

一种优化的手写字符自动分割算法

黄一琦¹, 郑佳春², 曹长玉¹

(1. 集美大学航海学院, 福建 厦门 361021; 2. 集美大学信息工程学院, 福建 厦门 361021)

[摘要] 在手写字符自动识别时, 由于手写字符中存在字符大小间距不一、粘连、断点, 以及不连贯等情况, 给字符自动分割识别带来极大的困难。针对该问题, 提出了一种优化的手写字符自动分割算法。该方法依据滴水算法的原理, 结合 CFS (color filling segmentation) 做初步分割; 再根据分割字符的连续黑色像素点的宽度判断是否为粘连字符, 若为粘连字符, 则在分割字符图片 0.2 倍宽度与 0.8 倍宽度之间扫描黑色像素位置; 结合分割图片中间位置来确定滴水算法起始滴落点, 解决特殊情况下的起始滴落点的定位不精准问题。经手写字符识别实验结果表明, 优化后手写字符分割准确率比传统方法分割准确率提高了 11.6%, 且有良好的通用性, 可提高手写字符的单个识别率。

[关键词] 手写字符分割; 滴水算法; 粘连字符; 起始滴落点

[中图分类号] TP 312

An Optimized Automatic Segmentation Algorithm for Handwritten Characters

HUANG Yiqi¹, ZHENG Jiachun², CAO Changyu¹

(1. Navigation College, Jimei University, Xiamen 361021, China;

2. School of Information Engineering, Jimei University, Xiamen 361021, China)

Abstract: In the automatic recognition of handwritten characters, due to the differences in character size and spacing, adhesion, breakpoints, and inconsistencies in handwritten characters, it has brought great difficulties to the automatic segmentation of characters during recognition. Aiming at this problem, an optimized automatic segmentation algorithm for handwritten characters is proposed. This method is based on the principle of the drip algorithm, combined with color filling segmentation (CFS) for preliminary segmentation; and then judges whether the characters are sticky characters based on the width of the continuous black pixels of the segmented characters. The black pixel position is scanned between 0.8 times the width and combined with the middle position of the segmented picture to determine the initial drip point of the drip algorithm, which solves the inaccurate positioning of the initial drip point in special cases. The experimental results of handwritten characters show that the segmentation accuracy of optimized handwritten characters is 11.6% higher than that of traditional methods. After optimization, it has good generality and can improve the single recognition rate of handwritten characters.

Keywords: handwritten character segmentation; drip algorithm; sticky character; initial drip point

[收稿日期] 2019-09-29

[基金项目] 福建省科技计划重点项目 (2017H0028); 福建省自然科学基金项目 (2015J01265); 集美大学校基金项目 (2P2020042)

[作者简介] 黄一琦 (1997—), 男, 硕士生, 从事交通通信及物联网技术研究。通信作者: 郑佳春 (1965—), 男, 教授, 从事通信工程研究。E-mail: jchzheng@jmu.edu.cn

<http://xuebaobangong.jmu.edu.cn/zkb>

0 引言

传统的人工阅卷方式需要耗费大量的人力和时间, 为了让教师将更多的时间投入到重要的教学任务中, 出现了机器自动阅卷。试卷自动评阅的实质是手写体答案的自动识别。识别过程包括: 试卷答案预处理、答案字符分割、答案识别等阶段。其中, 字符分割是最为关键的, 它直接影响着字符识别的准确率和答案判定的可靠性。因此, 手写字符分割值得研究。

