

一种复杂表格识别和处理方法

张平 黄尚廉 潘保昌

(重庆大学光电精密机械研究所 重庆 630044)

摘要 本文提出了一种复杂表格识别和处理方法。该方法首先输入一张未填有用信息的空表格作背景信息表格,然后输入已填有用信息的同类表格作有用信息表格,对有用信息表格图象采用几何位置配准使两幅表格的背景信息重合,然后利用数字逻辑运算和智能相关处理技术以达到对残留背景信息的彻底清除,保留后填入的有用信息供识别处理。此方法对表格形式不限,适用于背景信息和有用信息采用同色或异色笔填写。方法仅需第一次输入一张空表格后,就可连续处理同类型表格。文中给出了表格处理结果。

关键词 表格处理,文件识别,数字逻辑运算,智能相关处理

1 引言

在数字、文字识别领域中大部分情况是待识别和统计的数字、文字(称有用信息)要填入各种各样表格中。由于所填内容与表格原有信息(原有信息可能是各种线段组成的框架、文字、数字和符号,统称背景信息)混合于一幅图象中,如不正确加以区分,必将给识别有用信息带来一定困难,以致不能正确有效地进行有用信息的识别和各种处理。

现有技术的一种方法是:各种表格和背景信息采用一种颜色绘制或印刷,而使用者需填入的内容用黑色或与背景信息相异的笔填写,表格背景颜色与有用信息颜色为非同一颜色。表格输入各种表格处理机时,采用附加光学处理装置,自动滤出背景信息,只保留有用信息直接供识别机进行识别处理。

但这种方法有以下不足:(1)表格采用异色笔绘制,不便印制和复印,成本相对昂贵;(2)信息输入时需附加光学滤色装置;(3)人们日常事务中绝大多数都是采用黑色表格,改变表格颜色,实现办公自动化,用户难以接受。

本文正是为了克服上述表格处理方法不足而提出的一种复杂表格识别和处理方法。该方法能快速有效地处理各种表格,并根据需要快速自动地滤出各种框架、文字、数字和符号等背景信息,保留后填入的有用信息。此方法适用于有用信息和背景信息采用异色或同色笔填写,即该方法能处理黑色表格,是文件表格识别不可缺少的预处理。

2 处理方法

本方法可分为以下几个步骤,即(1)表格信息输入,(2)背景表格与信息表格图象

1993-03-17收到,1993-08-03定稿

张平 男,1963年生,讲师,现从事图象处理,模式识别方面的研究工作。

黄尚廉 男,1937年生,教授,博士生导师,长期从事光电技术及系统方面的研究。

潘保昌 男,1949年生,博士,教授,长期从事人工视觉,模式识别方面的研究。

表 2

年 龄	总 人 口		占总人口的百分比	年 龄	总 人 口		占总人口的百分比
	合 计	其中: 男			合 计	其中: 男	
72	800634	52859	89	033	230	989	
73	107573	48992	90	030	147	658	
74	121310	46734	91	029	74	371	
75	80936	39844	92	026	51	246	
67	383764	33661	93	022	40	194	
68	121310	28664	94	019	27	133	
69	70841	25333	95	017	17	78	
70	121310	21292	96	015	16	90	
71	112906	18127	97	013	3	47	
76	107573	14749	98	011	8	40	
77	95837	11883	99	009	6	30	
78	112906	9373	100	007	2	5	

用线段跟踪算法^[1-3]找出背景表格图象 I_X 的外框架 2 个最远边缘点坐标值 (x_{00}, y_{00}) , (x_{01}, y_{01}) 和有用信息表格图象 I_Y 对应的 2 个外框架最远边缘点坐标值 (u_{00}, v_{00}) , (u_{01}, v_{01}) 。边缘点的选取方法为寻找表格外框架上、下、左、右中任一方向直线最远间距的横竖线段 2 个交叉点边缘点坐标值。那么两幅图象的角度差:

$$\varphi = \arctg \left(\frac{v_{01} - v_{00}}{u_{01} - u_{00}} \right) - \arctg \left(\frac{y_{01} - y_{00}}{x_{01} - x_{00}} \right). \quad (1)$$

对信息表格图象 $I_Y(u, v)$ 中的每点 (u, v) , 首先以 (u_{00}, v_{00}) 为中心旋转 φ 角度后, 然后再平移至 (x_{00}, y_{00}) , 生成几何校正后新图象用 $I_Z(x, y)$ 表示。

