

# Decision Engines Data Science Position Challenge

Simon Olsen

## Approach

Thank you for the exciting challenge. Please see the attached script for my solution.

The three invoices differed quite a bit in their imperfections, which required different approaches for obtaining a best possible result for OCR.

I started by doing a manual analysis of the three invoices. My major concerns were related to the grayscale backgrounds in the forms for Gibson and Merck's invoices. My initial thought was to use opening (erode, then dilate) to try and remove the specks of the grayscaled backgrounds, while still retaining the important textual information. I knew OpenCV had such functionality ready, therefore I spent some time reading through their documentation on image processing, focusing on Smoothing and Morphological Transformations.

For Merck's invoice, opening seemed to be a good choice due to the high density of dots in the grayscale background and the small, but quite bold text.

For Gibson's invoice, however, an initial erosion would cause drastic information loss, especially in the contact information on the top left, due to its weakness and "white holes" within the text. Performing a closing before an opening would also not yield beneficial results, due to the grayscale background interfering with the text clarity. Smoothing, however, could efficiently remove the grayscale, without disrupting the text too much.

Upon experimenting I found a need to handle these two invoices quite differently. While an initial smoothing of Gibson's interface gave good results, it was too hard on the small text in Merck's invoice. Thus, I decided to implement a script with some options as to how the files should be processed.

Now, a simple opening on Merck's invoice did not give a perfect result. Some specks from the grayscale were still retained. An increase in the erosion's intensity was efficient in removal of the grayscale, but this also removed too much information from the foreground text. Upon closer inspection I found that after the initial erosion, much of these specks were single pixels that stood alone. In an attempt to further clear the grayscale background while still retaining the clarity of the text, I implemented a method for removing the pixels that stood alone, and ran this before the following dilation. This process is run by passing the argument `--hard-open` (or `-ho`) when running the script. The pixel removal method can be found in the script defined as `remove_isolated_pixels`.

The original GE Healthcare invoice was, in my opinion, quite good. Though, the top right box had a quite noisy black background, and some of the text was a bit thin. A simple smoothing reduced the noise drastically, while closing and dilating enhanced the text's clarity.

## The Script

*Please see the email attachments for the script. (ocr\_preprocessor.py)*

The script contains the following optional arguments:

### **--input, -i: Input Destination**

Path to file, or containing folder for process. If given path is a folder, all containing TIF's will be processed.

### **--out, -o: Output Destination**

Path to folder where the optimized files should be stored.

### **--blur, -b: Blur**

If provided performs median blur on the input. This helps remove unwanted grayscale at the cost of text clarity.

### **--hard-open, -ho: Hard Open**

If given performs a hard open on the input. This helps remove grayscale but at a high cost on weak and unclear text with much "holes" within.

### **--open, -op: Open**

If given performs a normal open on the input.

### **--close, -c: Close**

If given performs a normal close on the input.

### **--erode, -e: Erode**

If given performs a normal erosion on the input.

### **--dilate, -di: Dilate**

If given performs a normal dilation on the input.

### **--debug, -d: Debug**

If given prints debugging information, and shows one large and one small preview of the processed image.

## Obtaining the Results from the Script

Below follows the commands to use with the script to obtain the attached results for each of the single invoices. If no output argument is passed, the result will be stored in a subfolder named "results" with the same filename.

The instructions use relative paths and assume that the current directory contains both the script and each of the original invoices.

The scripts also assume the following filenames for the invoices:

- *merck.tif*
- *gibson.tif*
- *ge.tif*

### **Merck**

Only the extended opening (*hard\_open*) is run.

```
>> python .\ocr_preprocessor.py -i merck.tif -ho
```

### **Gibson**

First runs median blurring to clear background noise followed by a dilation and closing in an attempt to fill in gaps in the weakly visible text.

```
>> python .\ocr_preprocessor.py -i gibson.tif -b -di -c
```

### **GE Healthcare**

A similar procedure to the one used on the Gibson invoice gave a good result. Swapping the order of the closing and dilation also further increased the clarity of the text.

```
>> python .\ocr_preprocessor.py -d -i ge.tif -b -c -di
```

[illegible][illegible]



1040 Manchester Street  
Lexington, KY 40505  
Toll Free: 800 477 4763  
Phone: 859.254.9557  
Web: [www.gibsonbioscience.com](http://www.gibsonbioscience.com)

