

A simple detector for crowd counting using OpenCV and Python

Martí Gelabert Gómez

November 27, 2022

Contents

1	Introduction	1
2	General Procedure	2
3	algorithms	2
3.1	Annotation	2
3.2	Background Removal	2
3.3	Binarization	4
3.4	Dilation	5
3.5	Find contours	5
4	Results and Analysis	5
5	Is a different approach capable of improving our results?	7
6	Conclusions	7
A	Appendices	9
A	Output	9
B	Images related to the basic Algorithm	13
C	Complementary Images	16

1 Introduction

The assignment consists in **counting** the number of people that appear on the given images. There was not an established way to approach this problem. In this case I have taken the decision of using the computer vision algorithms seen in class and not rely on deep learning or any kind of artificial intelligence algorithm. This is because not by the amount of data needed for training, due to there are really good datasets out there, like MS COCO, but because of the effort of building a good performing detector architecture. Use an already set framework as YOLO or Mask-RCNN would be the way to accomplish the task, but I would be thinking about it as "cheating". Therefore, this way, using the content seen in class, I will be more cautious about my decisions and It will be much easier to justify them.

In the following document we will be focusing on the process taken and the algorithms used, their implementation and their performance.

2 General Procedure

The program takes the following steps :

1. Annotation of the data
2. Import images as black and white.
3. Apply Adaptative Histogram equalization to the images.
4. Subtract the background to the images using the image with the background.
5. Apply a thresholding algorithm to binarize the image.
6. Apply a dilation operation into the binarized images to expand the whites.
7. Use a contour algorithm to extract the different regions containing persons.
8. Obtain the bounding boxes from these regions.
9. Count them and compare the number of detections to the real quantity.

3 algorithms

In this section we will discuss the algorithms selected and the output we obtain from them in the application of our problem.

3.1 Annotation

For the annotation of the dataset it was used an online annotation tool, and some basic criteria was applied in order to decide how to label and if a person would be annotated or not:

- The annotation would be points in the image.
- The annotations should be placed as near the head possible.
- If too many people where too close it may be annotated as just one person.
- In case of occlusion just annotate the most clear figure.
- Detections in the shore could be useful, so if the head was visibly recognizable it would be annotated.
- Some sections of the image will be masked, and those regions will not be annotated.

The ground truth is saved as a csv file and the program loads it when it is needed.

3.2 Background Removal

Background removal is a technique that allows to remove the background from the image, this way, the output will be only the foreground highlighted in the form of a binarized image as it can be illustrated on the figure 1.



(a) Output of background subtraction using
cv2.threshold(image, 100, 255, cv2.THRESH_BINARY)

(b) Original 1660320000.jpg

Figure 1: Reference images

In the following sections will be explained how the preprocessing of this technique was prepared and for what they are used for. They are ordered by

CLAHE

In the case of our problem, first we have imported all the images in black and white integer values using `cv2.imread(img, 0)`. Then for a better extraction, it has been applied an **algorithm of CLAHE** seen previously on the course to try to uniformize the illumination of the images. With a more uniform images, we will be able to distinguish more precisely our foreground and artifact like shadows casted by objects on the images may have less impact on the **binarization process**.

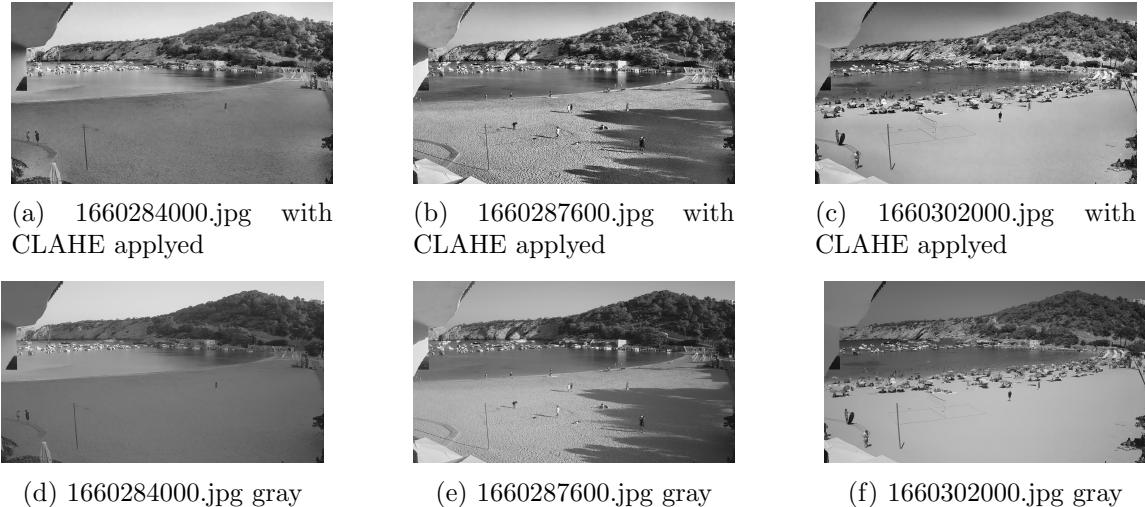


Figure 2: Images from Gelabert folder where the first row presents a more uniform illumination across frames compared with the second row images.

