



DISSERTAÇÃO DE MESTRADO

**AUTOMATED NON-CONTACT HEART
RATE MEASUREMENT USING
CONVENTIONAL VIDEO CAMERAS**

Gustavo Luiz Sandri

Brasília, Fevereiro de 2016

UNIVERSIDADE DE BRASÍLIA

FACULDADE DE TECNOLOGIA

UNIVERSIDADE DE BRASÍLIA
Faculdade de Tecnologia

DISSERTAÇÃO DE MESTRADO

**AUTOMATED NON-CONTACT HEART
RATE MEASUREMENT USING
CONVENTIONAL VIDEO CAMERAS**

Autor

Gustavo Luiz Sandri

Prof. Ricardo Lopes de Queiroz, Ph.D.

Orientador

Prof. Eduardo Peixoto Fernandes da Silva, Ph.D.

Coorientador


**UNIVERSIDADE DE BRASÍLIA
FACULDADE DE TECNOLOGIA
DEPARTAMENTO DE ENGENHARIA ELÉTRICA**

**AUTOMATED NON-CONTACT HEART RATE MEASUREMENT
USING CONVENTIONAL VIDEO CAMERAS**

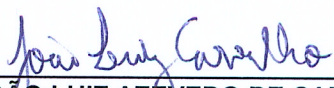
GUSTAVO LUIZ SANDRI

DISSERTAÇÃO DE MESTRADO SUBMETIDA AO DEPARTAMENTO DE ENGENHARIA ELÉTRICA DA FACULDADE DE TECNOLOGIA DA UNIVERSIDADE DE BRASÍLIA, COMO PARTE DOS REQUISITOS NECESSÁRIOS PARA A OBTENÇÃO DO GRAU DE MESTRE.

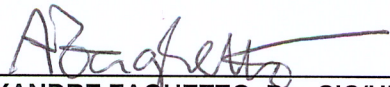
APROVADA POR:



RICARDO LOPES DE QUEIROZ, Dr., CIC/UNB
(ORIENTADOR)



JOÃO LUIZ AZEVEDO DE CARVALHO, Dr., ENE/UNB
(EXAMINADOR INTERNO)



ALEXANDRE ZAGHETTO, Dr., CIC/UNB
(EXAMINADOR EXTERNO)

Brasília, 16 de fevereiro de 2016.

FICHA CATALOGRÁFICA

SANDRI, GUSTAVO LUIZ

AUTOMATED NON-CONTACT HEART RATE MEASUREMENT USING CONVENTIONAL VIDEO CAMERAS [Distrito Federal] 2016.

xvi, 71 p., 210 x 297 mm (ENE/FT/UnB, Mestre, Engenharia Elétrica, 2016).

Dissertação de Mestrado - Universidade de Brasília, Faculdade de Tecnologia.

Departamento de Engenharia Elétrica

1. Pulse detection

2. Video processing

3. Heart rate

4. Photoplethysmography

I. ENE/FT/UnB

II. Título (série)

REFERÊNCIA BIBLIOGRÁFICA

SANDRI, G.L. (2016). *AUTOMATED NON-CONTACT HEART RATE MEASUREMENT USING CONVENTIONAL VIDEO CAMERAS*. Dissertação de Mestrado, Departamento de Engenharia Elétrica, Universidade de Brasília, Brasília, DF, 71 p.

CESSÃO DE DIREITOS

AUTOR: Gustavo Luiz Sandri

TÍTULO: AUTOMATED NON-CONTACT HEART RATE MEASUREMENT
USING CONVENTIONAL VIDEO CAMERAS.

GRAU: Mestre em Engenharia de Sistemas Eletrônicos e Automação ANO: 2016

É concedida à Universidade de Brasília permissão para reproduzir cópias desta Dissertação de Mestrado e para emprestar ou vender tais cópias somente para propósitos acadêmicos e científicos. Os autores reservam outros direitos de publicação e nenhuma parte dessa Dissertação de Mestrado pode ser reproduzida sem autorização por escrito dos autores.

Gustavo Luiz Sandri

Depto. de Engenharia Elétrica (ENE) - FT

Universidade de Brasília (UnB)

Campus Darcy Ribeiro

CEP 70919-970 - Brasília - DF - Brasil

Acknowledgments

I do not want to just acknowledge the help of those who contributed for this work, but also dedicate it to them.

I acknowledge, first, Eduardo Peixoto and Ricardo Queiroz, my tutors, who answered my prayers and gifted me with a work which I am proud of. I dedicate this work as an offering and sign of my gratitude for the patience, help, suggestions and all the time they devoted to me.

I want to extend my thanks to my parents, who supported me unconditionally and proudly told all their friends about the work I was developing, even when they wouldn't understand it in its plenitude.

Marcelo Campitelli, Ariane Alvez and Josua Carreño, my colleges and friends, are also subject to my thanks for their friendship and help during the evolution of this work, reading my manuscript and pointing what to improve. I'm thankful for their patience listening to me while I babbled all the things trapped in my head.

This work is also dedicated to Jonathan Alis and Wellington, the kind hearts with whom I shared a room and experiences at UnB. I'll will be eternally thankful for them having mercy on my soul and letting me borrow their computers to perform thousands simulations.

Least, but not less important, I want to thanks Geovany Borges for lending me the camera employed here and all volunteers who participated on my experiments. Without them this manuscript would never have reached this state.

Gustavo Luiz Sandri

ABSTRACT

As the blood flows through the body of an individual, it changes the way that light is irradiated by the skin, because blood absorbs light differently than the remaining tissues. This subtle variation can be captured by a camera and be used to monitor the heart activity of a person.

The signal captured by the camera is a wave that represents the changes in skin tone along time. The frequency of this wave is the same as the frequency by which the heart beats. Therefore, the signal captured by the camera could be used to estimate a person's heart rate. This remote measurement of cardiac pulse provides more comfort as it avoids the use of electrodes or others devices attached to the body. It also allows the monitoring of a person in a canceled way to be employed in lie detectors, for example.

In this work we propose two algorithms for non-contact heart rate estimation using conventional cameras under uncontrolled illumination. The first proposed algorithm is a simple approach that uses a face detector to identify the face of the person being monitored and extract the signal generated by the changes in the skin tone due to the blood flow. This algorithm employs an adaptive filter to boost the energy of the interest signal against noise. We show that this algorithm works very well for videos with little movement.

The second algorithm we propose is an improvement of the first one to make it more robust to movements. We modify the approach used to define the region of interest. In this algorithm we employ a skin detector to eliminate pixels from the background, divide the frames in micro-regions that are tracked using an optical flow algorithm to compensate for movements and we apply a clustering algorithm to automatically select the best micro-regions to use for heart rate estimation. We also propose a temporal and spatial filtering scheme to reduce noise introduced by the optical flow algorithm.

We compared the results of our algorithms to an off-the-shelf fingertip pulse oximeter and showed that they can work well under challenging situations.

RESUMO

Conforme o sangue flui através do corpo de um indivíduo, ele muda a forma como a luz é irradiada pela pele, pois o sangue absorve luz de forma diferente dos outros tecidos. Essa sutil variação pode ser capturada por uma câmera e ser usada para monitorar a atividade cardíaca de uma pessoa.

O sinal capturado pela câmera é uma onda que representa as variações de tonalidade da pele ao longo do tempo. A frequência dessa onda é a mesma frequência na qual o coração bate. Portanto, o sinal capturado pela câmera pode ser usado para estimar a taxa cardíaca de uma pessoa. Medir o pulso cardíaco remotamente traz mais conforto pois evita o uso de eletrodos. Também permite o monitoramento de uma pessoa de forma oculta para ser empregado em um detector de mentira, por exemplo.

Neste trabalho nós propomos dois algoritmos para a estimação da taxa cardíaca sem contato usando câmeras convencionais sob iluminação não controlada. O primeiro algoritmo proposto é um método simples que emprega um detector de face que identifica a face da pessoa sendo monitorada e extrai o sinal gerado pelas mudanças no tom da pele devido ao fluxo sanguíneo. Este algoritmo emprega um filtro adaptativo para aumentar a energia do sinal de interesse em relação ao ruído. Nós mostramos que este algoritmo funciona muito bem para vídeos com pouco movimento.

O segundo algoritmo que propomos é uma melhora do primeiro para torná-lo mais robusto a movimentos. Nós modificamos o método usado para definir a região de interesse. Neste algoritmo é utilizado um detector de pele para eliminar *pixels* do plano de fundo do vídeo, os frames dos vídeos são divididos em micro-regiões que são rastreados com um algoritmo de fluxo ótico para compensar os movimentos e um algoritmo de clusterização é aplicado para selecionar automaticamente as melhores micro-regiões para efetuar a estimação da taxa cardíaca. Propomos também um esquema de filtragem temporal e espacial para reduzir o ruído introduzido pelo algoritmo de fluxo ótico.

Comparamos os resultados dos nossos algoritmos com um oxímetro de dedo comercial e mostramos que eles funcionam bem para situações desafiadoras.

CONTENTS

1	INTRODUCTION	1
1.1	GOALS	2
1.2	CONTRIBUTIONS	2
1.3	PRESENTATION OF THE MANUSCRIPT	3
2	LITERATURE REVIEW	4
2.1	PHOTOPLETHYSMOGRAPHY	4
2.2	SKIN DETECTION	6
2.2.1	COLOR SPACES	7
2.2.2	COLOR THRESHOLDING	9
2.2.3	HISTOGRAM BASED APPROACH	10
2.2.4	GAUSSIAN MODEL	11
2.2.5	ARTIFICIAL NEURAL NETWORKS	11
2.3	PULSE DETECTION	12
2.3.1	PIXEL BASED	12
2.3.2	ALGORITHM OF POH <i>et. al.</i>	15
2.3.3	MICROMOVEMENTS	16
3	HEART RATE ESTIMATION WITH NOISE FILTERING	18
3.1	METHOD	18
3.2	PRELIMINARY RESULTS	22
4	HEART RATE ESTIMATION THROUGH MICRO-REGION TRACKING	25
4.1	SKIN DETECTION	26
4.2	MICRO-REGIONS	27
4.3	MICRO-REGION TRACKING	29
4.4	CLUSTERING	37
4.4.1	DISTANCE METRIC	37
4.4.2	ALGORITHM	39
5	RESULTS	42
5.1	DATABASE	42
5.2	EVALUATION OF HR-FD	44
5.2.1	STEADY VIDEOS	44
5.2.2	VIDEOS WITH MOVEMENT	46
5.2.3	SYNTHETIC DATA	47
5.2.4	CONCLUSION	49
5.3	EVALUATION OF HR-MRT	49

5.3.1	BLOCK DURATION	49
5.3.2	SKIN DETECTION	50
5.3.3	POINT TRACKING FILTERING	51
5.3.4	HR-MRT PERFORMANCE.....	52
5.4	ANALYSIS OF THE AUTOMATIC ROI SELECTION.....	54
5.5	SUMMARY	56
6	CONCLUSIONS.....	60
6.1	SUMMARY OF CONTRIBUTIONS	60
6.2	FUTURE WORK	61
	REFERENCES	62
	APPENDICES.....	68
I	AFFINE TRANSFORMATION MATRIX	69
II	RANDOM ESTIMATION OF THE HEART RATE.....	71

LIST OF FIGURES

2.1	Millimolar absorptivity spectra of hemoglobin in the visible range.....	5
2.2	Measurement of PPG signal using infrared light in contact with finger.	5
2.3	Schematic comparing PPG and ECG signals.	6
2.4	Region of interest, as employed in the literature, for heart rate detection.....	12
2.5	Schematic of the algorithm employed by Poh <i>et al.</i>	15
2.6	Signal windowing.	15
3.1	Schematic of the signal processing by HR-NF.	18
3.2	Distribution of the α_r , α_g and α_b values estimated using ICA	19
3.3	Low-pass filter employed to compute the adaptive filter mask.	20
3.4	Derivative filter.	21
3.5	Adaptive filtering.	21
3.6	Performance of the HR detection for three different videos	23
3.7	Effect of the adaptive filtering on the signal	24
4.1	Block diagram for the ROI definition in HR-MRT.	25
4.2	Skin detection algorithm.	26
4.3	Watershed segmentation method.....	27
4.4	Image smoothing with bilateral filtering.	29
4.5	Segmentation of an image using watershed.	30
4.6	Image segmentation in micro-regions	30
4.7	Optical flow estimation scheme.	32
4.8	Feature tracking filtering.....	33
4.9	Clustering algorithm.	40
5.1	Average of volunteers face from the first frame of the videos.	42
5.2	Oximeter data filtering.	43
5.3	HR detection performance of Poh versus HR-NF.....	44
5.4	Discrete Fourier transform of the noise in real data.	47
5.5	Discrete Fourier transform of the simulated noise.....	47
5.6	Evaluation of algorithm performance to integrated Gaussian noise.	48
5.7	Performance of HR-MRT for different block duration.	50
5.8	Performance of HR-MRT for different skin detection strategies.....	51
5.9	Performance of HR-MRT for different affine transform settings.	52
5.10	Performance of HR-MRT for different polynomial orders.	53
5.11	Performance of HR-MRT compared to HR-NF and Poh.	53
5.12	Percentage of time that the regions were chosen as ROI by HR-MRT.....	55

5.13	Performance of HR-MRT for the videos with movement, for each volunteer individually.....	56
5.14	HR detection for volunteer 18 - Steady video.	57
5.15	HR detection for volunteer 05 - Steady video.	58
5.16	HR detection for volunteer 05 - Video with movement.	58
5.17	HR detection for volunteer 13 - Video with movement.	59
5.18	HR detection for volunteer 18 - Video with movement.	59
II.1	Joint probability distribution.....	71

LIST OF TABLES

2.1	Examples of skin discrimination using explicit thresholding.	10
2.2	Algorithms employed in the literature for heart rate detection with video.....	14
5.1	Evaluation of the contribution of the adaptive filter, derivative filter and BSS strategy on HR-NF for the steady videos.	45
5.2	Average contributions of the adaptive filter, derivative filter and BSS strategy on HR-NF for steady videos.	45
5.3	Evaluation of the contribution of the adaptive filter, derivative filter and BSS strategy on HR-NF for the videos with movement.	46
5.4	Average contributions of the adaptive filter, derivative filter and BSS strategy on HR-NF for videos with movement.	46
5.5	SNR values where which algorithm reached the given percentage of correct estimations.	48
5.6	Parameters employed for HR-MRT.	49
5.7	Performance of the algorithms of Poh, HR-NF and HR-MRT.....	57

LIST OF SYMBOLS

Symbols

AC	As in Alternating Current, represents the signal component of frequency different than zero
DC	As in Direct Current, represents the signal component of frequency zero
f_s	Sampling frequency of the input video, in frames per second
T	Window duration
ΔT	Window time increment
$I[i]$	i -th frame of the input video
N_{FT}	Number of elements used to compute the DFT
N_T	Number of frames within T seconds
$N_{\Delta T}$	Number of frames within ΔT seconds
$p[j]$	j -th estimated Heart Rate
$x_r[i];$ $x_g[i];$ $x_b[i]$	The i -th red, green and blue traces, respectively, computed by the average value of the pixels inside the ROI, computed from the corresponding color channel
$y_r[j, k];$ $y_g[j, k];$ $y_b[j, k]$	The k -th normalized red, green and blue traces, respectively, on the j -th time window
$Z[j, v]$	Discrete Fourier Transform of the normalized trace on the j -th time window
$Z_f[j, v]$	Filtered Discrete Fourier Transform of the normalized trace on the j -th time window by the adaptive filter
$ \epsilon $	Absolute error between estimate heart rate and oximeter reading
σ_1, σ_2	Parameters of the bilateral filter
σ_s	Parameter of the similarity measure

Acronyms

BCG	Ballistocardiograph
BPM	Beats per minute
BSS	Blind source separation
CIE	<i>Commission Internationale de l'Eclairage</i>
DFT	Discrete Fourier transform
DRMF	Discriminative response map fitting
ECG	Electrocardiogram
FFT	Fast Fourier transform
FPS	Frames per second
HR	Heart rate
ICA	Independent component analysis
MLP	Multi layer perceptron
PCA	Principal component analysis
PPG	Photoplethysmography or photoplethysmographic
RMS	Root mean squared
ROI	Region of interest
SNR	Signal to noise ratio
SOM	Self organizing map
SpO ₂	Arterial oxygen saturation
STFT	Short time Fourier transform
HR-MRT	Video heart rate estimation through micro-region tracking
HR-NF	Video heart rate estimation with noise filtering
USB	Universal serial bus

Color Spaces

R	Red
G	Green
B	Blue
H	Hue
T	Tint
S	Saturation
V	Value
L	Lightness
I	Intensity
Y	Luma
Cr	Chrominance red
Cg	Chrominance green
Cb	Chrominance blue

1 INTRODUCTION

The heart contracts rhythmically to drive nutrients and oxygen necessary for life to our cells. The heart rate, that is the frequency by which the heart is beating, is a measure that can be used to verify a person's health and emotions.

As the heart beats, several characteristics of the body changes along the beating and can be employed to detect someones heart rate, such as electrical activity, blood pressure and light absorption by the skin. Electrocardiography (ECG) is the most precise and traditional method for detecting heart rate and some anomalies [1, 2]. It records the electrical activity due to the heart muscle depolarization in each heartbeat, using electrodes placed on the patient skin. One disadvantage of this technique is its need to use electrodes and equipment specially build for this purpose, making it harder for the everyday usage.

An alternative for the ECG is to use photoplethysmography (PPG) [3], a technique that uses light-based technology to capture the changes in skin light absorption as the blood flows. This is possible because the blood absorbs light differently than the other tissues of the skin. Hence, as the density of blood beneath the skin changes, as a result from the heartbeats, the total light absorption of the skin changes accordingly. Thus, methods for heart rate measurement that employ light as a signal measures the frequency by which the light emitted by the skin changes and attribute this frequency as the heart rate.

To capture the PPG signal, one can use a photodiode to produce a controlled light, place it near the skin and capture, with a photoreceptor device, either the light that traversed the skin or the one that was reflected by it. These devices are commonly called pulse oximeters [4]. In a less controlled environment, it is possible to avoid the use of a photodiode and employ the ambient light as a light source. The disadvantage is that, as the ambient light is hard or, sometimes, impossible to control, we do not know precisely which wavelengths and magnitudes are actually used. The advantage is that, without the used of a controlled light source, we can perform non-contact heart rate estimation using conventional cameras that are widely available. In this fashion, we trade precision in heart rate estimation for comfort, as we avoid the use of electrodes and/or sources of lights placed on the skin, and the possibility of non-contact estimation that enable us to monitor a patient for long time periods or even to monitor someone in a concealed way.

For heart rate estimation using videos, most methods focus their attention on the face of the person being monitored, because it has been shown that the face is a part of the body where the changes in skin absorption is sufficiently high to be perceived by a camera, even under natural light [5]. This is due to the fact that the skin in the face is more vascularized near the surface than in the rest of the body [6]. From a video, a region of interest comprising the face of the person being monitored is defined and the mean value of red, green and blue inside the region of interest is computed. These mean values are commonly called as red, green and blue traces. The heart

rate is then estimated examining how this traces change with time.

1.1 GOALS

This work aims to develop an algorithm to estimate the heart rate of human beings based on videos captured from their face under ambient light of indoors environments. The algorithm should be robust to noise and movements and capable of performing the estimation without supervision.

1.2 CONTRIBUTIONS

A number of algorithms have been proposed to estimate the heart rate using a video camera [5, 7–16]. Our work builds on top of these algorithms, modifying some key aspects in order to improve the overall reliability of the system. The main contributions of this work are as follows:

- An adaptive filter in order to make the algorithm perform better when facing noise;
- The way how the red, green and blue traces are combined, aiming to increase the number of correct estimated heart rate. Most algorithms in the literature employ blind source separation [11, 12, 14, 15, 17] or only the green channel [5, 8, 9, 13]. In this work we proposed a fixed mixture of the three signals, avoiding explicit blind separation. We justify why this approach is similar to blind separation of the signals, while avoiding artifacts that could be introduced by the blind source separator when it is not able to correct separate the sources;
- A method to compensate for movements while capturing the red, green and blue traces from the face. We divided the first frame of a time block of the video in micro-regions and tracked some points of them to determine how the micro-regions evolved with time. The information of the tracking points are temporally filtered using a polynomial model and spatial filtered using an affine transformation;
- A clustering algorithm to automatically select the best micro-regions to be used for heart rate estimation. Thus, the algorithm chooses the region of interest automatically having an arbitrary shape.

We also showed where are the best regions on the human face for heart rate estimation. Also, the proposed algorithm for point tracking filtering is new and could be used for other purposes.

1.3 PRESENTATION OF THE MANUSCRIPT

Chapter 2 presents a review of the literature. We explain the photoplethysmographic signal, its bases and why we are able to capture it and estimate a person's heart rate using a conventional camera. We review skin detection algorithms that we employ in our work to define the region of interest and we present an overview of other algorithms in the literature that estimate the heart rate using conventional cameras.

Chapters 3 and 4 disclose the proposed algorithms. In Chapter 3 we describe the first proposed algorithm, introducing an adaptive and a derivative filter and how the signal from the three color channels are combined to enhance the signal to noise ratio. We also present evidences that this algorithm do not perform very well in the presence of movement. Therefore, Chapter 4 describes the second proposed algorithm, an improvement of the first to make it robust to movements. In this algorithm we divide the frames in micro-regions and use point tracking and filtering to compensate for movements. We also present the clustering algorithm used to select micro-regions for heart rate estimation, automatically defining the region of interest.

Finally, Chapter 5 presents the results obtained with our algorithms, compared to the literature, for real and synthetic signals. Chapter 6 presents the conclusion and the future work.

2 LITERATURE REVIEW

This chapter addresses the main topics found in the literature related to our work. A complete review is out of the scope of this manuscript, hence, we concentrate our attention to explain the ideas upon which this work is based. Further literature is also found in subsequent chapters.

Section 2.1 is an introduction to the biological concepts that explain how one can estimate the heart rate using video cameras. Section 2.2 is dedicated to explain skin detection algorithms that are commonly used for heart rate detection based on the fact that the interest regions of the videos is the subject skin. In section 2.3 we present the work of other authors in the domain of heart rate detection using conventional cameras.

2.1 PHOTOPLETHYSMOGRAPHY

The signal of interest in our work is the wave produced by blood flowing through the body, known as the photoplethysmographic (PPG) waveform [18]. With each heart beat, a pressure pulse radiates out to the peripheral circulatory system causing significant change in the arterial and capillary diameters [3, 19], thus increasing the volume of blood beneath the skin. As pressure decreases, arteries return to their normal size.

