

COL 783

Assignment 4

Rajbir Malik
2017CS10416

November 3, 2019

Robust Form Processing

Overview

The assignment required us to process some forms in different natural settings, and trying to extract meaningful information from the forms.

Being much open ended, we had the opportunity to try various procedures to get the task done. I was, after many attempts, able to come up with a "mostly" automated process to get the required information.

I have presented this approach and explained the procedure in the upcoming report. I've also discussed another approach that was reasonable to some scale, but couldn't be scaled, by me, to the desired level.

Image Alignment

I worked upon the lines of a standard form scanner, to get all the input images aligned correctly. Having created a template-default form, I used it as reference to get the ideal warp-perspective for the input image, with the help of descriptor-matching.

The exact pseudo-code for the same has been described below...

```
def image_alignment(_input_image, _template):
    - Create ORB Descriptor
    - Get Key-Points and Descriptors
    - Match the features
    - Filter the matches, to get the ideal ones for mapping
    - Using keypoints, discover ideal homography
    - Use the homography, and warp the input image
    - Return the warped image
```

Results

The results for alignment procedure were quite-impressive, achieving as good as 99% accuracy for **printouts** and **scanned** and reasonable results for **booklets**. Perhaps, if I could find ideal templates for booklets as well (which I couldn't), then there too a decent result may have been procured.

Further, I present the results of alignment on all the categories (2 of each kind).

Another Approach

I also tried with a much simpler approach earlier, which involved finding the lines using **Canny-Edge Detector** and then using **Hough-Transform** to get the angle at which maximum no. of lines appeared.

This approached, though naive seemed to work sensibly for **scanned**-images, but failed badly for both **printouts** and **booklets**. So, I was compelled to drop the idea.

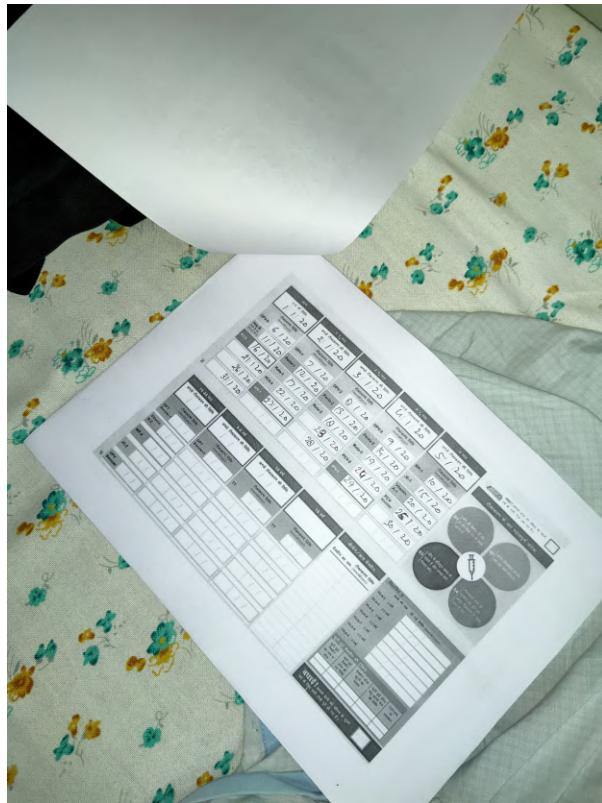


Figure 1: Printouts₁

५ वर्ष	१-२ वर्ष	२-३ वर्ष	३-५ वर्ष	५ वर्ष
५ वर्ष की रिपोर्ट 1/1/20	अपने दोस्राम ली रिपोर्ट 2/1/20	अपने तीसराम ली रिपोर्ट 3/1/20	अपने चौथाम ली रिपोर्ट 4/1/20	अपने पांचवाम ली रिपोर्ट 5/1/20
मासिक रिपोर्ट (अपने दोस्राम ली रिपोर्ट)				
OPV-4 6/1/20	OPV-1 7/1/20	OPV-2 8/1/20	OPV-3 9/1/20	MRI-1 10/1/20
Hep B Nex ५ वर्ष 1/1/20	Penta-1 12/1/20	Penta-2 13/1/20	J-E-1 15/1/20	Vitamin A+ 20/1/20
BCG 1/1/20	Rota-1 7/1/20	Rota-2 8/1/20	Rota-3 19/1/20	PCV booster 25/1/20
2/1/20	PCV1 22/1/20	28/1/20	PCV2 24/1/20	PCV booster 30/1/20
3/1/20	IPV-1 27/1/20	28/1/20	IPV-2 29/1/20	
4/1/20				
5/1/20				
6/1/20				
7/1/20				
8/1/20				
9/1/20				
10/1/20				
11/1/20				
12/1/20				
13/1/20				
14/1/20				
15/1/20				
16/1/20				
17/1/20				
18/1/20				
19/1/20				
20/1/20				
21/1/20				
22/1/20				
23/1/20				
24/1/20				
25/1/20				
26/1/20				
27/1/20				
28/1/20				
29/1/20				
30/1/20				
31/1/20				
32/1/20				
33/1/20				
34/1/20				
35/1/20				
36/1/20				
37/1/20				
38/1/20				
39/1/20				
40/1/20				
41/1/20				
42/1/20				
43/1/20				
44/1/20				
45/1/20				
46/1/20				
47/1/20				
48/1/20				
49/1/20				
50/1/20				
51/1/20				
52/1/20				
53/1/20				
54/1/20				
55/1/20				
56/1/20				
57/1/20				
58/1/20				
59/1/20				
60/1/20				
61/1/20				
62/1/20				
63/1/20				
64/1/20				
65/1/20				
66/1/20				
67/1/20				
68/1/20				
69/1/20				
70/1/20				
71/1/20				
72/1/20				
73/1/20				
74/1/20				
75/1/20				
76/1/20				
77/1/20				
78/1/20				
79/1/20				
80/1/20				
81/1/20				
82/1/20				
83/1/20				
84/1/20				
85/1/20				
86/1/20				
87/1/20				
88/1/20				
89/1/20				
90/1/20				
91/1/20				
92/1/20				
93/1/20				
94/1/20				
95/1/20				
96/1/20				
97/1/20				
98/1/20				
99/1/20				
100/1/20				
101/1/20				
102/1/20				
103/1/20				
104/1/20				
105/1/20				
106/1/20				
107/1/20				
108/1/20				
109/1/20				
110/1/20				
111/1/20				
112/1/20				
113/1/20				
114/1/20				
115/1/20				
116/1/20				
117/1/20				
118/1/20				
119/1/20				
120/1/20				
121/1/20				
122/1/20				
123/1/20				
124/1/20				
125/1/20				
126/1/20				
127/1/20				
128/1/20				
129/1/20				
130/1/20				
131/1/20				
132/1/20				
133/1/20				
134/1/20				
135/1/20				
136/1/20				
137/1/20				
138/1/20				
139/1/20				
140/1/20				
141/1/20				
142/1/20				
143/1/20				
144/1/20				
145/1/20				
146/1/20				
147/1/20				
148/1/20				
149/1/20				
150/1/20				
151/1/20				
152/1/20				
153/1/20				
154/1/20				
155/1/20				
156/1/20				
157/1/20				
158/1/20				
159/1/20				
160/1/20				
161/1/20				
162/1/20				
163/1/20				
164/1/20				
165/1/20				
166/1/20				
167/1/20				
168/1/20				
169/1/20				
170/1/20				
171/1/20				
172/1/20				
173/1/20				
174/1/20				
175/1/20				
176/1/20				
177/1/20				
178/1/20				
179/1/20				
180/1/20				
181/1/20				
182/1/20				
183/1/20				
184/1/20				
185/1/20				
186/1/20				
187/1/20				
188/1/20				
189/1/20				
190/1/20				
191/1/20				
192/1/20				
193/1/20				
194/1/20				
195/1/20				
196/1/20				
197/1/20				
198/1/20				
199/1/20				
200/1/20				
201/1/20				
202/1/20				
203/1/20				
204/1/20				
205/1/20				
206/1/20				
207/1/20				
208/1/20				
209/1/20				
210/1/20				
211/1/20				
212/1/20				
213/1/20				
214/1/20				
215/1/20				
216/1/20				
217/1/20				
218/1/20				
219/1/20				
220/1/20				
221/1/20				
222/1/20				
223/1/20				
224/1/20				
225/1/20				
226/1/20				
227/1/20				
228/1/20				
229/1/20				
230/1/20				
231/1/20				
232/1/20				
233/1/20				
234/1/20				
235/1/20				
236/1/20				
237/1/20				
238/1/20				
239/1/20				
240/1/20				
241/1/20				
242/1/20				
243/1/20				
244/1/20				
245/1/20				
246/1/20				
247/1/20				
248/1/20				
249/1/20				
250/1/20				
251/1/20				
252/1/20				
253/1/20				
254/1/20				
255/1/20				
256/1/20				
257/1/20				
258/1/20				
259/1/20				
260/1/20				
261/1/20				
262/1/20				
263/1/20				
264/1/20				
265/1/20				
266/1/20				
267/1/20				
268/1/20				
269/1/20				
270/1/20				
271/1/20				
272/1/20				
273/1/20				
274/1/20				
275/1/20				
276/1/20				
277/1/20				
278/1/20				
279/1/20				
280/1/20				
281/1/20				
282/1/20				
283/1/20				
284/1/20				
285/1/20				
286/1/20				
287/1/20				
288/1/20				
289/1/20				
290/1/20				
291/1/20				
292/1/20				
293/1/20				
294/1/20				
295/1/20				
296/1/20				
297/1/20				
298/1/20				
299/1/20				
300/1/20				
301/1/20				
302/1/20				
303/1/20				
304/1/20				
305/1/20				
306/1/20				
307/1/20				
308/1/20				
309/1/20				
310/1/20				
311/1/20				
312/1/20				
313/1/20				
314/1/20				
315/1/20				
316/1/20				
317/1/20				
318/1/20				
319/1/20				
320/1/20				
321/1/20				
322/1/20				
323/1/20				
324/1/20				
325/1/20				
326/1/20				
327/1/20				
328/1/20				
329/1/20				
330/1/20				
331/1/20				
332/1/20				
333/1/20				
334/1/20				
335/1/20				
336/1/20				
337/1/20				
338/1/20				
339/1/20				
340/1/20				
341/1/20				
342/1/20				
343/1/20				
344/1/20				
345/1/20				
346/1/20				

