

## EBTool使用手册

本工具采用的是字符分离再识别方向的OCR技术，至于其它更牛逼的OCR方法这里暂且不谈。

功能:

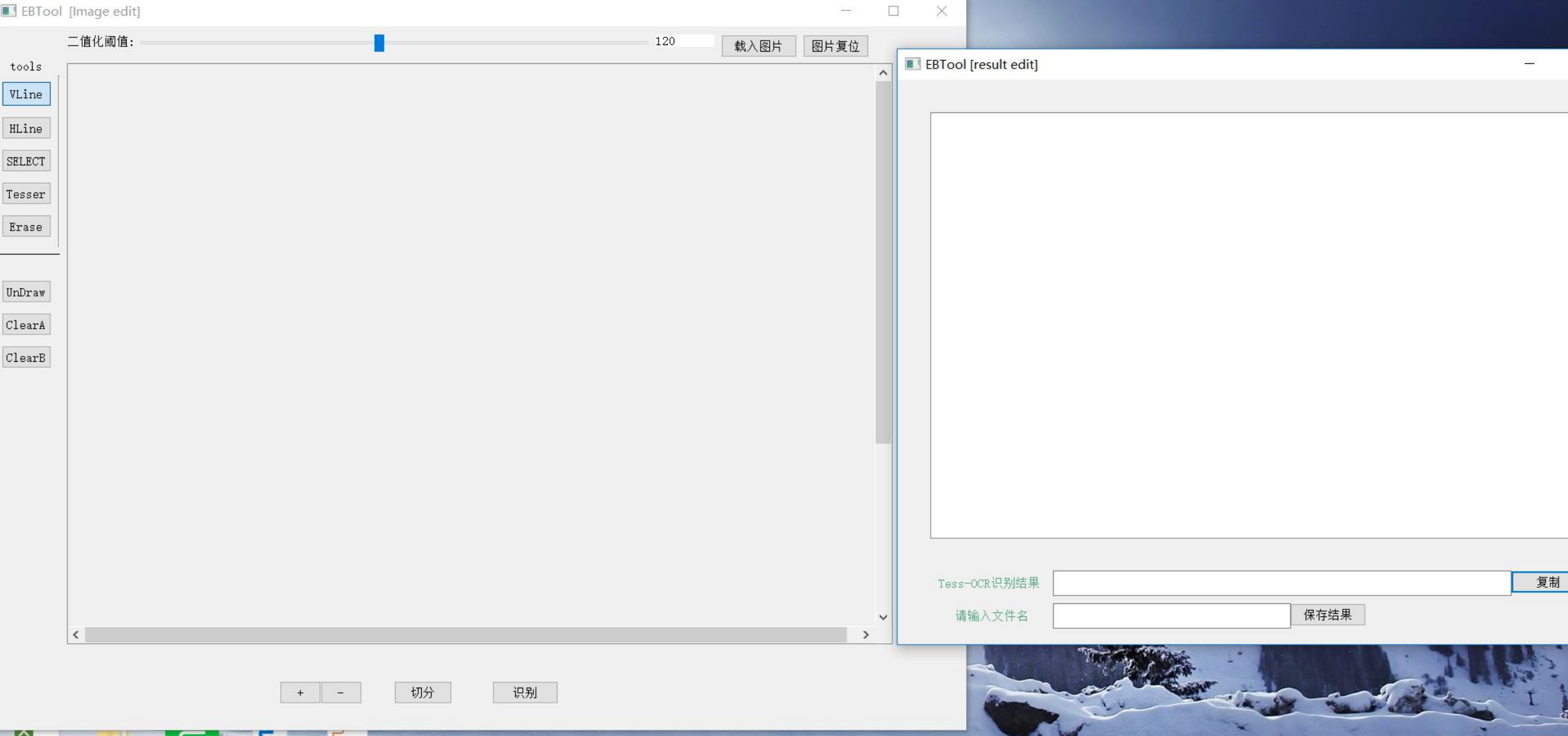
辅助提取图片中的表格内容。其它方法无法提取的pdf，可以通过ghostscript转图片后再通过此工具提取。--技术有限、无法开发出复杂图片环境下全自动化高识别率的工具，只能辅助工具了。

主要方式:

本着低成本原则，采用开源的tesseract对中文进行识别，自建模型对数字识别（需训练样本少，且准确率高）。

基于OCR识别时复杂环境下准确率不能100%问题，以及表格非标的问题，提供了识别结果人工编辑校正功能。

**使用:**



二值化阈值:

120

载入图片

图片复位

tools

VLine

画水平线

HLine

画竖直线

SELECT

扣出选择区域并放大

Tesser

tesseract-ocr识别选择区域内容

Erase

擦除选择区域内所包含的 vline,hline

UnDraw

取消最近一次所画的线

ClearA

清空所有线

ClearB

清空点击切分后自动生成的线

图片二值化阈值设置

载入图片

重新载入原图

自动切分图片

获取ocr结果

图片放大

图片缩小

切分

识别

+

-

add col  
del col  
add row  
del row  
add head col  
add head row  
copy  
paste  
copy from tess

右击，弹出菜单栏

add head col/row 解决添加时无法添加第一列/行问题

把识别结果直接载入光标所在单元格  
(与paste区别是，paste碰到\n,\t会换行换列，这里内容只在单元格)