Blood possesses a great amount of red cells, responsible for oxygen transportation. These cells contain hemoglobin, that absorbs light differently from others structures of epithelial tissue [20]. Skin tissue has a relatively low absorptivity for wavelengths in the visible and near-infrared light spectrum (400–2000 nm). This characteristic is mainly due to skin pigments (particularly melanin) and water [21]. On the other hand, the absorption for the hemoglobin varies depending whether it is loaded or not with oxygen or carbon monoxide [20]. Two of these variants are of interest in this work: the oxyhemoglobin and de-oxyhemoglobin.

The oxyhemoglobin (HbO_2) is formed when oxygen binds to hemoglobin in red blood cells during physiological respiration, while de-oxyhemoglobin (Hb), or deoxygenated hemoglobin, is the form of hemoglobin without the bound oxygen [20]. Figure 2.1 shows their absorption spectra. Oxyhemoglobin has significantly lower absorption for green light at 560 nm and a higher absorption for blue light at 480 nm.

As blood flows, the density of oxyhemoglobin and de-oxyhemoglobin underneath the skin changes periodically. The amount of red cells near the surface of skin also changes, which alters the average distance that light must travel before being reflected. Light, while traveling through the skin tissue, interacts with it, resulting primarily in reflection and absorption, although scattering, transmission and fluorescence may also be present [22]. This affects the way that skin radiates back the environment light.

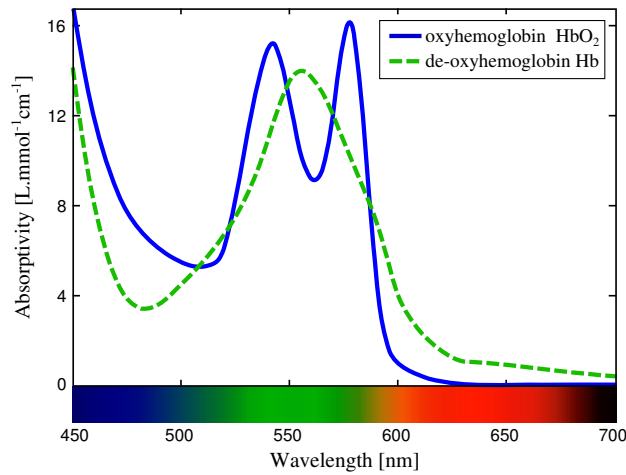


Figure 2.1: Millimolar absorptivity spectra of hemoglobin in the visible range. Adapted from Zijlstra *et al.* [20]

The most traditional way to measure the photoplethysmographic signal is to use an infrared light source (usually a photodiode) to illuminate the skin and a photodetector placed either in the opposite side, to capture the transmitted light, or on the same side, to capture the reflected light [18] (Figure 2.2). The main peripheral sites where PPG can be obtained are fingers' tissue pads, ears and toes where there is a high degree of superficial vasculature [23].

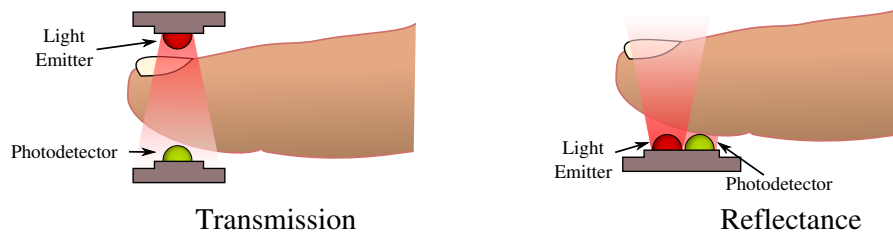


Figure 2.2: Measurement of PPG signal using infrared light in contact with finger. Either the transmitted or reflected light is used, captured by a photodetector.

The changes on the skin absorption due to blood circulation is a phenomenon that has been known for a long time [4]. It was first described by Alrick Hertzman in 1937 [24]. He believed, based on his observations, that the origin of the PPG signal was linked to blood volume changes. Therefore, he named it as "photoelectric plethysmograph". The term "plethysmos" derives from the Greek word for fullness and expressed his belief that he was measuring the fullness of the tissue when he measured the amount of light absorption. Later researches demonstrated that he was not far on his assumption, with results for the PPG being close to the more traditional strain gauge plethysmograph [25], a method that provides a quantitative measure of the bloodflow.

Electrocardiogram (ECG) and PPG signals share some characteristics, as both of them originate from the heart beats. For example, they have the same fundamental frequency (the heart rate) and the systole and diastole are visible on both of them [18]. However, their waveform differs significantly, as it can be seen in Figure 2.3, which is related to the way how each wave travels through the body and how it is measured. Also, one can notice a phase difference between them due to the fact that electric waves travel much faster than pressure waves, originating different

delays.

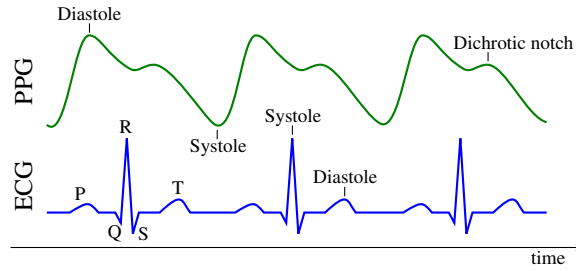


Figure 2.3: Schematic comparing PPG and ECG signals.

More recently, it has been shown that the PPG signal can be estimated from videos captured with a commercial camera filming the face of the subject being monitored [9, 12, 14, 26–28]. Conventional cameras divide the input light into three channels using a color filter [29]. As the PPG signal results in a periodic variation of the reflected light spectrum, this signal can be appreciated in all three color channels multiplied by a constant that is intrinsic to the light wavelengths gathered on each channel of the camera. Therefore, each channel will receive an attenuated version of the PPG signal, added to noise coming from movement artifacts, sensor temperature (depending on the sensor quality), among other noise sources.

Wu *et al.* [28] used conventional cameras and Eulerian Video Magnification to make the small changes in the skin color, due to blood circulation, visible to the naked eye. The aim of this work was only concentrated on the amplification of the skin color change. They did not try to detect the subject pulse, but their works shows the feasibility of contactless estimation of the heart rate using conventional cameras. Later, they presented Phase-Based Video Motion Processing [30], a modification of the previous algorithm, dedicated to amplify exclusively the motion. This new algorithm can be employed in the previous algorithm to make it more robust to subtle movements, resulting in a better outcome. But this approach only works when the movements are of low to very low amplitudes.

Among the image processing tools, skin detection stands out, as it is used by many techniques. A detailed description of skin detection techniques is presented next, and a more thorough review of works in heart rate detection using cameras is given in Section 2.3.

2.2 SKIN DETECTION

Skin detection is an important tool in heart rate estimation through videos as it can segment skin regions — where the PPG signal is likely to be found — from the rest of the scene. But it is also applied to a large range of other applications, including face detection [31], content-based image retrieval [32], and nudity detection [33].

The best performance for skin detection is acquired using the visual and non-visual spectrum of light, for-instance infrared [34, 35]. However, the majority of researches concentrate their effort

only in visible light, given the fact that non-visible spectrum information is usually not available (or too expensive to capture) and the wide use of conventional cameras that only capture red, green and blue wavelengths.

Therefore, they must cope with problems that arise with this limitation, mainly illumination, ethnic group and camera characteristics, that influence the skin color [36]. Illumination is a problem because skin detectors don't assume a controlled environment and should deal with indoor and outdoor scenes, shadows, highlights, etc. Skin color varies from white, yellow, reddish to dark accordingly to the subject ethnicity [37]. Finally, the camera, even under the same illumination, may present different resulting images depending on sensor sensitivity and color filter. There are also some other minor factors that may influence the detection, such as makeup, motion blur, etc.

The simplest skin detectors are the ones called pixel-based detectors. They classify each pixel as skin or non-skin independently from its neighborhood [37]. In this case, the crucial issue is to choose the color space more suitable to solve the problem. Surprisingly, many papers related to skin detection do not provide a strict justification of their choice [37]. Just some few works were devoted to do a comparative analysis of different color spaces used for skin detection [38–41].

We dedicate Section 2.2.1 to review the most employed color spaces for skin detection. A summary of the main techniques used for skin detection is presented in Sections 2.2.2 to 2.2.5.

2.2.1 Color Spaces

The performance of a given color space for skin detection is related to how well they can separate the skin pixels from non-skin pixels and is measured by the classification error obtained on a test data [36]. It has been observed that skin colors vary more in intensity (or luminance) than in chrominance [42]. Therefore, it has become a common practice to drop the luminance component in order to add robustness to the classification when facing different illumination and ethnicity.

We present below the main color spaces exploited for skin detection.

2.2.1.1 RGB

RGB originated from cathode ray tube display applications and it is most commonly used for storing and representing digital images, since the human eye naturally captures the three primary colors used for image representation (red, green and blue), which gave the name to this color space [43]. The RGB representation was standardized in 1930 by the *Commission Internationale de l'Eclairage* (CIE), using the primary colors of red (700.0 nm), green (546.1 nm) and blue (435.8 nm) [29].

However, these three channels are highly correlated and their dependencies with illumination make this space not a favorable choice [37]. Nevertheless, some researchers have successfully

employed this color space for the purpose of skin detection, avoiding the conversion of the colors to another space, such as Brand and Mason [44] and Jones and Rehg [45].

To overcome this problem we can perform a simple normalization on the R , G and B components, resulting in the normalized RGB (or simply rgb):

$$r = \frac{R}{R + G + B}, \quad g = \frac{G}{R + G + B} \quad \text{and} \quad b = \frac{B}{R + G + B}, \quad (2.1)$$

where R , G and B are the levels of red, green and blue, respectively. We can reduce the space dimensionality dropping the third component, as it becomes redundant after the normalization.

It has been reported that the normalized RGB space is more robust to skin color changes due to lighting and ethnicity and the clusters in rgb space present a lower variance than the corresponding cluster in the RGB space [42, 46].

The ratio between the colors in the RGB space is also used to detect the presence of skin. It was observed that skin contains a significant high level of red, independent of ethnicity [47]. Therefore, the R/G ratio has been used for skin detection [47]. Others ratios (R/B and G/B) were tested for skin detection by Brand and Mason [44].

2.2.1.2 Perceptual Color Spaces

Developed in the 1970s for computer graphics applications, hue-saturation color spaces were introduced to numerically represent the artistic idea of tint, saturation and tone [48]. The dominant color (red, yellow, purple, etc.) is represented by the Hue (or Tint) while saturation represents how pure the color is, starting with gray when the saturation is zero and going up to pure color when the saturation is maximum. The third component (Intensity, Value or Lightness) measures how bright the color is.

They cannot be directly described by the RGB space, but many non-linear transformations were proposed to relate them with the RGB space, for example [37]:

$$H = \cos^{-1} \left(\frac{1}{2} \frac{(R - G) + (R - B)}{\sqrt{(R - G)^2 + (R - B)(G - B)}} \right) \quad (2.2)$$

$$S = 1 - 3 \frac{\min(R, G, B)}{R + G + B} \quad (2.3)$$

$$V = \frac{R + G + B}{3} \quad (2.4)$$

These transformation are invariant to highlights at white lights, ambient light and orientation relative to light sources, making it a good choice for skin detection [29]. Another color space similar to the hue-saturation is the normalized chrominance-luminance TSL (tint, saturation, lightness) space, that is a transformation of the normalized RGB into more intuitive values, close to hue and saturation in their meaning.

2.2.1.3 Chrominance Based Spaces

YCbCr is an orthogonal color space, commonly used for image compression, that uses statistical independent components aiming to reduce the redundancy of RGB [29]. Color is represented by *luma* (Y), which is the luminance computed from a weighted sum of RGB values, and two chrominances: chrominance blue (Cb) and chrominance red (Cr) that are computed by subtracting *luma* from red and blue components.

$$Y = 0.299R + 0.587G + 0.114B \quad (2.5)$$

$$C_b = B - Y \quad (2.6)$$

$$C_r = R - Y \quad (2.7)$$

The explicit separation between luminance and chrominance and its simplicity makes YCbCr one of the most popular choices for skin detection. Other similar color spaces, such as YCgCr, YIQ, YUV and YES, derived from YCbCr, are also employed.

2.2.1.4 Psychophysical Based Spaces

The RGB space is not a perceptually uniform color space, meaning that the same amount of variation in the values of red, green and blue, applied at two different levels, will not be perceived the same way by the human brain. Therefore, the CIE, based on psychophysical experiments, introduced the CIELAB and CIELUV, that try to match the characteristics of human visual system [49]. The price for better perceptual uniformity is complex transformation functions from and to RGB space.

2.2.2 Color Thresholding

Given a color space, the components of skin pixels from different individuals tend to cluster in a small region [36]. Hence, one simple method to discriminate skin pixels is to explicitly define the boundaries of the skin cluster. One or more color spaces can be chosen and a pixel is classified as skin only if its parameters fall within all predetermined ranges.

The obvious simplicity of this method and its easy implementation, that leads to fast computation, has attracted many researchers [50–53]. The main difficulty of this approach is the need to empirically select the decision boundaries, which is strongly dependent on the chosen color space. So, both a good color space and good decision boundaries must be found. Also, the cluster of skin pixels overlaps the cluster of non skin pixels as there are several colors in nature very similar to that of skin. This contributes to make skin pixels discrimination a hard task.

Table 2.1 presents examples of thresholds employed for skin detection as proposed by some authors.

Table 2.1: Examples of skin discrimination using explicit thresholding.

Authors		Color Space	Boundaries
Peer <i>et. al.</i>	[50]	RGB	$R > 90, G > 40, B > 20, R > G + 15, R > B,$ $\max(R, G, B) - \min(R, G, B) > 15$
Tsekeridou and Pitas	[51]	HSV	$V \geq 40, 0.2 < S < 0.6,$ $0^\circ < H < 25^\circ$ or $335^\circ < H < 360^\circ$
Chai and Ngan	[52]	YCbCr	$77 \leq Cb \leq 127, 133 \leq Cr \leq 173$
Khan <i>et. al.</i>	[53]	CIELAB	$2 \leq A \leq 14, 0.7 \leq B \leq 18$

2.2.3 Histogram Based Approach

This method is used to produce a probabilistic classifier. It employs a database of previously segmented image pixels into two groups: skin and non-skin; and computes a 3D or 2D color histogram. In the case of 2D histograms, the brightness component is dropped out. The pixel components are quantized into a number of histogram bins that stores the number of times the given bin color occurred in the test data [37].

Two histograms are computed, one for skin pixels and another for non-skin pixels. These histograms are then normalized by their total sum for producing an estimation of the $P(c|skin)$ and $P(c|\overline{skin})$ (where c is a color column vector) that represent the probability of occurrence of c , given that it is a skin pixel and a non-skin pixel, respectively.

Finally, with the estimated probabilities, a lookup table is computed assigning that a bin is a skin pixel whenever its probability of occurrence is higher in the skin class than in the non-skin class.

For a more complete representation, one can use a Bayes classifier, that classifies a pixel as skin when

$$\frac{P(skin|c)}{P(\overline{skin}|c)} > \Theta, \quad (2.8)$$

for a given detection threshold Θ . $P(skin|c)$ and $P(\overline{skin}|c)$ can be computed from $P(c|skin)$ and $P(c|\overline{skin})$, resulting in

$$\frac{P(c|skin)P(skin)}{P(c|\overline{skin})P(\overline{skin})} > \Theta, \quad (2.9)$$

where $P(skin)$ and $P(\overline{skin})$ are estimated from the database.

2.2.4 Gaussian Model

The cluster formed by skin pixels in the chosen color space can be modeled by a multivariate Gaussian distribution [42], defined as:

$$P(c|skin) = \frac{1}{2\pi\sqrt{|C|}} \exp\left(-\frac{1}{2}(c - \mu)^T C^{-1}(c - \mu)\right), \quad (2.10)$$

where μ is the mean value and C the covariance matrix of the cluster. Let c_i be the parameters of i -th skin pixel in the database and N the total number of skin pixels, then:

$$\mu = \frac{1}{N} \sum_i c_i \quad (2.11)$$

$$C = \frac{1}{N-1} \sum_i (c_i - \mu)(c_i - \mu)^T \quad (2.12)$$

After estimating the parameters μ and C , $P(skin|c)$ is used to calculate the likelihood of pixel c to be a skin pixel and the classification is performed by means of thresholding this value.

To represent a more complex-shaped cluster, a weighted sum (also called mixture) of Gaussians can be employed in the form [54]:

$$P(c|skin) = \sum_j \frac{w_j}{2\pi\sqrt{|C_j|}} \exp\left(-\frac{1}{2}(c - \mu_j)^T C_j^{-1}(c - \mu_j)\right), \quad (2.13)$$

where w_j is the weight attributed to the j -th Gaussian and μ_j and C_j are its parameters.

2.2.5 Artificial Neural Networks

Artificial intelligence approaches have been successfully applied to skin detection due to their ability to represent complex non-linear relationship between inputs and outputs [55]. The two most used algorithms are the self organizing map (SOM), devised by Kohonen, that is based on a competitive algorithm, and the multi layer perceptron (MLP), a simple feed forward network [37]. Their main difference is that SOM is an unsupervised algorithm, *i.e.*, the input training data is not labeled, while the MLP is a supervised algorithm.

In the MLP, the weights in each neuron are updated iteratively using a training set for which we know the correct outcome. At each iteration of the training process, all inputs are processed in a feed-forward fashion and the output is compared with the expected result. Through a gradient descent technique, the error is back-propagated to adjust the neurons weights.

The SOM, on the other hand, constructs a topological structure that represents high dimensional data in a fashion that allows relative distances to be preserved. During the training phase, the weights of each neuron are updated to become similar to the input data, and neighboring neurons will have similar weights, thus creating a clustering.

Neural networks have shown good performance, as they can easily generalize complex structures [56], but the network must be tuned in order to obtain an optimal performance.

2.3 PULSE DETECTION

Detection of heart rate using cameras has the advantage of being a non invasive technique. For example, it enables one to monitor, in comfort, a patient for a long time or to monitor a subject in a concealed manner, where knowing that he is being monitored wouldn't be desired, such as in a lie detector. However, when dealing with heart rate detection with cameras one must take into account other sources of errors, such as movements, uncontrolled illumination, etc.

There are two main approaches to estimate the heart rate of a subject using conventional cameras without contact. The first is the pixel based approach where the estimation is based on how skin pixel change their tone with time [8,9,14]. The second one is based on micromovements of the body, that are detected using feature tracking [57].

2.3.1 Pixel Based

In the pixel based approaches, a conventional camera films the subject being monitored and the PPG signal is captured by three color channels, red, green and blue. One or multiple regions of interest (ROI) are defined and the average value for the red, green and blue, within the ROI, is computed over time. These three signals are then employed to estimate the heart rate, either in time or frequency domain. Table 2.2 is a summary of pixel based algorithms employed in the literature.

Most algorithms, such as those of Poh *et al.* [14, 15], Kwon *et al.* [17] and Purshe *et al.* [12], uses a cascade classifier to detect the subject face as they have shown that the PPG signal is relatively easy to detect on the face skin. Capdevila *et al.* [5] have shown that the forehead, cheeks and chin are the best regions on the face to measure the PPG signal. The ROI was defined as a rectangle comprising most of the subjects face (Figure 2.4).

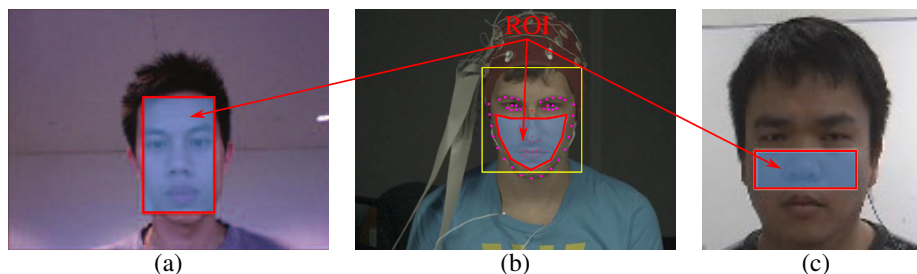


Figure 2.4: ROI, represented in red, as defined by (a) Poh *et al.* [14], (b) Li *et al.* [9] and (c) Yu *et al.* [11]. Adapted.

Li *et al.* [9] employed Discriminative Response Mat Fitting (DRMF) [58], a more robust face fitting algorithm, that estimates a set of parameters describing the position of the eyes, mouth,

chin, nose and eyebrows. They defined the ROI as the region comprising mainly the cheeks, as shown in Figure 2.4.

From the average values computed on the ROI, one for each color channel, researchers either use only the green channel (based on the fact that the PPG wave is more strongly present in this channel [16]) or a combination of the three, by means of Independent Component Analysis (ICA), Principal Component Analysis (PCA), or with fixed weight, linearly mixing them in an attempt to maximize the Signal to Noise Ratio (SNR) for the PPG signal. ICA is the preferred approach among researchers. Band-passing the signal to eliminate those frequencies that do not correspond to the heart rate (HR) is also a common practice.

Xu *et al.* [10], on the other hand, computed the discrete derivative of the logarithmic ratio between red and green traces. The use of the logarithmic function is explained making use of the Lambert–Berr law, where the signal received by each color channel can, approximately, be said to be proportional to $e^{-(\nu(\lambda)\rho(t)+A_0(\lambda))}$, where $\nu(\lambda)$ is the absorptivity of hemoglobin multiplied by the mean path that light travels before being reflected and $A_0(\lambda)$ the absorbance of other tissues in the skin, respectively, for a given wavelength λ . The concentration of hemoglobin is given by $\rho(t)$, and vary with time t . The discrete time derivative is used to eliminate DC components and to attenuate low frequency noises.

The signal is converted to the frequency domain using Discrete Fourier Transform (DFT) [12, 14–17], Short Time Fourier Transform (STFT) [11] or Welch periodogram [9, 13] (a method used to estimate the power spectra [59]), and the frequency corresponding to the peak of maximum power is attributed as the heart rate (HR) frequency. The search of the HR is limited to a range of values where they expect the pulse frequency to be.

Some works also use an approach in the time domain. Couderc *et al.* [8], for example, first applies a band-pass filter to the signal to eliminate those frequencies that do not correspond to the pulse and cubic interpolate the signal. Then, they search for the position of the peaks and valleys and the HR frequency is computed as the inverse of the signal period. The problem with this approach is that it degrades rapidly with noise.

Table 2.2: Summary of algorithms employed in the literature for HR detection with video.