Figure 2: Aligned

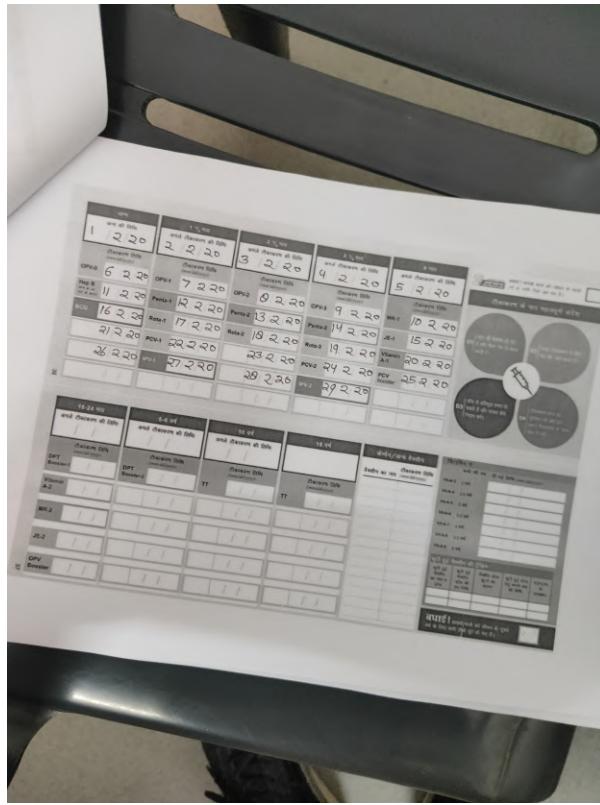


Figure 3: Printouts₂

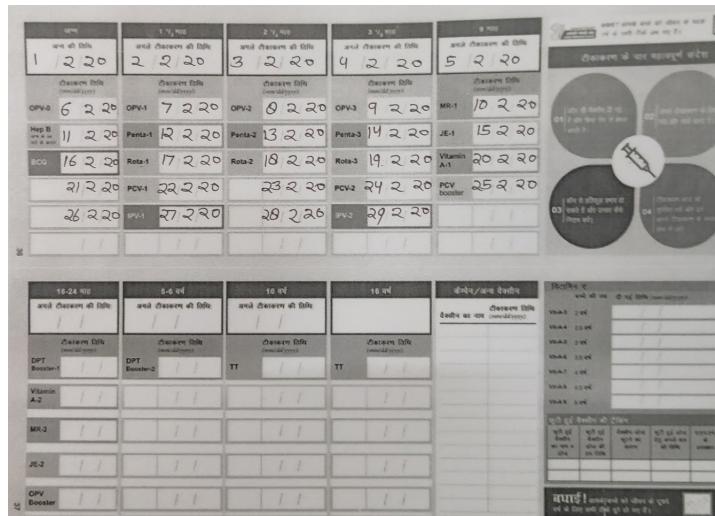


Figure 4: Aligned

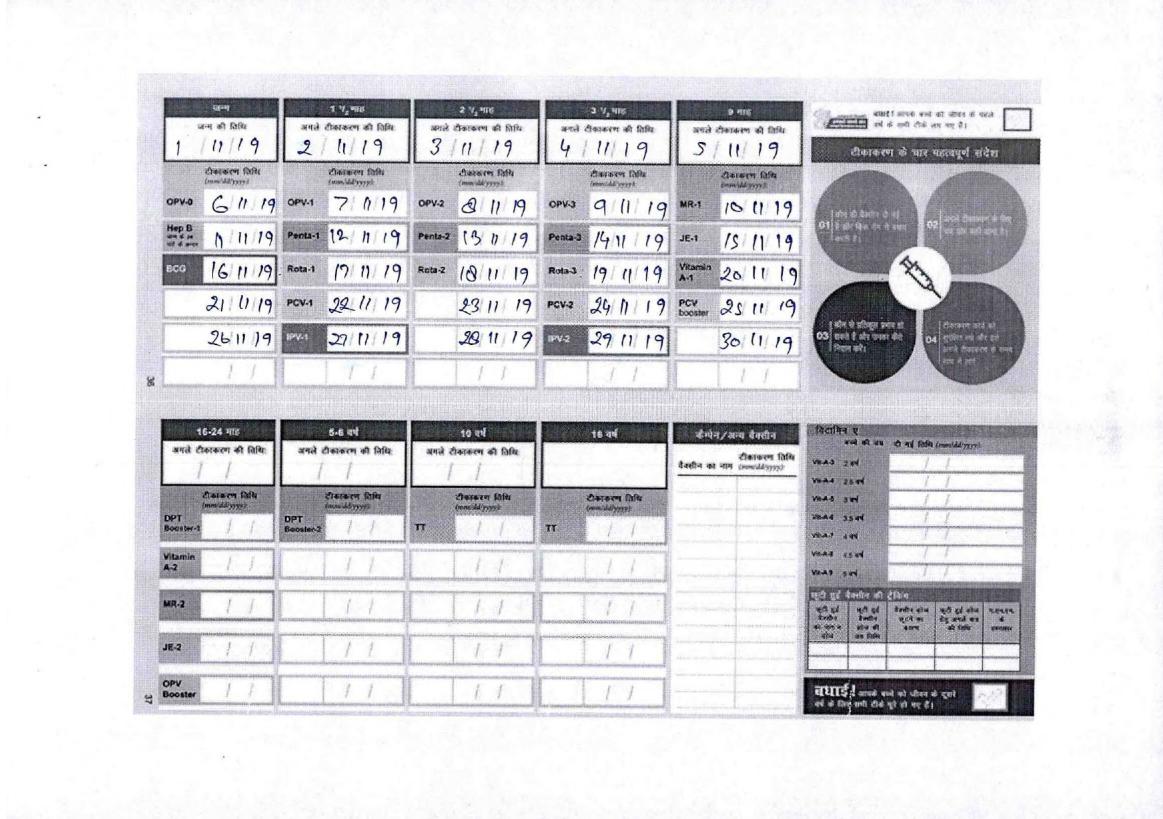


Figure 5: Scanned1

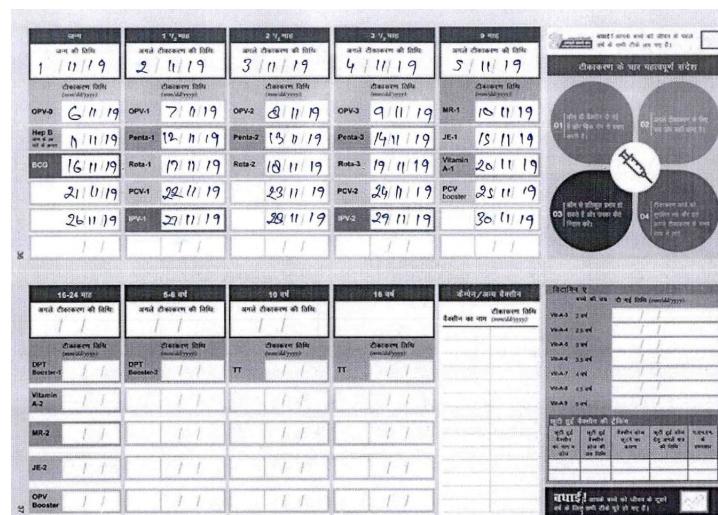


Figure 6: Aligned