目前国内外对字符分割的研究已经取得一些成果: 文献[1]提出将图片上的字符分割开后, 再结合传统的机器学习相关算法, 可以得到一个较好的识别效果; 文献[2]提出一种基于投影的竖直分割的方法, 它解决了大部分情况下非粘连字符的分割, 但在粘连字符的分割上效果不理想; 文献[3]提出一种利用中轴点作为字符分割点的方法, 其中中轴点表示在背景中两个分离的字符像素之间的背景中心像素点, 该算法针对那些工整字符之间宽度一致的不粘连的字符具有良好的分割效果; 文献[4]提出一种基于局部极小值和最小投影值的字符分割方法, 但它对交叉重叠和扭曲字符进行分割效果较差; 文献[5]根据平均宽度的相关倍数得到一个预设宽度来分割图像, 每个特征对应着不同的结果, 然后找出最佳片段来作为最终结果; 文献[6]提出了一种滴水算法分割手写数字字符, 它解决了一些基本粘连字符的分割问题; 文献[7]采用上下轮廓差投影法大体确定字符间的坐标位置, 然后再利用水滴算法进行字符切割; 文献[8]认为在验证码字符粘连处字符像素的密度较大, 通过SOM聚类来找到这些像素密集的地方, 但是这种假设依赖于字符笔画的宽度与字符粘连的程度, 并不适用于所有粘连的情况。

现有的分割方法在粘连手写字符间宽度差较大、字符笔画重叠粘连、字符不连贯且含有断点等特殊情况的分割效果不够理想。针对这些问题, 本文在滴水算法研究基础上, 提出一种优化的通用性分割方法。

1 传统字符分割算法存在的问题

1.1 字符分割存在的普遍问题

对于字符分割问题, 竖直投影分割是较早得到应用的方法, 它对于不粘连的字符分割效果还是较理想的。由于手写字符大多比较不工整, 常出现书写连笔或笔画搭在另一个字符上的情况, 导致字符粘连在一起, 此时用竖直投影分割法会出现字符分割不完整、错误分割、笔画断裂等现象, 进而导致识别率不理想。而滴水算法可以解决上面竖直投影分割法对于简单粘连字符分割出现的问题。从样本中取一张字符“48”(见图1a), 两种分割方法的路径及结果如图1所示。

从图1可以看出, 滴水算法分割粘连字符的效果明显好于竖直投影分割效果。因为滴水算法的分割路径是沿着右边字符“8”的边缘轮廓向下的, 避免了竖直投影分割中出现的破坏字符完整性的情形。

关于一些特殊的粘连形式, 传统分割方法的效果一般。在竖直投影时, 由于上下交错导致投影得到的像素值大于0从而误判为粘连情况, 但实际上两个字符之间没有真正粘连在一起, 这种