Authors	ROI	Camera	Signal employed ¹
Ufuk Bal [7]	face ² (Viola Jones + skin detection)	Webcam; 640x480; 30 FPS	fixed mixture of R, G and B, denoised and band-pass filtered
Couderc <i>et al.</i> [8]	face ³ (entire frame)	Microsoft LifeCam Cinema; 2MP; 15 or 30 FPS	green channel band-pass filtered and interpolated
Li <i>et al.</i> [9]	face ⁴ (DRMF)	IPAD (built-in iSight camera); 640x480; 30 FPS	green channel band-pass filtered + Welch periodogram
Xu <i>et al.</i> [10]	manual selection ³	Conventional camera and smartphones (using a dermatoscope to amplify the skin 20x)	derivative of $\log(R/G)$ band-pass filtered + DFT
Balakrishnan <i>et al.</i> [57]	face ⁴ (Viola Jones)	Panasonic Lumix GF2; 1280x720; 30 FPS	micro-movements (vertical) band-pass filtered + PCA
Yu <i>et al.</i> [11]	face ²	conventional camera; 720x576; 25 FPS	ICA + STFT
Capdevila <i>et al.</i> [5]	face ³ (manual)	Canon Ixus 80is; 640x480;	green channel band-pass filtered and interpolated
Kwon <i>et al.</i> [17]	face ² (Viola Jones)	smartphone (iPhone 4); 640x480; 30 FPS	ICA + DFT
Pursche <i>et al.</i> [12]	face	webcam; 1.3 MP; 30 FPS	ICA + DFT
Bolkhovskiy <i>et al.</i> [13]	right index finger ³ (covering the lens)	smartphones; 30 or 20 FPS	green channel interpolated + Welch periodogram
Poh <i>et al.</i> [14, 15]	face ² (Viola Jones)	Macbook Pro (built-in iSight camera); 640x480; 15 FPS	ICA + DFT
Verkruyssen <i>et al.</i> [16]	manual selection ³	Canon Powershot A560; 640x480 or 320x240; 15 or 30 FPS	R, G, and B channels (individually) + DFT

¹Describes how the red, green and blue traces computed on the ROI are used to estimate de HR. Most of them remove the DC component from the traces. Thus this step was omitted.

² Reset at each frame

³ Fixed region over time

⁴ Updated through point tracking

2.3.2 Algorithm of Poh *et. al*

The work we developed here is based on the algorithm of Poh *et al.* Many of the related works in this field follow a similar structure (or are based in this algorithm) as shown in Table 2.2. Therefore, we explain it in more details here. In the following of this manuscript, we refer to this algorithm as Poh, for simplicity.

Figure 2.5 presents a schematic of their algorithm.



Figure 2.5: Schematic of the algorithm employed by Poh *et al.*

Let $I[i]$ be the i -th frame of the input video. For each frame, a cascade of boost classifier that uses 14 Haar-like features trained with positive and negative examples of frontal faces, based on the work of Viola and Jones [60] and Lienhart and Maydt [61], is used to detect the position of the face. For each region detected as face, the algorithm returns the x and y coordinates along with the height and width of a square that describes the position of the face. The ROI is defined as the rectangle centered on the square found by the algorithm, with 60% of its width and 100% of its height. Whenever the algorithm is not able to find the ROI for a given frame, the ROI from the previous frame is employed. The average value of the red, green and blue channel of the pixels inside the ROI is stored on the signals $x_r[i]$, $x_g[i]$ and $x_b[i]$, respectively.

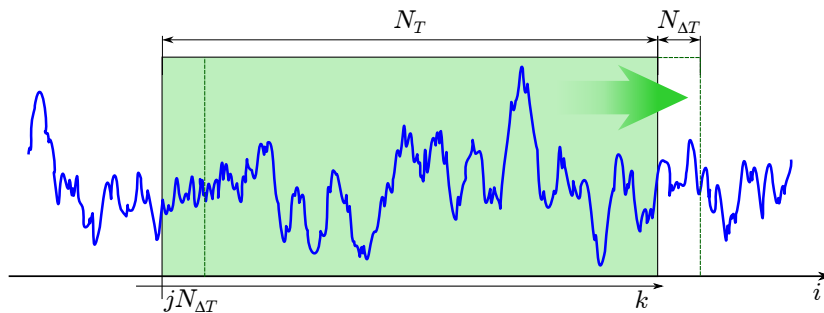


Figure 2.6: Signal windowing.

The normalization phase is executed in a window of duration T ($T = 30$ s in their algorithm), as shown in Figure 2.6. The window moves with ΔT seconds of increment ($\Delta T = 1$ s). The normalized traces are computed removing the mean (DC component) from $x_c[i]$ (where c denotes one the three color channels) and adjusting its amplitude in order to obtain a signal with unitary variance, as follows:

$$y_c[j, k] = \frac{x_c[jN_{\Delta T} + k] - \mu_c[j]}{\sigma_c[j]}, \quad 0 \leq k < N_T. \quad (2.14)$$

The integer $j \geq 0$ is the window number and defines the window starting point. k is the relative position inside the window. We define N_T and $N_{\Delta T}$ as the number of frames comprised in T and

ΔT seconds, respectively. The mean and variance (μ and σ , respectively) at the j -th window are given by:

$$\mu_c[j] = \frac{1}{M} \sum_{k=0}^{N_T-1} x_c[jN_{\Delta T} + k] \quad (2.15)$$

and

$$\sigma_c[j]^2 = \frac{1}{M} \sum_{k=0}^{N_T-1} (x_c[jN_{\Delta T} + k] - \mu_c[j])^2. \quad (2.16)$$

For a fixed j , $y_c[j, k]$ is a vector of duration T over which we try to determine the subject's HR. As j progresses, the window moves ΔT seconds along the signal, with an overlap of $(T - \Delta T)/T$.

Independent Component Analysis, based on the work of Cardoso [62], is used to separate the PPG signal from other noise components that are comprised in the $y_c[j, k]$ traces. This algorithm uses fourth-order cumulant tensors to automatically define weights for a linear mixture of the three channels in order to provide the best SNR, resulting in $z[j, k]$.

Finally, the Discrete Fourier Transform (DFT) is applied over $z[j, k]$ to obtain its spectrum $Z[j, v]$ and a peak detector determines the component of highest power within the range between 45 to 240 Beats per Minute (BPM), that correspond to the HR they expect to find for an adult individual. To account for noise, if the absolute difference between the current estimated pulse rate is above 12 BPM from the last computed value, the algorithm reject the actual estimation and searches for the next highest power component that meet this constraint. If no frequency peaks meet this criteria, then the current pulse frequency is retained.

Some works also have shown that arterial oxygen saturation (SpO_2) [63] can also be contactless estimated using cameras by means of the light reflected by the skin [64, 65]. However, they require a device capable of capturing light in the infrared spectrum and a controlled illumination in order to provide good estimation, since SpO_2 is computed measuring the amount of light absorbed by oxyhemoglobin and de-oxyhemoglobin.

2.3.3 Micromovements

Balakrishnan *et al.* [57] also proposed an approach to estimate HR using videos, but instead of trying to capture the PPG signal, they focus their attention on the ballistocardiograph (BCG) signal, that correspond to subtle movements of the body due to the heart beats.

The head is subject to movement in most axis and can be considered, for small amplitude movements, an inverted pendulum. As the blood enters and leaves the head, propelled by the heart, micro-movements appears.

The ROI is defined as being the head with the eyes removed (found with Viola Jones face detector). Head movements are estimated using feature tracking and Principal Component Analysis is used to give more robustness against noise.

This approach is successful for detecting the HR and can be employed even when no skin is visible on the video (when the subject is using a mask, for example). However, it suffers many others limitations because respiration, posture changes and voluntary or involuntary movements, that have higher amplitude, dominate the trajectory of the tracking points. This makes it harder for the tracker to perceive the BCG signal, negatively affecting the estimation.

3 HEART RATE ESTIMATION WITH NOISE FILTERING

This method is based upon Poh *et al.* work [14, 15]. Our contribution was to improve the robustness to noise by adding an adaptive and derivative filter to boost the energy of the PPG signal over noise. We also propose and employ another approach to mix the information from the red, green and blue channels, avoiding the use of ICA that adaptively determine the weight for each color channel, in an attempt to enhance the SNR.

In this manuscript, we refer to this proposed method as “video heart rate estimation with noise filtering” (HR-NF), that employs the face of the person being monitored as region of interest, from where the PPG signal will be extracted.

We present our method in Section 3.1. Some preliminary results are given in Section 3.2.

3.1 METHOD

Figure 3.1 presents a block diagram of the processing employed to estimate the HR from the video in this approach. This is very similar to the schematic for Poh *et al.* given in Figure 2.5, except for those blocks in red that were added or significantly modified.

$I[i]$ represents the i -th frame, that was acquired with a constant sampling rate, and $p[j]$ is the j -th HR detected. Signals x_r, x_g and x_b are computed from the ROI and normalized to the signals y_r, y_g and y_b following the same strategy, with a moving window with $T = 30$ s and $\Delta T = 0.5$ s. The HR is calculate at each 0.5 s and need a sequence of frames to be estimated.

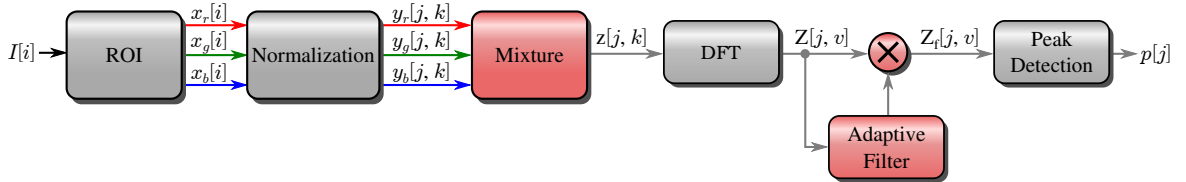


Figure 3.1: Schematic of the signal processing by HR-NF.

For the mixture stage, on the other hand, where y_r, y_g and y_b are linearly combined to provide a better SNR for the PPG signal, we removed the ICA to blindly separate the sources from the signal processing. In the ICA approach, three constants, α_r, α_g and α_b , are estimated to mix the signals as:

$$z[j, k] = \alpha_r y_r[j, k] + \alpha_g y_g[j, k] + \alpha_b y_b[j, k]. \quad (3.1)$$

However, we noticed in our work that even though the α values vary greatly from one video to another, or even from one signal site to another, their ratios are relatively constant, being almost

independent from skin tone and scene illumination. Therefore, we normalized the α values, making them comparable between each other, using the constraints $\alpha_g > 0$ and $|\alpha_r| + |\alpha_g| + |\alpha_b| = 1$. Figure 3.2 presents the distribution obtained for these values in our database, where we can observe their tendency to concentrate in a given region.

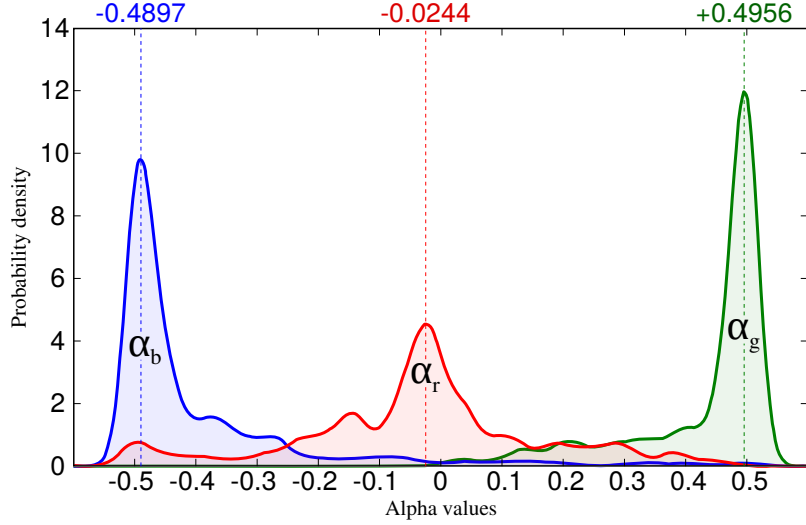


Figure 3.2: Distribution of the α_r , α_g and α_b values estimated using ICA, after normalization, in our database. The α values are represented by their corresponding colors

Due to this tendency, we decided to assume these values as constants, given by their position of highest probability density, reducing the algorithm complexity as no blind source separation, which is computationally expensive, must be performed. Also, when the SNR is low or the noise characteristics are similar to the PPG signal, the ICA is not capable to precisely estimate the weights to employ in order to obtain a good SNR. This can lead to errors and instabilities on the estimation. As in our algorithm the weight of each channel is a constant, we are not susceptible to such errors.

The values employed are $\alpha_r = -0.0244$, $\alpha_g = +0.4956$ and $\alpha_b = -0.4897$, based on the distribution shown in Figure 3.2. These values are in accordance with the literature because the green channel is the one presenting the highest energy for the PPG signal [9, 16]. But they are dependent on the camera parameters, such as sensor sensitivity, color filter, color and gamma correction. Thus, they must be determined for each camera.

After the mixture, for a fixed j , the signal $z[j, k]$ is zero padded to contain a total of $N_{FT} = 2^{14}$ elements. The magnitude of its DFT is computed, resulting in $Z[j, v]$, where v represents the frequency. We retain only those frequencies in the range going from 30 to 240 BPM that correspond to the values that we consider acceptable for a HR measurement for an adult individual. That is,

$$30 \frac{N_{FT}}{f_s} \leq v \leq 240 \frac{N_{FT}}{f_s}, \quad (3.2)$$

where f_s is the signal sampling frequency, given in frames per second.

An adaptive filter is then applied onto $Z[j, v]$, multiplying it by a mask $M[j, v]$ (on frequency

domain), resulting in the signal

$$Z_f[j, v] = Z[j, v]M[j, v]. \quad (3.3)$$

The mask aim is to amplify $Z[j, v]$ for those frequencies that have a higher probability of being the HR frequency and attenuate others, reducing the effect of noise in the estimation. The mask is defined supposing that the HR varies slowly with time and that within $\Delta T = 0.5s$ it should be almost the same. Therefore, the previous signals, $Z[j - 1, v]$ and $Z[j - 2, v]$, provide a good estimation of where the HR peak should be.

In order to accommodate small variations of the HR, we first apply a convolution with a low-pass filter to $Z[j - 1, v]$ and $Z[j - 2, v]$ (the convolution is made in the frequency domain). This filter will horizontally stretch the peaks and is almost equivalent to applying an apodization function (or window function) [66] on the time domain. For the low-pass filter, we chose $T_2[v]$, a triangular function with bandwidth of 2 BPM, as shown in in Figure 3.3, that roughly allows a ± 1 BPM variation.

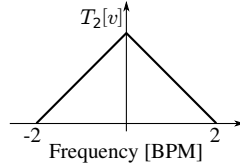


Figure 3.3: Low-pass filter employed to compute the adaptive filter mask.

As $Z[j - 2, v]$ is further away in time to the j -th signal than $Z[j - 1, v]$, we apply the same low-pass filter again in order to accommodate higher HR variations. The mask is given by a combination of them as

$$M[j, v] = (T_2[v] * Z[j - 1, v]) (T_2[v] * T_2[v] * Z[j - 2, v]), \quad (3.4)$$

where $*$ denotes convolution.

Also, it was depicted by Xu *et. al.* [10] that using the derivative of the traces improves the HR detection performance, as the noise tends to be more intense for low frequencies. Computing the derivative of a signal is equivalent to applying a high-pass filter that amplifies the signal proportional to its frequency. Hence, we decided to include a high-pass filter on the definition of the mask. The filter employed is shown in Figure 3.4. It applies a gain of 0.2 for frequencies inferior to 20 BPM and a gain of 1 for frequencies superior to 150 BPM. Between 20 and 150 BPM it applies a gain that varies linearly with frequency. This filter has a behavior similar to a derivative filter, but we restricted its actuation between 20 and 150 BPM to avoid over-attenuation of low frequencies, that could destroy useful information, and over-amplification of high frequencies, that could boost high frequency noise.

For simplification, we refer to this filter as derivative filter ($F_d[v]$). It is applied in cooperation

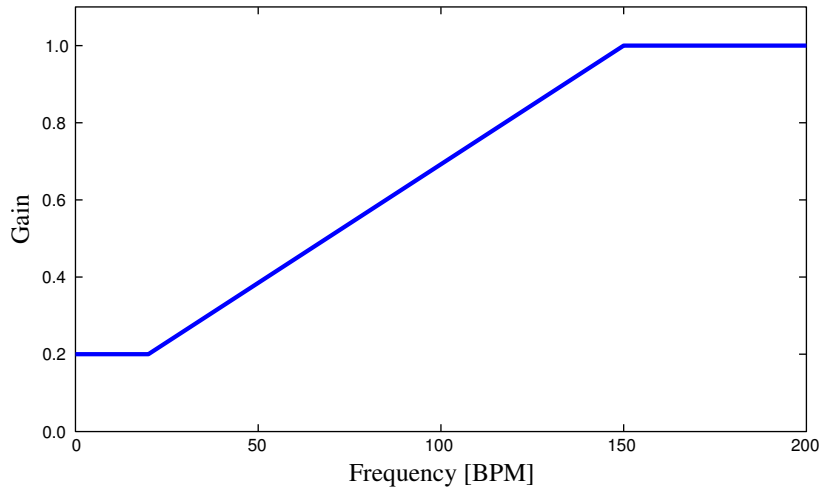


Figure 3.4: Derivative filter.

with the adaptive filter to further improve the mask, that is now expressed as:

$$M[j, v] = F_d[v] (T_2[v] * Z[j - 1, v]) (T_2[v] * T_2[v] * Z[j - 2, v]). \quad (3.5)$$

Figure 3.5 shows the application of the mask on a real signal. Note that the HR peak becomes more prominent on the filtered signal. This can be explained by two effects: (i) the peaks corresponding to the HR usually have a high amplitude and their amplitude and frequency changes slowly with time; (ii) those peaks corresponding to noise usually have a lower amplitude and their amplitude and frequency change faster than the HR with time. Hence, for regions where the signal is composed by noise the mask will apply an attenuation.

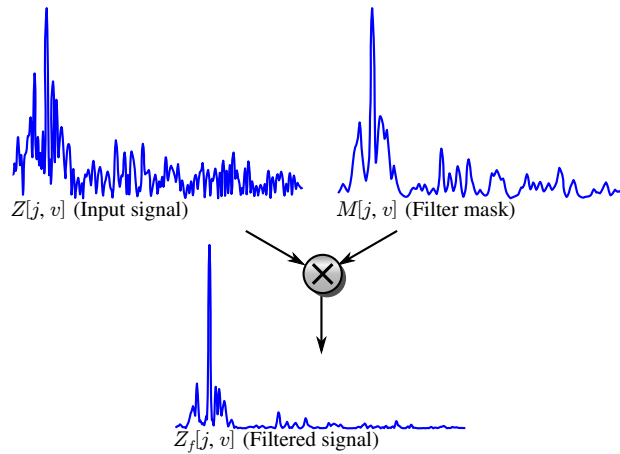


Figure 3.5: Adaptive filtering. The signal is multiplied by a mask in an attempt to attenuate noise and make the HR peak more prominent

If the assumptions made for the noise are not valid, that is, its frequency and amplitude do not vary with time, the amplitude of the mask on the corresponding frequency will be proportional the noise peak amplitude. Thus, the multiplication by the mask will have an effect similar to square the input signal at that frequency. As squaring a signal do not alter the ordering of the

peaks amplitude, the signal filtered by means of multiplication by this mask will have no negative effect, in terms of estimation, for noise components that vary slowly with time and will attenuate the noise components that variate faster.

Finally, from the filtered signal, the position of the peak of higher amplitude is found and its frequency is the j -th estimated HR frequency $p[j]$ for the subject. If the absolute difference between $p[j]$ and $p[j - 1]$ is higher than 12 BPM, the algorithm searches between the four highest peaks for the one that is the nearest to $p[j - 1]$, with an absolute difference inferior to 12 BPM. If none of them meet this constraint, the frequency of the peak of higher amplitude is used as estimation.

3.2 PRELIMINARY RESULTS

Figure 3.6 presents the HR estimated with the proposed method in comparison with that of Poh [14] for three different videos. For the ground truth we used an finger oximeter to monitor the subject HR at the same time that the video was captured.

In the first two videos, I and II, the subject were asked to remain still, moving as little as possible. In video III, the subject was talking and moving freely while the video was captured. It can be seen that for the video I the performance of our method and that of Poh were very similar. Indeed, they outcome the same values and their curves are overlapped. For video II, between 20 and 33s, the algorithm of Poh diverged from the HR captured from the oximeter while our approach remained very similar to it. We attribute this gain in performance due the adaptive filtering employed that improves the SNR (see Figure 3.7).

In the third video, on the other hand, both algorithms diverged from the values measured from the oximeter probably due to artifacts introduced from movement.

The effect of the adaptive filtering can be observed in Figure 3.7, that shows the signal used to feed the peak detector in both algorithms, $Z[j, v]$ for Poh and $Z_f[j, v]$ for our method. The filter was able to improve the SNR on all instants, except for the last video, as it can be seen. For video I, the SNR was already high before the filtering, therefore both algorithms were able to correctly estimate the HR. The same is valid for video II, as it can be seen in (b), except for the interval between 20 to 33s, where Poh algorithm failed due to noise. As the adaptive filtering removed most of this noise, the performance of our algorithm was little influenced by it.

However, video III shows that both algorithms are not robust to movements of higher amplitude and further improving is necessary concerning this case.

