertical_lines = []
7
8 if linesP is not None:
9     for i in range(0, len(linesP)):
10         l = linesP[i][0]
11         if (is_vertical(l)):
12             vertical_lines.append(l)
13
14         elif (is_horizontal(l)):
15             horizontal_lines.append(l)
16 for i, line in enumerate(horizontal_lines):
17     cv.line(cImage, (line[0], line[1]), (line[2], line[3]), (0,255,0), 3, cv.L
18
19 for i, line in enumerate(vertical_lines):
20     cv.line(cImage, (line[0], line[1]), (line[2], line[3]), (0,0,255), 3, cv.L
21
22 cv.imshow("with_line", cImage)
23 cv.waitKey(0)
24 cv.destroyWindow("with_line") #close the window
25

```

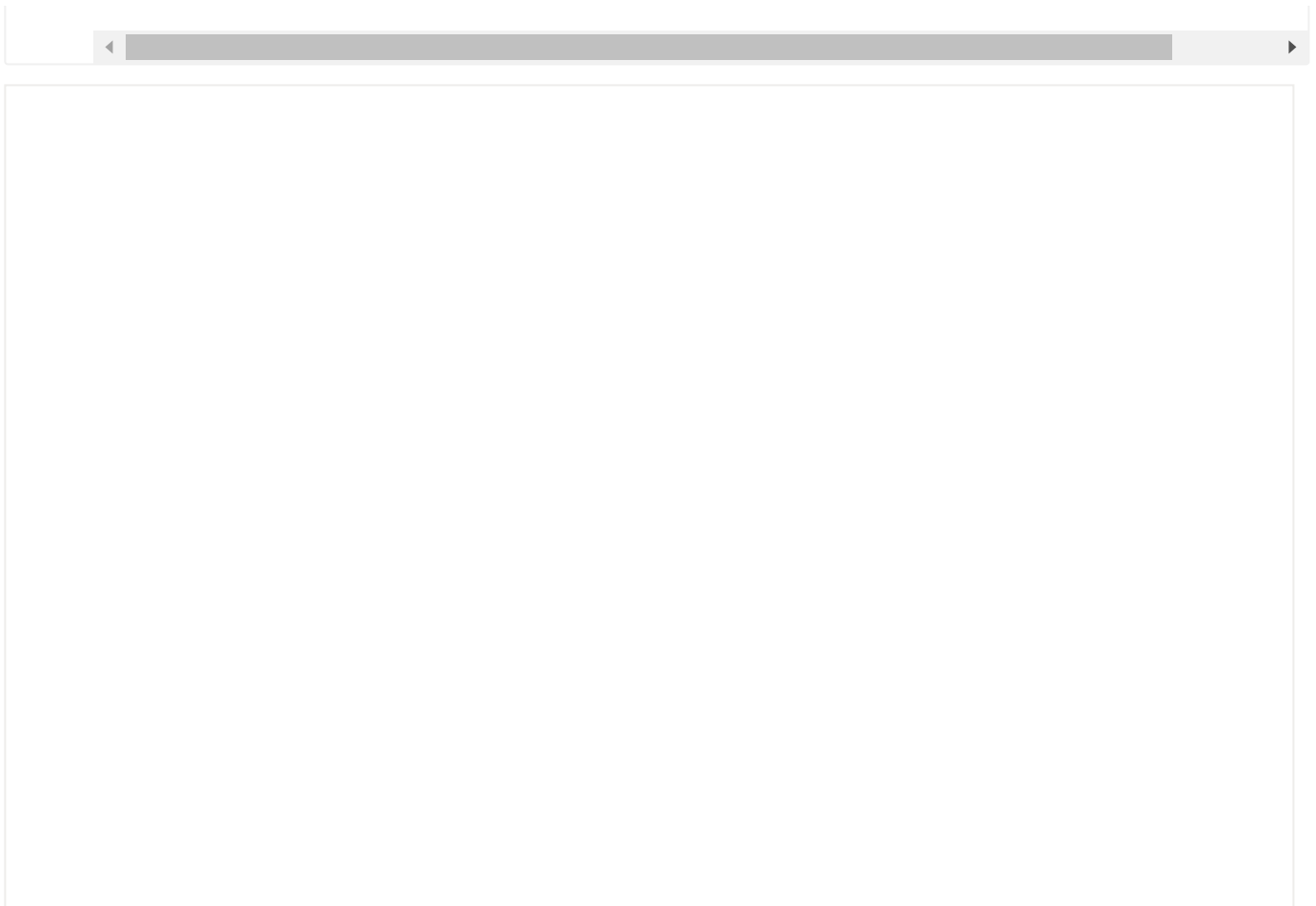


图4. 霍夫线变换结果—没有重叠滤波器

重叠滤波器

检测到的线如上图所示。但是，霍夫线变换结果中有一些重叠的线。较粗的线由多个相同位置，长度不同的线组成。为了消除此重叠线，我们定义了一个重叠过滤器。

最初，基于分类索引对线进行分类，水平线的 y_1 和垂直线的 x_1 。如果下一行的间隔小于一定距离，则将其视为与上一行相同的行。

