

维吾尔文笔迹鉴别预处理及边缘提取方法研究

沈 洁¹, 卡米力·木依丁¹, 张祖莲²

(1. 新疆大学 信息科学与工程学院, 新疆 乌鲁木齐 830046;

2. 新疆气象局 新疆兴农网信息中心, 新疆 乌鲁木齐 830002)

摘要:笔迹鉴别是一种行为特征的识别方法,笔迹容易获取并且具有唯一性,因此,基于笔迹的身份验证在安全等领域有广阔的应用前景。笔迹鉴别可分为在线、离线两种形式,鉴别方法有文本依存和文本独立两大类。文中结合维吾尔语自身的文字特点,主要针对离线的文本无关的维吾尔文笔迹鉴别中预处理和笔迹的边缘提取技术进行细致的研究。该研究为维吾尔文计算机笔迹鉴别系统的最终实现提供了关键的技术支持,推动了少数民族语言文字笔迹鉴别自动化的进程。

关键词:笔迹鉴别;文本无关;维吾尔文;预处理;边缘提取

中图分类号:TP391

文献标识码:A

文章编号:1673-629X(2012)04-0065-04

Pre-Processing Andedge Extraction Research on Uyghur Writer Identification

SHEN Jie¹, KAMIL · Moydi¹, ZHANG Zu-lian²

(1. College of Information Science and Engineering, Xinjiang University, Urumqi 830046, China;

2. Information Center of Xinjiang Development Agriculture Net, Xinjiang Weather Bureau, Urumqi 830002, China)

Abstract: Writer identification is a method of behavioral characteristics recognition. Handwriting is unique and easy to obtain, and therefore, handwriting-based authentication has broad application prospects in security field. Writer identification can be divided into online and offline two forms. There are two kinds of identification methods, they are text-dependent and text-independent. Do the meticulous research on pre-processing and handwriting edge extraction techniques in offline text-independent Uyghur writer identification based on the characteristics of the Uyghur itself. This research provides key technical support for the ultimate realization of the Uyghur computer writer identification system, pushes the automated process of the minority language writer identification.

Key words: writer identification; text-independent; Uyghur; pre-processing; edge extraction

0 引言

笔迹鉴别也称为笔迹识别,是一种通过分析和对比手写笔迹的书写风格来判断书写人身份的技术^[1]。随着生物识别技术的发展,笔迹鉴别逐渐成为图像处理和模式识别领域一个非常活跃的研究课题。

笔迹的优点很多,例如唯一性、稳定性,并且笔迹可以采集,获取笔迹没有任何侵犯性^[2]。因此,进行快速、有效、可靠的笔迹鉴别是很有社会价值的,并具有很大的实用意义。从身份识别的角度,基于生物特征的身份鉴别技术的发展为人们提供了一种更加方便

和可靠的解决方案。传统的笔迹鉴别采用人工的方法,容易引入人的感情因素,影响鉴定效果的真实性^[3]。利用计算机辅助进行笔迹鉴定,可以提高鉴定速度与准确性,为更为客观的人工鉴定提供有利的帮助。因此,笔迹鉴别的自动化和智能化,使笔迹鉴别结果更客观,成为该领域中的重要研究目标。随着计算机和网络技术的发展与普及,笔迹鉴别技术的应用领域变得越来越广,突破了原有的应用范畴。可以说凡是需要安全保密和身份鉴别的地方,例如金融、保险、公安司法部门的刑事调查和法庭审判领域都用着笔迹鉴别^[4]。

现今笔迹鉴别技术研究大部分基于英文和中文,然而少数民族语言文字自动鉴别技术研究尚不成熟。本研究正是弥补我国少数民族语言文字鉴别体系中的这一不足,展开对维吾尔文的笔迹鉴别技术的研究。汉字是典型的方块结构字符,笔迹中的字符之间一般

收稿日期:2011-09-02;修回日期:2011-12-08

基金项目:国家自然科学基金资助项目(61065001);新疆少数民族科技人才特殊培养计划(201023116)