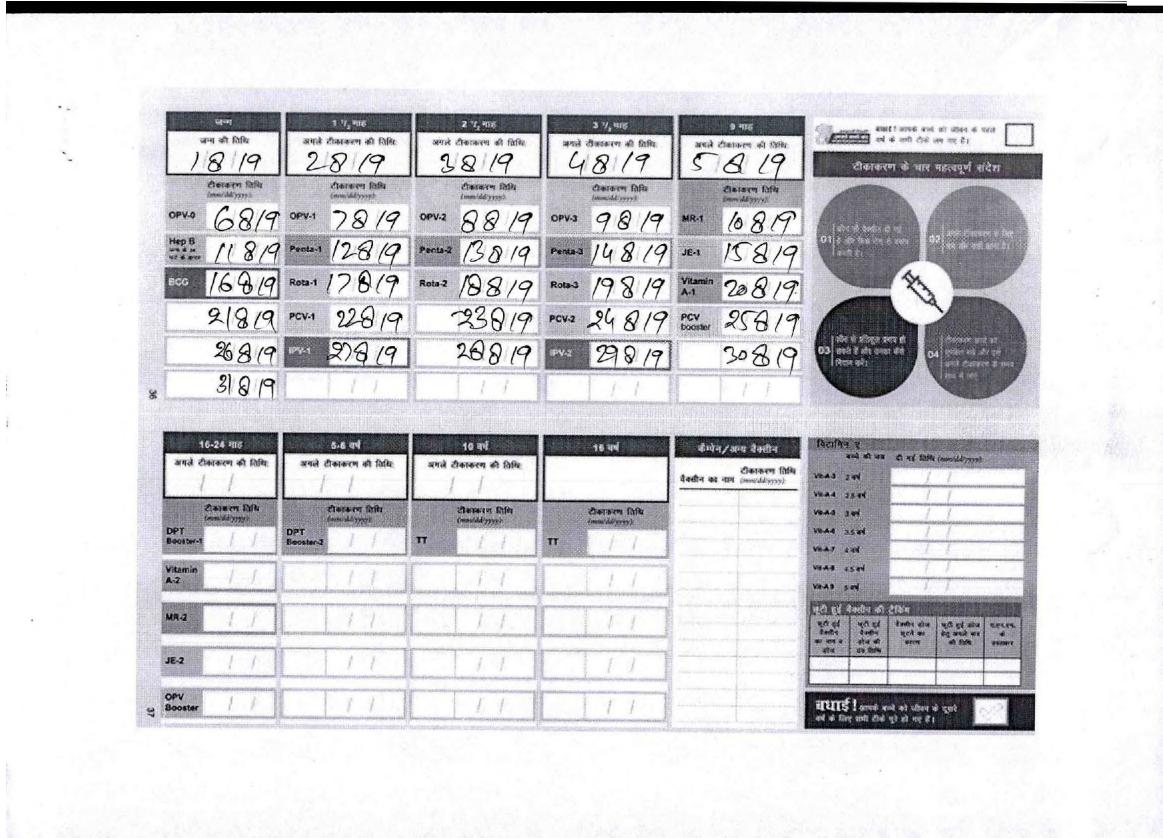


Figure 7: Scanned2

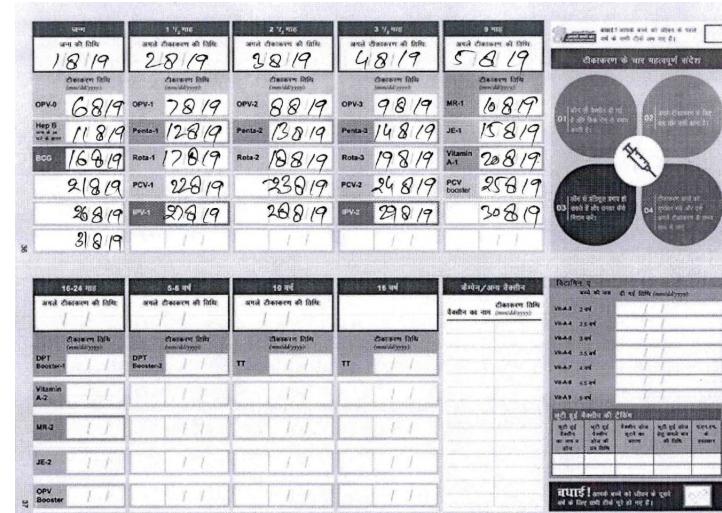


Figure 8: Aligned

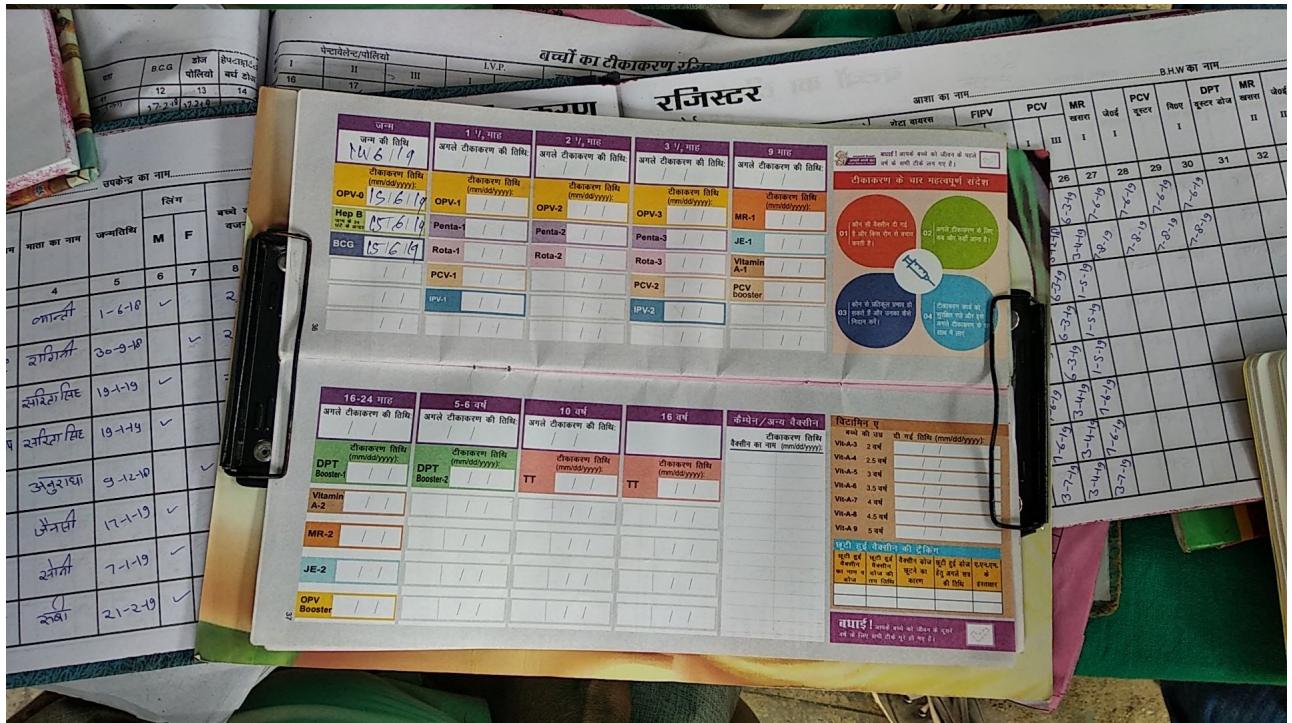


Figure 9: Booklet1

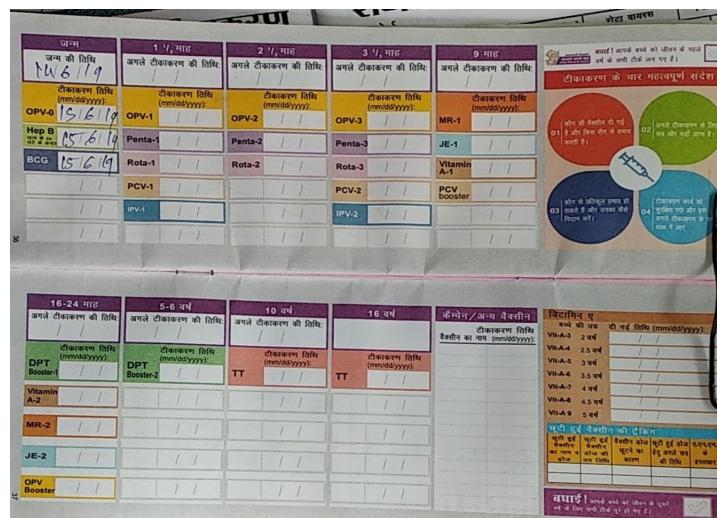


Figure 10: Aligned