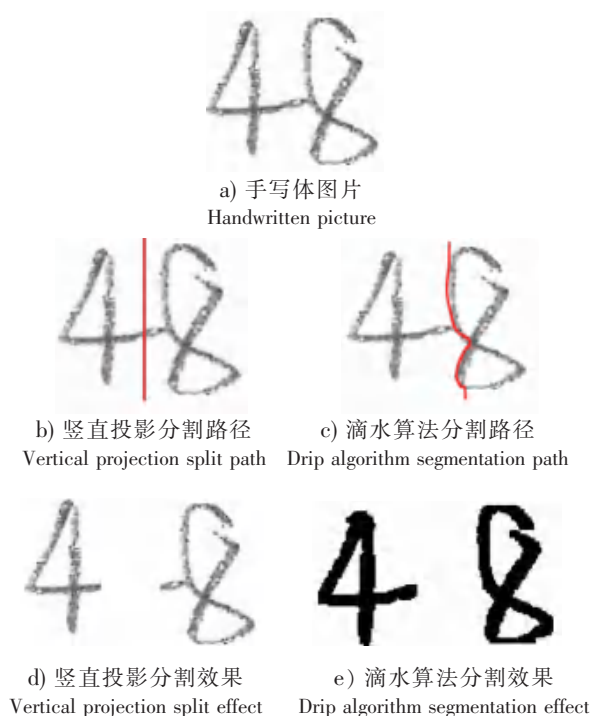


图1 两种方法分割效果对比

Fig.1 Comparison of two methods of segmentation effects

情况一般被称为投影粘连^[9],如图2所示。使用竖直投影分割的话会导致字符分割断裂,而传统滴水算法分割在此情况下的效果比前者更差。

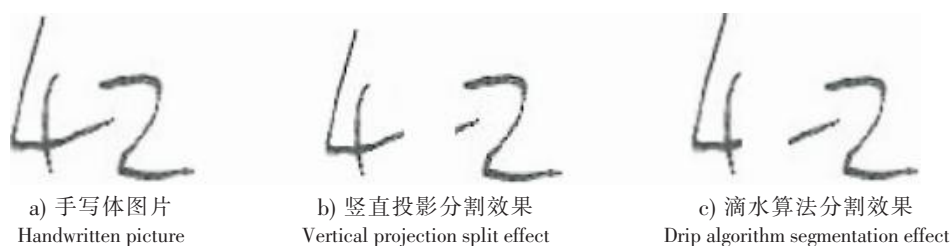


图2 投影粘连字符及分割效果

Fig.2 Projection adhesion characters and segmentation effect

1.2 特殊情况字符的分割问题

在图3a所示的“12”这个字符中可以看到,左边“1”和右边“2”的宽度值差距较大并局部粘连在一起。此情况下利用传统滴水算法分割,可能会将其识别为一个字符;也有可能因为起始点位置选取的不准确,导致字符从右边“2”开始分割,造成字符错误分割的现象,其效果如图3所示。



图3 特殊情况字符及分割效果

Fig.3 Special case characters and segmentation effects

在图4所示的“20”这个字符中可以看到,左边“2”和右边“0”局部粘连在一起,且字符“0”上出现两个的断点。这种笔画不连续且含断点的粘连字符,使用之前的分割方法导致将其错误分割为3个字符,其效果如图4所示。针对此类情况,需要加一些字符宽度以及高度阈值条件进行判断从而合并字符,即当小的连通区域的宽度范围包含在大的连通区域范围之内,则将小的连通区域与大的连通区域合并。

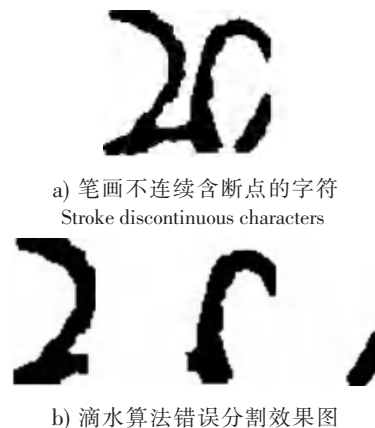


图4 特殊情况字符及分割效果
Fig.4 Special case characters and segmentation effects

2 字符分割算法的优化

2.1 滴水算法

滴水算法(drop fall algorithm)基本思想是模拟水滴从高处向低处滴落的过程来对粘连字符进行分割。在重力作用下,水滴从字符串之间的顶部向下滴落,当水滴遇到字符时只能沿字符轮廓向下滴落或者水平左右滚动。当水滴陷在轮廓的凹陷处时,则进行渗透处理,即从轮廓最低点渗漏到字符笔画中,然后穿透笔画继续滴落,最终水滴所经过的轨迹就是分割路径^[7]。这种方法可以解决竖直投影分割字符所带来的字符断裂、过分割等问题。滴水算法的影响因素主要有:起始滴落点、移动规则和方向不同。水滴的周围状况如图5所示。其中,d0表示水滴当前的位置,水滴的下一滴落位置由d1、d2、d3、d4、d5这5个周围像

d5	d0	d4
d1	d2	d3

图5 水滴及周围像素
Fig.5 Drop point near pixels

素点的情况决定。

2.2 CFS 分割

CFS 算法是由 Yan 等^[11]提出, 通过使用颜色填充字符块的方式, 将字符标记成颜色不同的区块, 这样就可以将没有粘连的字符分割出来, 因此称之为 CFS (color filling segmentation) 分割法。其主要工作原理是: 将字符图片从左到右, 从上到下进行扫描, 当扫描到第一个黑色像素点时, 以它为中心在其相邻的八个方向继续探测, 若存在新黑像素点, 就认为这是一个字符区块。然后以这个新像素点为新的中心点继续进行上述步骤, 直到不再探测到新黑像素点时就认为这一字符区块结束, 并用颜色填充探测到的字符区块^[12]。接着在区块外新的像素点继续重复以上流程, 直到所有字符区块都被检测出来, 这样就会得到若干字符区块, 然后根据填充颜色的不同进行分割, 图 6 为示例图。传统的连通域分割是利用像素点之间的连通性, 把不同连通域的字符块分割开来。只要字符之间不粘连, 即使字符存在倾斜扭曲, 其分割效果都不受影响。