保存结果

识别结果复制到剪切板

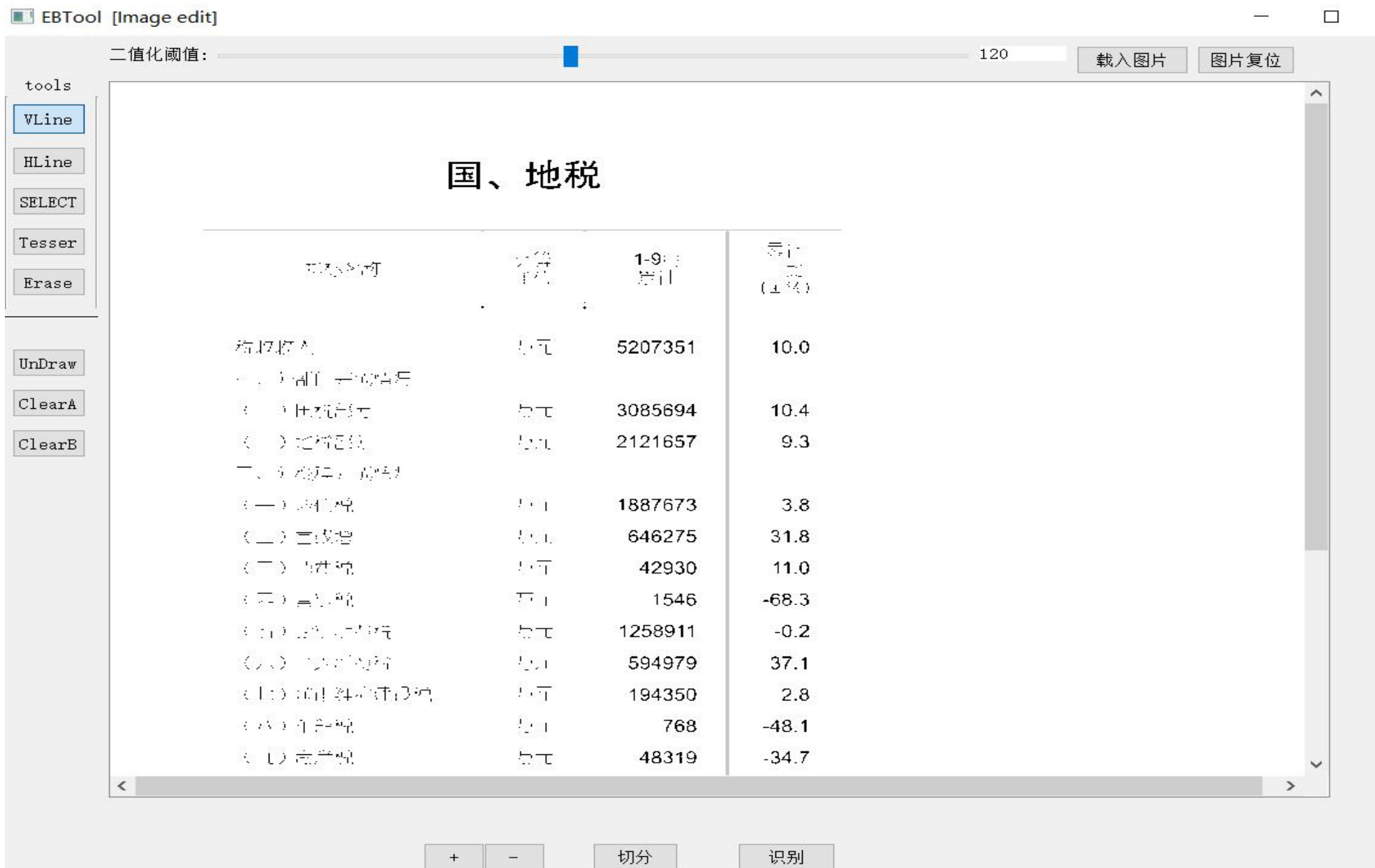
Tess-OCR识别结果

请输入文件名

保存结果

复制

## 使用demo: 载入图片



## 调整二值化阈值

EBTool [Image edit]