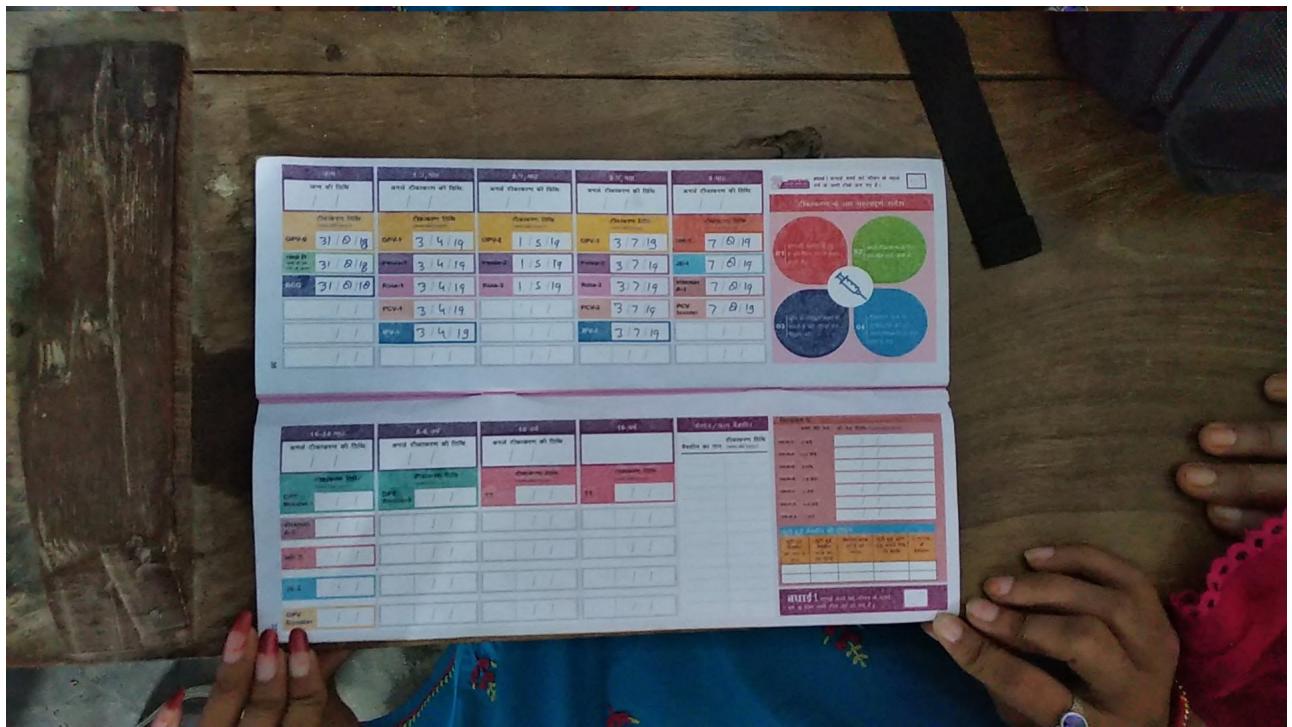


Figure 11: Booklet2



Figure 12: Aligned

Image Segmentation

I tried, various approaches for segmenting the text fields from the complete image. Of these various attempts, I obtained decent results via two approaches, described below.