图 6 CFS 的分割示例图

Fig.6 CFS segmentation example diagram

2.3 起始点的设计

通常传统滴水算法的起始滴落点是以从上到下从左到右扫描手写体图片中每一行像素点, 然后使第一个满足像素分布为 $(\dots 0 * 1 \dots 10 \dots)$ 的白像素点 (*) 作为起始滴落点, 其中 0, 1 分别代表黑色、白色像素点^[13]。可以看出传统滴落起始点的选择具有一般性, 对于特殊情况, 这种选择方法得到的分割效果不太理想。按照传统滴水算法规则, 在遇到字符笔画凹陷时, 很大可能就以此为起始点, 进而造成手写体字符分割断裂。并且水滴在字符上移动时还会由于字符轮廓的不平滑, 而出现错误分割现象。综上可知, 起始滴落点的选择对于滴水算法分割的效果是至关重要的。

一般来说, 字符的粘连点会出现在竖直投影直方图中的极小值处。因此, 本文结合两种方法来确定最佳起始滴落点, 即利用竖直投影法辅助局部扫描分析起始滴落点, 这样便能在粘连字符中选择较准确的起始点, 起始点重新设计的流程如图 7 所示。

本研究的局部扫描模块是指根据分割字符的宽度 W , 取粘连字符 $0.2W$ 与 $0.8W$ 之间为扫描区域 I 。选取这个区域是为了消去粘连字符两端的开始、结尾笔画的影响, 即把起始点选择区域集中在粘连字符居中的位置。对区域 I 按水平方向从左到右、从上到下依次扫描像素点, 选出最初满足像素分布为 $(\dots 0 * 1 \dots 10 \dots)$ 的一行。先将所有满足条件的白像素点 (*) 位置坐标存入列表 number 中, 接着将竖直投影法中存在的极小值点的位置坐标按大小存入列表 number 中, 最后参照字符中间位置坐标, 选取列表 number 中最右边的白像素点 $start_x_n$ 作为最优滴落起始点。

本研究根据得到的最优滴落起始点进行滴水算法分割粘连字符, 具有良好的分割效果。图 8 为起始滴落点的正确与错误选择的对比示意图, 不同起始滴落点的分割效果有明显的差异。

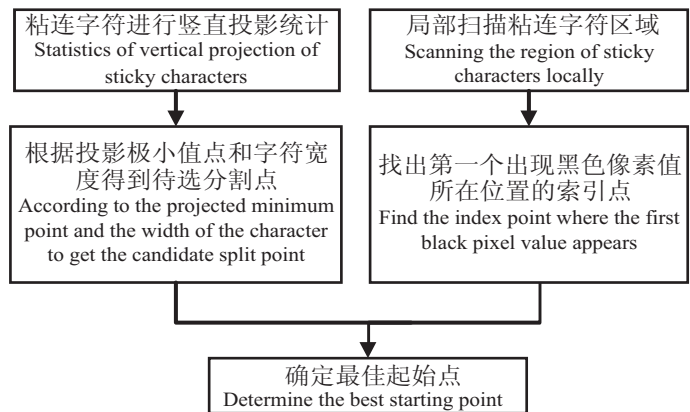


图 7 最佳起始滴落点的选择流程

Fig. 7 Best starting drip point selection process

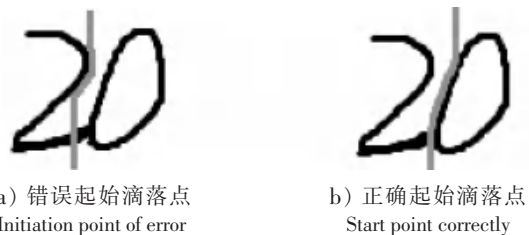


图 8 正确与错误起始滴落点

Fig.8 Correct and wrong start drop point

2.4 通用性优化设计

根据 2.1 分析可知: 手写字符出现粘连时, 不能使用竖直投影算法分割; 字符之间紧密粘连且字符有断裂部分, 不能直接使用连通域算法分割; 有些粘连字符之间宽度差较大, 起始滴落点选择不理想, 不能直接使用滴水算法分割。本文结合以上三种分割算法, 设计出一种通用性更强的方法, 分割出单个字符, 并且对样本字符的宽度、高度和像素点数进行统计, 在分割时用来作为判断字符块中包含的字符个数的依据。流程图见图 9。