二值化阈值:  186

载入图片 图片复位

tools

- VLine
- HLine
- SELECT
- Tesser
- Erase
- UnDraw
- ClearA
- ClearB

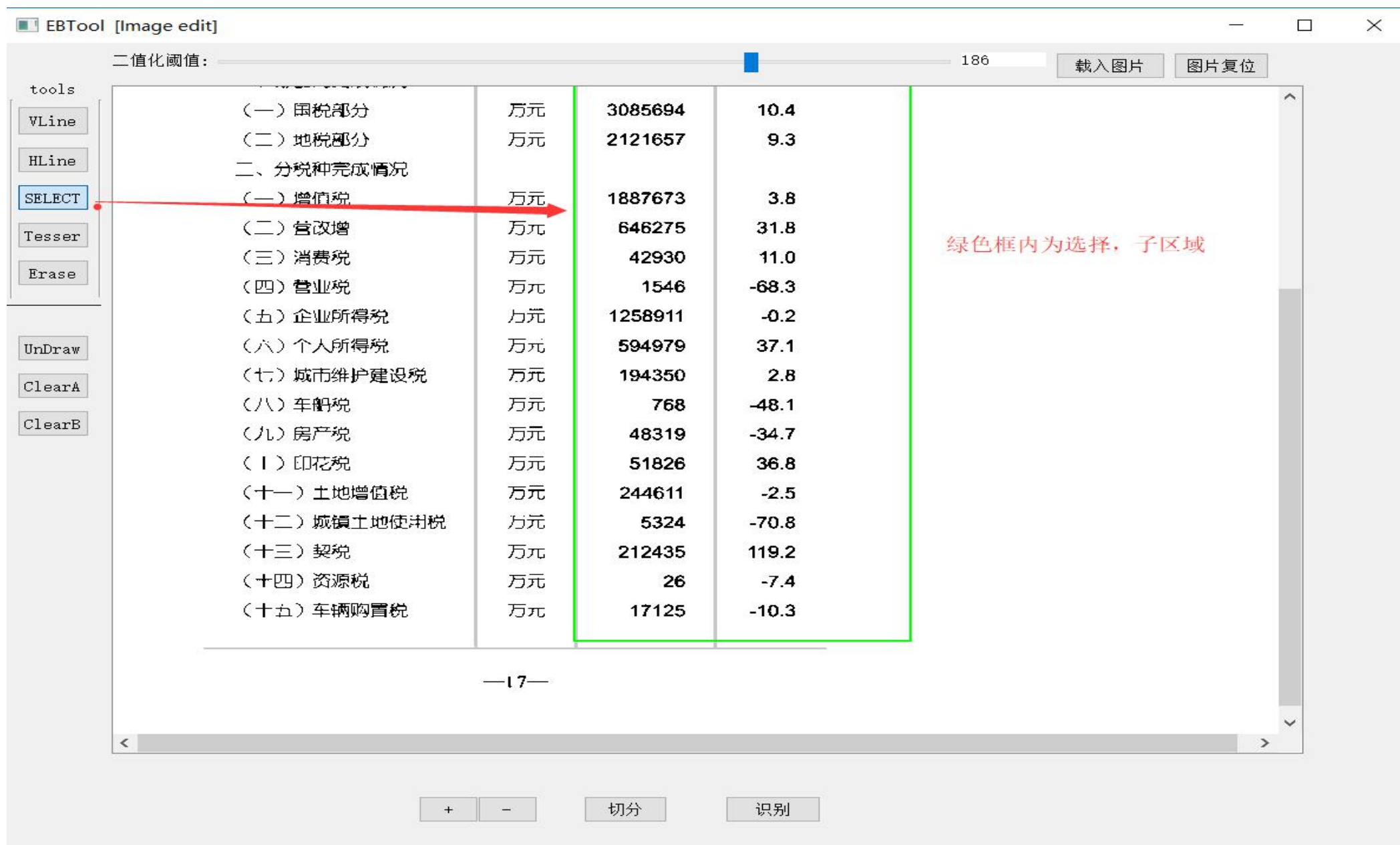
### 国、地税

指标名称	计算单位	1-9月累计	累计同比(±%)
税收收入	万元	5207351	10.0
一、分部门完成情况			
（一）国税部分	万元	3085694	10.4
（二）地税部分	万元	2121657	9.3
二、分税种完成情况			
（一）增值税	万元	1887673	3.8
（二）营改增	万元	646275	31.8
（三）消费税	万元	42930	11.0
（四）营业税	万元	1546	-68.3
（五）企业所得税	万元	1258911	-0.2
（六）个人所得税	万元	594979	37.1
（七）城市维护建设税	万元	194350	2.8
（八）车船税	万元	768	-48.1
（九）房产税	万元	48319	-34.7