Contour Map

I implemented **Canny-Edge Detection** on the aligned image, and carved out boxes of desired aspect-ratio to get the boxes. This approach did work, but wasn't complete, i.e. many boxes were left out.

Predefined-Map

Since, the image alignment procedure, tries to warp most of the boxes in the same position, majority of boxes have fixed coordinates. I, thus, worked on a single image and through various measures (contouring) came up with the ideal mask for that image, but that mask can now be applied to all the images. This made segmentation very fast as well as majorly complete! The mask is presented below.

Also, the results for both the procedures are presented below.

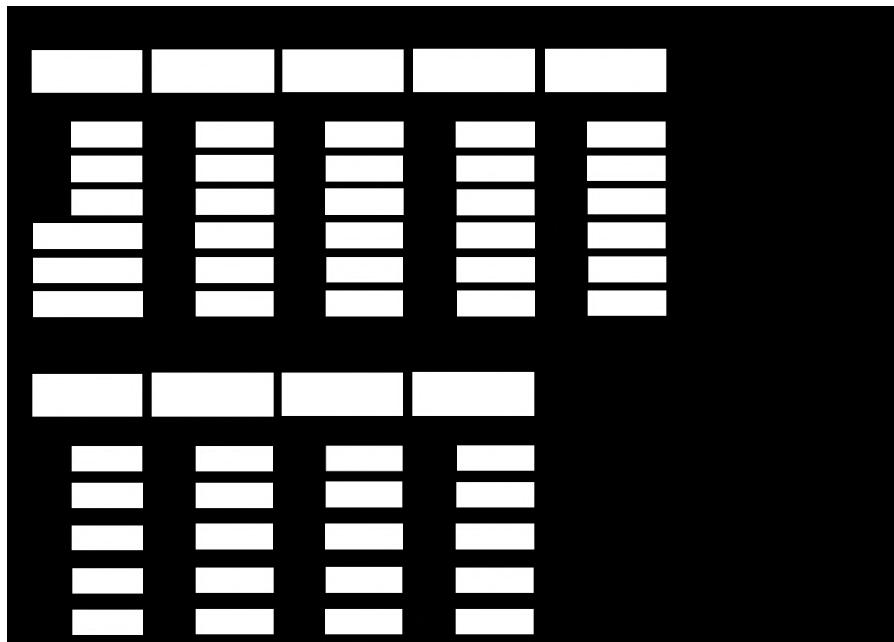


Figure 13: Mask

Contour Map

Figure 14: Scanned1

Figure 15: Scanned2

Contour Map

प्राप्ति दिनी	प्राप्ति दिनी	प्राप्ति दिनी	प्राप्ति दिनी	प्राप्ति दिनी
01/01/19	02/01/19	03/01/19	04/01/19	05/01/19
प्राप्ति दिनी प्राप्ति दिनी प्राप्ति दिनी	प्राप्ति दिनी प्राप्ति दिनी प्राप्ति दिनी	प्राप्ति दिनी प्राप्ति दिनी प्राप्ति दिनी	प्राप्ति दिनी प्राप्ति दिनी प्राप्ति दिनी	प्राप्ति दिनी प्राप्ति दिनी प्राप्ति दिनी
OPV-1 Hep B BCG	OPV-1 Penta-1 Rota-1	OPV-2 Penta-2 Rota-2	OPV-3 Penta-3 Rota-3	MRI-1 JE-1 Vitamin A-1
06/01/19 11/01/19 16/01/19 21/01/19 26/01/19 31/01/19	07/01/19 12/01/19 17/01/19 22/01/19 27/01/19 28/01/19	08/01/19 13/01/19 18/01/19 23/01/19 28/01/19	09/01/19 14/01/19 19/01/19 24/01/19 29/01/19	10/01/19 15/01/19 20/01/19 25/01/19 30/01/19
PCV-1 IPV-1	PCV-2 IPV-2	PCV-3 IPV-3	PCV-4 IPV-4	PCV-5 IPV-5
16-24 वर्षीय वर्षीय दिनी	3-6 वर्षीय वर्षीय दिनी	1 वर्षीय वर्षीय दिनी	1 वर्षीय वर्षीय दिनी	वर्षीय/अन्य वर्षीय वर्षीय दिनी
DPT DPT Booster-1	DPT DPT Booster-2	TT	TT	वर्षीय वर्षीय दिनी
Vitamin A-2				
MRI-2				
JE-2				
OPV Booster				