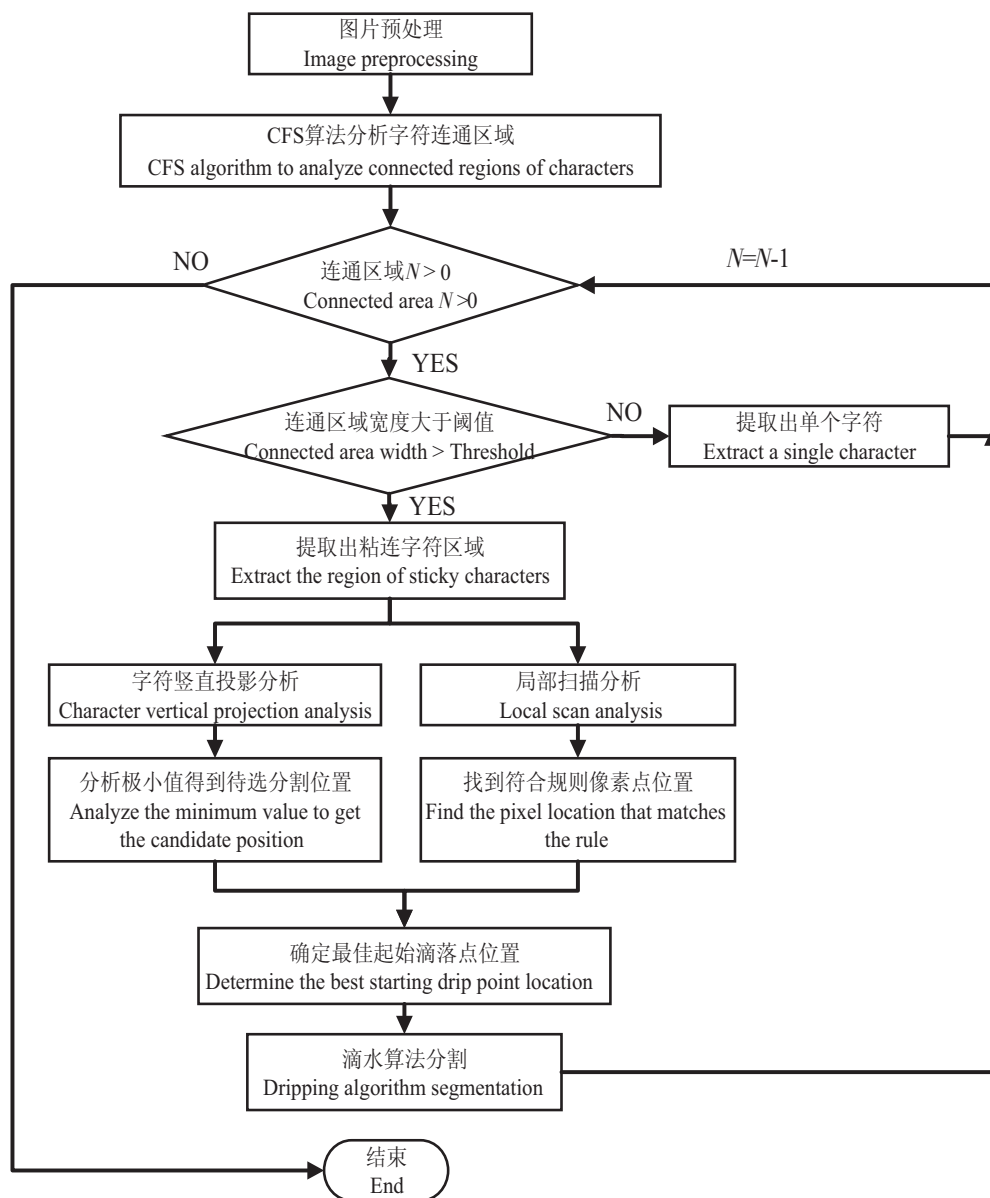


图 9 通用性字符分割方法流程

Fig. 9 Flow of universal character segmentation method

本文通用性分割方法设计主要有以下三个模块:

1) 手写字符输入模块。该模块首先通过 CFS 分割, 它主要有两个作用: 一是当字符为非粘连字符时, 直接进行字符分割; 二是当字符为非粘连字符时, 它可以用来判断字符区块是否含有粘连字符的情况, 若存在则进入下一步分析。所以, 该方法通用于字符不粘连和字符粘连的情况。

2) 滴落起始点选择模块。该模块会对粘连字符进行局部扫描, 并将符合条件的像素点与竖直投影得到的极小值像素点进行分析, 以便确定最佳的起始滴落点。

3) 字符分割与优化模块。该模块根据最佳起始滴落点来进行滴水算法分割粘连字符, 将字符分割之后的结果送入阈值判断中, 经再次判断以避免字符的分割遗漏。相当于自适应的优化分割, 从而对分割结果进行最后的确认。

本分割方法根据粘连情况将三种分割方法依次使用, 并不是一性分割出所有粘连字符, 而是根据粘连程度逐步进行字符分割。

3 实验与分析

3.1 实验环境及数据

本文实验在 Ubuntu18.04 系统进行, 在 PyCharm 开发环境下使用 Python3.6 编写, 测试代码。实验主要使用 30 张小学数学试卷答案, 取其中字符不粘连、字符轻微粘连、字符特殊粘连各 500 个, 总计 1500 个字符。手写体答案在字体大小、书写格式、清晰度上等有所不同。本研究还采用 500 个手写英文字符和网络验证码字符样本作为实验验证集, 验证其通用性效果。

3.2 实验结果

针对传统分割方法出现的分割错误情况, 本研究改进的分割方法与传统分割方法的效果对比如图 10 所示。改进算法对 30 张试卷答案手写字符实验样本进行分割测试, 实验部分分割效果如图 11 ~ 12 所示。