```
1 def overlapping_filter(lines, sorting_index):
2     filtered_lines = []
3
4     lines = sorted(lines, key=lambda lines: lines[sorting_index])
5     separation = 5
6     for i in range(len(lines)):
7         l_curr = lines[i]
8         if(i>0):
9             l_prev = lines[i-1]
10            if ( (l_curr[sorting_index] - l_prev[sorting_index]) > separati
11                filtered_lines.append(l_curr)
12            else:
```

```
13         filtered_lines.append(l_curr)
14
15     return filtered_lines
```

实现重叠滤镜并在图像上添加文本，现在代码应如下所示：

```
1 horizontal_lines = []
2 vertical_lines = []
3
4 if linesP is not None:
5     for i in range(0, len(linesP)):
6         l = linesP[i][0]
7         if (is_vertical(l)):
8             vertical_lines.append(l)
9
10        elif (is_horizontal(l)):
11            horizontal_lines.append(l)
12    horizontal_lines = overlapping_filter(horizontal_lines, 1)
13    vertical_lines = overlapping_filter(vertical_lines, 0)
14    for i, line in enumerate(horizontal_lines):
15        cv.line(cImage, (line[0], line[1]), (line[2], line[3]), (0,255,0), 3, cv.LINE_AA)
16        cv.putText(cImage, str(i) + "h", (line[0] + 5, line[1]), cv.FONT_HERSHEY_SIMPLEX, 1, (0,0,255))
17    for i, line in enumerate(vertical_lines):
18        cv.line(cImage, (line[0], line[1]), (line[2], line[3]), (0,0,255), 3, cv.LINE_AA)
19        cv.putText(cImage, str(i) + "v", (line[0], line[1] + 5), cv.FONT_HERSHEY_SIMPLEX, 1, (0,255,0))
20
21 cv.imshow("with_line", cImage)
22 cv.waitKey(0)
23 cv.destroyAllWindows() #close the window
24
```

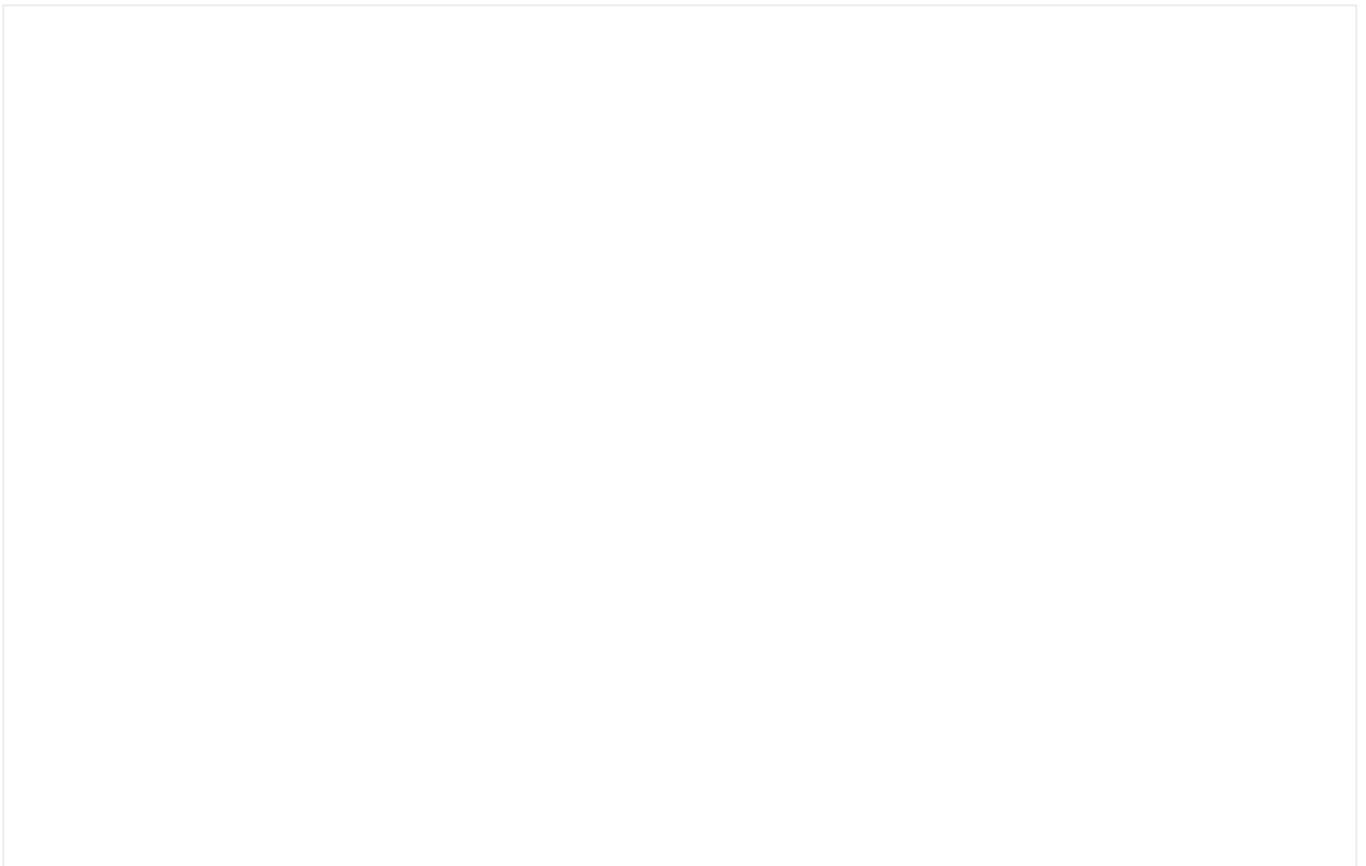



图5. 霍夫线变换结果一带重叠滤波器

有了这个代码，就不会提取出重叠的行。此外，我们还将在图像中写入水平和垂直线的索引，这将有利于ROI的选择。

ROI选择

首先，我们需要定义列数和行数。这里我们只对第二行第十四行以及所有列中的数据感兴趣。对于列，我们定义了一个名为关键字的列表，将其用于字典关键字。