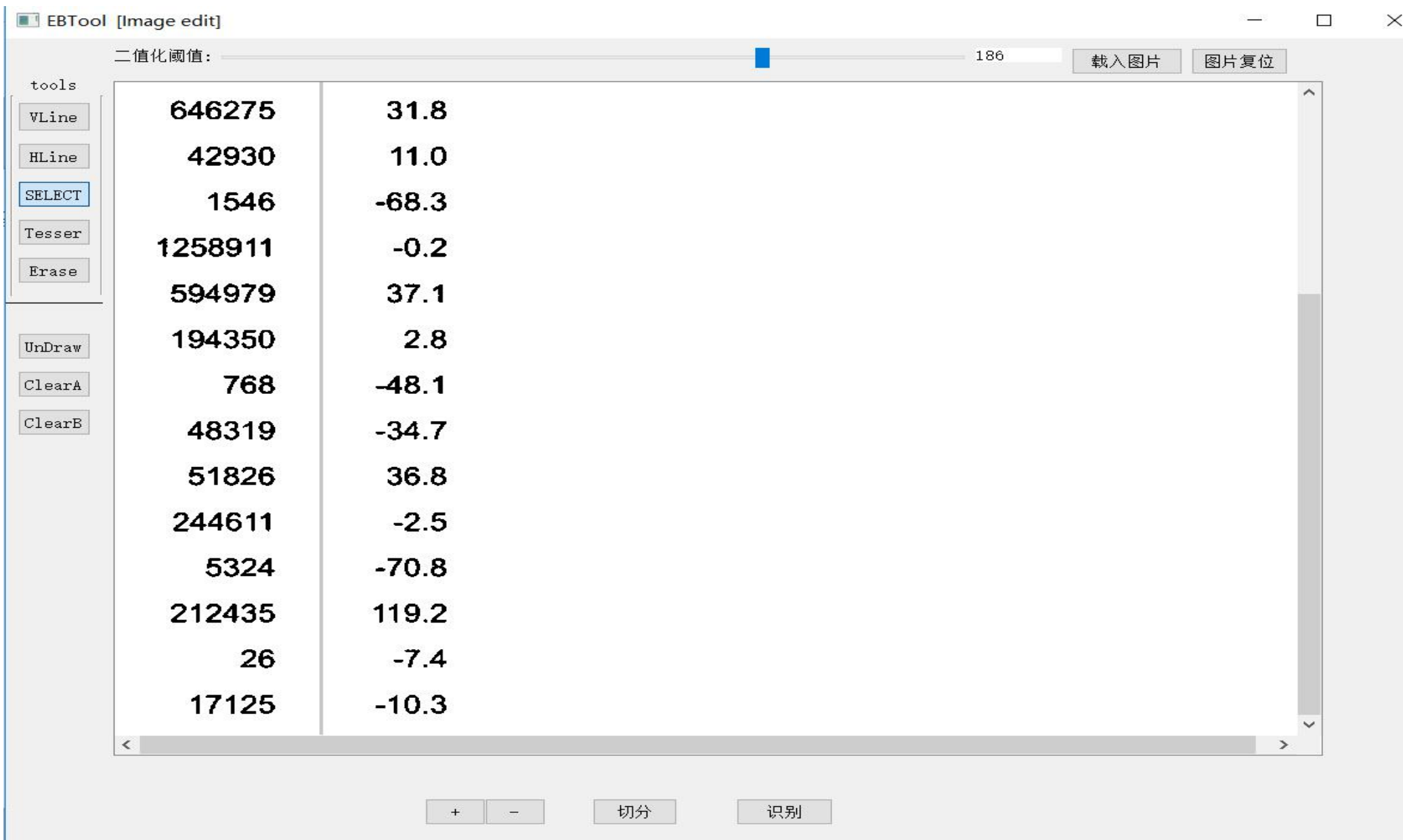
阈值调大，图片清晰了

+ - 切分 识别

选择子区域：点击select后，光标移动至画图区，按下鼠标左键拖动画图。

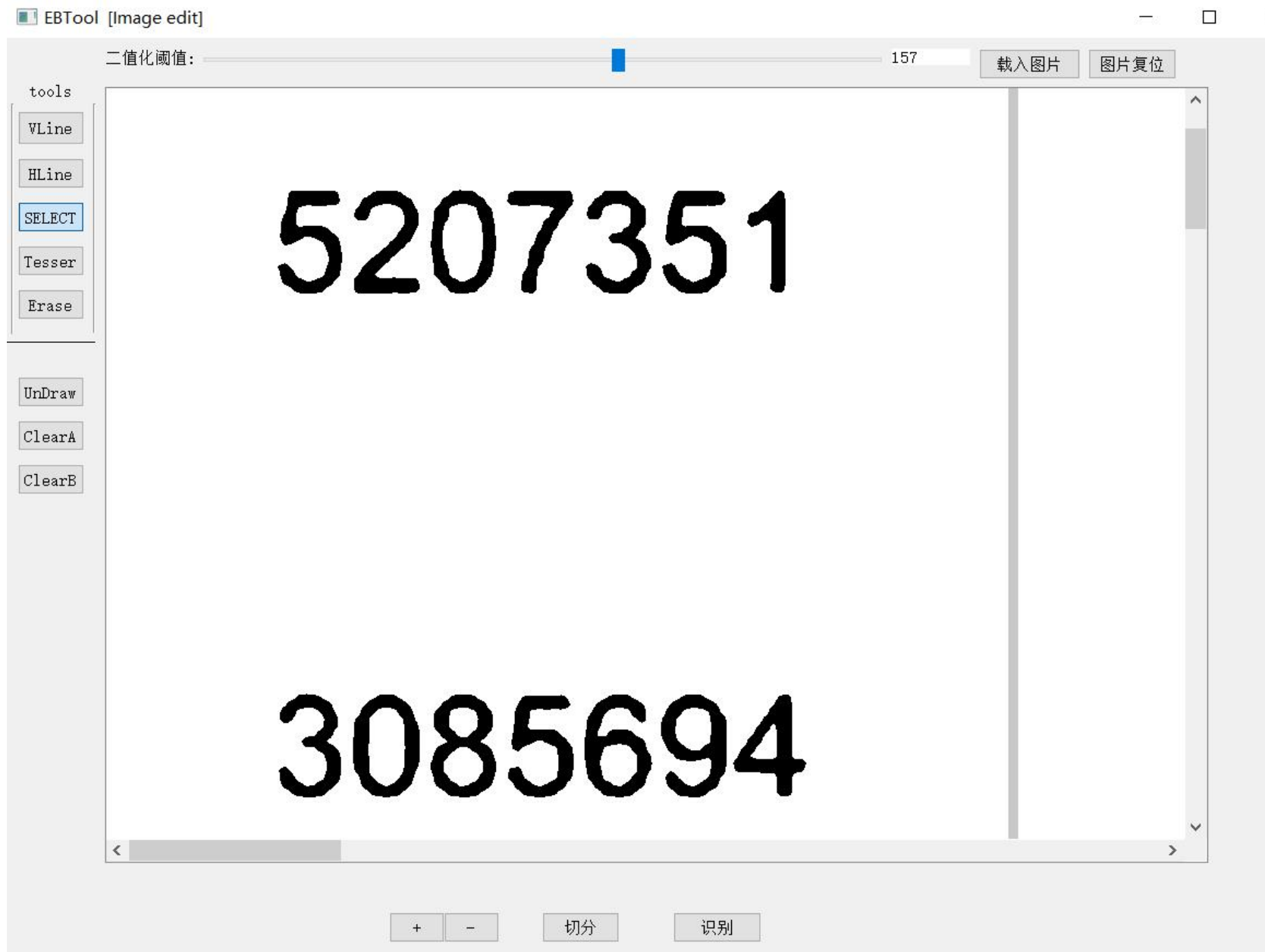


## 跳转到选择的子区域



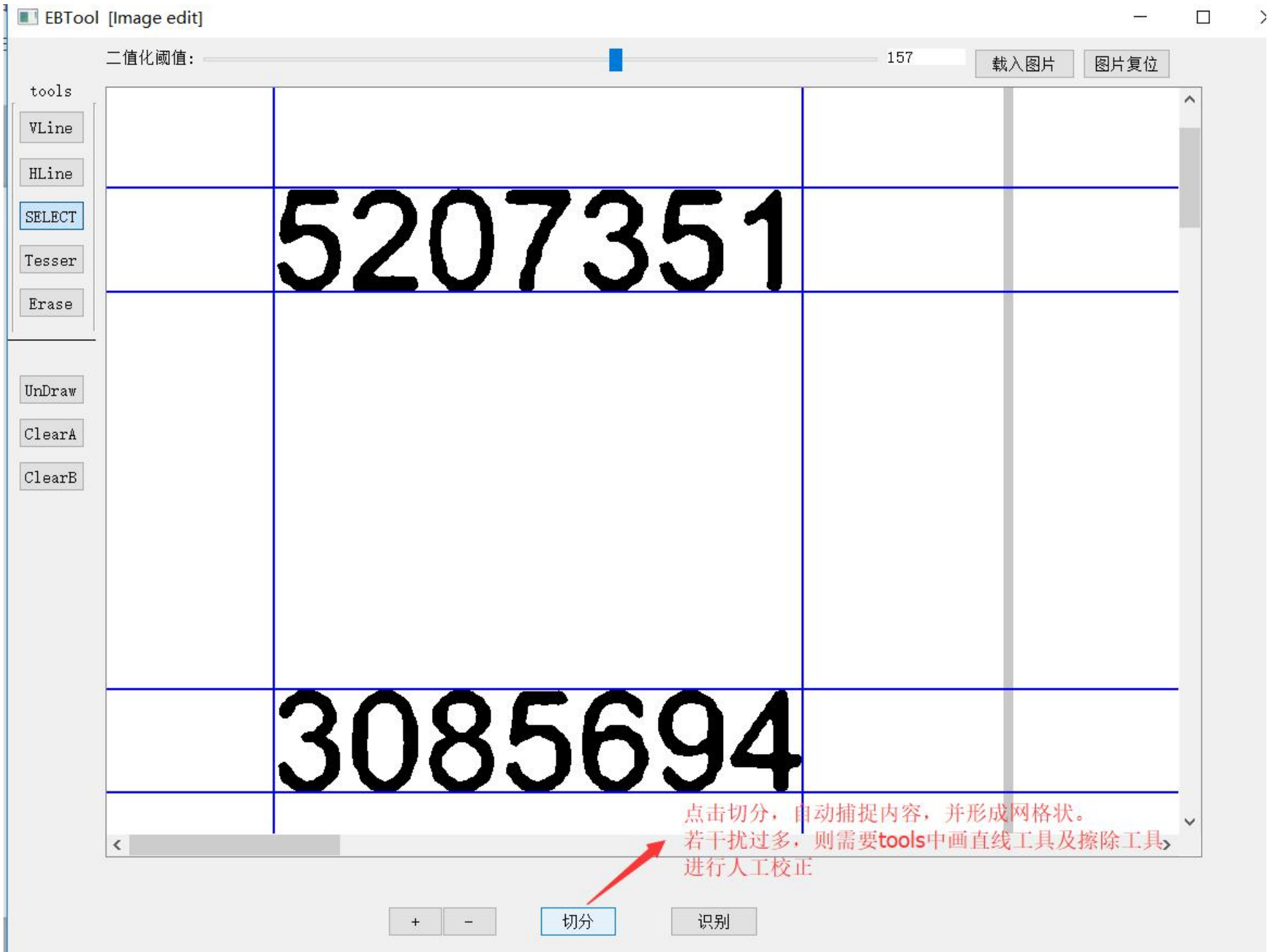


放大图片(+),调整阈值,使数字尽量不黏连。(一般 图片放越大越好, 阈值160左右, 字体看起来不带毛边且圆润即可)



调整后字体特征明显，且不黏连

识别时单元格形式和直线(绿, 蓝)网格类似(灰色直线忽略)



二值化阈值:

157

载入图片

图片复位

tools

VLine

HLine

SELECT

Tesseract

Erase

UnDraw

ClearA

ClearB

5207351

3085694

点击识别后，右边生成识别结果

EBTool [result edit]

	1	2
1	5207351	10.0
2	3085694	10.4
3	2121657	9.3
4	1887673	3.8
5	646275	31.8
6	42930	11.0
7	1546	-68.3
8	1258911	-0.2
9	594979	37.1
10	194350	2.8
11	768	-48.1
12	48319	-34.7
13	51826	36.8

Tess-OCR识别结果

请输入文件名

保存结果

+

-

切分

识别

二值化阈值:

157

载入图片

图片复位

tools

VLine

HLine

SELECT

Tesser

Erase

UnDraw

ClearA

ClearB

## 国、地税

指标名称	计算单位	1-9月累计	累计可比(±%)
税收收入	万元	5207351	10.0
一、分部门完成情况			
（一）国税部分	万元	3085694	10.4