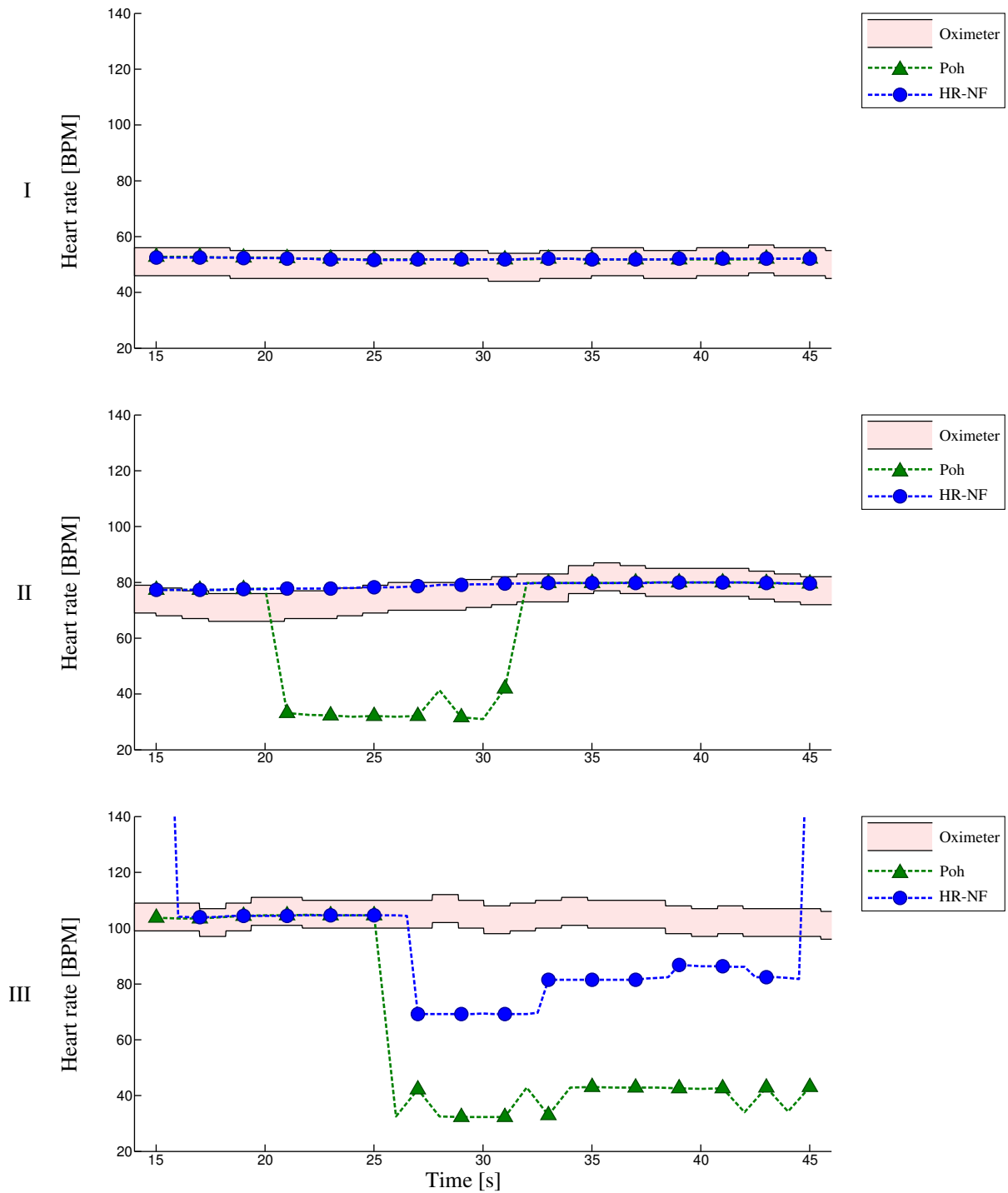


Figure 3.6: Performance of the HR detection for three different videos. In I and II we asked the subject to stay as still as possible and in III the subject is freely talking and moving as the video is captured.

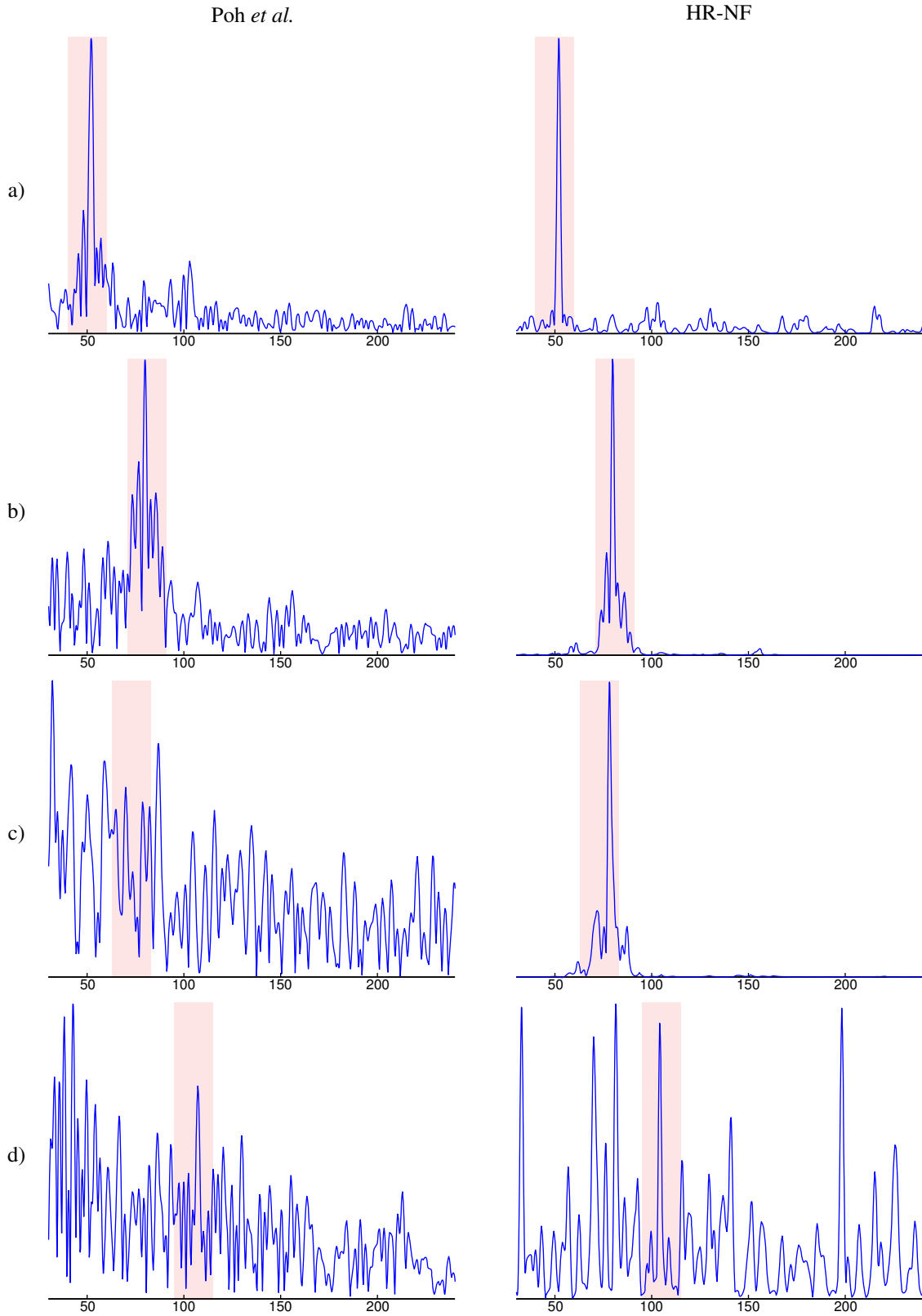


Figure 3.7: Signal fed to the peak detector as processed by the algorithm of *Poh et al.* and the proposed method. In the abscissa axis is represented the frequency, in BPM, and on the ordinate axis the magnitude. These signals were extracted from the same videos used in Figure 3.6: a) I at 22s; b) II at 38s; c) II at 25s and d) III at 34s. The colored rectangle correspond to the oximeter reading at that time.

4 HEART RATE ESTIMATION THROUGH MICRO-REGION TRACKING

The method presented in this chapter, denominated “video heart rate estimation through micro-region tracking” (HR-MRT), emerged as a modification of the HR-NF method in an attempt to produce an algorithm more robust to movements, avoiding the problem found in Figure 3.6-III.

Movement affects the ROI. The face deforms while someone speaks and the movement of the lips, eyelid, eyebrows, etc, are all captured by the values of x_r , x_g and x_b (the average value for red, green and blue inside the ROI) and will introduce artifacts that makes the task of HR estimation more difficult.

If these artifacts are of small amplitude or if their frequency are different than the range of frequencies where we search for the HR, they will have little influence on the estimation. That is the case for videos I and II in Figure 3.6, where most of the artifacts are due to blinking and other small movements, but not for video III, where the movements have higher amplitudes.

Therefore, we modify the ROI to make it more robust to movements. Figure 4.1 is a schematic of how the ROI is defined in this method. We divide the video in blocks with the same duration, segment the first frame of each block in micro-regions and extract some features from them to track. The tracking of these features enable us to estimate the parameters of an affine transformation that describe how the micro-regions move with time. We then extract the red, green and blue traces for each micro-region, normalize and mixture them with the same procedure employed for HR-NF. From the DFT for each micro-region, we feed a clustering algorithm that will define the ROI as a combination of micro-regions where the PPG signal is visible. The output of the clustering algorithm are the red, green and blue traces. The subsequent steps to estimate the HR were omitted, since they are the same employed for HR-NF, in Figure 3.1.

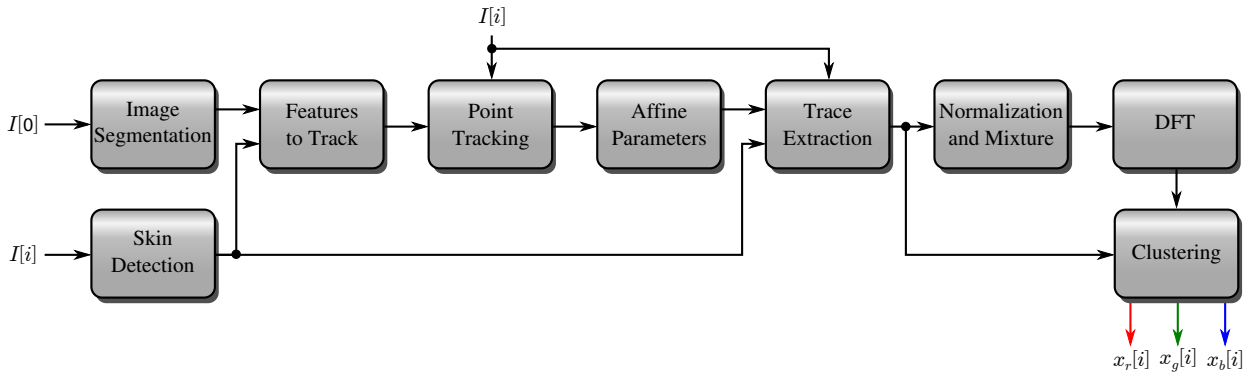


Figure 4.1: Block diagram for the ROI definition in HR-MRT.

Section 4.1 describes the skin detection algorithm employed to ignore non-skin pixels from the micro-regions. Section 4.2 shows how we divide the image in micro-regions to take movements

into account and Section 4.3 how we use tracking to compensate for movements. Section 4.4 explain the clustering algorithm used to define the ROI.

4.1 SKIN DETECTION

As the PPG signal is only present in pixels that correspond to the skin of the person being monitored, an algorithm capable of correctly discriminate skin pixels is a good tool to eliminate those pixels that could add noise to the red, green and blue traces computed on the ROI. Hence, we decided to employ a skin detector to ignore non-skin pixels, following three strategies.

In the first strategy, we created a look-up table, using the histogram based approach described in Section 2.2.3. We used the database provided by Jones and Rehg [45]¹, which is composed of images obtained from the internet, for which we know the ground truth. The database is comprised of 3789 images containing skin from people in different backgrounds and 6187 images where no skin is present. The color space employed is the YCbCr, as it is the preferred among researchers in the area. We do not drop the luminance component and the histograms have 256^3 bins. We used a cubic filter of size 3^3 to smoth the histogram, before creating the look-up table . As it can be seen in Figure 4.2-(b), this algorithm results in a great quantity of false positives because, in real life images, the clusters for skin and non-skin pixels overlap. Therefore, in the second strategy we combine this skin detector with the Viola-Jones face detector to eliminate those pixels that do not correspond to the face. Only pixels in the rectangle with 100% of the height and 80% of the width of the region found by the face detector are kept.

In the third strategy, we manually selected the skin pixels to observe the performance of our algorithm when we eliminate its dependence on the performance of the skin detector. The manual discrimination is made only on the first frame of a video block.

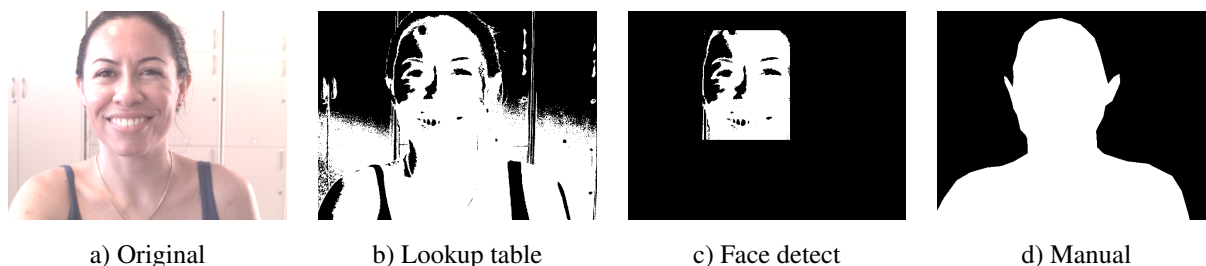


Figure 4.2: The image in a) was provided to the skin detection algorithm using b) histogram based approach employing a lookup table, c) the histogram approach combined with Viola-Jones face detector and d) manual detection.

¹<https://drive.google.com/folderview?id=0Bz-X0E2bqx9YcW9GaEM0OS0xb28&usp=sharing>

4.2 MICRO-REGIONS

The micro-regions are segments with uniform color, found by a segmentation algorithm in the first frame of each block of the video. We used the watershed method, an idea introduced by Beucher and Lantuéjoul [67] based in a topological representation of an image, to segment the frame. This method was chosen because of its simplicity and also because it divides the image in several segments and we can indirectly control the size and number of segments.

Prior to the segmentation with watershed, we usually compute the gradient modulus of the image. The gradient is more intense for regions of transition, like the edges between objects and less intense elsewhere. It is a good representation of the image borders.

The regions of high gradient can be seen as mountains separating valleys (low gradient). Referring to Figure 4.3, if we let a drop of water fall in a given position, it will flow down the mountains formed by the gradient until it arise to a local minimum. As more water drops arrive they start to form a puddle. We can intuitively interpret the watershed segmentation by means of these drops. Those pixels that flow to the same puddle belong to the same region (they are in the same catchment basin of that minimum) and the gradient is used as a base to separate one region from another.

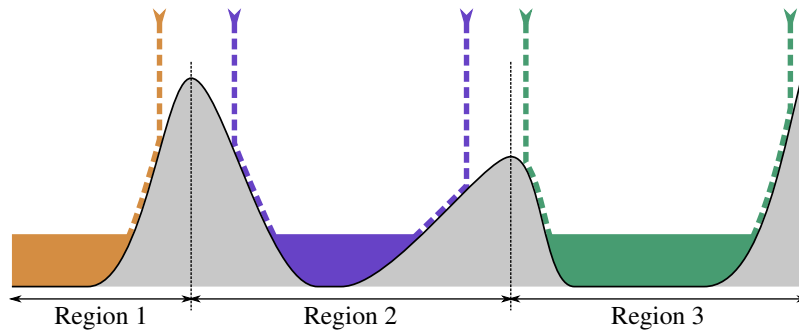


Figure 4.3: Two dimensional representation of the watershed segmentation method. A region correspond to all positions for which a drop of water, flowing down the hills, would fall in the same puddle.

The algorithm for the watershed segmentation can be given in three steps:

1. Create a group containing all non-labeled pixels of the image. Initially, none of them is labeled, so this group will contain the entire image;
2. From the group of non-labeled pixels, extract those of minimal altitude and attribute to them the label of an adjacent labeled pixel. If there is no adjacent labeled pixel, attribute a new label to it.
3. Repeat step 2 until there is no more non-labeled pixels left.

One disadvantage of the watershed method is its tendency to over-segmentation. Thus, it is a standard practice to smooth the image before applying it, eliminating or attenuating gradients

of low amplitude in textured portions of the image. In our work, we employed a bilateral filter, a non-linear edge preserving low-pass filter that was first proposed by Tomasi and Manduch [68]. This filter smooths pixels by averaging them with neighboring pixels, similar to isotropic filters, but adds a penalty based on their color difference. Let $I(x, y)$ be the original image that we want to filter, then the filtered image is given by

$$I_{BF}(x, y) = \sum_{(i,j) \in \Theta} I(x+i, y+j) D_{ij} P\{\Delta I_{ij}(x, y)\}, \quad (4.1)$$

$$\Delta I_{ij}(x, y) = I(x+i, y+j) - I(x, y),$$

where Θ defines the neighborhood. D_{ij} is a distance penalty function that reduces the weight of a pixel based on how far it is from the pixel being filtered and $P\{\Delta I_{ij}(x, y)\}$ introduces a penalty based on how different the pixels are. We suppose that pixels that are in different sides of the border will have a significant difference in color and those in the same side will be rather similar. Therefore, the second penalty tries to adjust the filtering in order to take the borders into account. These functions are subject to the constraint

$$\sum_{(i,j) \in \Theta} D_{ij} P\{\Delta I_{ij}(x, y)\} = 1. \quad (4.2)$$

In our work, we used a Gaussian function to express these penalties in a rectangular neighborhood. The filtering is then expressed as

$$I_{BF}(x, y) = \left(\sum_{i=-N}^N \sum_{j=-N}^N I(x+i, y+j) W_{ij}(x, y) \right) / \left(\sum_{i=-N}^N \sum_{j=-N}^N W_{ij}(x, y) \right), \quad (4.3)$$

where, in this case, N defines the size of the neighborhood and

$$W_{ij}(x, y) = \exp\left(-\frac{i^2 + j^2}{2\sigma_1^2}\right) \exp\left(-\frac{|\Delta I_{ij}(x, y)|^2}{2\sigma_2^2}\right), \quad (4.4)$$

is the weight attribute to pixel at position $(x+i, y+j)$ to compose the filtered pixel at (x, y) .

Equation (4.4) has two components. The first one is the distance penalty. The second term introduces the non-linearity and anisotropy on the filtering trying to take in account the borders.

Figure 4.4 shows the application of this filter in the image of Lena with noise added. The filter is capable of reducing the noise without corrupting the borders. However, one disadvantage of this algorithm is its susceptibility to create a carton-like effect on the filtered image.

From the filtered image we compute its luminance, $Y(x, y)$, the same way as the Y component of the YCbCr color space (see Section 2.2.1.3). The gradient is given by the horizontal and vertical



Figure 4.4: Image smoothing with bilateral filtering: a) is the original image with a dimension of 512x512 pixels and in b) a Gaussian noise was added to the image in a). c) shows the result of the filtering with $\sigma_1 = 3$, $\sigma_2 = 0.15$ and $N = 8$. Figure d) is a section of Lena's arm showing in close up the before and after filtering.

derivative of the luminance. To compute the horizontal and vertical derivative, we use kernels

$$K_h = \begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix} \quad \text{and} \quad K_v = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ +1 & +2 & +1 \end{bmatrix}, \quad (4.5)$$

respectively, and the squared gradient is given by

$$G(x, y)^2 = (Y(x, y) * K_h)^2 + (Y(x, y) * K_v)^2. \quad (4.6)$$

We further smooth $G(x, y)^2$ with a rectangular kernel to attenuate the gradient on textured regions and avoid over segmentation and we apply the watershed method to segment the image (see Figure 4.5).

Finally, after segmentation, we ignore those micro-regions that contain less than 80% of skin pixels.

Figure 4.6 presents a summary of the image segmentation in micro-regions using the first frame of the video block. The image is first low pass filtered with the bilateral filter, than we compute its gradient. The gradient is further smoothed and we apply the watershed algorithm that segments the image in small regions respecting the edges.

4.3 MICRO-REGION TRACKING

Each micro-region is delimited by the pixels in its border. Hence, if we can track how this pixels position evolved with time we can describe how the micro-region moved and deformed in the video sequence. In this fashion, we would be correcting the artifacts introduced by movement in the video. Although, tracking the pixels on the borders of the micro-regions is not a easy task because some of them may be in low textured regions where tracking algorithms fail to precisely



Figure 4.5: Segmentation of the image in a) using watershed with different kernel sizes to smooth the gradient. The bilateral filtering was performed using $\sigma_1 = 3$, $\sigma_2 = 0.06$ and $N = 5$.

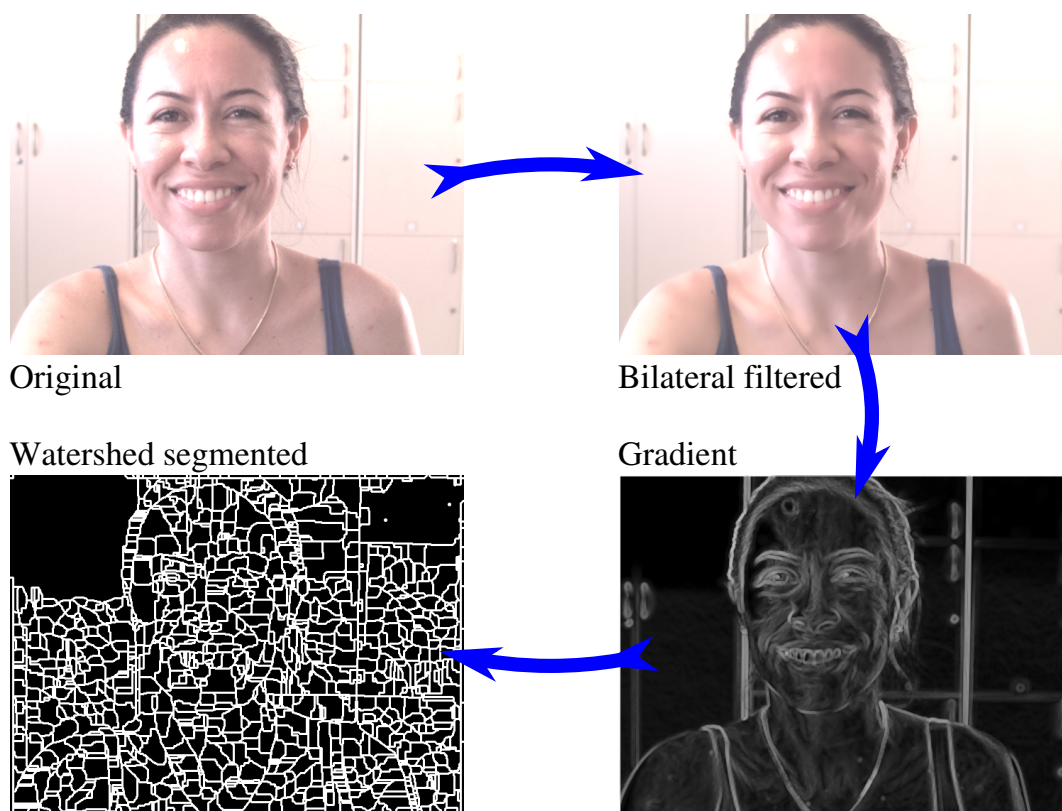


Figure 4.6: Image segmentation in micro-regions

estimate the pixel disparity and the amount of pixels to track would slow down the algorithm.

Therefore, for each segmented micro-region, we select a set of easy to track points that are inside or on the border of the micro-region (up to 12 points for each micro-region). We assume that the pixels inside the region are subject to the same movement as the border pixels. Thus, we use these points to estimate how the micro-region moved and deformed with time. The set of points is selected using the algorithm of Shi and Tomasi [69], as implemented in OpenCV 2.4. This algorithm finds the most prominent corners within the micro-region. This set of points are then tracked with the Lucas–Kanade algorithm as implemented by Yves [70].