作者简介:沈 洁(1986-),女,硕士研究生,主要研究方向为模式识别;卡米力·木依丁,副教授,硕士生导师,主要研究方向为模式识别、信息检索。

有明显间隔,字符间的连写少。而维吾尔语属阿尔泰语系突厥语族,维吾尔文采用阿拉伯文字母,由 32 个字母组成,在形态结构上属于粘着语类型,其字符形式复杂,字符间连写非常频繁。笔迹鉴别的整个过程中,预处理是一个很重要的环节,预处理的目的是抑制无用信息增强有用信息,直接关系到最终的识别准确率^[5]。维吾尔文的特殊性给笔迹图像的预处理带来了很大的困难。因此需要在已有的汉字笔迹图像的预处理技术的基础上进一步研究针对维吾尔文笔迹图像预处理的方法。

文中主要研究维吾尔文计算机笔迹鉴别的前期关键技术与方法。整个过程包括笔迹图像的预处理和笔迹图像边缘检测。针对维吾尔文文字的特殊性,为了保存较为完整的笔迹信息,文中的预处理工作主要包括笔迹图像采集、去噪、灰度化和二值化。

1 预处理

1.1 笔迹图像采集

文中采集了 80 名学生的笔迹样本。本实验中书写的笔迹均使用钢笔、圆珠笔、碳素笔等普通的书写工具。首先扫描收集的 80 名学生的笔迹样本,图像的扫描精度通常有 200dpi、300dpi、600dpi 等。扫描精度的高低直接影响图像的清晰度,过高的清晰度又会促使数据量增大,进而增加内存需求和处理时间,最终会影响鉴别速度。综合以上各因素,本研究采用 300dpi 的扫描精度,扫描类型设置为彩色(24 位),扫描后所得的笔迹图像被保存成位图(BMP)格式存入计算机中^[6],如图 1 所示。

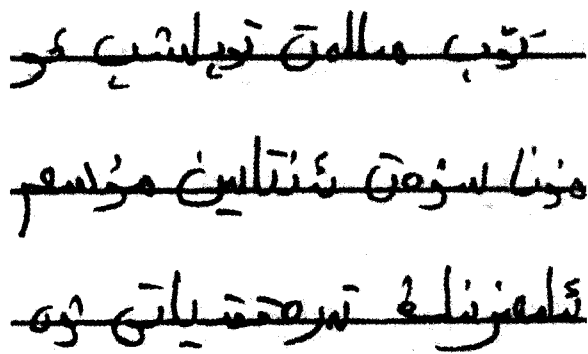


图 1 笔迹样本

1.2 纸张背景噪声和网格线的去除

笔迹鉴别即运用图像处理技术来分析不同笔迹书写人的书写风格,笔迹的书写过程不受书写工具和纸张类型的限制。在实际运用图像处理技术分析笔迹样本图像中,常常会因纸张类型、背景颜色、污迹、格线以及其它杂色等影响分析结果,因此,必须要消除这些干扰因素。

文中在预处理中设计了一个图像颜色取色器(如图 2 所示),通过图像取色器获取文本的颜色,把与文本颜色不一致的其他颜色置为背景色。首先将鼠标放在文字的笔道上,通过单击鼠标来取出鼠标指针处的颜色值。经过大量实验设定一个阈值,当图像各个墨迹点处的 R、G、B 三基色的值与获得的文本颜色的三基色值的差别均小于设定的阈值时,将该墨迹点的颜色置为黑色,否则认定为干扰色,将该墨迹点的颜色置为背景色。该方法能有效地消除其它颜色的干扰,并且不会出现文字断笔等现象,从而较为完整的保留了笔迹书写者的文本信息。通过使用图像颜色取色器获取笔迹颜色的方法,能够在干扰色较多的情况下,大大简化操作人员的工作。

实际上,笔划往往是很细的,维吾尔文字的连写又相当复杂,鼠标的点击位置往往不准确,为了有效获取笔迹的颜色,在取色器中设计了一个对话框,对话框中有一个专门显示鼠标点击处颜色的矩形区域,可以通过手工修改的方式来解决颜色偏差的问题,颜色的变化会显示在矩形区域内。笔迹去噪效果如图 3 所示。