（二）地税部分	万元	2121657	9.3
二、分税种完成情况			
（一）增值税	万元	1887673	3.8
（二）营改增	万元	646275	31.8
（三）消费税	万元	42930	11.0
（四）营业税	万元	1546	-68.3
（五）企业所得税	万元	1258911	-0.2
（六）个人所得税	万元	594979	37.1
（七）城市维护建设税	万元	194350	2.8
（八）车船税	万元	768	-48.1
（九）房产税	万元	48319	-34.7

点击图片复位

+

-

切分

识别

tools

VLine

HLine

SELECT

Tesser

Erase

UnDraw

ClearA

ClearB

## 国、地税

指标名称	计算单位	1-9月累计	累计同比(± %)
税收收入	万元	5207351	10.0
一、分部门完成情况			
（一）国税部分	万元	3085694	10.4
（二）地税部分	万元	2121657	9.3
二、分税种完成情况			
（一）增值税	万元	1887673	3.8
（二）营改增	万元	646275	31.8
（三）消费税	万元	42930	11.0
（四）营业税	万元	1546	-68.3
（五）企业所得税	万元	1258911	-0.2

调整阈值及大小

编辑结果编辑工具，添加横列

国、地税

指标名称	计算单位	1-9月累计	累计同比(± %)
税收收入	万元	5207351	10.0
一、分部门完成情况			
（一）国税部分	万元	3085694	10.4
（二）地税部分	万元	2121657	9.3
二、分税种完成情况			
（一）增值税	万元	1887673	3.8
（二）营改增	万元	646275	31.8
（三）消费税	万元	42930	11.0
（四）营业税	万元	1546	-68.3
（五）企业所得税	万元	1258911	-0.2
（六）个人所得税	万元	594979	37.1

EBTool [result edit]

