

CS 229 Project report: Extracting vital signs from video

D.Deriso, N. Banerjee, A. Fallou

December 9, 2013

Introduction

Cardiovascular health is the *sin qua non* of human life. Early detection of cardiovascular disease is of paramount importance in public health. This project aims to develop a method to visualize the perfusion of blood through the skin via pulse oximetry. Pulse oximetry is a technique that exploits the fact that oxygenated and deoxygenated hemoglobin changes the color of red blood cells. The technique maps these changes in rgb color of the visible skin to the invisible presence of oxygenated vs deoxygenated blood in the local vasculature underneath the skin.

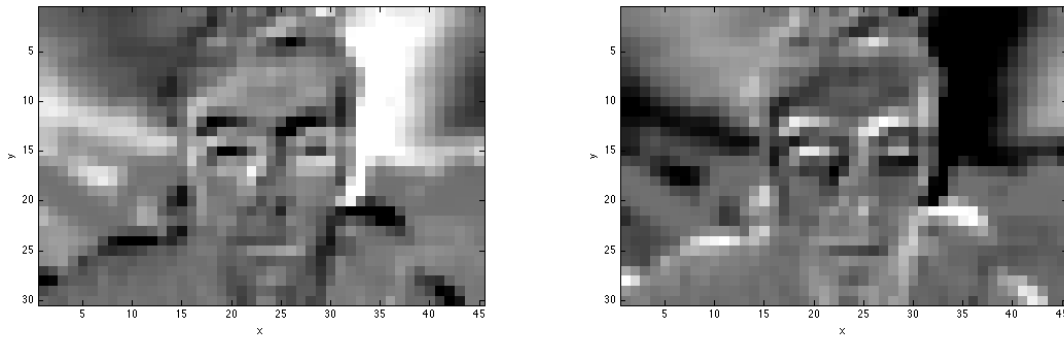
Previous studies have shown that video obtained from an ordinary webcam can be used to visualize perfusion by selectively amplifying temporal frequencies in video ¹. A study by the MIT CSAIL showed that this technique can also be used to infer heart rate from the person being taped. The present project aims to extend this work to detect the relative changes in oxygenated vs deoxygenated blood and reconstruct the pulse oximeter waveform from an ordinary webcam video.

1 Methods

Previous work on Eulerian Video Magnification (EVM), has led to the following processing pipeline, which we later modify for our purpose. The video can undergo different treatments, but the central goal is to amplify spatial or temporal changes that are normally invisible to the naked eye. The process can be summarized as:

- Separate the video into distinct spatial frequency bands by performing a 2D spatial fourier decomposition.
- For each spatial frequency band, blur and downsample several times using Gaussian or Laplacian Pyramids. The first step preserves spatial features (e.g. high frequencies such as edges) through this destructive process.

¹<http://people.csail.mit.edu/mrub/vidmag/>



(a) At a given time

(b) Half a heartbeat period later

Figure 1: Red channel of the EVM output, before recombining with the original video.

- Amplify a pre-selected temporal frequency band.
- Recombine spatial frequencies from step 1, and add the amplified video to the original video.

Figure 1 presents an example for the output of the next-to-last-step.

Strikingly, these two frames showed that the whole video undergoes a periodic color change with a frequency that seems equal to the heartbeat. This did not happen with the data provided in the paper, where in a similar video only the person's face changes color. To extract features relevant to pulse oximetry, we seek out a weighted combination of periodic changes in color space that best reconstruct our training pulse oximetry signal. We therefore undergo the following process:

- Extract each pixels time course from the video.
- For each pixel, separate the colors into R, G, B bands and perform a fourier decomposition where frequency is grouped into n bins.
- Learn weights via linear regression for frequencies within each color band to reproduce the simultaneously recorded pulse oximeter data in the training data set.
- Amplify the temporal frequencies according to the weights of the regression.
- Combine the amplified video with the original video.

2 Preprocessing

2(a) Video data

We needed to find a way to convert a video into features. Our feature set initially started out vaguely; we wanted to incorporate the frequency data inherent within each pixel in the

video and output frequency data about a pulse oximeter waveform. The three color channel intensity values through time, $I_R(t)$, $I_G(t)$, $I_B(t)$ were our starting points. We converted the video into a four dimensional matrix of pixel location x and y, color channel and time and performed the descontruction mentioned in the previous section.

2(b) Pulse Oximeter Data

Pulse oximeter data was collected using an Arduino and a Pulse Oximeter connected to a laptop. A program was written that started recording the pulse ox data and simultaneously started taking pictures every 50-80 ms on a webcam. This gave the pulse ox data corresponding to the image/video data. Thus we had many images taken over small time intervals with a corresponding pulse oximeter values.

Further processing had to be done because the pictures were not taken uniformly in time and so computing a fourier decomposition of the pixels taken from each image would not be correct. An interpolation was needed for images (to create a set of images with uniform time between them) and also for the pulse oximeter data to corroborate with the images. In order to accurately train our model, we also needed to bin the pulse oximeter values and be able to recreate the signal from the binned values so that they remained true to life.

3 Results

Our initial computations on the pixel-intensity data revealed, unsurprisingly, that pixel intensity variation over time for a reasonably still video was not large. Often a single color channel for a pixel would have a range of 3 (when it could vary from 0-255) over the whole video, even if that pixel was contained in the face. The DFT of the data would produce large peaks at 0 Hz and would rapidly diminish in magnitude. Because of this, we have not had any significant linear regression results. A simple solution is to subtract a mean intensity value of our pixel, so as to increase the relative intensity variations over time. Ultimately, we may have to resort to using EVM-enhanced videos if our training data has signals that are too weak to pick up.

4 Conclusions

5 Further Work

References