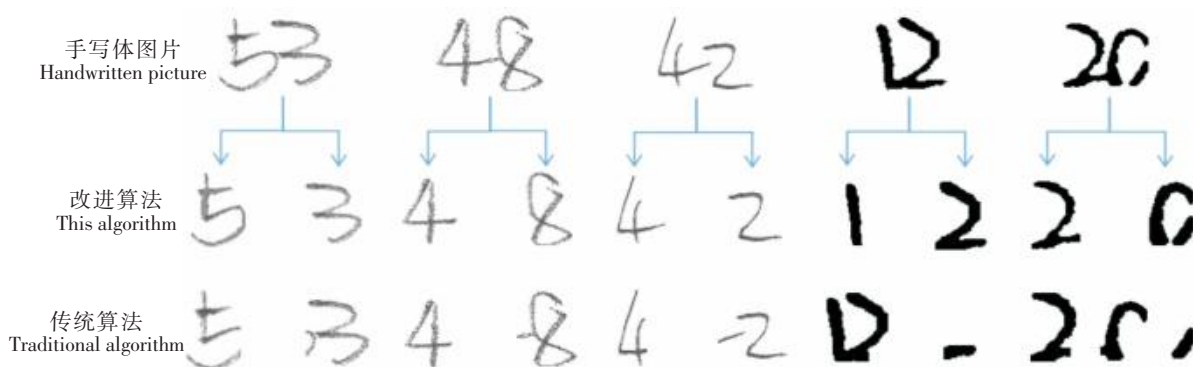


图 10 特殊类型字符示例及分割结果

Fig.10 Examples of special types of characters and segmentation results

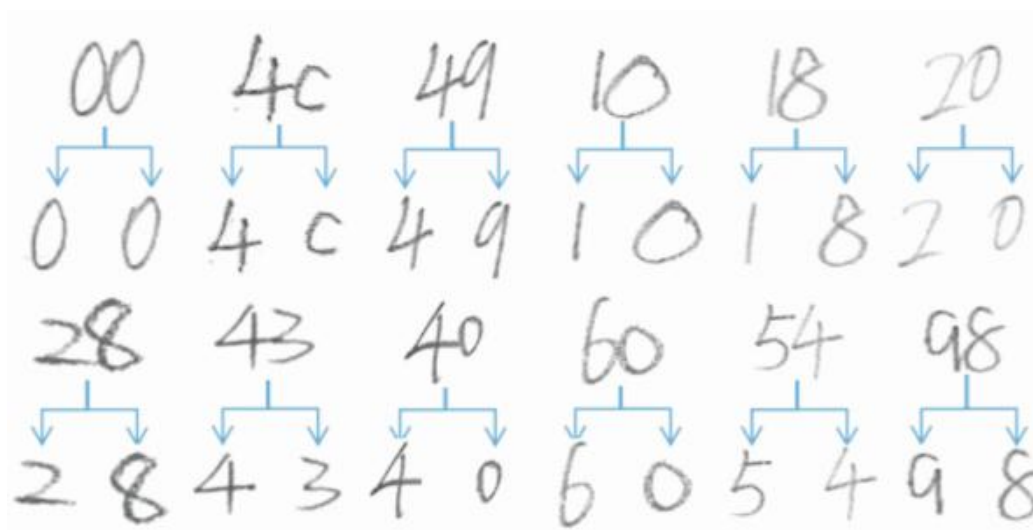


图 11 字符分割正确结果

Fig.11 Character segmentation results correctly

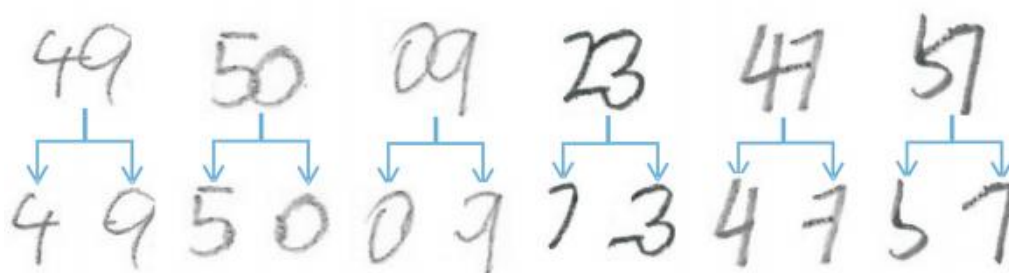


图 12 字符分割不正确结果

Fig.12 Character segmentation results incorrectly

字符分割准确率计算公式为: 字符分割正确率 (%) = 正确分割字符个数/字符总个数 $\times 100\%$ 。实验结果如表 1 所示。对于 1500 个样本数据集, 改进算法的准确率达到 92.4%, 传统滴水算法为 80.8%, CFS 和竖直投影算法都为 33%。相对于传统滴水算法, 改进算法准确率提升了 11.6%。

表 1 字符分割结果对比

Tab. 1 Comparison of character segmentation results

算法 Algorithm	不粘连 Non-stick (500)		轻微粘连 Light stick (500)		特殊粘连 Special adhesion (500)	
	正确分割数 Correct number of divisions	准确率 Accuracy/%	正确分割数 Correct number of divisions	准确率 Accuracy/%	正确分割数 Correct number of divisions	准确率 Accuracy/%
竖直投影 Vertical projection	500	100	80	16	0	0
CFS	500	100	80	16	0	0
滴水算法 Dripping algorithm	500	100	426	85.2	286	57.2
改进算法 This algorithm	500	100	468	93.6	418	83.6

本文对验证集样本进行字符分割测试, 验证集中粘连字的分割效果如图 13 所示。从图 13 中看出改进算法的分割效果是比较理想的, 说明改进算法在粘连字符分割上具有良好的通用性。



图 13 实验验证集分割效果

Fig.13 Experimental verification set segmentation effect

4 结束语

在实验样本中, 很多手写字符是粘连在一起的, 尤其是投影粘连这种情况。竖直投影分割法对无粘连字符但存在倾斜的字符, 以及粘连字符的分割大多会造成字符分割断裂, 即把完整的字符分割为