Background Image

For the selection of our background image, in a starting point the **average image** was computed as the main way for accomplish foreground extraction, but the resulting image generally was really noisy

and with ghosting effect very present on it, which decreased the performance of the algorithm as it can bee seen on the Figure 8. Even with a gaussian blur applyed to try to counter the ghosty figures, they are still really present on the image. The image averaging is a good idea generally for noise reduction and in this case to try to obtain a neutral image were the background is predominant, but there is not enought images to obtain a more sharp image to really have a good background. Therefore, we **end up** using the image 1660284000.jpg with a gaussian blur applied as our background, this way we obtain better results on the **binarization process**.

Gaussian Blur

With an image already selected for the image background, it is applyed a gaussian filtering with a kernel of size (7,7), this way we can obtain a more standard image, and get a more smuthered image that will allow the algorithm to contrast slightly better our foreground.

The ending size used for gaussian blur was obtained by trial and error experimentation, and could be changed to get different ending behaveours.

Substraction

Once we have selected our background and applied the preprocessing, we should subtract it to every image, this way, the result we end up with is the foreground (in this case the persons and some residuals) highlighted as we can see in the figure 4.

3.3 Binarization

Once the substraction has been aplyed to all images, the resulting foreground will contain hopefully our persons. All those whity areas should be further treated using a binarization progress, allowing a solid division between the background and our foreground. This is acomplish by the process of binaryzation. Depending of the process selected, the resulting image will transform into a binary black and white image from a gray scale one. In this assignment there has been some experimentation using OTSU and the binary thresholding function available on cv2 library, but for the sake of the performance, the last one has been ultimately selected.

The problem with using OTSU is that for each images it obtains a threshold selection automatically, and this should not be a problem, and sometimes could perform better than a fixed one. But for the samples tested there are a lot of artifacts related to the shadow casting of the image, as it can be seen in the figure 9. The problems are seen on the bottom of the image, where there are an amount of shadows casted by the sand crests which can be a bit annoying for the next steps.

By using a set thresholding we are loosing some flexibility and a fixed values could not be the wiser idea, but in this case using a tigh threshold may be the best decision as may reduce the number of false detections considerably. In this case the values selected are setted in the instruction `cv2.threshold(subtracted, 100, 255, cv2.THRESH_BINARY)`. The resulting image can be shown in figure 5.

Masking

For a better filtrage of our data, we will be using a binary mask to get rid of centain areas of the image where we can get non-persons in the foreground (i.e the quay part) and those areas where there is to much information to get clear data about it. We can see the last case in the umbrella section of the

beach and the pavement area, where there are too many persons overlapping each others. The mask applied can be seen in the figure 6.

3.4 Dilation

Now, with the image binarized we will apply a dilation to the image. The dilation will expand the white color on a binarized image, obtaining as results the highlighted areas in form of white chunks as it can be seen on figure 7. With these chunks the countour detection will obtain better results because they are more clear shapes. This could merge figures into one, so the parameters of the kernel size may be tinkered with to obtain the best performance possible.

3.5 Find contours

Once a binarized image is obtained, there should be a process of contour detection. Here we use the dilated images over the function that openCV provides `cv2.findContours()`. From this functions, we can extract the bounding boxes that contain each of the object contours. In the figure 3 appears the image with some bounding boxes applied.

4 Results and Analysis

For convenience, we will quantify our detection accuracy using the following assumptions:

- If detection contains more than one person it will count as only one detection.
- In the case we would have a massive region we will discharge that detection.
- For checking the dimensionality of the bounding box, we would assume that width or height higher than a third of the image will not be acceptable.
- Regions minuscule will also be discarded.
- Some labels of the ground truth could be double checked.

The reasons of the anterior list, are bounded to the complexity of the problem. We are not restrictive enough to lose a lot of information. Allowing more flexibility in order to have more valid detections let us obtain a better performance overall. Also, some cases are unlikely to happen really often. In the tables 1 and 2 is shown how the algorithm performs.