Figure 16: Printout1

Figure 17: Printout2

Contour Map



Figure 18: Booklet1

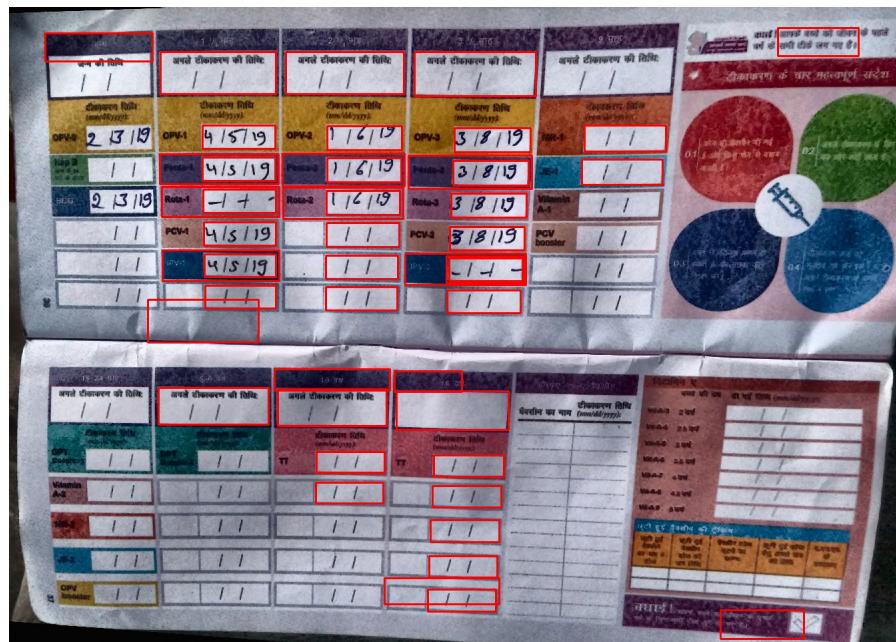


Figure 19: Booklet2

Masked

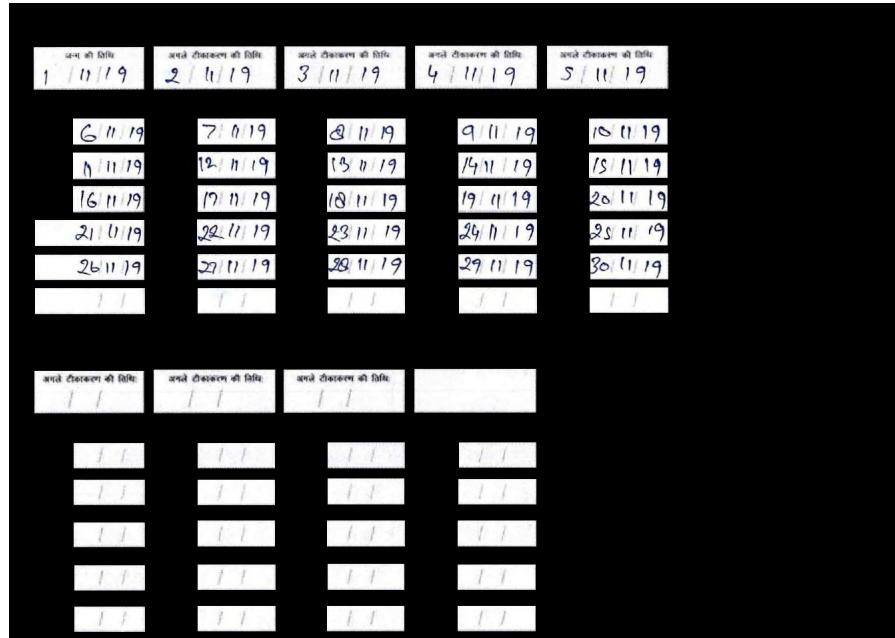


Figure 20: Scanned1



Figure 21: Scanned2

Masked



Figure 22: Printout1



Figure 23: Printout2

Masked

Figure 24: Booklet1

अन्य की दिनी	अपने दोस्राम की दिनी	अपने दोस्राम की दिनी	अपने दोस्राम की दिनी	अपने दोस्राम की दिनी
31/0/19	3/4/19	1/5/19	3/7/19	7/8/19
31/0/19	3/4/19	1/5/19	3/7/19	7/8/19
31/0/19	3/4/19	1/5/19	3/7/19	7/8/19
1/1	3/4/19	1/1	3/7/19	7/8/19
1/1	3/4/19	1/1	3/7/19	7/8/19
1/1	1/1	1/1	1/1	1/1
				
				
				
				

Figure 25: Booklet2

Character detection

Having segmented the image, using the latter approach (masking), it was fairly easy to do character detection using thresholding and then contouring the remaining characters.

Note: The thresholding was binary, after adaptive histogram-equalisation was applied to the aligned image. This allowed for faster and easier detection.

Bonus: It was observed that median filtering the images led to removal of pre-printed slashes and thus they didn't appear in the detected character-set.

The results, including the detected characters and thresholding results are presented further.