```
1  ## set keywords
2  keywords = ['no', 'kabupaten', 'kb_otg', 'kl_otg', 'sm_otg', 'ks_otg', 'not_cv
3              'kb_odp', 'kl_odp', 'sm_odp', 'ks_odp', 'not_cvd_odp', 'death_odp
4              'kb_pdp', 'kl_pdp', 'sm_pdp', 'ks_pdp', 'not_cvd_pdp', 'death_pdp
5              'positif', 'sembuh', 'meninggal']
6
7  dict_kabupaten = {}
8      for keyword in keywords:
9          dict_kabupaten[keyword] = []
10
11  ## set counter for image indexing
12  counter = 0
13
```

```

14  ## set line index
15  first_line_index = 1
16  last_line_index = 14
17

```

然后，要选择ROI，我们定义了一个函数，该函数将图像（水平线和垂直线都作为输入）以及线索索引作为边框。此函数返回裁剪的图像及其在图像全局坐标中的位置和大小

```

1  def get_cropped_image(image, x, y, w, h):
2      cropped_image = image[ y:y+h , x:x+w ]
3      return cropped_image
4  def get_ROI(image, horizontal, vertical, left_line_index, right_line_index, top_line_index, bottom_line_index):
5      x1 = vertical[left_line_index][2] + offset
6      y1 = horizontal[top_line_index][3] + offset
7      x2 = vertical[right_line_index][2] - offset
8      y2 = horizontal[bottom_line_index][3] - offset
9
10     w = x2 - x1
11     h = y2 - y1
12
13     cropped_image = get_cropped_image(image, x1, y1, w, h)
14
15     return cropped_image, (x1, y1, w, h)

```

裁剪的图像将用于下一个任务，即文本提取。返回的第二个参数将用于绘制ROI的边界框

文字提取

现在，我们定义了ROI功能。我们可以继续提取结果。我们可以通过遍历单元格来读取列中的所有数据。列数由关键字的长度指定，而行数则由定义。

首先，让我们定义一个函数来绘制文本和周围的框，并定义另一个函数来提取文本。

```

1  import pytesseract
2  pytesseract.pytesseract.tesseract_cmd = r'C:\Program Files (x86)\Tesseract-OCR\
3  def draw_text(src, x, y, w, h, text):

```