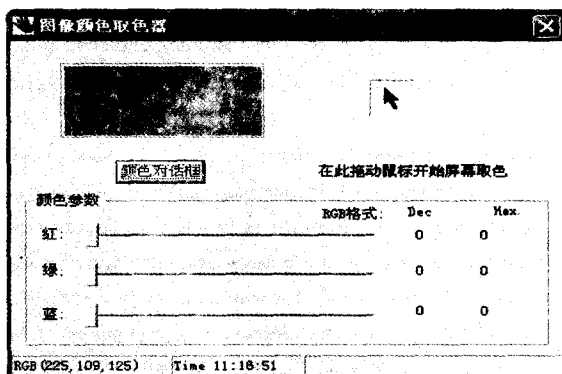


图 2 图像颜色取色器

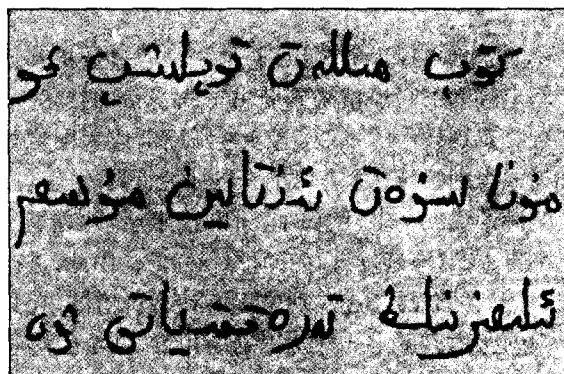


图 3 笔迹样本去噪(背景格线去除)

1.3 笔迹图像的灰度化

本研究实现的笔迹鉴别是建立在二值图像的基础上,首先要进行笔迹图像的灰度化操作,然后将灰度图像转化为二值图像。将彩色图像转化为灰度图像的方法有最大值法、平均值法和加权平均法等^[7],最大值法是将彩色图像中的三个分量的亮度的最大值作为灰度

图的灰度值,平均值法将彩色图像中的三个分量的亮度求平均得到一个灰度图,而加权平均法是根据特殊指标对 R、G、B 分别赋权值。加权平均值法将参考 WR, WG, WB 的不同取值,根据重要性及其它指标,将三个分量以不同的权值进行加权平均。由于人眼对绿色的敏感最高,对蓝色敏感最低,因此,使用加权平均法最终能得到较合适的灰度图像。通过实验和理论的证明可以得出,当 $f_R = 0.299$, $f_G = 0.587$, $f_B = 0.114$ 时,将得出最合理的灰度图像,即

$$f(i,j) = 0.299R(i,j) + 0.587G(i,j) + 0.114B(i,j) \quad (1)$$

$f(i,j)$ 为像素的灰度值。笔迹灰值化效果如图 4 所示。

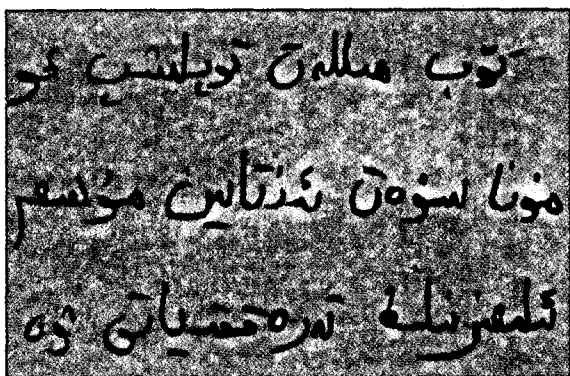


图 4 笔迹样本灰值化

1.4 笔迹图像黑白二值化

对图像进行二值化操作是为了去除扫描图像中的冗余信息,为笔迹鉴别的特征提取过程打下基础。二值图像受书写环境的影响较小^[8],我们认为经过预处理中二值化操作后得到的二值图像中的所有黑像素点都是笔迹书写人的书写墨迹。

黑白二值化的方法根据其运算的范围不同,可分为全局阈值法和局部比较法。全局阈值法(S. Watanabe^[9])根据图像的直方图或灰度的空间分布确定一阈值,并根据此阈值实现灰度文本图像到二值化文本图像的转化。局部比较法通过定义考察点的邻域,并由邻域计算模板来实现考察点灰度与邻域的比较。

局部比较法实现速度慢,并且针对维吾尔文文字字符间连写较多的特点,不能保证字符笔画连通性。因此文中采用全局阈值法,根据灰度直方图确定最佳阈值,当像素的灰度值大于阈值,则将此像素灰度值置为 1,当像素点的灰度值小于阈值,则将此像素的灰度值置为 0,最终实现二值化。笔迹样本二值化效果如图 5 所示。