Scanned1

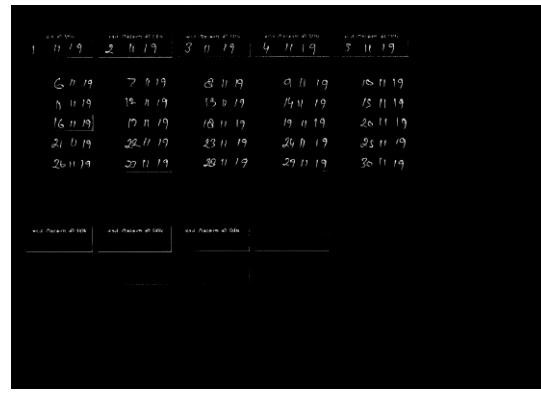


Figure 26: Thresholded Image

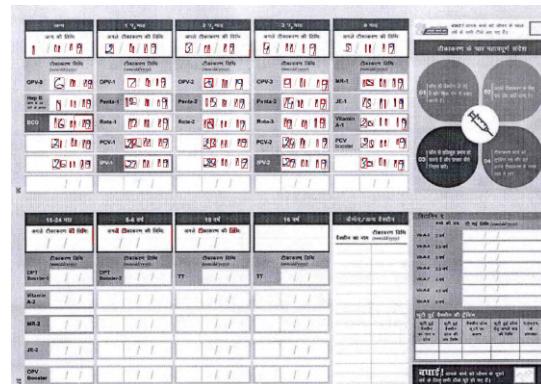


Figure 27: Character Detection

Scanned2

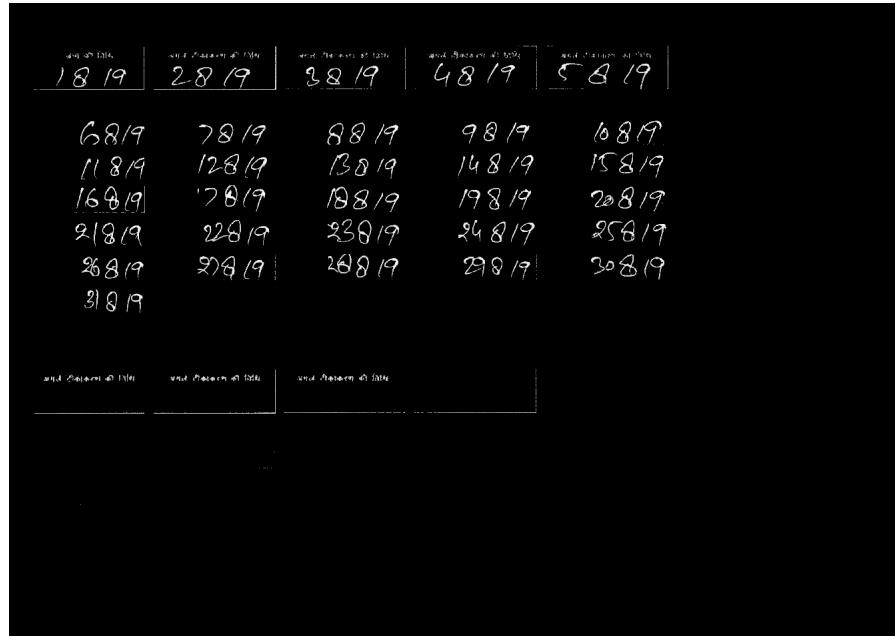


Figure 28: Thresholded Image

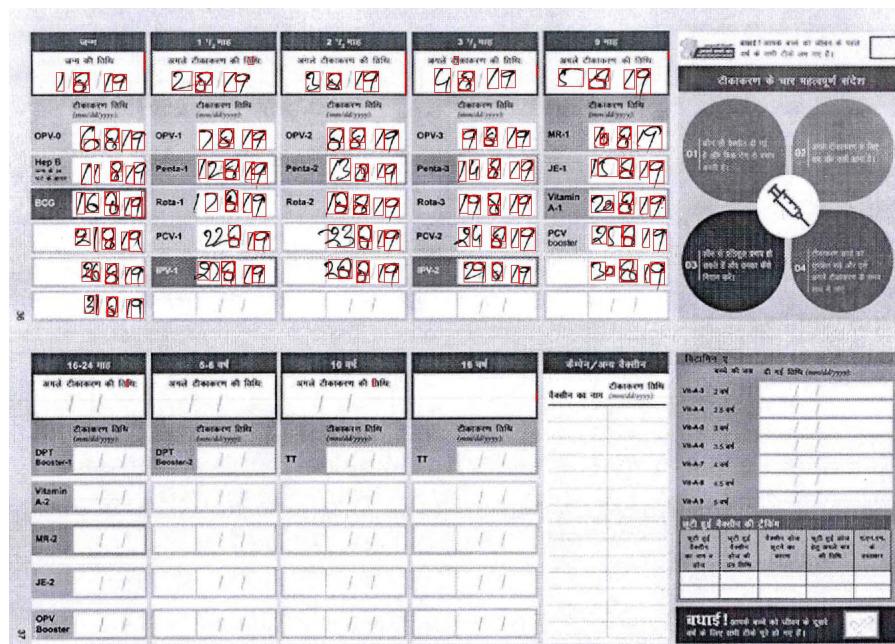


Figure 29: Character Detection

Printout1

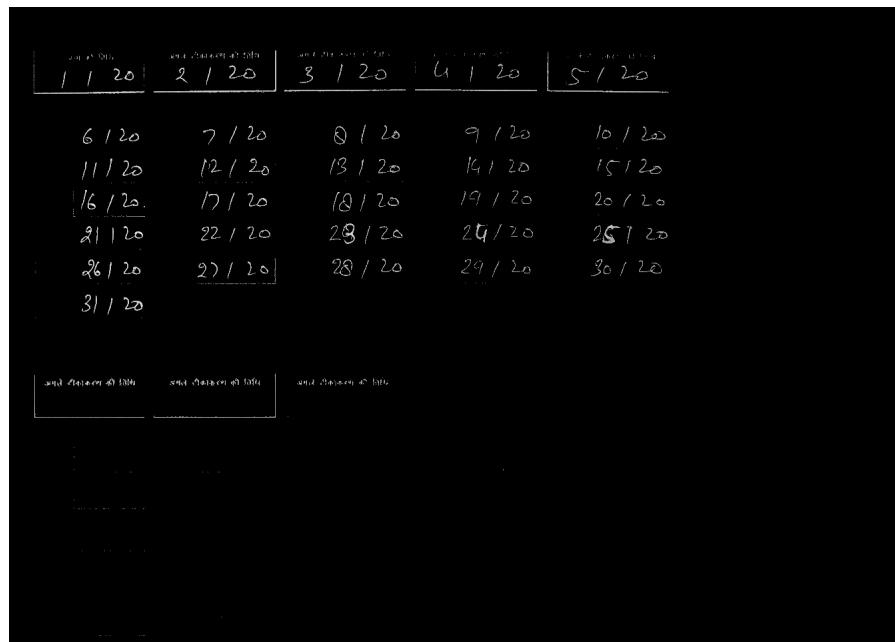


Figure 30: Thresholded Image

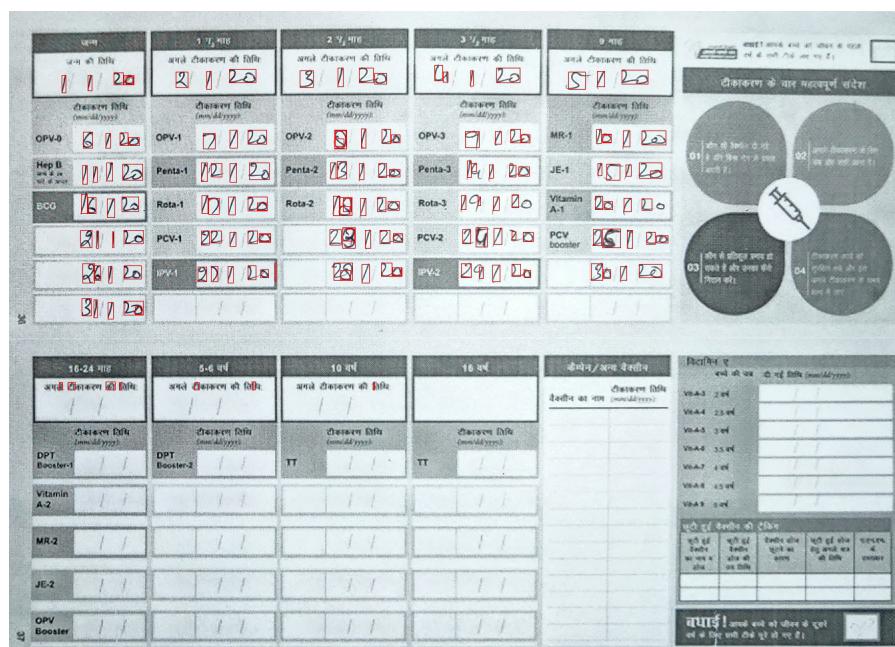


Figure 31: Character Detection

Printout2

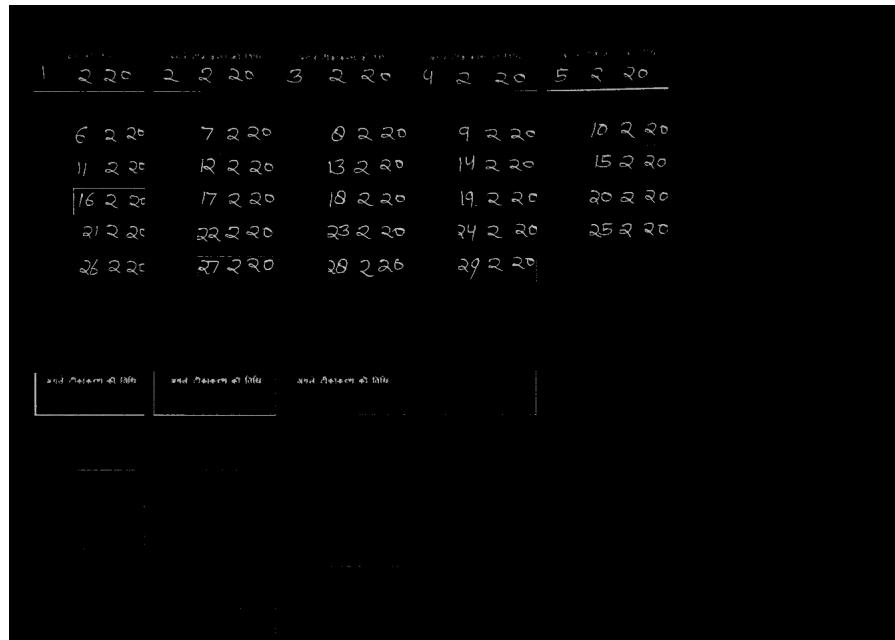


Figure 32: Thresholded Image

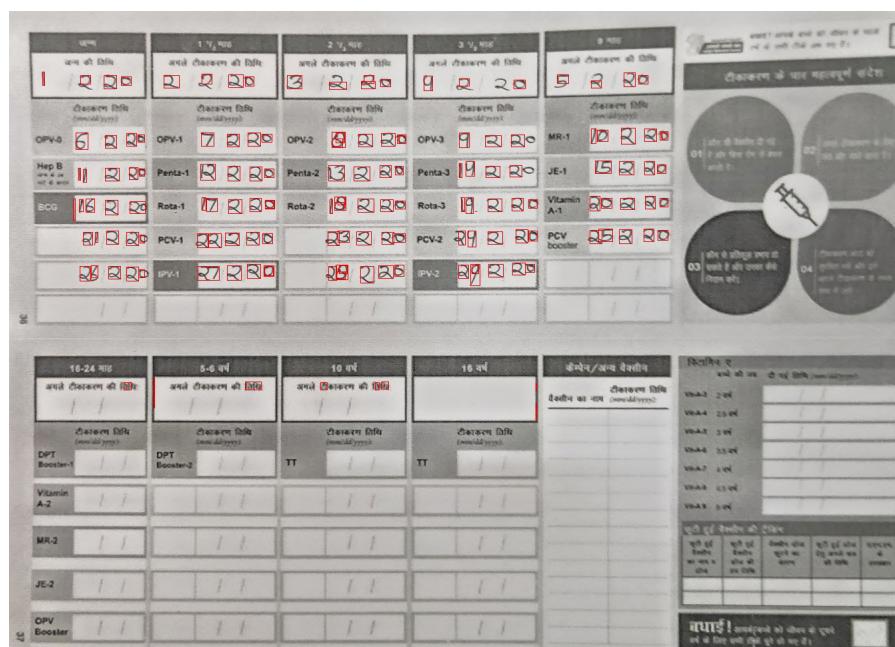


Figure 33: Character Detection