```

4     cFrame = np.copy(src)
5     cv.rectangle(cFrame, (x, y), (x+w, y+h), (255, 0, 0), 2)
6     cv.putText(cFrame, "text: " + text, (50, 50), cv.FONT_HERSHEY_SIMPLEX, 2,
7
8     return cFrame
9 def detect(cropped_frame, is_number = False):
10     if (is_number):
11         text = pytesseract.image_to_string(cropped_frame,
12                                             config = '-c tessedit_char_whitelist
13     else:
14         text = pytesseract.image_to_string(cropped_frame, config='--psm 10')
15
16     return text

```

将图像转换为黑白以获得更好的效果，让我们开始迭代！

```

1 counter = 0
2 print("Start detecting text...")
3 (thresh, bw) = cv.threshold(gray, 100, 255, cv.THRESH_BINARY)
4 for i in range(first_line_index, last_line_index):
5     for j, keyword in enumerate(keywords):
6         counter += 1
7
8         left_line_index = j
9         right_line_index = j+1
10        top_line_index = i
11        bottom_line_index = i+1
12
13        cropped_image, (x,y,w,h) = get_ROI(bw, horizontal, vertical, left_line
14
15        if (keywords[j]=='kabupaten'):
16            text = detect(cropped_image)
17            dict_kabupaten[keyword].append(text)
18
19        else:
20            text = detect(cropped_image, is_number=True)
21            dict_kabupaten[keyword].append(text)
22        image_with_text = draw_text(img, x, y, w, h, text)

```

问题解决

这是文本提取的结果！我们只选择了最后三列，因为它对某些文本给出了奇怪的结果，其余的很好，所以我不显示它。



图6. 检测到的文本一版本1

一些数字被检测为随机文本，即39个数据中的5个。这是由于最后三列与其余列不同。文本为白色时背景为黑色，会以某种方式影响文本提取的性能。



图7. 二进制图像

为了解决这个问题，让我们倒数最后三列。