The optical flow computed with Lucas–Kanade method may contain some noise. This noise should not affect short sequences of video, but for longer sequences it adds a drift to the motion estimation, which can lead to incorrect values for the optical flow, which can spoil the micro-region tracking.

To overcome this issue, one can filter the data to attenuate this noise in order to reduce its effect for longer sequences. But before that, we need to define a model to the data to decide how to filter it. In this work, we suppose that, for a short interval, the movement of the objects in the scene can be described by a polynomial function of a given order. Let $\mathbf{p}(t) = [x_t(t), y_t(t)]^T$ (\bullet^T means transpose) be the position of a pixel that we are tracking in time t , then we can express $\mathbf{p}(t)$ as

$$\mathbf{p}(t) = \sum_{n=0}^{N_{poly}} \mathbf{r}_n \frac{t^n}{n!}, \quad (4.7)$$

where $\mathbf{r}_n = [x_n, y_n]^T$ are constants and N_{poly} is the polynomial order. The higher the polynomial order, the better it can represent complex movements, but it becomes more sensible to noise. A good trade-off between data representation and noise cancellation is found for $N_{poly} = 3$. Thus, for simplicity of reading, we will consider only the case of third order polynomials models in this chapter. Results for other orders can be found following the same steps. A demonstration of this statement will be presented in the next chapter. Therefore we have

$$\mathbf{p}(t) = \begin{bmatrix} x_t(t) \\ y_t(t) \end{bmatrix} = \begin{bmatrix} x_0^n \\ y_0^n \end{bmatrix} + \begin{bmatrix} v_x^n \\ v_y^n \end{bmatrix} t + \begin{bmatrix} a_x \\ a_y \end{bmatrix} \frac{t^2}{2} + \begin{bmatrix} u_x \\ u_y \end{bmatrix} \frac{t^3}{6}, \quad (4.8)$$

or, in a more compact way,

$$\mathbf{p}(t) = \mathbf{p}_0 + \mathbf{v}t + \mathbf{a}\frac{t^2}{2} + \mathbf{u}\frac{t^3}{6}, \quad (4.9)$$

where \mathbf{p}_0 , \mathbf{v} , \mathbf{a} and \mathbf{u} are constants vectors that define the initial position, velocity, acceleration and jerk.

Let us now assume that we have two reference images, R_1 and R_0 , that are consecutive in the video sequence (R_1 being first), for which we know the position of the point we are tracking, \mathbf{p}^1 and \mathbf{p}^0 , respectively. We also have six other images in the sequence, I_1, I_2, I_3, I_4, I_5 and I_6 ,

for which we don't know the tracking point position yet. Their position will be estimated using Lucas–Kanade method.

The sequence of images, $[R_1, R_0, I_1, I_2, I_3, I_4, I_5, I_6]$, form a sub-sequence of the video, meaning that they are all consecutive. We suppose that Equation (4.9) is a good model to represent the movement of the tracking pixel in this sub-sequence between frame R_0 and I_6 .

Using the reference frames, we estimate the optical flow as schematized in Figure 4.7. We denote as $\mathbf{p}_n^0 = [x_n^0, y_n^0]^T$ the position estimated for frame I_n using information coming from reference R_0 and as \mathbf{p}_n^1 when coming from reference R_1 . From this estimated values, we find the constants in Equation (4.9) that minimizes the quadratic error. We assume that R_1 and R_0 correspond to $t = -1$ and 0, respectively, and I_n correspond to $t = n$.

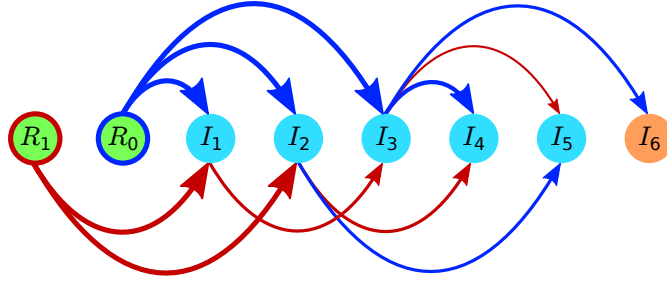


Figure 4.7: Optical flow estimation scheme. The arrows indicate from which to which frame the optical flow is estimated using the Lucas–Kanade method. Blue arrows represent that the flow was estimated using information coming from reference R_0 and red arrows from reference R_1 .

The constants can be found by linear algebra, solving the equation

$$\begin{bmatrix} 1 & 0 & 0^2/2 & 0^3/6 \\ 1 & 1 & 1^2/2 & 1^3/6 \\ 1 & 1 & 1^2/2 & 1^3/6 \\ 1 & 2 & 2^2/2 & 2^3/6 \\ 1 & 2 & 2^2/2 & 2^3/6 \\ 1 & 3 & 3^2/2 & 3^3/6 \\ 1 & 3 & 3^2/2 & 3^3/6 \\ 1 & 4 & 4^2/2 & 4^3/6 \\ 1 & 4 & 4^2/2 & 4^3/6 \\ 1 & 5 & 5^2/2 & 5^3/6 \\ 1 & 5 & 5^2/2 & 5^3/6 \\ 1 & 6 & 6^2/2 & 6^3/6 \end{bmatrix} \begin{bmatrix} x_0 & y_0 \\ v_x & v_y \\ a_x & a_y \\ u_x & u_y \end{bmatrix} = \begin{bmatrix} x^0 & y^0 \\ x_1^0 & y_1^0 \\ x_1^1 & y_1^1 \\ x_2^0 & y_2^0 \\ x_2^1 & y_2^1 \\ x_3^0 & y_3^0 \\ x_3^1 & y_3^1 \\ x_4^0 & y_4^0 \\ x_4^1 & y_4^1 \\ x_5^0 & y_5^0 \\ x_5^1 & y_5^1 \\ x_6^0 & y_6^0 \end{bmatrix}, \quad (4.10)$$

which can be more compactly written as

$$TA = P, \quad (4.11)$$

where T is the time matrix, A the parameters matrix and P the position matrix.

Equation (4.11) falls in the category of a very well known mathematical problem: that of linear least squares fitting for an overdetermined system of linear equations [71]. These problems are convex and have a closed-form solution that is unique and given by

$$A = (T^T T)^{-1} T^T P. \quad (4.12)$$

Now that we know the parameters, we can recalculate the position of the tracking pixel in frame I_n , $1 \leq n < 6$, resulting in

$$T'(T^T T)^{-1} T^T P = FP, \quad (4.13)$$

$$T' = \begin{bmatrix} 1 & 1 & 1^2/2 & 1^3/6 \\ 1 & 2 & 2^2/2 & 2^3/6 \\ 1 & 3 & 3^2/2 & 3^3/6 \\ 1 & 4 & 4^2/2 & 4^3/6 \\ 1 & 5 & 5^2/2 & 5^3/6 \end{bmatrix}. \quad (4.14)$$

$F = T'(T^T T)^{-1} T^T$ is therefore the filtering matrix that will eliminate from the set of points P the movements that do not obey the model given by Equation (4.9), represented in Figure 4.8, from where we can see that it corresponds to a low-pass filter, as expected. A higher (or lower) order model could also be employed, but as the order increases it starts to become unable to attenuate noise, as noise components can be represented by the high order terms of the model.

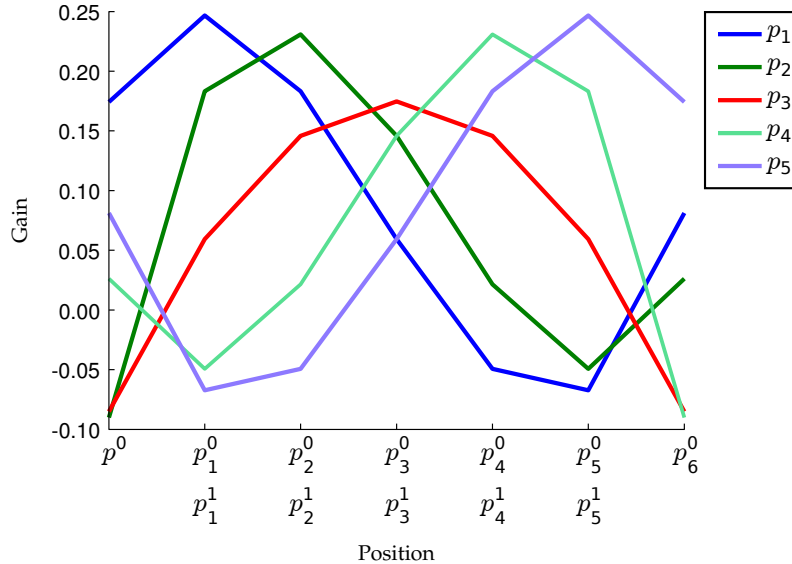


Figure 4.8: Feature tracking filtering. Shows the value of the filtering matrix F that multiply each component of the position matrix P .

Frame I_6 was used only to calculate the filtering parameters, but the filtered pixel position is not computed for this frame. We advance in the video sequence. I_4 and I_5 become now the

reference and we use the same approach to estimate the optical flow for $I_6, I_7, I_8, I_9, I_{10}$ and I_{11} . We keep advancing until the entire video sequence is covered.

Finally, after tracking how the set of points for a given micro-region evolved with time using the above algorithm, we can determine how the micro-region moved. However, to this point, we only know how a set of points inside the region evolved with time. To use the information of the set of tracking points we suppose that the k -th points in the border of the micro-region, given by $\mathbf{b}^k(t) = [x_b^k(t), y_b^k(t)]^T$ in time t , undergo an affine transformation of the form

$$\begin{bmatrix} x_b^k(t) \\ y_b^k(t) \end{bmatrix} = R(t) \begin{bmatrix} x_b^k(0) \\ y_b^k(0) \end{bmatrix} + \begin{bmatrix} d_x(t) \\ d_y(t) \end{bmatrix}, \quad (4.15)$$

where R is a 2×2 affine matrix transformation and d_x and d_y are the translations. We don't know these parameters, but we can estimate them using the set of tracking points with the assumption that the tracking points undergo the same transformation as the border points. Let $[x_t^n(t), y_t^n(t)]^T$ be the n -th tracking points in time t for the micro-region we are processing; then

$$\begin{bmatrix} x_t^n(t) \\ y_t^n(t) \end{bmatrix} = R(t) \begin{bmatrix} x_t^n(0) \\ y_t^n(0) \end{bmatrix} + \begin{bmatrix} d_x(t) \\ d_y(t) \end{bmatrix}. \quad (4.16)$$

We can eliminate the translation in Equation (4.16) by subtracting each side of the equation by the mean value of the tracking point position in time, $[\bar{x}_t(t), \bar{y}_t(t)]^T$,

$$\begin{bmatrix} x_t^n(t) \\ y_t^n(t) \end{bmatrix} - \begin{bmatrix} \bar{x}_t(t) \\ \bar{y}_t(t) \end{bmatrix} = R(t) \left(\begin{bmatrix} x_t^n(0) \\ y_t^n(0) \end{bmatrix} - \begin{bmatrix} \bar{x}_t(0) \\ \bar{y}_t(0) \end{bmatrix} \right) \quad (4.17)$$

Equation (4.17) can be divided in 5 simple transformations. Ignoring translation, for simplicity, we have:

$$\begin{aligned}
\text{Scaling} \quad & \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \alpha(t) \begin{bmatrix} x(0) \\ y(0) \end{bmatrix} \\
\text{Rotation} \quad & \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} \cos(\theta(t)) & -\sin(\theta(t)) \\ \sin(\theta(t)) & \cos(\theta(t)) \end{bmatrix} \begin{bmatrix} x(0) \\ y(0) \end{bmatrix} \\
\text{Vertical Shearing} \quad & \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \alpha_v(t) & 1 \end{bmatrix} \begin{bmatrix} x(0) \\ y(0) \end{bmatrix} \\
\text{Horizontal Shearing} \quad & \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} 1 & \alpha_h(t) \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x(0) \\ y(0) \end{bmatrix} \\
\text{Mirroring} \quad & \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} (-1)^{m(t)} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x(0) \\ y(0) \end{bmatrix}
\end{aligned}$$

α is the scaling factor (a positive value), θ the anticlockwise angle of rotation, α_h and α_v are the horizontal and vertical shearing and $m = \{0, 1\}$ is a integer that indicate if the points suffered horizontal mirroring or not.

If we restrict the affine transformation to only rigid transformation, ignoring shearing and mirroring, matrix R can be expressed as

$$R = \alpha \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} = \begin{bmatrix} a_{11} & -a_{21} \\ a_{21} & a_{11} \end{bmatrix}. \quad (4.18)$$

On the other hand, allowing shearing and mirroring,

$$R = \alpha \begin{bmatrix} (-1)^m & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} 1 & \alpha_h \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \alpha_v & 1 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}. \quad (4.19)$$

These two distinct transformation are referred, in this work, as rigid affine transform and full affine transform, respectively.

To solve for R , we use the procedure described by Lawson and Hanson [71]. We find the parameters that minimize the squared error. The solution is given in terms of sums

$$\begin{aligned}
C_{xx}(t) &= \sum_n [x_t^n(0) - \bar{x}_t(0)] \cdot [x_t^n(t) - \bar{x}_t(t)] \\
C_{xy}(t) &= \sum_n [x_t^n(0) - \bar{x}_t(0)] \cdot [y_t^n(t) - \bar{y}_t(t)] \\
C_{yx}(t) &= \sum_n [y_t^n(0) - \bar{y}_t(0)] \cdot [x_t^n(t) - \bar{x}_t(t)] \\
C_{yy}(t) &= \sum_n [y_t^n(0) - \bar{y}_t(0)] \cdot [y_t^n(t) - \bar{y}_t(t)] \\
D_{xx} &= \sum_n [x_t^n(0) - \bar{x}_t(0)]^2 \\
D_{xy} &= \sum_n [x_t^n(0) - \bar{x}_t(0)] \cdot [y_t^n(0) - \bar{y}_t(0)] \\
D_{yx} &= \sum_n [y_t^n(0) - \bar{y}_t(0)] \cdot [x_t^n(0) - \bar{x}_t(0)] = D_{xy} \\
D_{yy} &= \sum_n [y_t^n(0) - \bar{y}_t(0)]^2
\end{aligned}$$

Then, for Equation (4.18) we have

$$R = \frac{1}{D_{xx} + D_{yy}} \left(\begin{bmatrix} C_{xx} & -C_{xy} \\ C_{xy} & C_{xx} \end{bmatrix} + \begin{bmatrix} C_{yy} & C_{yx} \\ -C_{yx} & C_{yy} \end{bmatrix} \right) \quad (4.20)$$

and for Equation (4.19) we have

$$R = \begin{bmatrix} C_{xx} & C_{yx} \\ C_{xy} & C_{yy} \end{bmatrix} \begin{bmatrix} D_{xx} & D_{yx} \\ D_{xy} & D_{yy} \end{bmatrix}^{-1}. \quad (4.21)$$

We ignored the time dependency for simplicity of writing, but the reader must note that matrix R varies with time.

The translation can be found from Equation (4.17) as

$$\begin{bmatrix} d_x(t) \\ d_y(t) \end{bmatrix} = \begin{bmatrix} \bar{x}_t(t) \\ \bar{y}_t(t) \end{bmatrix} - R(t) \begin{bmatrix} \bar{x}_t(0) \\ \bar{y}_t(0) \end{bmatrix}. \quad (4.22)$$

Applying this affine transformation to the border of the micro-region, we can determine how it evolved with time from their initial position and where the micro-region is in time t . It can be shown² that if the micro-region has a total area of A_0 in $t = 0$ and undergo an affine transformation, its area in time t is given by $A_0 \cdot \det(R(t))$. We use this function as a metric to ignore those

²The demonstration was omitted here. The reader can refer to <http://www.mathopenref.com/coordpolygonarea.html> and <https://www.math.wisc.edu/~robbin/461dir/coordinateGeometry.pdf> for a rigorous mathematical demonstration

micro-regions for which its area became bigger than twice or smaller than half its initial size at some point as it can indicate that misleading occurred during point tracking or affine transformation estimation.

Finally, we computed the mean value of red, green and blue for each micro-region over time. We ignore those pixels that were not detected as skin-pixel from the traces computation. The traces are then normalized, following the same procedure as for the HR-NF, and we compute its DFT.

4.4 CLUSTERING

The traces extracted for each micro-region may contain different levels of noise. For example, when the skin detector does not provide a good segmentation, some micro-regions may be composed mainly by background pixels, where the PPG signal is not present. For the micro-regions composed mainly by true skin pixels, the energy of the PPG signal and the noise may vary due to the characteristics of the skin on the corresponding region, such as beard, vasculature, makeup, motion noise.

We want to average the traces found for the micro-regions in a single trace to be employed for HR estimation. To obtain a good SNR, we try to ignore those micro-regions where the PPG signal is not visible or where the noise energy obscures it. The clustering algorithm is used to resolve which micro-regions to use in order to maximize the SNR. It groups in a cluster those micro-regions that present similar traces, by comparing their Fourier transform. Based on the assumption that most micro-regions contain the PPG signal, we select the cluster of micro-regions with the highest number of elements and ignore the remaining clusters. In this fashion, the algorithm automatically selects the ROI.

Section 4.4.1 presents a new distance metric used to compute the similarity between the traces of micro-regions and Section 4.4.2 the clustering algorithm *per se*.

4.4.1 Distance Metric

For each micro-region, we mix their red, green and blue traces, using the weights given in Section 3.1, and we compute its DFT. Those micro-regions corresponding to the skin of the individual being monitored will have a Fourier transform composed of the PPG signal plus additive noise and those that do not correspond to the skin will be primarily composed by noise. Hence one can expect that the Fourier transform of some micro-regions will be similar between each other because they carry the same signal.

However, even though some micro-regions have similar Fourier transforms, their amplitude can vary significantly from one to another, as a result from shadows, changes in the characteristics of the skin, e.g. tone, texture, beard, blood circulation, etc. Therefore, traditional approaches to

compare functions, such as Euclidean distance, are not suitable in our case. Hence the necessity to create a new metric capable of comparing the similarity between Fourier transforms that is robust to the scale.

Let $F_1[v]$ and $F_2[v]$ be the magnitude of two discrete Fourier transforms that we want to compare, where v is the discrete frequency. They are said to be similar if they have peaks at the same frequencies, with amplitudes proportional to each other *i.e.*, exists a real constant $k > 0$ which yields

$$F_2[v] \approx k.F_1[v], \quad v_{min} \leq v \leq v_{max}, \quad (4.23)$$

where v_{min} and v_{max} define the range of frequencies where the two Fourier transforms are compared. This range is chosen as those frequencies where we expect the HR to be, from 30 to 240 BPM. Otherwise, when Equation (4.23) is not a good approximation, the functions are considered to be not similar.

One simple way to find the value of the constant k for equation (4.23) is to normalize $F_1[v]$ and $F_2[v]$ by their root mean squared amplitude (RMS amplitude),

$$A_1 = \sqrt{\frac{1}{N} \sum_i F_1[v]^2} \quad \text{and} \quad A_2 = \sqrt{\frac{1}{N} \sum_i F_2[v]^2}, \quad (4.24)$$

respectively, where N is the number of elements of F_1 and F_2 . Therefore, equation (4.23) becomes:

$$\frac{F_2[v]}{A_2} \approx \frac{F_1[v]}{A_1} \quad (4.25)$$

Nevertheless, the RMS amplitudes can be influenced by noise. So, for finer tuning, we can add a real constant $\alpha > 0$ onto equation (4.25), resulting in:

$$\frac{F_2[v]}{A_2} \approx \alpha \frac{F_1[v]}{A_1} \quad \Leftrightarrow \quad \alpha^{-1/2} \frac{F_2[v]}{A_2} \approx \alpha^{1/2} \frac{F_1[v]}{A_1}. \quad (4.26)$$

Thus, based in the Euclidean distance, we can define

$$D \{F_1, F_2\} = \sum_v \left(\alpha^{1/2} \frac{F_1[v]}{A_1} - \alpha^{-1/2} \frac{F_2[v]}{A_2} \right)^2 \quad (4.27)$$

as a distance between F_1 and F_2 which results in small values when equation (4.26) is a good approximation (for the cases where F_1 and F_2 are similar) and large values otherwise.

The regions of low amplitude in functions $F_1[v]$ and $F_2[v]$ suffer more from the noise influence than regions with high amplitude. So, we can add a weight w_v to equation (4.27) in order to give

more importance for the regions of high amplitude, resulting in

$$D \{F_1, F_2\} = \frac{\sum_v w_v \left(\alpha^{1/2} \frac{F_1[v]}{A_1} - \alpha^{-1/2} \frac{F_2[v]}{A_2} \right)^2}{\sum_v w_v} = \frac{1}{W} \left(\alpha \frac{V_{11}}{A_1^2} + \frac{1}{\alpha} \frac{V_{22}}{A_2^2} - 2 \frac{V_{12}}{A_1 A_2} \right), \quad (4.28)$$

$$V_{11} = \sum_v w_v F_1[v]^2, \quad V_{22} = \sum_v w_v F_2[v]^2, \quad V_{12} = \sum_v w_v F_1[v] F_2[v] \quad \text{and} \quad W = \sum_v w_v,$$

where w_v is defined as:

$$w_v = \max \left(\frac{F_1[v]}{A_1}, \frac{F_2[v]}{A_2} \right). \quad (4.29)$$

The value of $\alpha > 0$ is found as the one minimizing the distance metric $D \{F_1, F_2\}$ and can be calculated by deriving equation (4.28) as a function of α and making it equal to zero, leading to:

$$\alpha^2 = \frac{A_1^2 V_{22}}{A_2^2 V_{11}}. \quad (4.30)$$

Substituting the value of α given by equation (4.30) in equation (4.28) and simplifying, produces:

$$D \{F_1, F_2\} = 2 \frac{\sqrt{V_{11} V_{22}} - V_{12}}{A_1 A_2 W}. \quad (4.31)$$

A similarity metric yielding values between 0 and 1, where 1 means completely similar and 0 completely different, can be calculated using a Gaussian function as:

$$S \{F_1, F_2\} = e^{-\frac{D \{F_1, F_2\}}{2\sigma_s^2}}, \quad (4.32)$$

for a given σ_s .

4.4.2 Algorithm

Now that we have a distance metric capable of comparing the similarity between two Fourier transforms, we can use it as a base for a clustering algorithm. This algorithm should be capable of grouping together regions that have similar Fourier transforms, indicating that they carry the same signal.

The algorithm proposed in this work is a modification of the K-means algorithm [72]. The K-means algorithm is used to cluster a set of input vectors into k groups (or clusters). Each group is defined by its centroid and the vectors are assigned to the closest cluster. The k centroids are

initialized by a set of representative (or random selected) samples from the input set and each element of the input set is assigned to the closest centroid. After the elements are assigned, the centroids are updated. This process is repeated iteratively until the centroids converge (or when the differences between current centroids are below a threshold).