Results

In table 1, the column of *gt* represent the number of labels on each image, *detected* represent the number detections found in the image, and the column *matched* represent the number of detections where at least contained one label of our ground truth. As we can see, the tables show a normal scoring, having a balance in between the actual number of persons in each image and the detections, in figure 3 we can observe the output of the algorithm. But of course the general problem we will still have is the quantity of false positives we end up having in some images. Even the **mean squared error** with a value of 341.666667 tells how this algorithm is missing with a huge error, we are predicting with a really soft constrictions and allowing a lot of error on our detections. Although this is done entirely on purpose and in the practice this has turned out to be more beneficial for our images.

Table 1: Performance metrics using the proposed algorithm

files	precision	recall	f1 score	gt	detected	matched
1660309200.jpg	0.529	0.700	0.603	90	119	63
1660302000.jpg	0.446	0.563	0.498	103	130	58
1660294800.jpg	0.535	0.736	0.620	72	99	53
1660320000.jpg	0.770	0.770	0.770	135	135	104
1660287600.jpg	0.275	0.647	0.386	17	40	11
1660298400.jpg	0.663	0.594	0.627	106	95	63
1660305600.jpg	0.595	0.653	0.623	101	111	66
1660316400.jpg	0.888	0.856	0.872	139	134	119
1660291200.jpg	0.667	0.654	0.660	52	51	34

Table 2: Performance metrics overall using the proposed algorithm

Metric	Scoring
MSE	341.667
Macro-average precision	0.597
Macro-average recall	0.686
Macro-average F1	0.629



Figure 3: Output of foreground detector, as it can be seen there are some big detections due to contour detecting multiple forms as one shape, and some

Should be considered that the type of images we got and the placement of the ground truth can make a huge difference on our outcome, so this algorithm has been tuned to try to be the more general possible **only** with the images provided in my corresponding directory.

5 Is a different approach capable of improving our results?

For some of the test that were executed for this assignment, generally the effort was not useful. Ending up in situations where a precise modification to the parameters gave as output excellent results only in a subset of the images. Usually the over-tunning of some function parameters made really inconsistent outputs that provoked a bad performance.

In the following subsections it will be discussed some other paths to accomplish the task of the assignment, some of them had been implemented during the development stage if implemented or experimented with.

Color Spacing

For example, one of the approaches taken for background removal, was the use of colored images in the color space of LAB instead of gray scale. On the paper one may think that the use of color will improve our performance, because of the extra information we are just acquired from free. But in this case the process to obtain our foreground sections were in general less precise, giving as result a lot more objects as foreground, which is not ideal for our case, been a lot of those object just umbrellas and cast shadows.

The problem with shadows

In general, the sample images follow a consisting landscape without heavy cast shadows on the sand, but it is noticeable in some of them, and usually they will give you more than one headache.

We could try to counter some cases using a sharpening processes or even more filtering (i.e. median filtering), to avoid certain annoyance for shadows, but doesn't mean that it will be worth it to apply. We could have the situation again where is really useful for 1 case, but then the rest of images just got obliterated just for applying those transformations.

Annotation with bounding boxes

For a reason of time, the annotation was given only by dots on the image, but it would be a good idea to use bounding boxes and utilize more "sophisticated" methods as **intersection over union**, to allow us the execution of other metrics to improve our way of evaluating our results.

6 Conclusions

As we could be expecting, the performance is not good compared to a neural network. There are some cases where we can be completely sure about classifying incorrectly our detection because of how the ground truth is placed, some detections are actually correct but slightly shifted and without a manual intervention it is safer to just leave it as 'incorrect'. In general, the amount of false positives is consistent and overall generally concentrates in areas where there is a brusque lighting change like

cast shadows or where there is a lot of information compacted like in the shores.

The use of **background subtraction** to try to extract completely our foreground is a tricky task, and it can be performed following different strategies, for example, using an averaged image could allow to have an image with the solid background on it, but it works better when we dispose of a lot of samples. With just the images from the 'Gelabert' folder the output generated was not clear enough (even with some processing) to accomplish the objective, and in the end just using the empties image was the one given the best performance. We can still spice up our work if we start working with more complex color spaces, like LAB, or using more complex post-processing conditions for the sake of improving our results, but in some cases, working with simpler ideas give the best results.

Of course, we have to take in count how flexible the annotation and the accuracy measurements where, but as it can be shown on the outputs, we are capable of at least detecting areas where people concentrates, and even single persons in some images. Therefore, for the use of these techniques it should be enough.

A Output











B Images related to the basic Algorithm



Figure 4: Image obtained by applying `cv2.subtract(background, image)`



Figure 5: Binarized image using a fix thresholding method



Figure 6: Mask image



Figure 7: Dilation applied to the binarized image

C Complementary Images



Figure 8: Average image with a gaussian blur applied with a kernel of (11,11)



Figure 9: Binarized image using OTSU and the empty beach image