

CS 229 Project report: Extracting vital signs from video

D.Deriso, N. Banerjee, A. Fallou

December 12, 2013

Abstract

In both developing and developed countries, reducing the cost of medical care is a primary goal of science and government. In this project we seek to find and extract information from a video of a human that tells us the pulse rate and the oxygen level of the blood. We therefore aim to create a virtual pulse oximeter: the ultimate non-invasive, equipment-free medical diagnostics tool, which could be deployed to anyone with video recording capabilities. Features were chosen to be related to the three color channel intensity values, with the idea that changing color of the video would relate to blood flow around the body. Extensive pre-processing was required on both the video data and the pulse oximeter data to enable training. Early results showed that feature selection was vital in reducing the mean-squared error of the output, but plenty of further work can be done.

1 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.

Our goal was to develop a process to obtain the SpO2 signal from an ordinary webcam. However, this imaging modality (webcam) is not designed to detect SpO2, and thus our signal of interest is weak and noisy. Furthermore, ordinary webcams have limited spatial, temporal, and color resolution, that may be unable to detect the subtle changes in color necessary for accurate pulse oximetry. Finally, factors outside the webcams limitations, including between-subject variations in skin pigmentation and lighting conditions, introduce noise to the signal acquisition process. To overcome these obstacles, we focused on engineering rich features that could be used in standard learning models. The learning

task was therefore to learn which combination of features best represented the weak pulse oximetry signal.

A paper [1] by Wu et al detailed a process known as 'Eulerian Video Magnification'(EVM) whereby previously imperceptible signals were amplified such that they became visible to the naked eye. For example by amplifying the red color channel in a 0.6-1.1Hz band, EVM allows seeing the periodic blood flow on the face of a person in the frame. The processing steps are summarized below in the methods.

However, potential issues appeared when performing EVM on our own videos. We had time-periodic artefacts showing color fluctuations in unexpected regions. For example, part of the ceiling in the video would show a pulsating red color change at the frequency of the person's heart rate. Our pursuit of optimal features led us to extract the most relevant aspects of their pipeline, which included frequency selection and parallel treatments of different color channels. Moreover, the EVM process stipulates that the parameters for the color amplification had to be pre-determined, and we believed that would not play out well with our goal of being robust to lighting conditions and camera settings. We therefore wanted to relax the process such that those parameters did not need to be set manually.

The EVM code provided convincing evidence that there exists information in a video from which a pulse signal can be extracted. We therefore decided to base our feature extraction techniques on those used in EVM. Further down the line, we hope to incorporate EVM more substantially into the project. We aim to map out the blood flow across the whole human body by amplifying the right color frequencies. Furthermore, we could use the spatial amplification of the EVM to show the physical motion of the skin as blood travels under it; the potential impact of early identification of poor blood circulation is large.

2 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.
- 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.

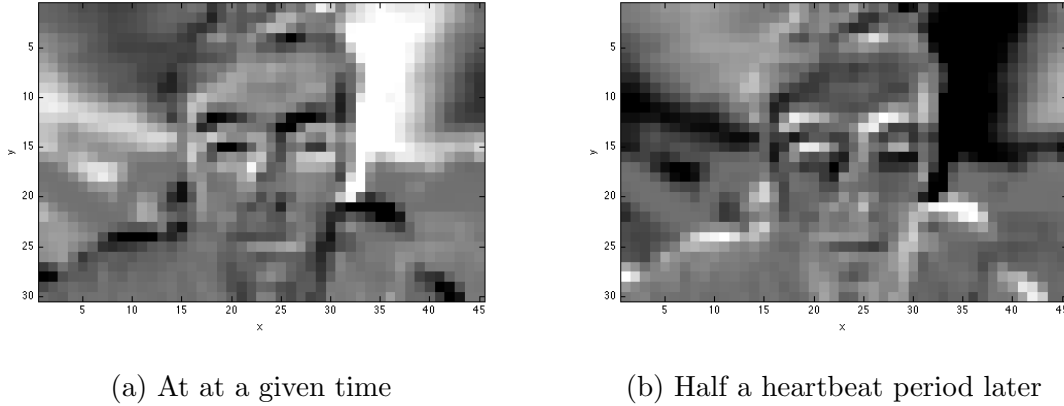


Figure 1: Red channel of the EVM output, before recombining with the original video.

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.

3 Preprocessing

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 deconstruction mentioned in the previous section.



Figure 2: Obtaining ground truth estimates from pulse oximeter values per image

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 (Figure 2).

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.