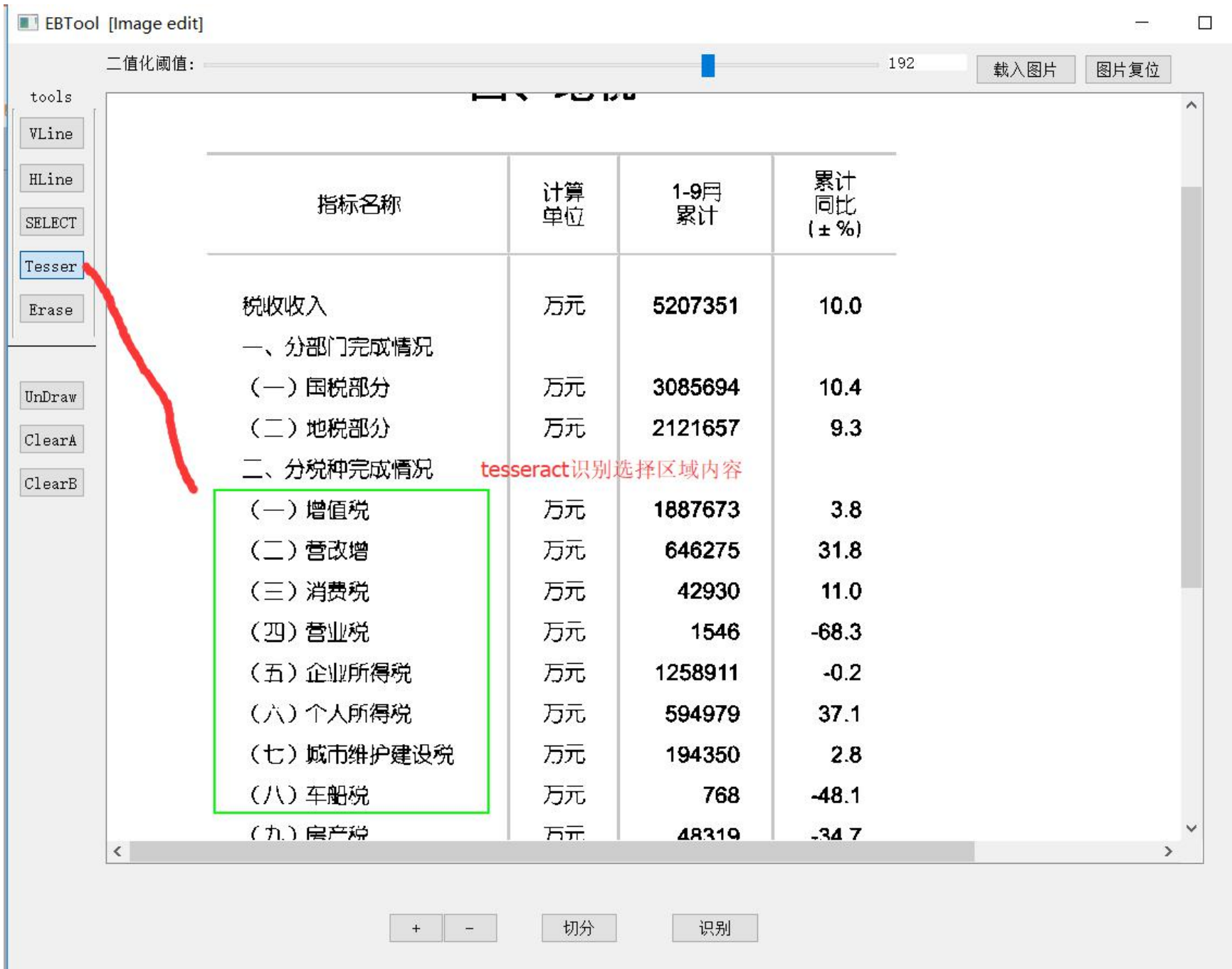
	1	2	3	4
1				
2			5207351	10.0
3			3085694	10.4
4			2121657	9.3
5			1887673	3.8
6			646275	31.8
7			42930	11.0
8			1546	-68.3
9			1258911	-0.2
10			594979	37.1
11			194350	2.8
12			768	-48.1
13			48319	-34.7

Tess-OCR识别结果

复制

请输入文件名

保存结果



迟  
tesseract接口识别较慢，可能有一定的延



## 二、分税种完成情况

(一) 增值税	万元	1887673	3.8
(二) 营改增	万元	646275	31.8
(三) 消费税	万元	42930	11.0
(四) 营业税	万元	1546	-68.3
(五) 企业所得税	万元	1258911	-0.2
(六) 个人所得税	万元	594979	37.1
(七) 城市维护建设税	万元	194350	2.8
(八) 车船税	万元	768	-48.1
(九) 房产税	万元	48319	-34.7
(十) 印花税	万元	51826	36.8

4		2121657	9.3
5		1887673	3.8
6		646275	31.8
7		42930	11.0
8		1546	-68.3
9		1258911	-0.2
10		594979	37.1
11		194350	2.8
12		768	-48.1
13		48319	-34.7

识别结果

Tess-OCR识别结果

莒血说 <丑> 企山所得税 <八> 个人所信碗 <亡> 贼币维沪建设税 <八> 车船皖

请输入文件名

保存结果

+

-

切分

识别



EBTool [result edit]

	1	2	3	4
1				
2			5207351	10.0
3			3085694	10.4
4			2121657	9.3
5			1887673	3.8
6			646275	31.8
7			42930	11.0
8			1546	-68.3
9			1258911	-0.2
10			594979	37.1
11			194350	2.8
12			768	-48.1
13			48319	-34.7

点击复制，内容复制到剪切板

Tess-OCR识别结果 莒皿说 <丑> 企山所得税 <八> 个人所信碗 <亡> 贼币维沪建设税 <八> 车船皖 复制

请输入文件名  保存结果

EBTool [result edit]

选择单元格，右击，弹出菜单栏，选择paste

	1	2	3	4
2			5207351	10.0
3			3085694	10.4
4			2121657	9.3
5			1887673	3.8
6			646275	31.8
7			42930	11.0
8			1546	-68.3
9			1258911	-0.2
10			594979	37.1
11			194350	2.8
12			768	-48.1
13			48319	-34.7

- add col
- del col
- add row
- del row
- add head col
- add head row
- copy
- paste
- copy from tess

Tess-OCR识别结果 莒皿说 <丑> 企山所得税 <八> 个人所信碗 <亡> 贼币维沪建设税 <八> 车船皖 复制

请输入文件名  保存结果

选择单元格不要双击进入单元格编辑模式，不然相  
copy from tess 操作

	1	2	3	4
4			2121657	9.3
5	〈一〉墙但税		1887673	3.8
6	〈二〉莒改墙		646275	31.8
7	〈三〉消费说		42930	11.0
8	〈四〉莒皿说		1546	-68.3
9	〈丑〉企山所得...		1258911	-0.2
10	〈八〉个人所信碗		594979	37.1
11	〈亡〉贼币维沪...		194350	2.8
12	〈八〉车船皖		768	-48.1
13			48319	-34.7
14			51826	36.8
15			244611	-2.5
16			5324	-70.8

Tess-OCR识别结果 莒皿说 〈丑〉企山所得税 〈八〉个人所信碗 〈亡〉贼币维沪建设税 〈八〉车船皖 复制

请输入文件名

保存结果

结果，需要人工校正。  
目前只支持竖直批量处理  
tesseract结果。横向请逐  
个操作。

该单元格复制粘贴类似  
excel内形式。可与通过复  
制粘贴excel交互。

EBTool [result edit]