对 $I_Y(u, v)$ 中每点阵 (u, v) 计算

$$t_x = u - u_{00}, \quad (2)$$

$$t_y = v - v_{00}; \quad (2a)$$

$$u_1 = t_x \times \cos(\varphi) - t_y \times \sin(\varphi), \quad (3)$$

$$v_1 = t_x \times \sin(\varphi) + t_y \times \cos(\varphi); \quad (3a)$$

$$x = u_1 + x_{00}, \quad (4)$$

$$y = v_1 + y_{00}. \quad (4a)$$

对(3)和(3a)式中, u_1, v_1 计算进行取整运算, 使计算新的象素点落在图象点阵中。

信息表格图象 $I_Y(u, v)$ 经几何校正后的图象 $I_Z(x, y)$ 其背景信息在空间位置上与背景信息表格图象 $I_X(x, y)$ 大部分重合。

2.3 数字逻辑运算处理

对两幅配准的背景表格图象 I_X 和信息表格图象 I_Z , 其背景信息在空间域位置上是重合的, 用数字逻辑处理方法可消除大部分背景信息, 其逻辑运算关系见表 3。

表 3 数字逻辑运算关系表

区 域	信息表格数值	背景表格数值	逻辑运算结果图象
背景信息	1	1	0
有用和部分 残留背景信息	1	0	1
背景表格 残留信息	0	1	0
无信息区域	0	0	0

表 3 中第一栏: 信息表格数据和背景表格数据在空间位置上都为 1, 只能是背景信息数据, 处理结果为 0, 以此消除大部分背景信息。

第二栏: 信息表格数据为 1, 背景表格数据为 0, 此集合包括了全部后填入的有用信息和由于两表格匹配误差, 而信息表格保留的小部分残留的背景数据, 处理结果为 1。

第三栏: 背景表格数据为 1, 而信息表格数据为 0, 这也是失匹产生的, 处理结果为 0。

第四栏: 为两表格空白区域。

经过上述逻辑运算, 只有第二栏保留了全部有用信息数据, 但同时也混进了小部分背景信息数据。逻辑运算后所得结果图象数据为 $IW(x, y)$ 。其数学表达式为

$$\left. \begin{aligned} & \forall (x, y) \in IZ, \\ & \{ [IX(x, y) = 0] \wedge [IZ(x, y) = 0] \} \rightarrow IW(x, y) = 0, \\ & \{ [IX(x, y) = 1] \wedge [IZ(x, y) = 1] \} \rightarrow IW(x, y) = 0, \\ & \{ [IX(x, y) = 0] \wedge [IZ(x, y) = 1] \} \rightarrow IW(x, y) = 1, \\ & \{ [IX(x, y) = 1] \wedge [IZ(x, y) = 0] \} \rightarrow IW(x, y) = 0. \end{aligned} \right\} \quad (5)$$

上述数字逻辑运算可用数字电路实现, 以提高图象处理速度。

2.4 智能相关处理

$IW(x, y)$ 数据保留了全部后填入的有用信息, 但也留下了部分由于图象输入时 CCD 扫描误差、寻找 2 个边缘点坐标精度误差、成象透视投影失真以及几何位置校正误差而产生的残留背景信息。根据在 $IW(x, y)$ 中残留背景信息在空间位置上与 $IX(x, y)$ 中的背景信息是连通的, 即相关联这一特点, 采用 $IX(x, y)$ 背景图象数据往上、下、左、右 4 个方向分别位移 $\pm \Delta x, \pm \Delta y$ 生成的背景位移图象 $IX(x \pm \Delta x, y \pm \Delta y)$, 再分别与 $IW(x, y)$ 进行第 2.3 节所示迭代数字逻辑运算, 以彻底消除残留的背景信息, 其处理最终结果图象 IW 只剩下有用信息。位移量 $\Delta x, \Delta y$ 的计算可利用编程方法计算 IW 中残留背景信息的最大宽度而得到^[9]。此时的 IW 图象数据, 根据用户需要可存入有用信息存贮体或计算机内存。其运算流程图见图 2。