4 Results and Discussion

Initial results

At first we took 10 bins over a 6Hz interval of the Fourier transform of each pixel. We outputted a binary vector, with a 1 if the binned value exceeded a threshold and a 0 otherwise. The binning involved with our initial results provided a coarse model and was in fact a very poor representation of the data. In particular, for each pixel the binary vector consisted of $[1, 0, 0 \dots, 0]$ which is exactly the same distribution we would expect for noise. Thus our weight matrix contained predominantly zeroes, as we were training to a predominantly zero vector.

To make our predictions closer to the ground truth, we increased the bin size, although limitations in our computational power meant that we were restricted to maximum 24 bins over the 6Hz interval. This applied to both the pulse oximeter value and the pixel's Fourier transforms. We also included the phase terms and the pixel locations as further features (Figure 3).

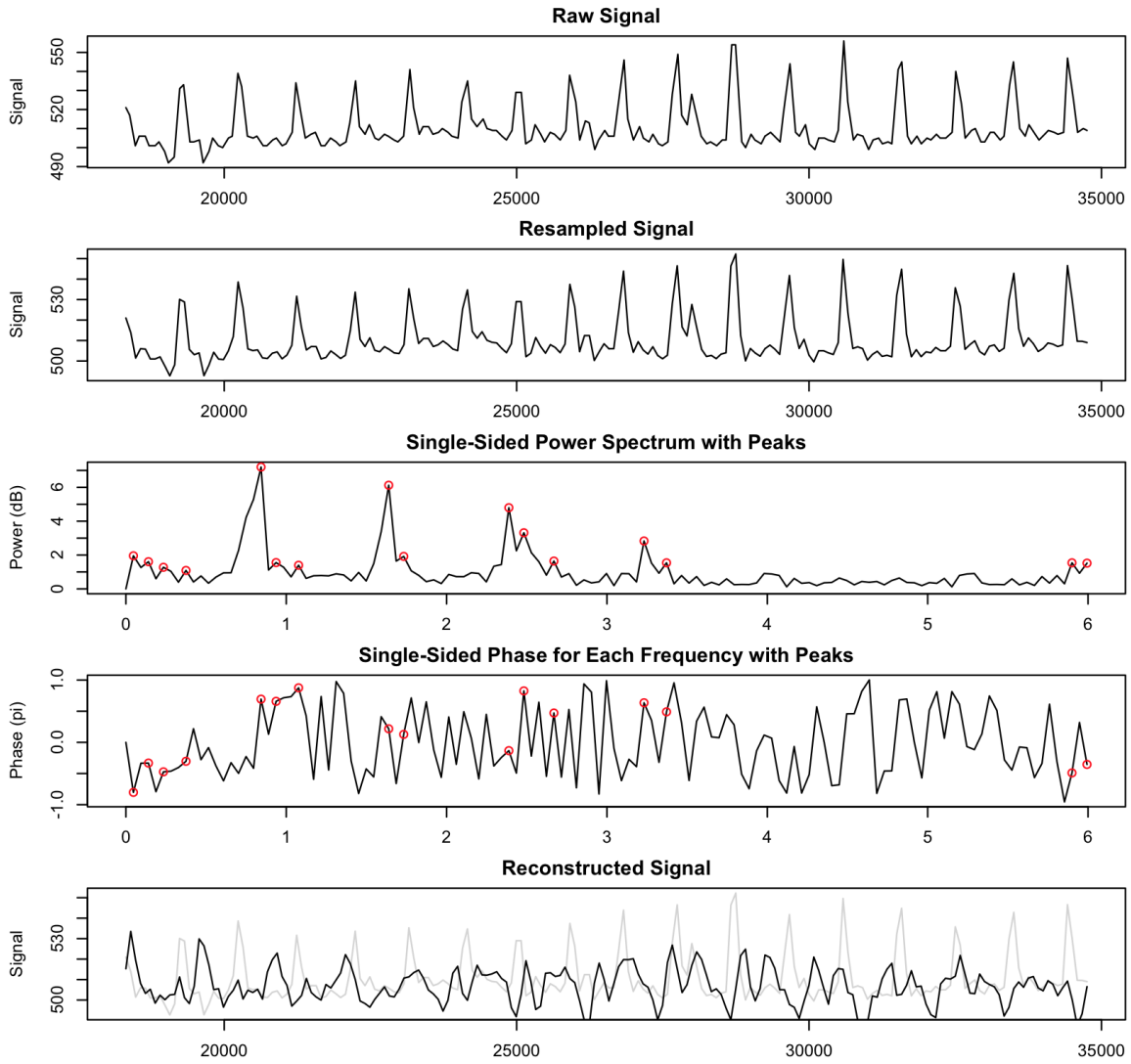


Figure 3: The pre-processing steps

Using the normal equations

Our feature matrix, X , which contained the features for every training example, gave a matrix $X^T X$ with a very high condition number. This meant the normal equations were unusable and so stochastic gradient descent was used.