2 笔迹图像的边缘提取

笔迹鉴别中前期的预处理是为了更好地实现笔迹

的特征提取,然而笔迹的特征提取需要在边缘图像上进行,因此图像的边缘检测也是一个很重要的环节。

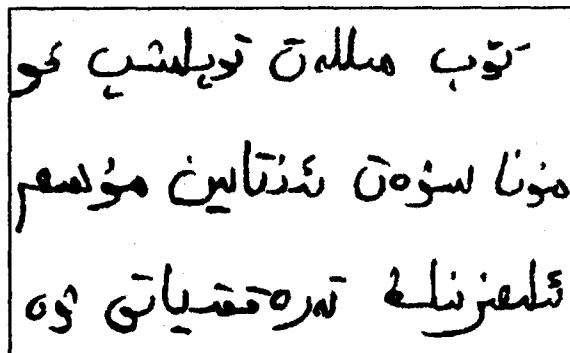


图 5 笔迹样本二值化

所谓边缘就是指图像局部亮度变化最显著的部分,它是检测图像局部显著变化的最基本的运算^[10]。使图像的轮廓更加突出的图像处理方法叫做边缘检测或边缘增强。边缘检测是一种重要的区域处理,它将突出图像的边缘,边缘以外图像区域通常将被削弱甚至被完全去掉,处理后的图像的亮度保持不变,像素值变化缓慢的区域变黑,而像素变化剧烈的区域被突出^[11]。

图像边缘检测算法^[12]很多,有基于求导 Roberts 算子、Prewitt 算子和 Sobel 算子检测梯度(数字图像中,梯度用来表示图像灰度值的显著变化)最大值法,还有检测二阶导数的过零点,以及利用多尺度小波进行边缘检测等方法。

Roberts 算子在 2×2 领域内采用对角线方向相邻两像素之差近似梯度幅值检测边缘。检测水平和垂直边缘的效果好于斜向边缘,定位精度高,但由于不包括平滑,所以对于噪声比较敏感。

Prewitt 算子是一种一阶微分算子的边缘检测,利用像素点上下、左右邻点的灰度差,在边缘处达到极值检测边缘,去掉部分伪边缘,对噪声具有平滑作用。其原理是在图像空间利用两个方向模板与图像进行邻域卷积来完成的,这两个方向模板一个检测水平边缘,一个检测垂直边缘。它对灰度渐变低噪声的图像有较好的检测效果,但是对于混合多复杂噪声的图像,处理效果就不理想了。

Sobel 是在 Prewitt 算子的基础上,在 3×3 邻域内做加权平均和差分运算,Sobel 算子很容易在空间上实现,Sobel 边缘检测器不但产生较好的边缘检测效果,而且受噪声的影响也比较小。该算子对噪声具有平滑作用,提供较为精确的边缘方向信息。与 Roberts 算子和 Prewitt 算子相比,Sobel 算子对于像素的位置的影响做了加权,因此效果更好。

Laplacian 算子是二阶微分算子。其具有各向同性,即与坐标轴方向无关,坐标轴旋转后梯度结果不

变。但是,其对噪声比较敏感,所以,图像一般先经过平滑处理,因为平滑处理也是用模板进行的,所以,通常的分割算法都是把 Laplacian 算子和平滑算子结合起来生成一个新的模板。

经过大量的实验比较,根据维吾尔文文字的特点,从笔迹样本的特征提取需要方面考虑,文中采用 Sobel 算子获取笔迹的边缘图像。Sobel 算子的实现公式为

$$G[i,j] = |f[i-1,j+1] + 2f[i,j+1] + f[i+1,j+1] - f[i-1,j-1] - 2f[i,j-1] - f[i+1,j-1]| + |f[i-1,j-1] + 2f[i-1,j] + f[i-1,j+1] - f[i+1,j-1] - 2f[i+1,j] - f[i+1,j+1]| \quad (2)$$

其中 $G[i,j]$ 表示处理后 (i,j) 点的灰度值, $f[i,j]$ 表示处理前该点的灰度值。

Sobel 算子实现笔迹图像边缘检测的方法如下:

首先对于数字图像,用一阶差分代替一阶微分,即

$$\Delta x f(x,y) = f(x,y) - f(x-1,y) \quad (3)$$

$\Delta y f(x,y) = f(x,y) - f(x,y-1) \quad (4)$

求梯度时对于平方和运算及开方运算,用两个分量的绝对值之和表示,即

$$G[f(x,y)] = |[\Delta x f(x,y)] + [\Delta y f(x,y)]| + |\Delta x f(x,y)| + |\Delta y f(x,y)| \quad (5)$$