两个部分,所以针对这种情况采用竖直投影分割法对手写体测试样本的分割成功率是最低的。CFS 分割对不粘连字符的分割效果理想,而对粘连严重的字符无法分割。传统滴水算法对于字符倾斜且没有明显粘连时分割效果理想;但在字符粘连扭曲复杂、宽度相差较大的两个字符粘连,以及字符中笔画不连续含断点等粘连情况下,它的分割效果不理想。从实验结果来看,本文提出对传统分割方法加以组合使用,找出滴水算法的最佳起始滴落点,按照新的起始滴落点的滴落轨迹分割粘连字符。该方法能够正确有效地分割粘连手写字符间宽度差较大、字符笔画重叠粘连、字符不连贯且含有断点等特殊情况,但在出现粘连严重、笔画重叠交叉严重的情况时,就会出现错误分割现象,导致分割准确率下降。整体来说,改进算法的分割准确率相较于传统分割算法提高了 11.6%,针对字符粘连和字符不粘连的情况都能有较理想的分割效果。

接下来的探索方向是利用 SOM 聚类结合字符骨架化分析来继续提高字符分割的准确率,以及通过卷积神经网络识别字符。

[参考文献]

- [1] BURSZEIN E, MARTIN M, MITCHELL J. Text-based CAPTCHA strengths and weaknesses [C] //Proceedings of the 18th ACM Conference on Computer and Communications Security. New York: ACM, 2017: 125-138. DOI:10.1145/2046707.2046724.
- [2] FUJISAW H, NAKANO Y, KURINO K. Segmentation methods for character recognition: from segmentation to document structure analysis [J] //Proceedings of the IEEE, 1992, 80(7): 1079-1092. DOI:10.1109/5.156471.
- [3] HUANG SY, LEE YK, BELL G, et al. An efficient segmentation algorithm for captchas with line clustering and character wrapping [J]. Multimedia Tools & Application, 2010, 48(2): 267-289. DOI:10.1007/s11042-009-0341-5.
- [4] 王璐,张荣,尹东,等.粘连字符的图片验证码识别[J].计算机工程与应用,2011,47(28):150-153.
- [5] CHEN J, LUO X, GUO Y, et al. A survey on breaking technique of text-based captcha [J]. Security & Communication Networks, 2017, 2017(1/2): 1-15. DOI:10.1155/2017/6898617.
- [6] CONGEDO G, D IMAURO G, IMPEDOV S, et al. Segmentation of numeric strings [C] //Proceedings of the 3rd International Conference on Document Analysis and Recognition, 1995: 1038-1041. DOI:10.1109/ICDAR.1995.602080.
- [7] 汪洋,许映秋,彭艳兵.基于 KNN 技术的校内网验证码识别[J].计算机与现代化,2017(2):93-97.
- [8] 简献忠,曹树建,郭强,等. SOM 聚类与 Voronoi 图在验证码字符分割中的应用[J].计算机应用研究,2015,32(9):2857-2861.
- [9] 尹龙,尹东,张荣,等.一种扭曲粘连字符验证码识别方法[J].模式识别与人工智能,2014,27(3):235-241.
- [10] YANG J Y, GUO J, JIANG W W. A novel drop-fall algorithm based on digital features for touching digit segmentation [C] //2016 IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON). Vancouver, BC: IEEE, 2016. DOI:10.1109/IEMCON.2016.7746350.
- [11] YAN J, AHMAD A S E. A low-cost attack on a Microsoft captcha. [C] // Proceedings of the 15th ACM Conference on Computer and Communications Security (CCS'08). New York: ACM, 543-554. DOI:10.1145/1455770.1455839.
- [12] WITTAYA J, THANARAT H C. Automatic detection and segmentation of text in low quality thai sign images [C] // Proceedings of the APCCAS 2006-2006 IEEE Asia Pacific Conference on Circuits and Systems. Singapore: IEEE, 2006. DOI:10.1109/APCCAS.2006.342256.
- [13] WANG X J, ZHENG K F, GUO J. Inertial and big drop fall algorithm [J]. International Journal of Information Technology, 2006, 12(4): 39-48.

(责任编辑 朱雪莲 英文审校 黄振坤)