One disadvantage of the K-means for our purpose is that we need to know, prior to clustering, how many clusters we want to form. The optimal number of clusters is video-dependent because the noise from different micro-regions tend to differ from one site to another.

The proposed algorithm is represented in Figure 4.9 and presents two major modifications:

1. We employ Equation (4.31) as a distance metric;
2. The algorithm automatically determines the optimal number of clusters.

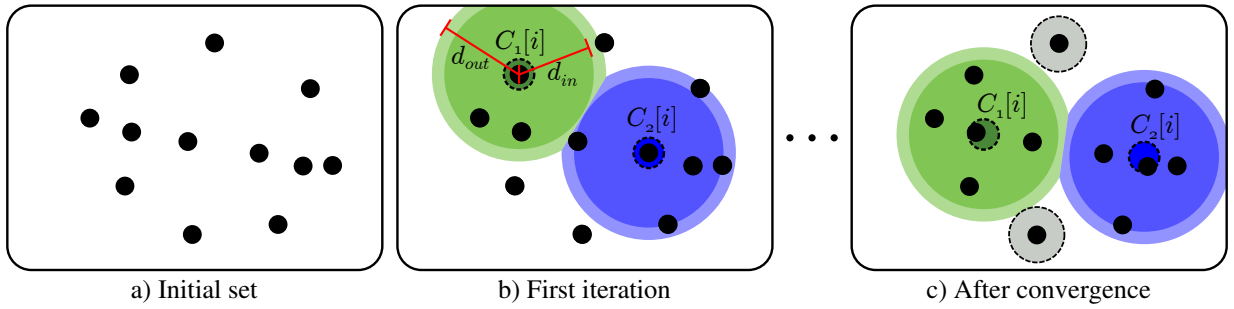


Figure 4.9: Clustering algorithm: a) Initial set of points that we want to cluster; b) Algorithm after first iteration. The points within one of the colored regions belong to the same cluster; c) Algorithm after convergence. The points inside the dashed regions belong to a special cluster where its elements are too different from others clusters.

The input of the algorithm is $\Upsilon = \{F_n[i]\}$, a set of vectors of the Fourier transform for each micro-region, in a total of N_R vectors. This set may also be seen as points in a nonlinear space, which have Equation (4.31) designed as a distance metric.

Differently from the K-means algorithm, the k -th cluster is entirely defined by the neighborhood around the vector $C_k[i]$, the cluster centroid, that is computed by the mean of the vectors inside the cluster after each iteration. That is, we do not need to know the centroid of other clusters to decide if a vector belongs or not to a given cluster.

In order for a vector to be included in the cluster, its distance to the cluster center must be less or equal to d_{in} . Once inside the cluster, the vector will only be excluded when, in the next update of the cluster center, its distance to the center is superior to d_{out} .

From the initial set Υ , 20% of vectors are randomly selected as cluster center. For all the remaining points, we calculate its distance to $C_k[i]$. If this distance is below or equal to d_{in} , the point is integrated to the cluster. Otherwise, if there are no cluster for which the distance is less or equal to d_{in} , the element remain unclassified.

After this step we calculate the mean value of all vectors within a cluster and this mean is used to update the cluster center. As it is possible that, in the random selection of elements as

cluster centers, we have separated elements that are in fact very similar, we compute the distance between all cluster centers to determine how similar they are. If this distance is less or equal to d_{out} , the clusters are aggregated together. After this step, for those clusters that have just one element, their vectors are aggregated together in a special cluster for those elements that are too different from others. Normally, this corresponds to elements composed primarily by noise.

It is possible that, after updating the cluster center, some elements within a cluster will actually no longer belong to the cluster. Therefore, we calculate the distance of all elements within a cluster to the cluster center. Those that have a distance higher than d_{out} are excluded and set as unclassified.

If we still have unclassified elements, we randomly select 20% of the unclassified elements to form new clusters and we repeat the previous steps till there are no more unclassified elements remaining. At this point, the algorithm is said to have achieved convergence.

For some cases it is possible that the algorithm never converges or converge just after a large number of iterations. Hence, we fixed a limit of 200 iterations and the algorithms stops, even without convergence, when it reaches this limit.

5 RESULTS

In this chapter we analysis the performance of our algorithms to real and synthetic data. Given the lack of video databases for the purpose of HR estimation, we created our own database to test the algorithms' performance.

Section 5.1 presents the database. The remaining sections evaluate the algorithms' performance.

5.1 DATABASE

The database is composed of video sequences with duration of 60 seconds. The videos are stored without compression, since the motion compensation employed in compression algorithms can eliminate the subtle variations of the skin tone that we employ for HR estimation.

We captured 2 different video sequences from 20 volunteers with an approximate duration of 65 seconds. The videos were cropped to an exact duration of 60 seconds. Among the volunteers, 15 were man and 5 were women, 4 used glasses and 6 had beards. Figure 5.1 sketch the average face of all volunteers as captured on the first frame of the videos.

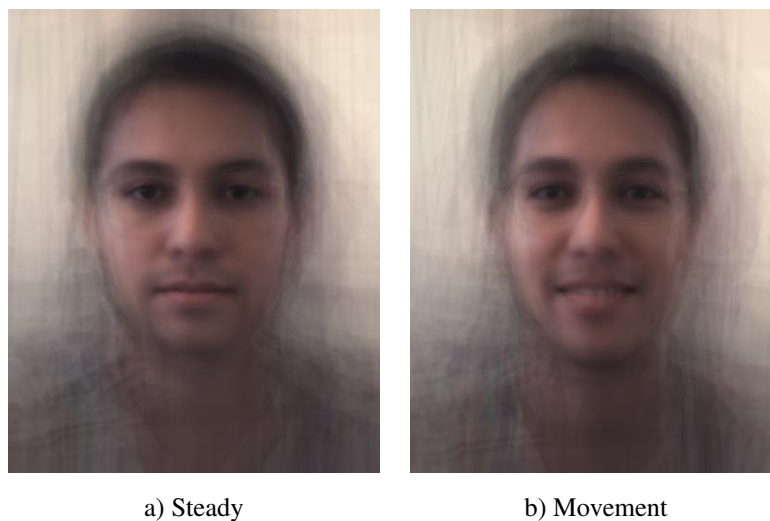


Figure 5.1: Average of volunteers face from the first frame of the videos. The images were aligned with respect to the eyes and mouth position.

The videos were captured using the camera Firefly MV FMVU-03MTC (Point Grey, Richmond, Canada), with a resolution of 0.3 megapixels (640×480 pixels) and a temporal sampling of 60 frames per second. Each frame is captured and stored without compression by a personal computer running an application that communicates with the camera through USB. The frames are stored as raw data.

During the acquisition process, when the application was not able to gather a frame from the camera buffer before the camera started to overwrite it, some frames were lost. As this error occurs with relative small frequency, our algorithms were able to estimate the HR even under such errors. The performance of our algorithms was evaluated with this database and the results are shown in subsequent sections.

The captured frames present a Bayer pattern [73] of type ‘RGGB’ with 8 bits per pixel. The frames are demosaiced in MATLAB[®] using the function `demosaic`¹.

The camera captured mainly the volunteer’s face under two conditions: for the first video, we asked the volunteer to stay as still as possible; for the second video, the volunteer was allowed to move freely and we interviewed them to encourage movements of the lips, face and hands. The audio was not recorded. The volunteers were not instructed on how to move on the second video in order to obtain a more natural reaction and avoid introducing a bias.

As the videos were captured, we monitored the volunteers HR using a commercial fingertip pulse oximeter that presents an accuracy of 2 BPM and a resolution of 1 BPM, as informed by the manufacturer. This measure is used as a reference (the ground-truth) to determine how well the algorithms perform. The oximeter data was filtered with a Gaussian filter with variance equal to 0, 7071 to eliminate its discontinuities due to the low precision, as shown in Figure 5.2.

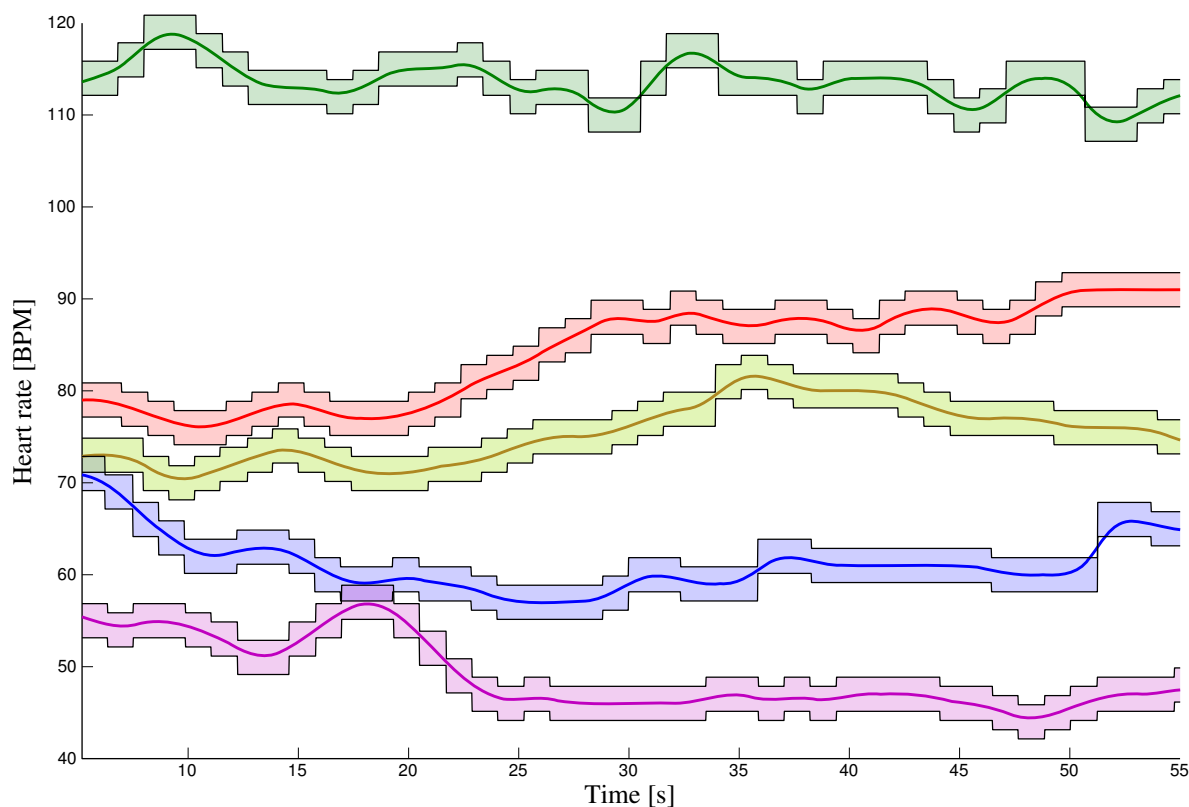


Figure 5.2: Oximeter data filtering. The colored regions correspond to the fingertip oximeter reading for 5 volunteers, within a range of ± 1.86 BPM. The lines show the filtered data for each volunteer.

¹<http://www.mathworks.com/help/images/ref/demosaic.html>

5.2 EVALUATION OF HR-FD

To evaluate the performance of the algorithms, we employ, in this work, bar plots similar to the one in Figure 5.3. These bar plots show the amount of time, as a percentage from the total time, that the algorithm stayed within a given absolute error range for all 20 volunteers. The error is considered to be the difference between the filtered oximeter reading and the estimated HR.

If the absolute error of the estimated HR, compared to the oximeter reading, is inferior or equal to 2 BPM, then the estimation is considered correct as it falls in the accuracy range of the oximeter. We also assume that errors inferior to 8 BPM are considered acceptable for HR estimations. On the other hand, if the error is superior to 11 BPM, we consider that the estimated HR is incorrect. Hence, the green regions correspond to correct estimations, red regions to incorrect estimations and yellow regions correspond to the transition between acceptable and incorrect.

In this section we evaluate the performance of HR-NF compared to the algorithm of Poh *et al* [14], referred as Poh algorithm, for simplicity. To perform the analysis we employ the steady videos and the videos with movement from our database, as well as synthetic data.

5.2.1 Steady Videos

Figure 5.3 shows the performance of the proposed algorithm compared to that of Poh. It can be noticed that the HR-NF performed better than Poh for the steady videos where Poh stayed 73.9% of the time with an error less or equal to 8 BPM, compared to 86.6% for HR-NF.

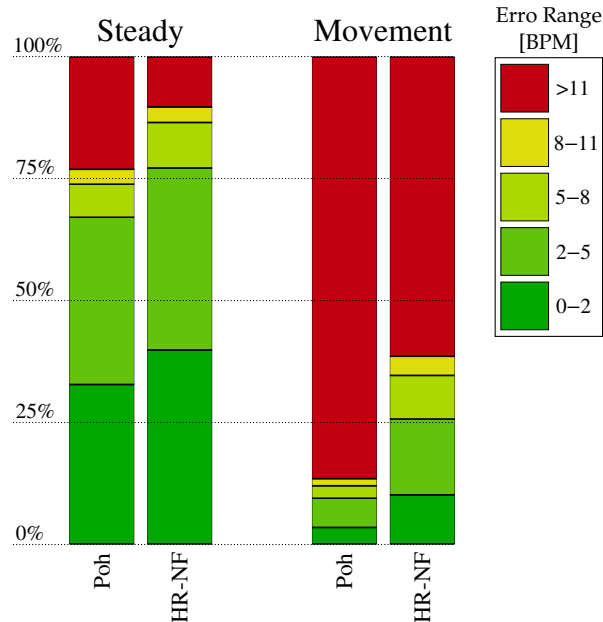


Figure 5.3: HR detection performance of Poh versus HR-NF.

This better performance is due to three factors: the use of the adaptive filter, the derivative filter and the fixed mixture of the traces to form the signal employed for HR estimation. As we apply

a fixed mixture for source separation, we are not susceptible to errors that could be introduced by the Independent Component Analysis algorithm that adaptively estimate the alpha values to apply to each color channel. Also, the adaptive and derivative filter are capable of increasing the SNR before heart rate estimation.

Table 5.1 portray the contribution of the adaptive and derivative filters, and the Blind Source Separation (BSS) strategy: fixed mixture of the red, green and blue channels (the strategy employed in HR-NF); adaptive estimation of the weights to apply to each color channel by means of ICA (the strategy employed by Poh).

Table 5.1: Evaluation of the contribution of the adaptive filter, derivative filter and BSS strategy on HR-NF for the steady videos.

Filter			Error range [BPM]				Rank	
Adaptive	Derivative	BSS	$ \epsilon \leq 2$	$ \epsilon \leq 5$	$ \epsilon \leq 8$	$ \epsilon > 11$		
No	No	Fixed	40.1%	79.4%	89.3%	07.4%	1 st	← Poh
		Adaptive	31.8%	65.0%	72.4%	25.7%	7 th	
	Yes	Fixed	38.1%	72.9%	81.1%	15.5%	5 th	
		Adaptive	28.6%	58.1%	64.4%	34.1%	8 th	
Yes	No	Fixed	40.8%	78.2%	87.3%	09.5%	2 nd	← HR-NF
		Adaptive	31.8%	66.1%	73.9%	22.8%	6 th	
	Yes	Fixed	39.8%	77.2%	86.6%	10.2%	3 rd	
		Adaptive	36.6%	73.0%	81.8%	14.0%	4 th	

The best performance was achieved when we used the fixed mixture for BSS and when no adaptive and derivative filters were employed. The HR-NF obtained the third place, but we can observe that its performance is very close to first and second place. The algorithm of Poh, on the other hand, obtained the seventh place.

It can be seen from the average contribution, depicted in Table 5.2, that the use of the adaptive filter and the fixed mixture improved the performance of the algorithm. The use of the filters did not provide the best performance for the steady videos since their function is to attenuate noise and the signal captured on this case comprises a good SNR as little motion artifact is present. However, they provided a performance close to the best, achieving the second and third position with only 2.0% and 2.7% less correct estimated HR than the first place, respectively.

Table 5.2: Average contributions of the adaptive filter, derivative filter and BSS strategy on HR-NF for steady videos. The average percentage of incorrect estimated HR ($|\epsilon| > 11$ BPM) and the average rank is shown in each cell.

Feature	Employed	
	Yes	No
Adaptive Filter	14.1% (3.75 th)	20.7% (5.25 th)
Derivative Filter	18.5% (5.00 th)	16.3% (4.00 th)
Fixed Mixture BSS	10.7% (2.75 th)	24.1% (6.25 th)

5.2.2 Videos with Movement

For the videos with movement, the difference between the performance of HR-NF and Poh is even higher than for the steady videos, as it can be seen in Figure 5.3, resulting in 11.9% correct estimations for Poh and 34.6% for HR-NF.

For these videos, the use of the adaptive and derivative filters becomes more important because they present a lower SNR, resulting in an improvement of the algorithm performance, as shown in Table 5.3. Indeed, the best performance is acquired when both filters are employed, with an advantage of 6.4% and 7.9% more correct HR estimations than the second and third place. The average contribution, shown in Table 5.4, enforces this statement. As the use of the adaptive and derivative filters provide a large advantage for the videos with movement and a performance close to the best for the steady videos, they are employed in our algorithm for HR estimation.

Table 5.3: Evaluation of the contribution of the adaptive filter, derivative filter and BSS strategy on HR-NF for the videos with movement.

Filter			Error range [BPM]				Rank	
Adaptive	Derivative	BSS	$ \epsilon \leq 2$	$ \epsilon \leq 5$	$ \epsilon \leq 8$	$ \epsilon > 11$		
No	No	Fixed	07.9%	20.3%	26.1%	71.7%	4 th	← Poh
		Adaptive	02.6%	07.5%	09.3%	89.7%	8 th	
	Yes	Fixed	06.8%	16.9%	23.0%	74.5%	5 th	
		Adaptive	04.0%	13.9%	18.4%	78.7%	6 th	
Yes	No	Fixed	07.9%	20.1%	28.1%	68.6%	2 nd	← HR-NF
		Adaptive	03.4%	10.4%	14.4%	83.9%	7 th	
	Yes	Fixed	10.1%	25.6%	34.5%	61.6%	1 st	
		Adaptive	06.7%	20.7%	26.6%	71.2%	3 rd	

Table 5.4: Average contributions of the adaptive filter, derivative filter and BSS strategy on HR-NF for the videos with movement. The average percentage of incorrect estimated HR ($|\epsilon| > 11$ BPM) and the average rank is shown in each cell.

Feature	Employed	
	Yes	No
Adaptive Filter	71.3% (3.25 th)	78.6% (5.75 th)
Derivative Filter	71.5% (3.75 th)	78.5% (5.25 th)
Fixed Mixture BSS	69.1% (3.00 th)	80.9% (6.00 th)

As it can be seen, the use of the adaptive BSS is the main factor that reduces the performance of the algorithm. The application of the fixed mixture, combined with the adaptive and derivative filters, provided a better robustness against noise.

5.2.3 Synthetic Data

Since it is hard to evaluate the noise performance of the algorithms for a real scenario, as signal and noise energy are unknown, we created a simulated experiment to investigate their performance against noise.

From our database, we can observe that the noise presents higher amplitudes for low frequencies, and the noise amplitudes tend to decay inversely with frequency, as shown in Figure 5.4. The low frequency characteristic of the noise was also observed by Xu *et al.* [10].

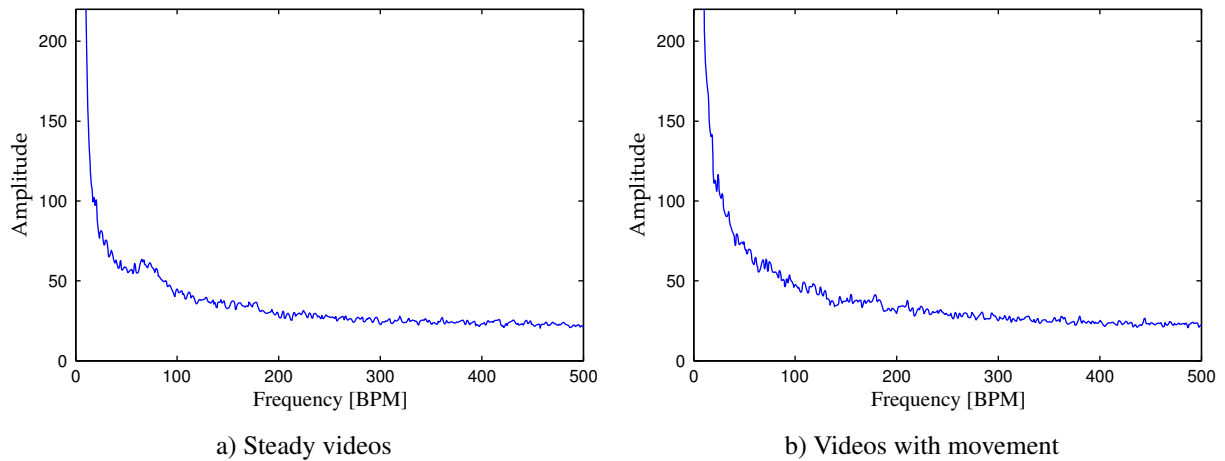


Figure 5.4: Discrete Fourier transform of the noise on real data. The curves depicts the average value of the DFT obtained for all volunteers. The frequencies near the value obtained by the pulse oximeter (± 10 BPM), for each volunteer, were not used to compute the average in order to eliminate the PPG signal from the resulting curves.

Therefore, to obtain a simulated noise more befitting to the real data, we integrate the Gaussian noise in time to amplify the low frequencies and attenuate the high frequencies, applying a gain inversely proportional to the frequency. We also remove the DC component of the noise and we multiply it by a constant in order to obtain the desired SNR. The frequency characteristic of the Gaussian noise and the integrated Gaussian noise are shown in Figure 5.5.

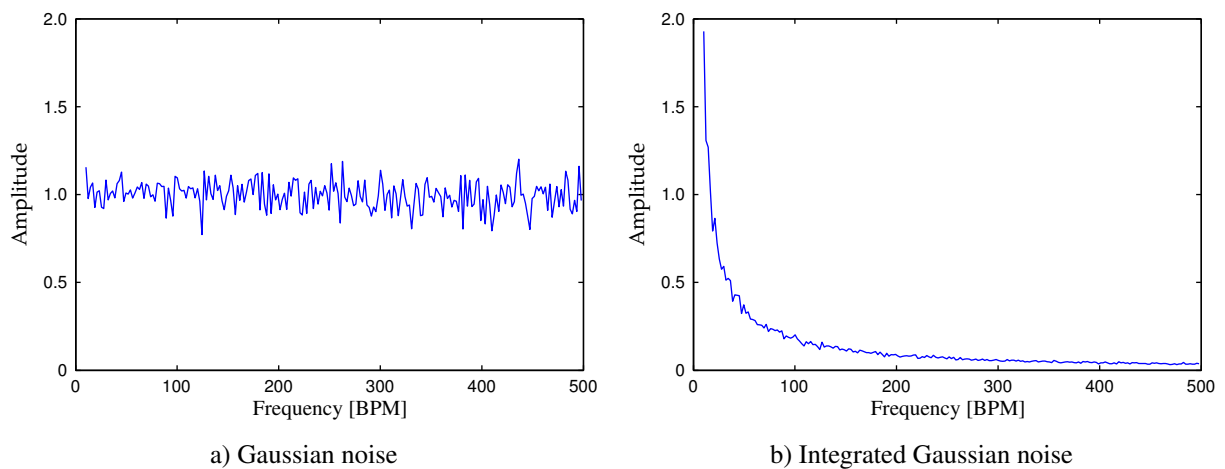


Figure 5.5: Discrete Fourier transform of the simulated noise.

