

OCR 文本识别

基于传统OCR与神经网络结合

复旦大学计算机科学与技术
16307130335
方煊杰

1

想法来源

- Cs 扫描仪

上传文件(作业)很方便

- Pdf 论文提取

PDF 转化为word是个十分常见的功能，有时候做自然语言处理并不想总是去用adobe reader 做预处理，而是希望能够在代码中一起实现

2

1

步骤

- 拍摄图片预处理
 - 边缘提取，透射变换，锐化增强
- 字符定位与切割
 - 投影算法
- 生成训练数据集
- 卷积神经网络

3

拍摄图片预处理

边缘提取

原图



灰度



强度变化

$$R = G = B$$

$$Gray(i, j) = 0.299 * R(i, j) + 0.578 * G(i, j) + 0.114 * B(i, j)$$

模糊



高斯降噪

核: 5 X 5

中间点“取”周围点的平均值，就会变成1。在数值上，这是一种“平滑化”

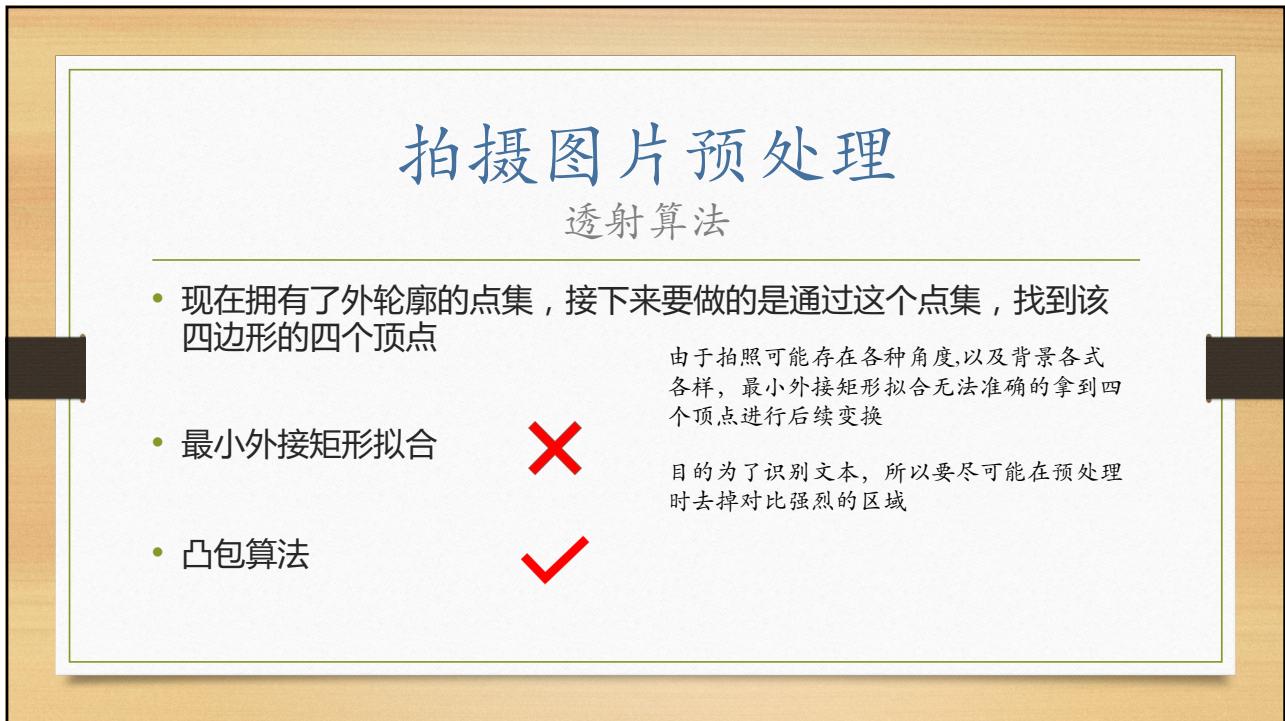
适当大小的核

4

2



5

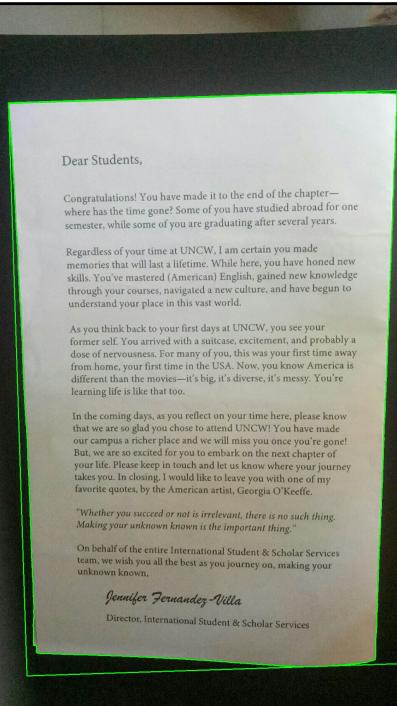


6

拍摄图片预处理

透射算法

- 现在拥有了外轮廓的点集，接下来要做的是通过这四边形的四个顶点
- 最小外接矩形拟合
- 凸包算法



7



8

拍摄图片预处理

透射算法

Dear Students,

Congratulations! You have made it to the end of the chapter—where has the time gone? Some of you have studied abroad for one semester, while some of you are graduating after several years.

Regardless of your time at UNCW, I am certain you made memories that will last a lifetime. While here, you have formed new skills, learned about different cultures, traveled to new places, worked through your courses, navigated a new culture, and have begun to understand your place in this vast world.

As you think back to your first days at UNCW, you are most likely to remember the excitement, and probably a dose of nervousness. For many of you, this was your first time away from home, your first time in the USA. Now, you know America is different, you are more confident, and you are beginning to understand what learning life is like that took place here.

In the coming days, as you reflect on your time here, please know that we are so glad you chose to attend UNCW. You have made our community a better place, and we will miss you when you leave. But, we are so excited for you to embark on the next chapter of your life. We hope you will continue to learn and grow as your journey takes you. In closing, I would like to leave you with one of my favorite quotes, by the American artist, Georgia O'Keeffe:

"Whether you succeed or not is irrelevant; there is no such thing. Making your unknown known is the important thing."

On behalf of the International Student & Scholar Services team, we wish you all the best as you journey on, making your unknown known.

Jennifer Hernandez-Villa
Director, International Student & Scholar Services

仿射变换是透视变换的特例

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} + \begin{bmatrix} X' = \frac{X}{Z} \\ Y' = \frac{Y}{Z} \\ Z' = \frac{Z}{Z} \end{bmatrix}$$

$$\begin{cases} X' = a_{11}x + a_{12}y + a_{13} \\ a_{21}x + a_{22}y + a_{23} \\ a_{31}x + a_{32}y + a_{33} \\ Z' = 1 \end{cases}$$