# Invoice



LX-94747

Acct. No. 109526.5 Date 1/15/2018 Invoice # LX-94747

Packing Slip # / Proforma #

LX-54576

Page 1 of 1

COUNTRY OF ORIGIN: USA

## Bill To

ATTN: ACCOUNTS PAYABLE  
MCKESSON MEDICAL SURGICAL INC  
PO BOX 4059  
DOCUMENT PROCESSING  
DANVILLE IL 61834

## Ship To

ATTN: RECEIVING / LAB  
MCKESSON MEDICAL SURGICAL INC KANSAS  
1405 NORTH CHOUTEAU  
KANSAS CITY PC 003  
KANSAS CITY MO 64120

PO # 21615548	Ship Date 1/15/2018	Shipping Method LEX-FedEx 2Day®	Tracking # 789346692527	
Ship Note	Ship Note Ref Numbers	Comments SHIPPING 01/15/2018		
Catalog #	Quantity	Description	Unit Price	Amount
60060	2	ID-AE KIT	107.90	215.80
Thank you for your business!		Subtotal		215.80
		Other (Royalty)		0.00
		Shipping Cost (1 FX-FedEx 2Day®)		50.14
		Total		265.94

\*\*\*\*\* EACH LOT OF THE PRODUCTS IN THIS SHIPMENT MEETS APPLICABLE NCCLS (M22) STANDARDS. \*\*\*\*\*  
VISIT OUR WEB PAGE AT <http://gibsonlabs.com>.

## Remittance Advice

Invoice # LX-94747

Customer 109526.5 MCKESSON MEDICAL SURGICAL INC-LEX : MCKESSON MEDICAL SURGICAL INC KANSAS-LEX

Amount Due 265.94

Amount Remitted \_\_\_\_\_



# GE HEALTHCARE

DBA: GE Medical Systems Information Technologies, Inc

FEDERAL ID# 39-1046671

REMIT INVOICE NUMBER: 2861892

INVOICE DATE: 09-JAN-18

CUSTOMER ACCT: 128009

GE REFERENCE#: 80797369

CUSTOMER PO#: 21732386

**AMOUNT DUE: 204.00 (US DOLLARS)****DUE DATE: 08-FEB-18**

**Remit to:** GE Medical Systems Information Technologies, Inc  
**US MAIL:** Attn: Accounts Receivable \* 5517 Collections Center  
 Drive \* CHICAGO IL 60693  
**Wire/EFT Information:** ABA 021000021 ACCOUNT 304183911  
 If Wire/EFT, please email remittance advice to: Remit.Healthcare@ge.com

**SOLD TO:**

MCKESSON MEDICAL SURGICAL INC  
 ACCOUNTS PAYABLE  
 PO BOX 4059  
 DANVILLE, IL 61834-4059

**SHIP TO:**

MCKESSON MEDICAL SURGICAL INC  
 7343 S HARDY DR STE 101  
 TEMPE AZ 85283-4479

<b>Payment Terms:</b> NET 30	<b>Contract #:</b>	<b>GE BA Number:</b>
<b>Cost Center:</b>	<b>FE Badge:</b>	<b>FE Name:</b>
<b>GE Sales Rep or FE:</b> MS SM Central Admin cc	<b>ProdLine:</b>	<b>SERV MANAGER :</b>

Inquiries regarding this Invoice should be directed to: 1-800-581-5600

QUANTITY	ITEM NUMBER	DESCRIPTION	NET UNIT PRICE	EXTENDED AMOUNT OF BILLING
2	002203	THIS INVOICE AMOUNT IS DUE AND PAYABLE NET 30 FROM INVOICE DATE.  DURA-CUF, ADULT, 2 TB SUBMIN, 23 - 33 CM, 5/ BOX Tax USAZ_EX_O2 @ 7.40	102.00	204.00          0.00

Please include the Invoice / Credit Memo number for proper credit:  2861892	<b>TOTAL</b>	204.00
	<b>Tax</b>	0.00
	<b>SHIPPING/HANDLING</b>	0.00
	<b>Total Amount</b>	204.00
PAST DUE INVOICES ARE SUBJECT TO A SERVICE CHARGE OF 1.5% PER MONTH, NOT TO EXCEED THE MAXIMUM RATE ALLOWED BY LAW. ALL ORDERS SUBJECT TO GE HEALTHCARE TERMS AND CONDITIONS.		Goods and services or reimbursements associated with the ordered products or services and provided under contract without separately identified charges constitute discounts or other reductions in price under applicable federal law. It is the customer's responsibility to disclose such discounts or other reductions in price in the manner required under state or federal program which provides reimbursement to the customer for or related to the products or services under the contract.