Using this noise, we obtain a different scenario. Results are shown in Figure 5.6, that depicts the percentage of time that the algorithm presented an absolute error inferior to 8 BPM for 100 simulations.

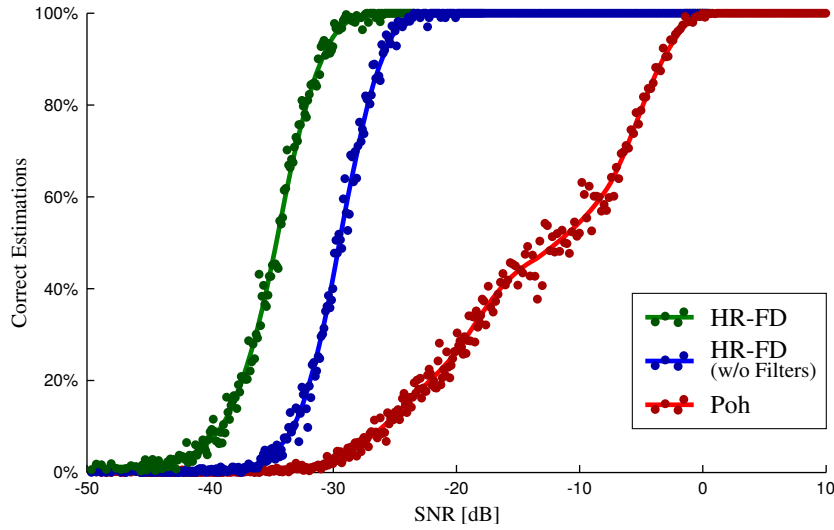


Figure 5.6: Evaluation of algorithm performance to integrated Gaussian noise. A synthetic signal composed by a sinusoidal wave of known frequency plus integrated Gaussian noise is fed to the algorithms. The percentage of time the estimation error is less or equal to 8 BPM is plotted for different SNR values. Each point on the curve is the average value obtained in 100 simulations.

Table 5.5 present the SNR values where each algorithm reached the given percentage of correct estimation. HR-NF reach 10% of correct estimations at -39.09 dB, 13.33 dB before the algorithm of Poh, and 95% at -30.18 dB, while Poh reach this percentage at -1.79 dB. This better performance is mainly due to the fixed mixture employed for BSS, that avoids the errors introduced by ICA at low SNR, because this algorithm needs a relative high value of SNR in order to correct estimate the weights to apply for each color channel. This can be observed by the performance of the algorithm employing only the fixed mixture (HR-NF w/o Filters) that presented a better performance than Poh.

Table 5.5: SNR values, in dB, where which algorithm reached the given percentage of correct estimations.

	10%	50%	95%
HR-NF	-39.09	-34.75	-30.18
HR-NF (w/o Filters)	-33.25	-29.36	-25.07
Poh	-25.76	-11.89	-01.79

Also, we can observe that the use of the adaptive and derivative filters contributed for better performance and that HR-NF reached a given percentage of correct estimations 5.4 dB, in average, before its versions that do not employ those filters.

5.2.4 Conclusion

HR-NF presented an improvement in performance, compared to the algorithm of Poh, both in the videos with and without movement and on the synthetic data. We noticed that this better performance is mainly due to the use of the fixed mixture employed for BSS. The adaptive and derivative filters also contributed and an enhancement in performance was observed with their use.

5.3 EVALUATION OF HR-MRT

Nonetheless, HR-NF and Poh do not show a satisfactory performance for videos with movement, as was already stated in Section 3.2, due to movement artifacts. This issue is addressed by HR-MRT, that employs a more complex trace extraction.

The HR-MRT is tuned by parameters that influence, for example, the frame segmentation in micro-regions, point tracking filtering, clustering algorithm, etc. Table 5.6 summarizes the parameters employed for HR-MRT. Sections 5.3.1, 5.3.2 and 5.3.3 disclose the choice of the main parameters. In Section 5.3.4 we compare the performance of HR-MRT to HR-NF and Poh.

Table 5.6: Parameters employed for HR-MRT.

Stage		Parameters
Frame Segmentation	Bilateral Filter	$\sigma_1 = 3$ (control the distance penalty) $\sigma_2 = 0.06$ (control the color difference penalty) Squared neighborhood of size 11×11
	Low-pass Filter	$N = 6$ (kernel size)
Skin Detection		Histogram based approach combined with Viola-Jones face detector
Block Duration		10 seconds
Point Tracking Filtering	Temporal Filter	Third order polynomials
	Spacial Filter	Full affine transformation
Clustering Algorithm		$d_{in} = -2 \ln(0.4)$ (similarity of 0.4 for $\sigma_s = 1$) $d_{out} = -2 \ln(0.42)$ (similarity of 0.42 for $\sigma_s = 1$)
Traces Prefilter		Employs the adaptive and derivative filters

5.3.1 Block Duration

One block corresponds to a set of frames upon which we perform point tracking in order to compensate for movements. The smaller the block duration, the less the algorithm will be affected by drift errors and occlusions on the video when, for example, a moving hand occludes part of the volunteer's face. On the other hand, blocks can introduce artifacts due to discontinuities in

the transition from one block to another. Therefore, very small blocks should also be avoided. In our work, the traces of each micro-region are normalized, removing its DC component and adjusting its amplitude, resulting in a signal with unitary variance. This normalization reduces the discontinuities between blocks, but it does not eliminate them.

Figure 5.7 shows the performance of the HR-MRT for two block durations: 60 and 10 seconds. The performance was better for the estimations using blocks of 10 seconds, compared to that of 60 seconds. This better performance is due to the point tracking that can get lost when applied for long sequences.

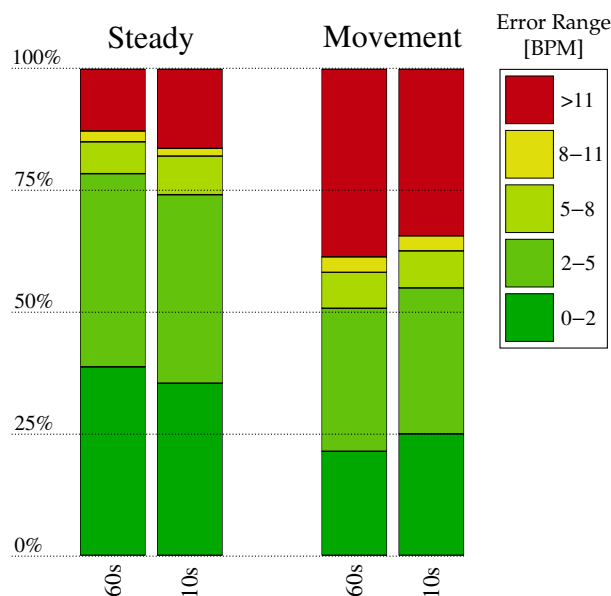


Figure 5.7: Performance of HR-MRT for different block duration.

The opposite behavior is found for the steady videos, where blocks of 60 seconds perform better than blocks of 10 seconds. As there is little movement, executing point tracking using blocks of 60 seconds is not very challenging and the algorithm offers, in general, a very good estimation while avoiding the artifacts introduced in the transitions from one block to another.

As the videos with movement present a more challenging scenario, we decided to employ blocks of 10 seconds. This block duration showed the best performance for videos with movement and a satisfactory performance for steady videos.

5.3.2 Skin Detection

Skin detection plays an important role in the definition of the ROI, eliminating the pixels from the background. Although, skin detection is a hard task, as the color of skin pixels overlap the color of non-skin pixels found in the nature. Therefore, we proposed three strategies, as described in Section 4.1:

- **Auto:** Automatic skin detection using a histogram based approach;

- **Viola-Jones:** The result of the automatic skin detection is combined with the Viola-Jones face detector to eliminate those pixels that are not on the volunteers face;
- **Manual:** The volunteer skin was manually selected on the first frame of the video.

Figure 5.8 shows the performance of the three strategies for two block duration.

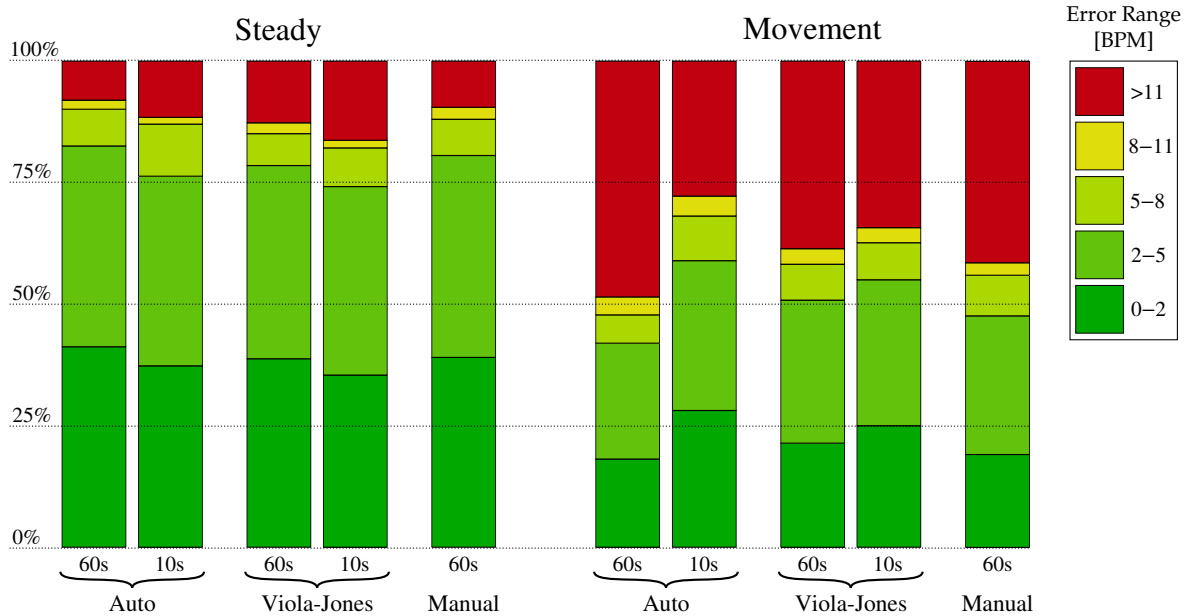


Figure 5.8: Performance of HR-MRT for different skin detection strategies.

The use of automatic skin detection presented the best performance when combined with blocks of 10 seconds, while for blocks of 60 seconds the automatic skin detection combined with the Viola-Jones face detector had the best performance. Manual detection was only tested for blocks of 60 seconds and did not offer the best performance because, as it can be seen in Figure 4.2, it is just a rough contour that includes eyes, eyebrows, teeth and parts of clothing that negatively influence the estimation.

5.3.3 Point Tracking Filtering

Points are tracked using eight successive frames of the video, as schematized in Figure 4.7. The output of the tracking algorithm is considered rather an estimation of the the actual optical flow and some noise may be added. To reduce noise, we exploit the temporal and spatial redundancy. The temporal redundancy is used by modeling the movement of the tracking points by polynomials. The position of the estimated tracking position is adjusted in order to fit to the model. The spatial redundancy is explored by means of an affine transformation, that combines up to 12 points within a small region (the micro-region).

The affine transformation is executed in two distinct ways:

- **Rigid:** Using only rigid transformation: translation, scaling and rotation;

- **Full:** Using a full search: translation, scaling, rotation, shearing and mirroring.

The performance of the different affine transform settings is presented in Figure 5.9. Their performance is, in fact, very similar and it is hard to decide which is one better. These results indicate that the rigid transformation can, actually, represent well the movements of the micro-regions. The full search would, in this case, result in parameters close to that for rigid transformation, what implies in similar performance.

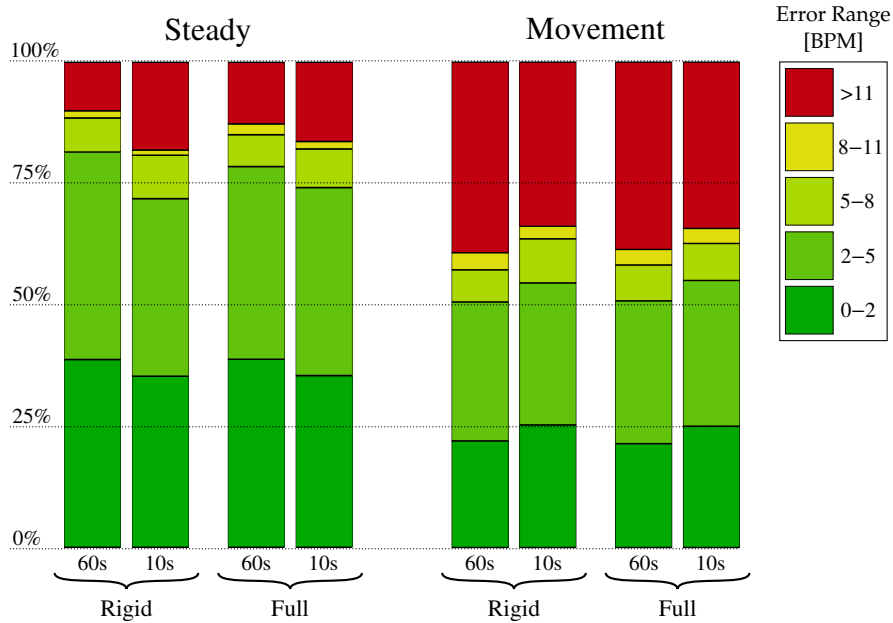


Figure 5.9: Performance of HR-MRT for different affine transform settings.

The polynomial used to model the movement varies in order. From the scheme used, it can be used polynomial of zeroth to sixth order. Figure 5.10 depicts the performance for polynomials from second to fourth order. Third order polynomials were the ones that presented the best performance among the three of them for videos with movement. The use polynomials of smaller order is not very well suited, because they can not satisfactorily model the movement of the tracking points. Higher order are also not a good choice, because, despite their ability to model complex movements, they do not offer good noise suppression as the noise can fit more easily to it.

5.3.4 HR-MRT Performance

Figure 5.11 presents a comparison between the algorithm of Poh, HR-NF and HR-MRT. The HR-MRT presents a significant gain in performance compared to Poh and HR-NF for the videos with movement, where the point tracking is used to reduce artifacts introduced by movement.

For the steady videos, the performance of HR-MRT is slightly inferior to that of HR-NF, but still superior to Poh, because of how the ROI is defined. For HR-MRT we employ a skin detector and point tracking. Therefore, it is possible that some micro-regions on the volunteers face, that could be used for HR detection, were rejected either by the skin detector or the point tracker

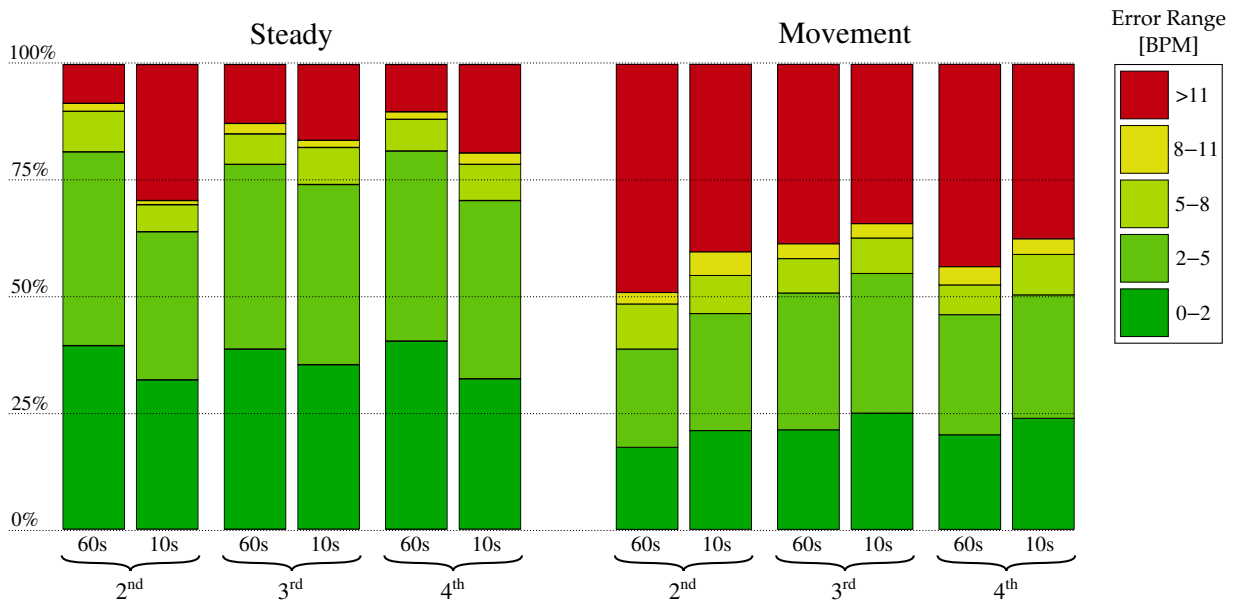


Figure 5.10: Performance of HR-MRT for different polynomial orders.

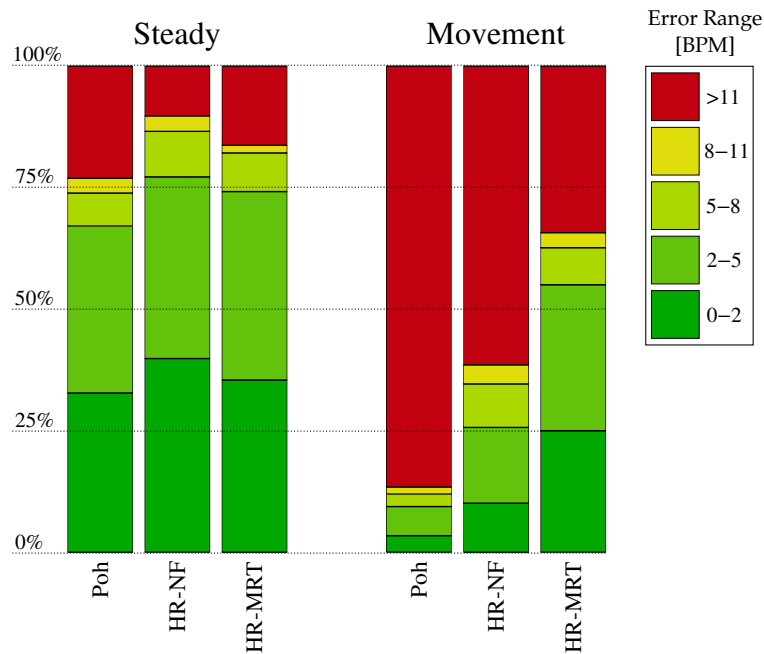


Figure 5.11: Performance of HR-MRT compared to HR-NF and Poh.

before the computation of the red, green and blue traces. Indeed, for the videos in our database, some frames were lost during data acquisition. The discontinuities introduced by these errors are not very well modeled by third order polynomials, which can introduce a drifting effect that may culminate in the rejection of the micro-region. This effect also occurs for the videos with movement, but most of the rejected regions on those videos are affected by noise. Therefore, eliminating them actually results in a better SNR.

In general, we can conclude that HR-NF is the best algorithm among the three shown in Figure 5.7 to be employed for videos with little movement. When movement is present, HR-

MRT provides better estimation and the use of small blocks is preferable.

5.4 ANALYSIS OF THE AUTOMATIC ROI SELECTION

HR-MRT divides the frames of the video in micro-regions. The micro-regions that do not meet the criteria imposed during skin detection, point tracking and clustering are eliminated. These criteria are:

1. **Skin detection:** those micro-regions that do not contain at least 80% of pixels detected as skin are rejected as the PPG signal can only be extracted from skin pixels;
2. **Point tracking:** the micro-regions for which the affine transform results in a area superior to twice or inferior to half of the initial area in a given time are rejected, as their estimated optical flow may contain errors;
3. **Clustering:** we keep only the cluster of micro-regions, grouped accordingly to their DFT, that presents the highest number of elements. The micro-regions of the other clusters are eliminated as they are likely to be affected by motion noise.

The remaining micro-regions are used for HR estimation and form the ROI. Figure 5.12 sketches the percentage of time that regions in the video were chosen to compose the ROI.

The forehead and cheeks were the most chosen, as they correspond to very well vascularized regions that facilitate the extraction of the PPG signal in the case of steady videos. For videos with movement, the algorithm prefers the forehead over the cheeks as the cheeks contain, in general, more motion noise in this case. We can also observe that eyes and mouth were rejected most of the time, mainly for videos with movement, as they introduce artifacts due to speaking and blinking. In the case of automatic skin detection without Viola–Jones face detector we can see that the region on the neck was also employed for HR detection.

The density of chosen regions were higher when using blocks of smaller duration, because fewer micro-regions are rejected during the point tracking stage. This is not crucial for steady videos, as a considerable number of micro-regions, in general, are kept after point tracking. On the other hand, for videos with movement, the number of rejected regions during point tracking increases, making the use of blocks of smaller duration a better choice. This effect can be seen in Figure 5.13 that depicts the number of micro-region employed for HR detection for each volunteer individually. Most of the micro-regions were eliminated during point tracking for blocks of 60 seconds and, for videos 02 and 05, for example, almost none remained for HR detection.