Sobel 梯度算子先做成加权平均,再微分,即

$$\Delta x f(x,y) = f(x-1,y+1) + 2f(x,y+1) + f(x+1,y+1) - f(x-1,y-1) - 2f(x,y-1) - f(x+1,y-1) \quad (6)$$

$$\Delta y f(x,y) = f(x-1,y-1) + 2f(x-1,y) + f(x-1,y+1) - f(x+1,y-1) - 2f(x+1,y) - f(x+1,y+1) \quad (7)$$

最后求得梯度及其方向,即

$$\text{梯度: } G[f(x,y)] = |\Delta x f(x,y)| + |\Delta y f(x,y)| \quad (8)$$

$$\text{方向: } \theta = \arctan(y/x) \quad (9)$$

Sobel 算子边缘检测后获取的笔迹边缘图像如图 6 所示。

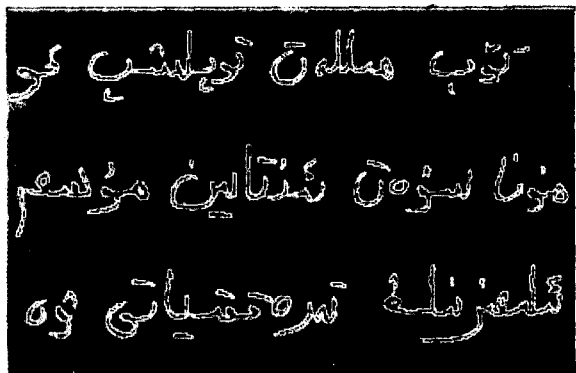


图 6 笔迹样本边缘检测

3 结束语

文中研究了维吾尔文笔迹鉴别中笔迹图像的预处理和笔迹边缘提取的方法,在图像处理算法的基础上进行了大量的实验与比较,提出了有效去除维吾尔文笔迹背景噪声的方法,并采用最合适的图像处理算法实现笔迹的二值化,最后比较并采用 Sobel 算子获取了笔迹的边缘图像。

该维吾尔文笔迹鉴别的预处理与边缘提取是在 VC++6.0 开发平台上实现的。在笔迹鉴别的整个过程中,预处理和边缘提取都是为了后面的笔迹特征提取打下基础,笔迹的特征是进行鉴别的关键内容,因此,特征提取是维吾尔文笔迹鉴别中的重点与难点,也是正在研究的内容与方向。

参考文献:

- [1] 张慧档,贺显曜. 基于小波变换和神经网络集成的笔迹鉴别方法[J]. 计算机应用研究, 2008, 25(3): 741-742.
- [2] 师宝山,张贵州. 笔迹鉴别预处理算法的设计与实现[J]. 电子器件, 2008, 31(4): 1357-1358.
- [3] 张德贤,郭小波,刘永平. 笔迹鉴别预处理与特征抽取技术研究[J]. 微计算机信息, 2006, 22(09S): 310-311.
- [4] 李. 莹. 汉字笔迹鉴别的算法研究[D]. 济南: 山东大学, 2007.
- [5] 王凤岭,刘连芳,蒋宗礼,等. 离线手写体笔迹鉴别方法研究[J]. 计算机工程与设计, 2006, 27(14): 2581-2582.
- [6] Bulacu M, Schomaker L, Vuurpijl L. Writer identification using edge-based directional features[C]//Proceedings of the 7th International Conference on Document Analysis and Recognition. Edinburgh, UK: IEEE, 2003: 937-941.
- [7] 白雪冰,张. 伟. 汉字手写笔迹鉴别预处理算法的研究[J]. 计算机工程与设计, 2009, 30(22): 5189-5190.
- [8] Bulacu M, Schomaker L. Text-independent writer identification and verification using textural and allographic features[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2007, 29(4): 701-717.
- [9] 沈. 聪. 基于改进的多通道 Gabor 小波变换的笔迹鉴别[D]. 北京: 北京工业大学, 2002.
- [10] 周长发. 精通 visual c++ 图像处理编程[M]. 第 2 版, 北京: 电子工业出版社, 2004: 229-290.
- [11] Bulacu M, Schomaker L, Brink A. Text-independent writer identification and verification on offline Arabic handwriting[C]//Proceedings of the 9th International Conference on Document Analysis and Recognition. Curitiba, Brazil: IEEE, 2007: 769-773.
- [12] 李. 昕,丁晓青,彭良瑞. 一种基于微结构特征的多文种文本无关笔迹鉴别方法[J]. 自动化学报, 2009, 35(9): 1200-1201.