	1	2	3	4
1				
2			207351	10.0
3			085694	10.4
4			121657	9.3
5	〈一〉墙但税		887673	3.8
6	〈二〉莒改墙		46275	31.8
7	〈三〉消费说		42930	11.0
8	〈四〉莒皿说		1546	-68.3
9	〈丑〉企山所得...		1258911	-0.2
10	〈八〉个人所信碗		594979	37.1
11	〈亡〉贼币维沪...		194350	2.8
12	〈八〉车船碗		768	-48.1
13			48319	-34.7

add col  
del col  
add row  
del row  
add head col  
add head row  
copy  
paste  
copy from tess

Tess-OCR识别结果 指标名枷 复制

请输入文件名 保存结果

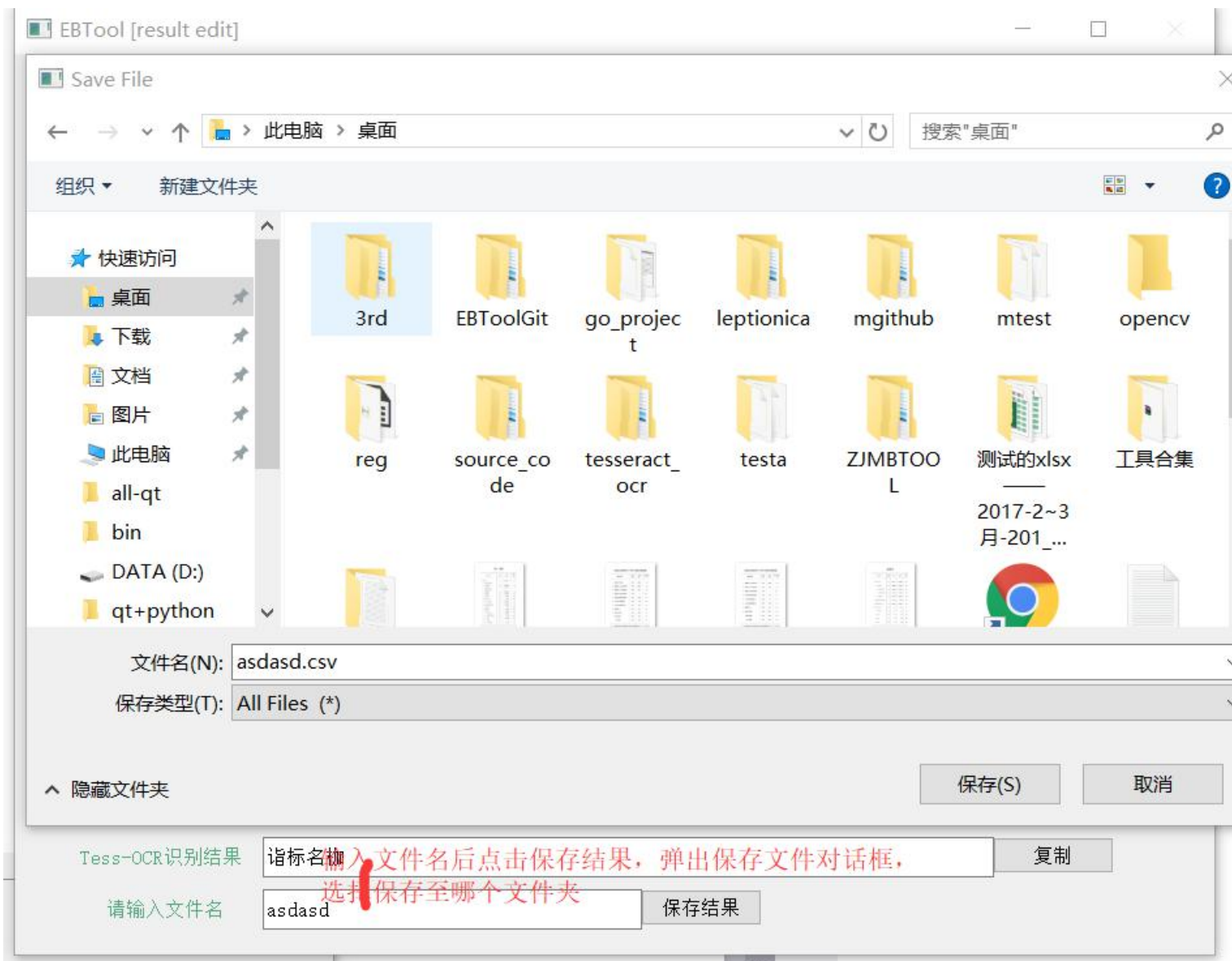
EBTool [result edit]

	1	2	3	4
1	指标名枷			
2			5207351	10.0
3			3085694	10.4
4			2121657	9.3
5	〈一〉墙但税		1887673	3.8
6	〈二〉莒改墙		646275	31.8
7	〈三〉消费说		42930	11.0
8	〈四〉莒皿说		1546	-68.3
9	〈丑〉企山所得...		1258911	-0.2
10	〈八〉个人所信碗		594979	37.1
11	〈亡〉贼币维沪...		194350	2.8
12	〈八〉车船碗		768	-48.1
13			48319	-34.7

Tess-OCR识别结果 指标名枷

请输入文件名 保存结果

直接从识别结果复制到单元格内(忽略换行，  
换列符号)



Microsoft Excel interface showing the '开始' (Home) tab. The ribbon includes options for '文件' (File), '开始' (Home), '插入' (Insert), '页面布局' (Layout), '公式' (Formulas), '数据' (Data), '审阅' (Review), '视图' (View), and '帮助' (Help). The '开始' tab is active, showing options for '剪贴板' (Clipboard), '字体' (Font), and '对齐方式' (Alignment). The font settings are set to '等线' (Dengxian) and size 11. The alignment settings are set to '左对齐' (Left Align).

The active cell is A1, containing the text '诣标名枷'. The formula bar shows the text '诣标名枷'.