```
1 def invert_area(image, x, y, w, h, display=False):
2     ones = np.copy(image)
3     ones = 1
4
5     image[ y:y+h , x:x+w ] = ones*255 - image[ y:y+h , x:x+w ]
6
7     if (display):
8         cv.imshow("inverted", image)
9         cv.waitKey(0)
10        cv.destroyAllWindows()
11    return image
12 left_line_index = 17
13 right_line_index = 20
14 top_line_index = 0
15 bottom_line_index = -1
16
17 cropped_image, (x, y, w, h) = get_ROI(img, horizontal, vertical, left_line_index, right_line_index, top_line_index, bottom_line_index)
18 gray = get_grayscale(cropped_image)
19 bw = get_binary(gray)
20 bw = invert_area(bw, x, y, w, h, display=True)
```

结果如下所示。

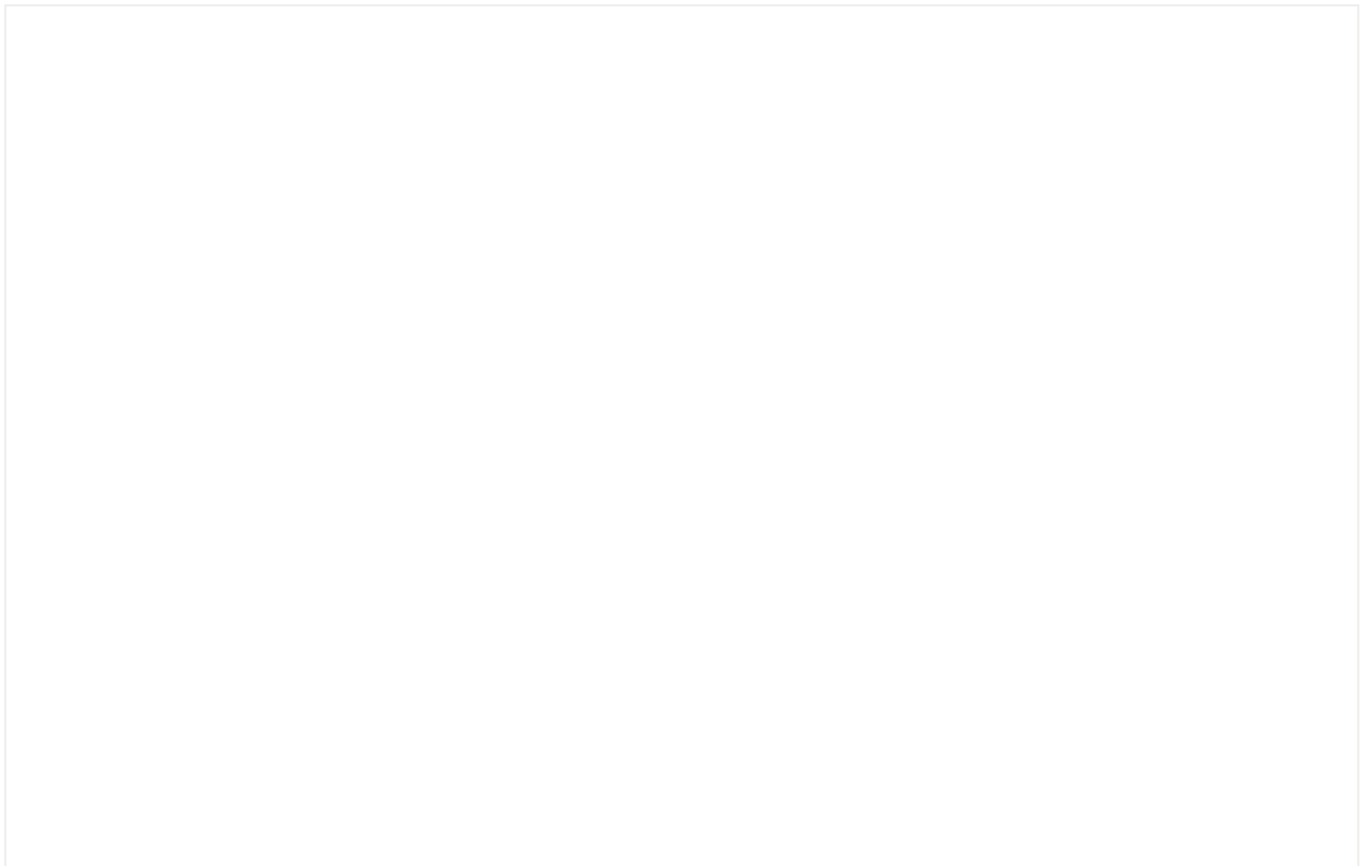


图8. 处理后的二进制图像

结果

反转图像后，重新执行步骤，这是最终结果！

算法成功检测到文本后，现在可以将其保存到Python对象（例如Dictionary或List）中。由于Tesseract训练数据中未包含某些地区名称（“Kabupaten / Kota”中的名称），因此无法准确检测到。但是，由于可以精确检测到地区的索引，因此这不会成为问题。文本提取可能无法检测到其他字体的文本，具体取决于所使用的字体，如果出现误解，例如将“5”检测为“8”，则可以进行诸如腐蚀膨胀之类的图像处理。

源代码：<https://github.com/fazlurnu/Text-Extraction-Table-Image>

努力分享优质的计算机视觉相关内容，欢迎关注：



AI算法与图像处理

考研逆袭985，非科班跨行AI，目前从事计算机视觉的工业和商业相关应用的工作。分...
248篇原创内容

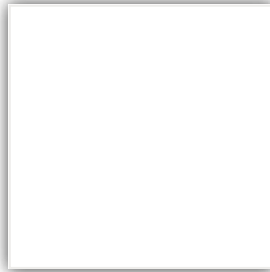
公众号

交流群

欢迎加入公众号读者群一起和同行交流，目前有**美颜、三维视觉、计算摄影、检测、分割、识别、NeRF、GAN、算法竞赛**等微信群

个人微信（如果没有备注不拉群！）

请注明：地区+学校/企业+研究方向+昵称



下载1：何恺明顶会分享

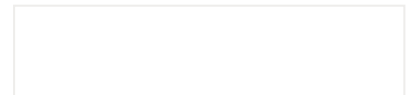
在「AI算法与图像处理」公众号后台回复：何恺明，即可下载。总共有6份PDF，涉及ResNet、Mask RCNN等经典工作的总结分析

下载2：终身受益的编程指南：Google编程风格指南

在「AI算法与图像处理」公众号后台回复：c++，即可下载。历经十年考验，最权威的编程规范！

下载3 CVPR2021

在「AI算法与图像处理」公众号后台回复：**CVPR**，即可下载1467篇CVPR 2020论文 和 CVPR 2021 最新论文



喜欢此内容的人还喜欢

基於python和OpenCV構建智能停車系統
小白學視覺

SpringBoot+flowable快速實現工作流, so easy!

小哈學Java

基於Python的OpenCV輪廓檢測聚類

小白學視覺