原图像里面的直线，经透视变换后仍为直线

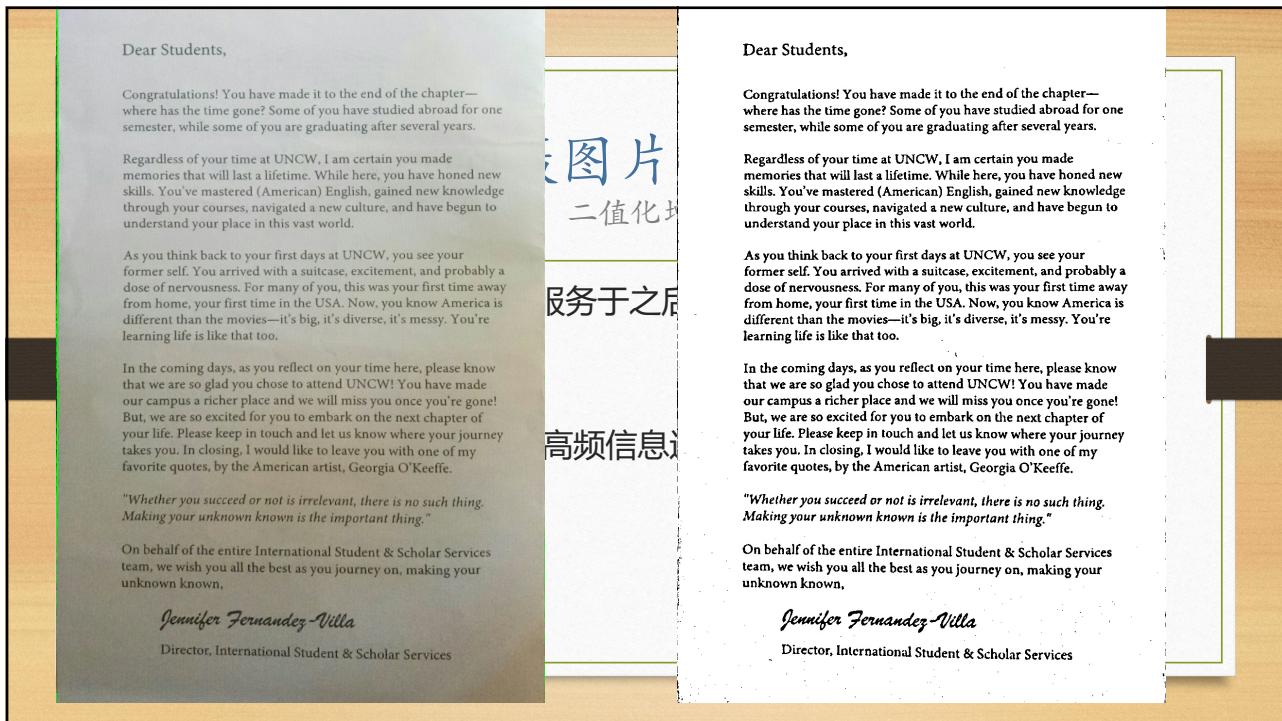
9

拍摄图片预处理

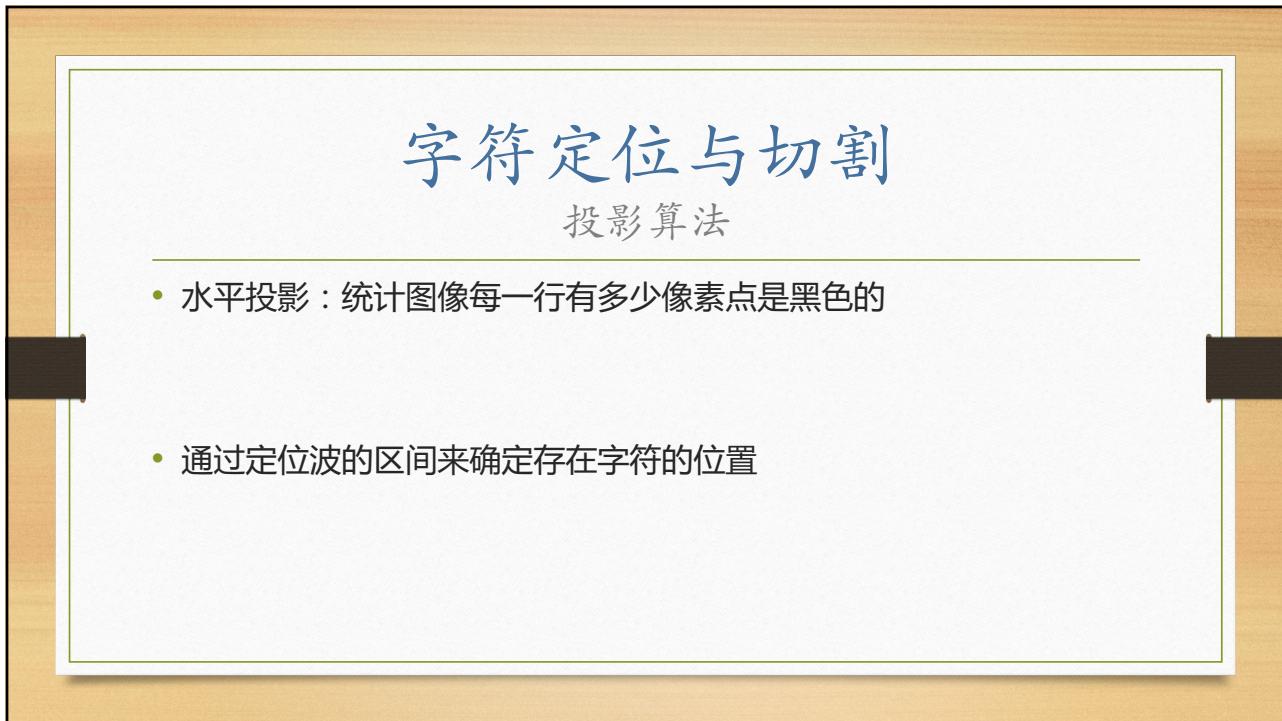
二值化增强

- 为了实现扫描效果，更好服务于之后的文本识别，还需要对图片进行二值化增强处理
- 图像锐化：高通滤波（让高频信息通过，过滤低频信息）边缘得到增益

10



11



12

字符定位与切割
投影算法

- 水平投影：统计图像每一行有多少像素点是黑色的

Dear Students,

- 通过定位波的区间来确定存在字符的位置

13

字符定位与切割
投影算法

- 将每一行再进行垂直投影，获得单个字符：（还需再次水平投影）

Dear Students,



14

生成训练字符

plt

- 使用plt 生成大量字符训练集

- 网上并无大量字符印刷数据集
- 从自身电脑中导出字体：[/Library/fonts](#)
- # 创建Font对象:
font = ImageFont.truetype('font_test/Geneva.dfont', 60)
创建Draw对象:
draw = ImageDraw.Draw(image)
输出文字:
draw.text((0, 0), string[k], font=font, fill=(0, 0, 0))