	A	B	C	D	E	F	G
1	诣标名枷						
2							
3			5207351	10			
4			3085694	10.4			
5			2121657	9.3			
6	〈一〉墙但税		1887673	3.8			
7	〈二〉莒改墙		646275	31.8			
8	〈三〉消费说		42930	11			
9	〈四〉莒皿说		1546	-68.3			
10	〈丑〉企山所得税		1258911	-0.2			
11	〈八〉个人所信碗		594979	37.1			
12	〈亡〉贼币维沪建设		194350	2.8			
13	〈八〉车船皖		768	-48.1			
14			48319	-34.7			
15			51826	36.8			
16			244611	-2.5			
17			5324	-70.8			
18			212435	119.2			
19			26	-7.4			
20			17125	-10.3			
21							
22							
23							
24							

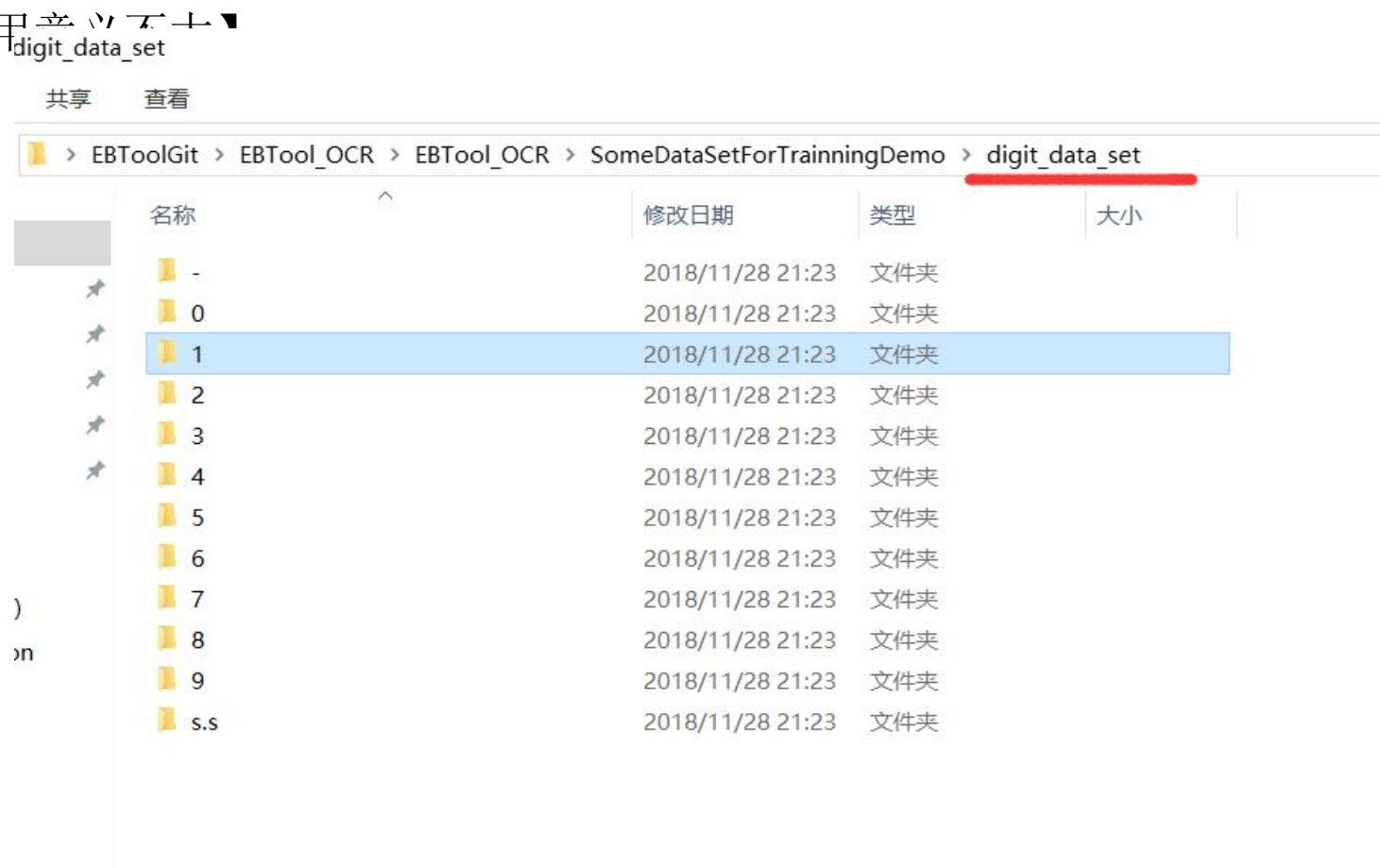
结果

TrainningTool使用：【很多参数写死，使用

打开cmd，  
cd 至训练结果需要保持的文件夹

在cmd输入指令  
路径/TrainningTool.exe 路径/数据集路径

项目  
中有  
部分  
数据  
集样  
本。  
位置  
如右  
图





## 程序配置

config.ini - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

TCP\_SERVER\_HOST=127.0.0.1

TCP\_SERVER\_PORT=8034

#MLAPI DEFAULT: 使用默认内置ANN, TCP:采用tcp方式与ML服务端交互进行识别

MLAPI=DEFAULT

# 如果采用DEFAULT模型, 则给出输出标签种类个数

MLAPI\_LABELS=12

#IMAGE\_SAVE\_DIR DEFAULT:不保存扣取的图片, 指定路径: 保存图片(采集数据集)到指定路径, 用于机器学习训练

#IMAGE\_SAVE\_DIR=xx/im

# xx/im 图片保存的位置, 当点击识别时, 识别之前会把需要识别的图片保存到该指定文件夹

IMAGE\_SAVE\_DIR=DEFAULT

|

TCP交互方式, 在EBTool\_OCR\mlServiceDemon, 有个python版的demo

## 数字识别模型的通用性

该工具在使用时，涉及图片放大缩小，以及阈值改变。

若图片不黏连，能切分出来，实际对数字识别模型影响有限，因为这些字体都非手写，特征明显。

非抱着0识别率的态度，怎样放大缩小以及阈值设置之后每个数字的特征还是会相当明显。故模型在年鉴识别等通用性较强，无需频繁更新模型。