Results

Features were extracted from each pixel over time, and included the pixels $[i,k]$ location, the FFT of the pixels time course binned to 4 buckets per frequencies 0-6Hz $[fq1, fq2, , fq24]$, and phase $[ph1, ph2, , ph24]$. Linear least squares regression models were trained on the following combination of features:

Pixel location	Phase	Frequency	MSE	Max cross correlation
✓	✓	✓	0.00884	0.0936
✓	✓	×	0.114	0.0943
✓	×	✓	0.00903	0.0936
✓	×	×	0.116	0.0943
×	✓	✓	0.0115	0.0936
×	✓	×	0.748	0.0967
×	×	✓	0.0117	0.0936

Table 1: Results of different feature selections with corresponding mean-squared error between predicted and input signal and maximum cross correlation

These results suggest that each of the features, frequency, phase, and pixel location; play a role in the prediction. Furthermore, the least error was obtained when all three features were used (MSE =0.00884), followed by just frequency and pixel (MSE=0.00903), and pixel location (MSE=0.0117). As expected, for a single video, with numerous pixels serving as training examples, the residuals were low.

5 Conclusions

- We built a program to enable training of the data: the program could record a series of images and pulse oximeter readings for those images.
- Using a binned Fourier Transform, we could efficiently reduce the video data to a finite dimensional feature space and perform regression on it.
- Increasing the variety of features included in the video reduced our mean squared error in the linear regression.

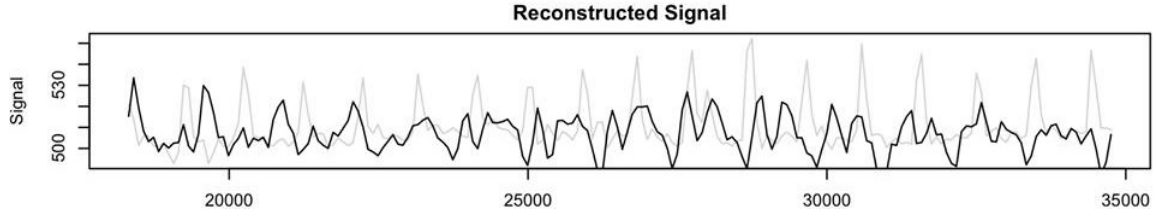


Figure 4: The reconstructed signal in bold and the actual signal in grey

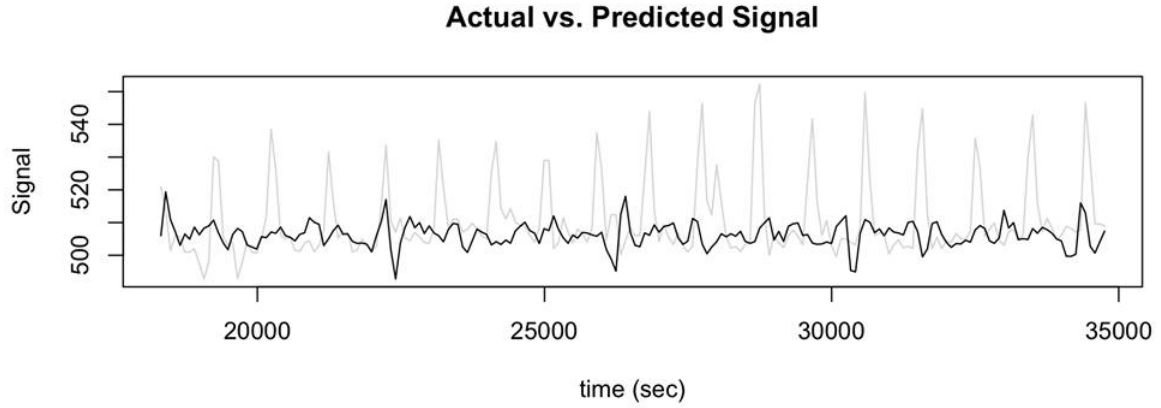


Figure 5: The predicted signal in bold and the actual signal in grey

6 Further Work

Given the time constraint on this project, we feel we have made great progress. However in order for the project to be considered successful, we have identified some areas where our methodology was non-ideal (due to time constraints or computational power).

- The main limitation on our method is the autoencoding process used for discretizing the FT of the pulse oximeter signal. As can be seen in Figure 4 if we tried to reconstruct the original pulse ox signal with our compressed version, the compression is a lossy one. Given this lossiness in our compression our predicted signal (Figure 5) has limits on its accuracy.

References