2.5 有用信息输出

经上述方法处理后的输出图象数据 $IW(x, y)$ 只保留了后填入的有用信息, 可以输出至各种文字、数字识别机进行自动识别, 也可反馈至有用信息存贮体存贮。有用信息的

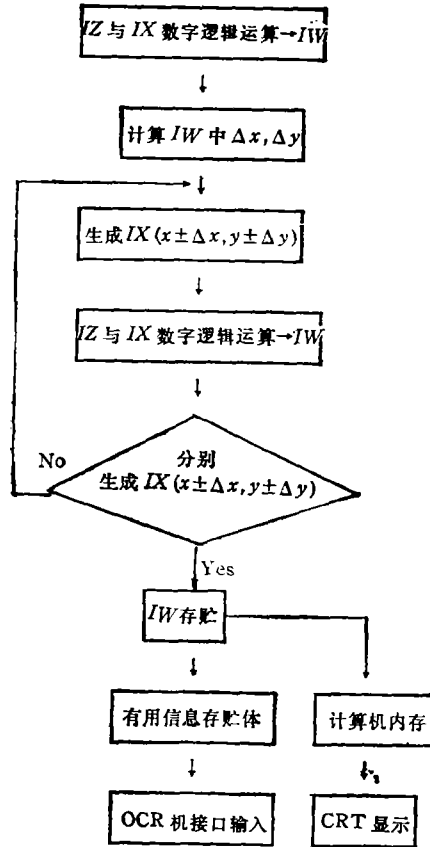


图 2 智能相关处理运算流程图

表 4

72	890634	52833	89	033	230	989
73	107573	48992	90	030	147	658
74	121310	46734	91	029	74	371
75	80956	39844	92	026	51	246
67	583764	33661	93	022	40	194
68	121310	28664	94	019	27	133
69	70841	25333	95	017	17	78
70	121310	21292	96	015	16	90
71	112906	18127	97	013	5	47
76	107573	14749	98	011	8	40
77	95837	11883	99	009	6	30
78	112906	9373	100	007	2	5

结果图象输出见表 4。

3 结论

文件版面处理是国内外新兴的领域,它涉及模式识别,图象处理,人工智能等学科.主

要研究文件结构分析,版面处理等方面。对它的深入研究必将推动办公自动化的蓬勃发展。本方法提出的文件表格自动清除背景的方法能处理包括黑色表格在内的各种复杂表格。此方法仅需第一次输入一张空表格后,即可连续处理同类型表格。我们完成了表格处理系统,并编制了整套软件,图象处理硬件正在研究之中^[5]。一幅表格在 IBM-PC386 机上平均处理速度为 15s。该方法的实现必将推广表格处理机和各种 OCR 机的广泛应用。

参 考 文 献

- [1] [日]田村秀行著,赫荣威译. 计算机图象处理技术. 北京: 北京师范大学出版社, 1988, 107—120.
- [2] 关直树. Method for Processing Hidden Lines of Figures. 国际出愿番号 PCT/JP88/00966, 国际公开番号 WO89/03095 (国际专利).
- [3] 小松浩一. Edge Information Extracting Device and Method Thereof. 国际出愿番号 PCT/JP91/00189, 国际公开番号 WO91/12585 (国际专利).
- [4] Yoh-Han Pao. Adaptive Pattern Recognition and Neural Network. Addison-Wesley Publishing Company Inc. 1989, 150—170.
- [5] 费旭东, 荆仁杰, 等. 电子学报, 1992, 20(4): 14—18.

A RECOGNITION AND PROCESSING METHOD OF COMPLEX FORM

Zhang Ping Huang Shanglian Pan Baochang
(Chongqing University, Chongqing 630044)

Abstract A recognition and processing method of document form background information is presented. First, an original form, called Background Information Form Image (BIFI), is scanned. The same kind of forms, which is filled by various characters that one wants to recognize and process, called Useful Information Form Image (UIFI), is stored. Then, a geometric remedy of UIFI is used for the best matching with the BIFI. A digital logical operation is carried out in order to erase the most of background information and an intelligent correlation technique is available to completely delete the remained background information. This method is suitable for processing various forms, including background and useful information printed or filled with same colour. It can continuously process same kind of UIFIs with the input of BIFI only once. The document forms are shown and satisfactory results are obtained.

Key words Form processing, Document recognition, Digital logical operation, Intelligent correlation processing