▼ fonts

- Apple Braille.ttf
- Apple Braille Outline 6 Dot.ttf
- Apple Braille Outline 8 Dot.ttf
- Apple Braille Pinpoint 6 Dot.ttf
- Apple Braille Pinpoint 8 Dot.ttf
- Apple Color Emoji.ttc
- Apple Symbols.ttf
- AppleSDGothicNeo.ttc
- AquaKana.ttf
- ArabicUDisplay.ttc
- ArabicUIText.ttc
- ArialHB.ttc
- Avenir.ttc
- Avenir Next.ttc
- Avenir Next Condensed.ttc
- Courier.ttf
- GeezaPro.ttc
- Geneva.dfont

15

生成训练字符

plt

- 字体
- 角度
- 颜色

68_95.jpg	68_96.jpg	68_97.jpg	68_98.jpg	68_99.jpg	98.jpg	99.jpg	100.jpg	101.jpg	102.jpg

68.jpg	69_0.jpg	69_1.jpg	69_2.jpg	69_3.jpg	103.jpg	104.jpg	105.jpg	106.jpg	107.jpg

69_4.jpg	69_5.jpg	69_6.jpg	69_7.jpg	69_8.jpg	108.jpg	109.jpg	110.jpg	111.jpg	112.jpg

69_9.jpg	69_10.jpg	69_11.jpg	69_12.jpg	69_13.jpg	113.jpg	114.jpg	115.jpg	116.jpg	117.jpg

69_14.jpg	69_15.jpg	69_16.jpg	69_17.jpg	69_18.jpg	118.jpg	119.jpg	120.jpg	121.jpg	122.jpg

93.jpg	94.jpg	95.jpg	96.jpg

16

卷积神经网络

cnn

- 网络定义

Conv(kernel size 8*8, stride [1, 1, 1, 1])*16 个 feature map→Relu→

Max_pooling(kernel size 2*2, stride [1, 2, 2, 1])→

Conv(kernel size 5*5, stride [1, 1, 1, 1])*32 个 feature

Max_pooling(kernel size 1*1, stride [1, 1, 1, 1])→

FullConnectedLayer→Dropout(keep_prob 0.5)→

Adam Optimizer, softmax_cross_entropy_with_logits

Test Loss: 0.000395, Acc: 0.994681

Test Loss: 0.000293, Acc: 0.996277

Test Loss: 0.000304, Acc: 0.995213

Test Loss: 0.000297, Acc: 0.995745

Test Loss: 0.000439, Acc: 0.992021

Test Loss: 0.000293, Acc: 0.995213

Test Loss: 0.000323, Acc: 0.995213

Test Loss: 0.000124, Acc: 0.998936

Test Loss: 0.000110, Acc: 0.998936

Test Loss: 0.000101, Acc: 0.999468

17

识别

- 演示

	数字&字母
Train set	12192
Validation set	371
Accuracy(validation)	98.6%
Test set	612
Accuracy(test)	97.1%
Time cost	840.2s

18

总结算法动机

- **选择透射而非仿射**：自然场景文字框可能存在由于拍摄角度不同而导致的变形，有可能出现梯形
- **选择凸包拟合而非最小外接矩形**：对于文本识别，需要尽可能的去除背景的干扰，保留“文本框”，以及由于角度变形可能导致变换后文字变形过大
- **需要三次投影**：第一次水平，得到一行文字；第二次在一行文字图像基础上垂直，得到单个字符，但此时单个字符竖直方向有空白，所以需要再次水平投影去除上下空白

19

问题及改进

- 针对场景较特殊，需要背景与所需文本框色差较大，否则矩形轮廓无法找到（需要针对不同图片调整opencv里的参数）
- 单独纯文字无法准确提取，这个需要使用（直线矫正）来进行预处理
- 识别率有待提高，可能还需要增加字体训练
- 当前并不支持手写体（外网新闻报道外国人写字五花八门特别是斜体粘连分割有很大麻烦）

20

10