Booklet1

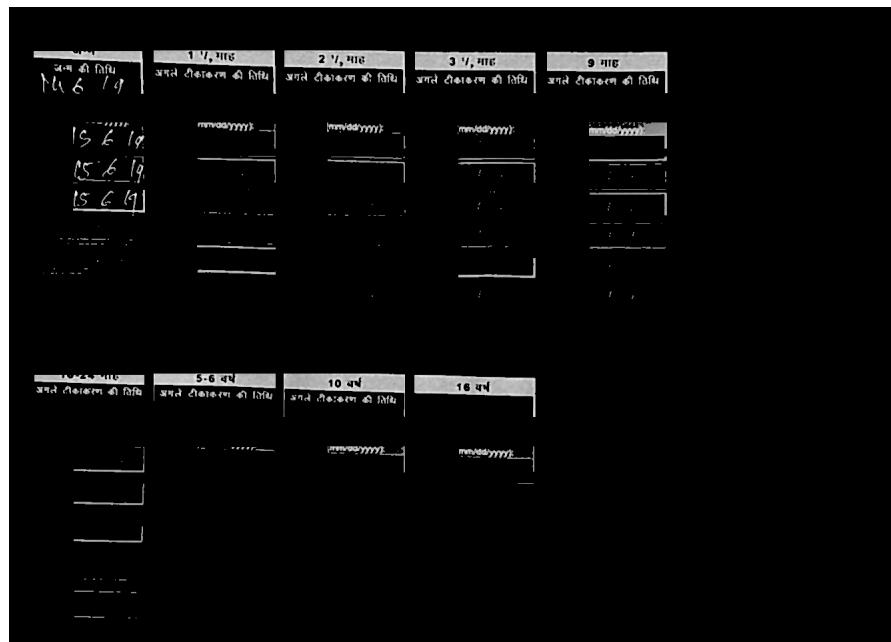


Figure 34: Thresholded Image



Figure 35: Character Detection

Booklet2

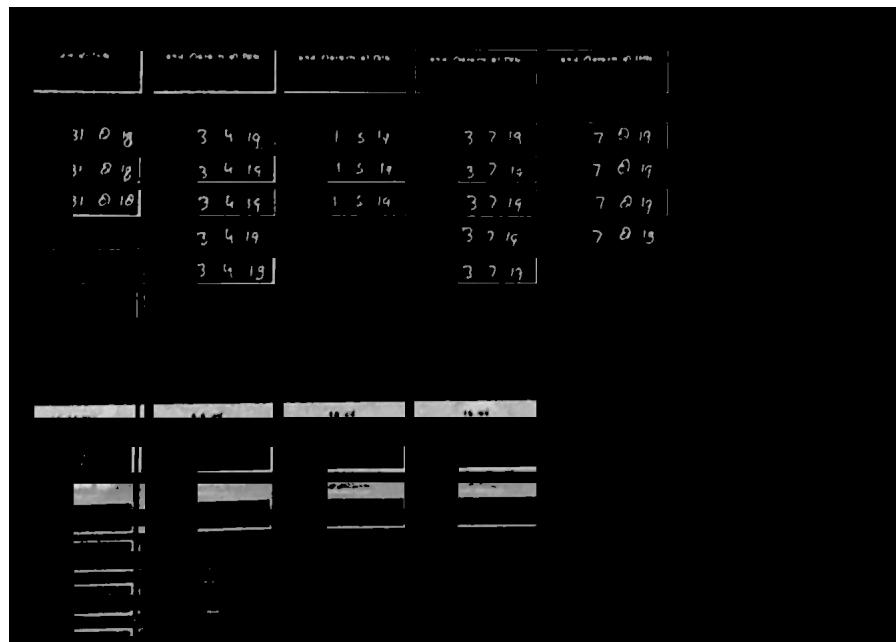


Figure 36: Thresholded Image

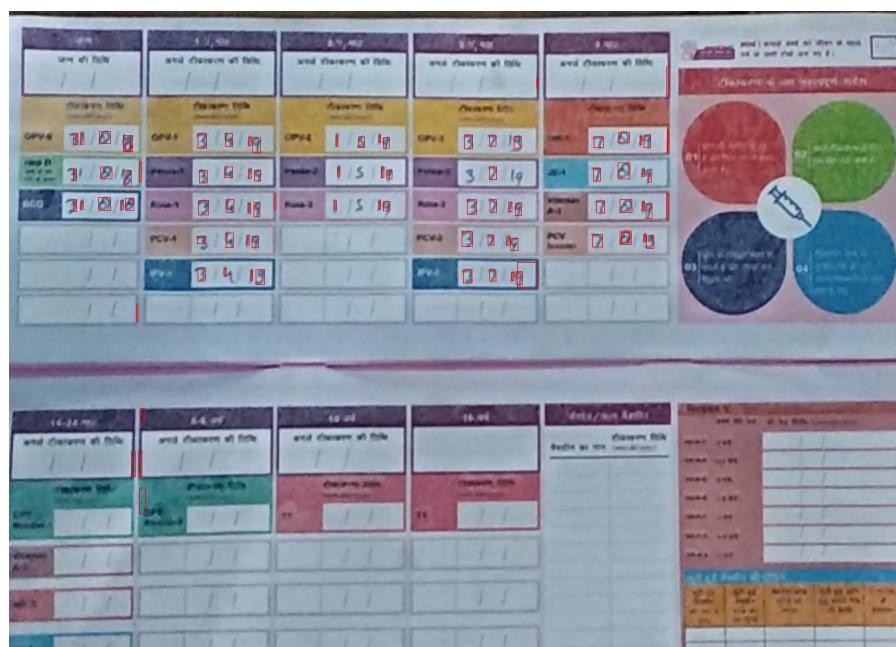


Figure 37: Character Detection

Summary

This assignment was very interesting, and I learnt a lot of new things. I really liked the idea of making this assignment open ended. As a direct practical application, we were very motivated to try all the possible things out and see how actually all methods we know of go around. What was discussed in class was very helpful.

Thanks and Regards

Rajbir