The volunteers were segregated into three groups, for better visualization, composed by those videos where: A) both tested block sizes presented poor performance; B) blocks of 10 and 60 seconds diverged greatly in performance; C) both block sizes presented good performance. It can be observed that, in general, the higher the number of micro-regions employed for HR detection,

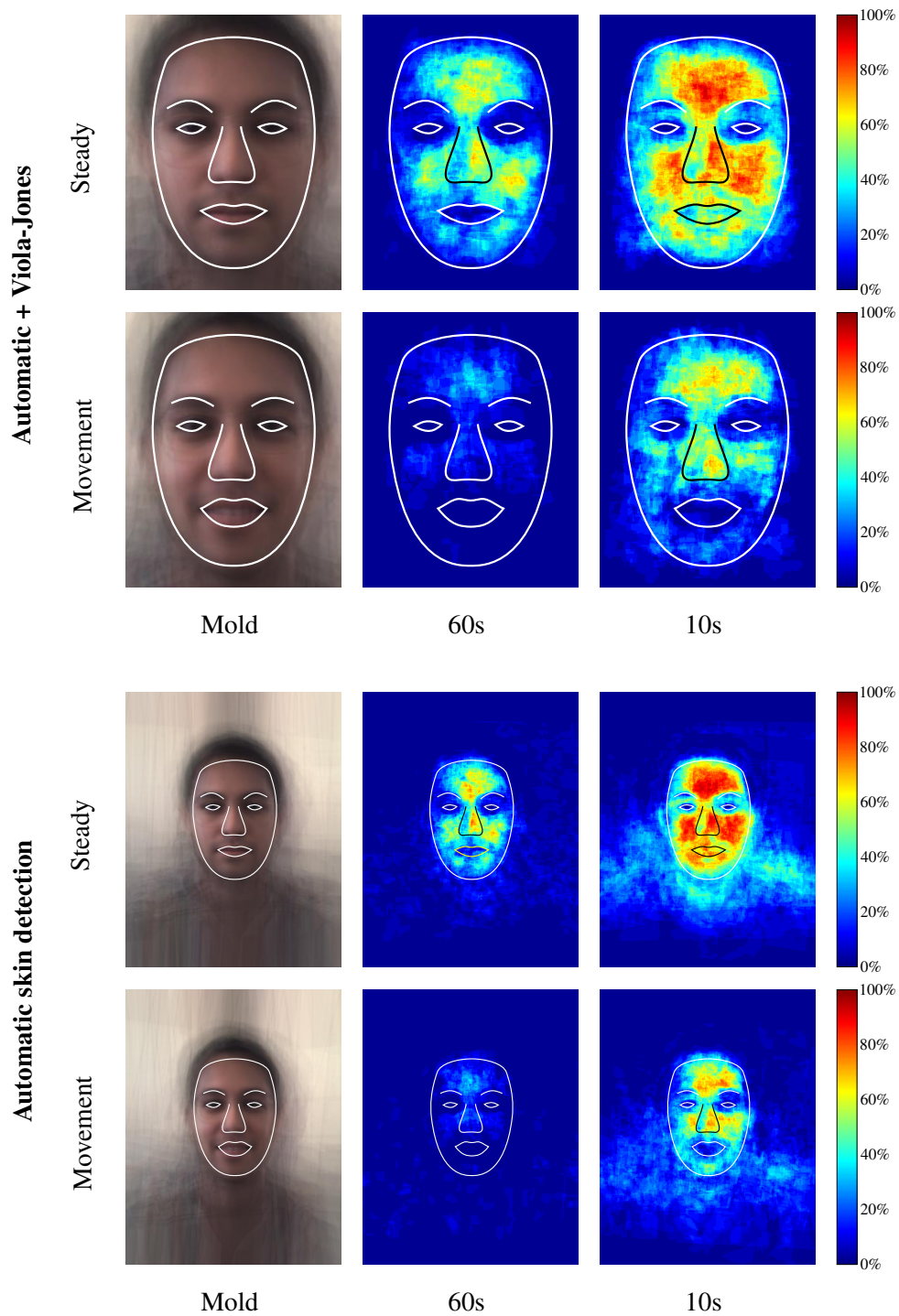


Figure 5.12: Percentage of time that the regions were chosen as ROI by HR-MRT. The mold is based on the average face of the volunteers, after aligning them, and depicts the position of face, eyes, mouth and nose. The second and third columns show the percentage of time that a given region was selected for HR detection for the block duration of 60 and 10 seconds.

the better the HR detection performance. We can also see that for videos with a poor performance, most of the micro-regions were eliminated during the point tracking stage.

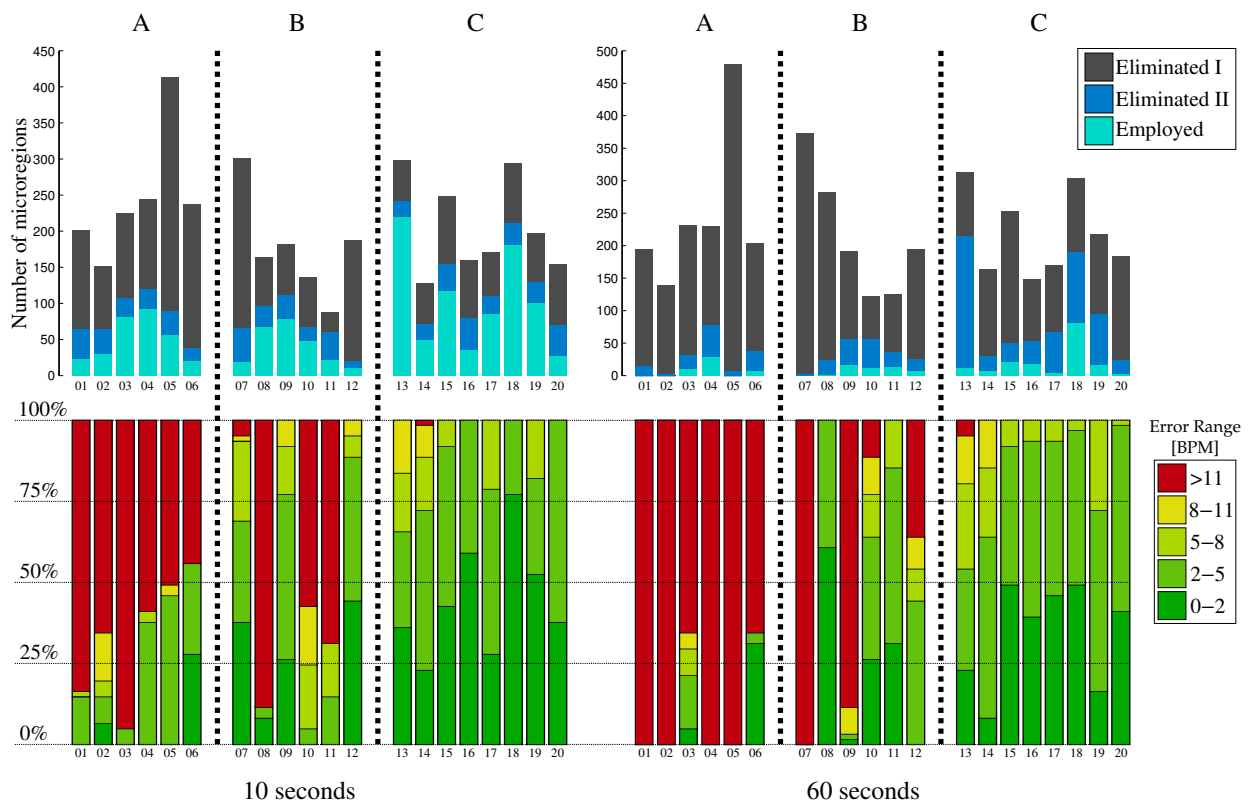


Figure 5.13: Performance of HR-MRT for the video with movement, for each volunteer individually. The first row depicts the number of micro-regions employed for HR estimation in light blue, for two blocks duration. In dark blue is represented the number of micro-regions eliminated by the clustering algorithm and in gray the number of micro-regions eliminated during point tracking. The total height of the bars show the total number of micro-regions initialized by the skin detector. The second row presents the performance separated in three groups: A, B and C.

5.5 SUMMARY

Table 5.7 summarizes the performance of the proposed algorithms, compared to that of Poh, when using skin detection combined with the Viola–Jones face detector. HR-NF, despite its simplicity, is the best suited for videos without movement, resulting in only 10.2% of incorrect HR estimations, less than half of the percentage found with Poh’s algorithm. For videos with movement, on the other hand, HR-MRT presented the best performance with 62.6% of correct estimations, almost twice that of HR-NF and 5.3 times than that of Poh.

Figures 5.14 and 5.15 show the HR estimation over time for two videos without movement. For most steady videos, all three algorithms are capable of correctly estimating the HR, as displayed in Figure 5.14. The proposed algorithms have better robustness against noise and are capable of maintaining a good performance for more challenging videos, such as the one in Figure 5.15.

Figures 5.16, 5.17 and 5.18 exhibit the HR estimation over time for three volunteers for the videos with movement. Figure 5.16 is an example where none of them were capable to offer satisfactory HR estimation. Nevertheless we can notice a better performance than the algorithm

Table 5.7: Performance of the algorithms of Poh, HR-NF and HR-MRT.

Video Group	Algorithm	Error range [BPM]			
		$ \epsilon \leq 2$	$ \epsilon \leq 5$	$ \epsilon \leq 8$	$ \epsilon > 11$
Steady	Poh	32.7%	67.1%	73.9%	23.1%
	HR-NF	39.8%	77.2%	86.6%	10.2%
	HR-MRT	35.2%	71.9%	80.8%	18.1%
Movement	Poh	03.4%	09.4%	11.9%	86.6%
	HR-NF	10.1%	25.7%	34.6%	61.5%
	HR-MRT	25.0%	55.0%	62.6%	34.3%

of Poh due to use of the fixed mixture and the filters employed.

Figure 5.17 depicts an example where all three algorithms offered satisfactory estimation. Notwithstanding, HR-MRT stayed closer to the ground-truth due to motion compensation. The effect of the motion compensation is more outstanding in Figure 5.18, where only HR-MRT was capable of correctly estimating the HR.

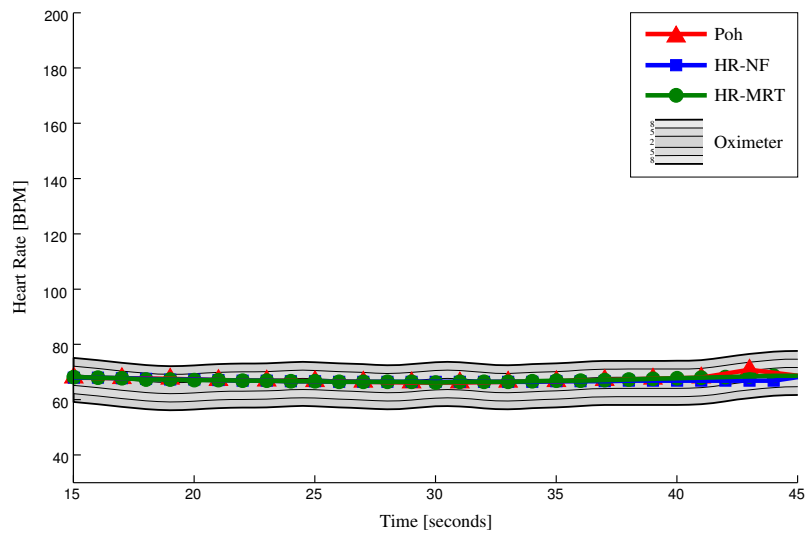


Figure 5.14: HR detection for volunteer 18 - Steady video. The oximeter reading is shown in the range of absolute errors of 2, 5 and 8 BPM.

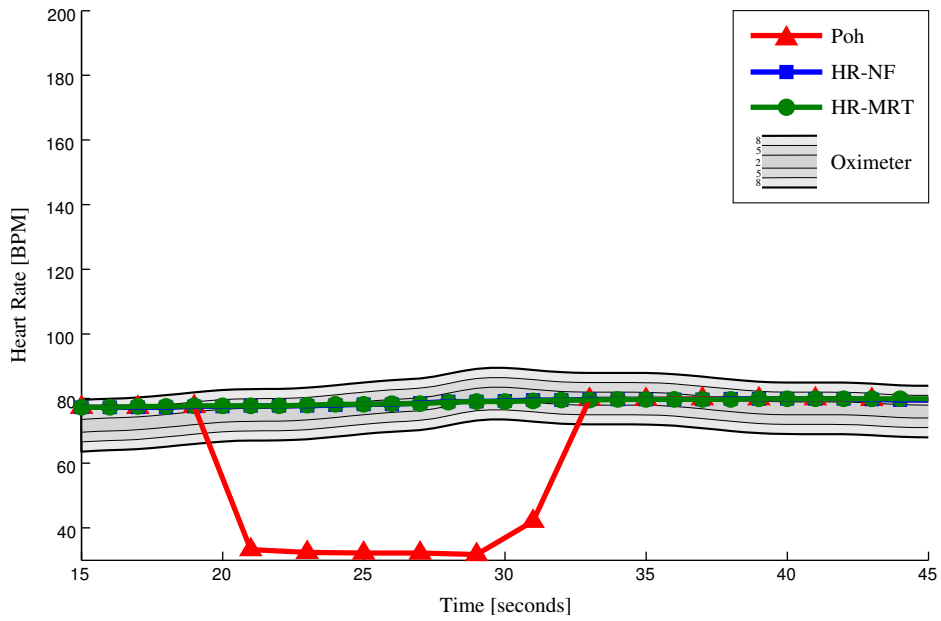


Figure 5.15: HR detection for volunteer 18 - Steady video. The oximeter reading is shown in the range of absolute errors of 2, 5 and 8 BPM.

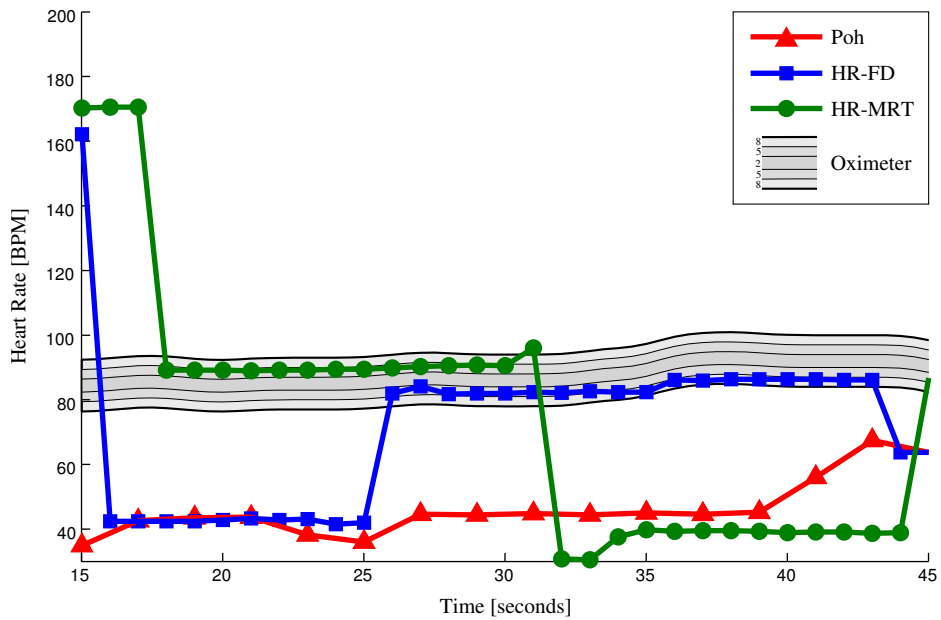


Figure 5.16: HR detection for volunteer 05 - Video with movement. The oximeter reading is shown in the range of absolute errors of 2, 5 and 8 BPM.

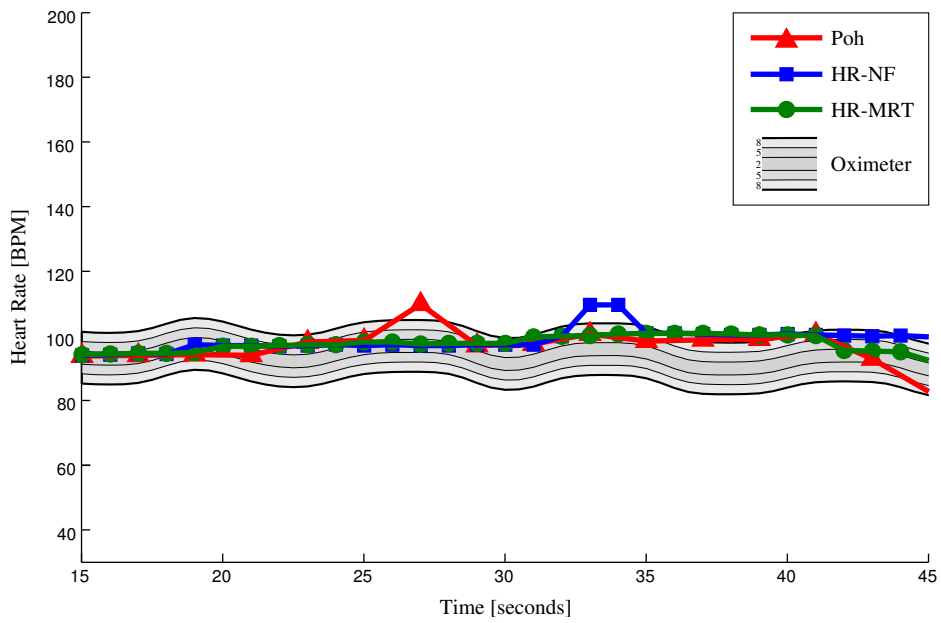


Figure 5.17: HR detection for volunteer 13 - Video with movement. The oximeter reading is shown in the range of absolute errors of 2, 5 and 8 BPM.

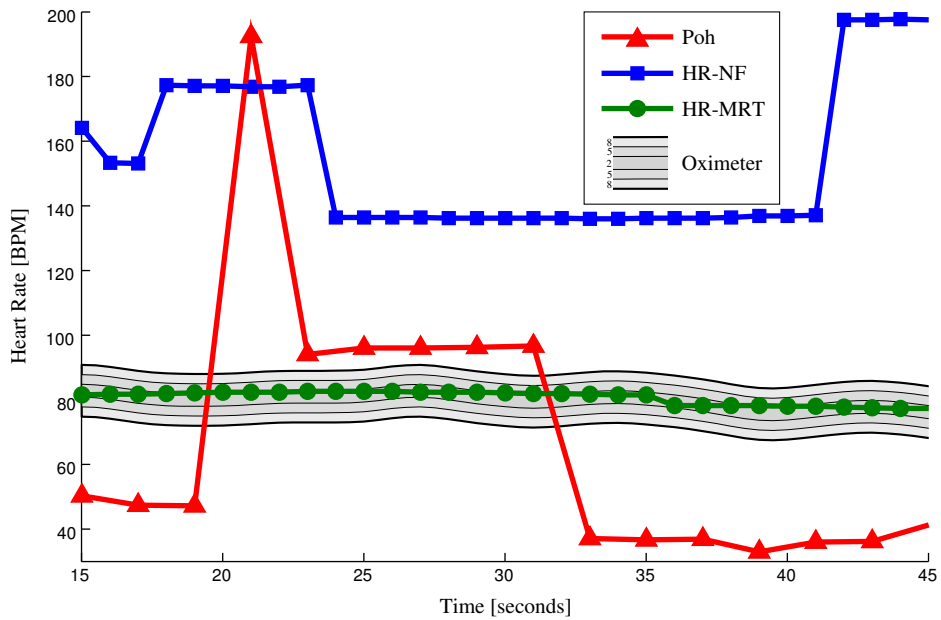


Figure 5.18: HR detection for volunteer 18 - Video with movement. The oximeter reading is shown in the range of absolute errors of 2, 5 and 8 BPM.

6 CONCLUSIONS

6.1 SUMMARY OF CONTRIBUTIONS

In this work we proposed two algorithms for heart rate (HR) estimation using videos of the human face under uncontrolled light in indoor environments. We compared our results to that of Poh [14] and observed substantial improvement.

The first algorithm, Video heart rate estimation through Face Detection (HR-NF), employs an adaptive filter that imposes a temporal coherence on the signal. This filter is based on the assumption that the heart rate varies slowly with time. A derivative filter was also employed to reduce the influence of low frequency noise. These filters boost the signal to noise ratio of the signal used for HR estimation and we showed that they are capable of reducing the number of incorrect estimated HR, particularly for noisy traces. Predominantly, a larger improvement was obtained by avoiding the use of Independent Component Analysis (ICA), commonly used in the literature. Without ICA, we eliminate the dependence on the performance of this algorithm and we are not susceptible to errors that could be introduced when it is not capable to correctly determine the mixture of the three traces. The proposed algorithm presents low complexity and, despite its simplicity, can correctly estimate the heart rate for videos with little movement.

The second algorithm, Video heart rate estimation through micro-region tracking (HR-MRT), further improved the performance by adding robustness to motion. This algorithm used a different Region of Interest (ROI) than the previous one. The video is divided in blocks and the first frame of the block is segmented in micro-regions using watershed segmentation. For each micro-region we select a set of tracking points that are used to compensate for movement. The optical flow of the tracking points is estimated using the Lucas–Kanade algorithm. The optical flow is spatial and temporal filtered to reduce the effect of noise.

For the temporal filter, we modeled the motion of the tracking points with a polynomial function. The order of the polynomial should be high enough to correctly accommodate the complex movements of the tracking points, but small enough to avoid data overfitting, which reduces the capacity of the algorithm to eliminate or attenuate noise. The results indicate that the use of third order polynomials presents the best trade-off between data representation and noise attenuation.

Spatial filtering is executed finding the affine transformation that represents the tracking points for each micro-region, minimizing the quadratic error, at a given time. This affine transformation is then applied to determine how the border of the micro-region evolved. We proposed two approaches for the affine transform: restrict the transformation to only rigid transforms and allow a full search. The results suggest that the full search is slightly better than the search restricted to rigid movements in terms of number of correct estimated HR for the case of videos with movement. However it is hard to decide which of the two approaches is better and we can conclude

that this choice is not crucial for HR estimation.

The HR-MRT algorithm also employs a clustering algorithm to decide which micro-regions to exploit for HR estimation, automatically defining the ROI. This algorithm is based on K-means and try to cluster those micro-regions that present a similar DFT. For the sake of automation, we use only the cluster that presents the largest number of elements, based on the assumption that most of the micro-regions will contain the photoplethimographic (PPG) signal and that the micro-regions composed primarily by noise do not cluster very well, as the noise varies greatly from one site to another. This assumption is dependent on the performance of the skin detection algorithm, as micro-regions that are incorrectly detected as skin do not contain the desired signal. The results show that performing skin detection using the histogram based approach, combined or not with the Viola-Jones face detector, presents a good performance.

6.2 FUTURE WORK

We noticed that, in general, the higher the number of micro-regions employed for HR detection, the better the algorithm performance is. We also noticed that most of the micro-regions are rejected during point tracking, especially for blocks of higher duration. Therefore, in a future work, special attention should be given to the optical flow algorithm.

The knowledge that, for steady videos, the forehead and the checks are the best regions on the face for HR estimation could be exploited to ameliorate the ROI, improving the performance for HR-NF and/or HR-MRT. For videos with movement, it was observed that the forehead is the best region for HR estimation, while the cheeks are rejected most of the time due to motion artifacts.

We could also improve the performance by pre-processing the traces obtained in each time block, to reduce the discontinuity artifacts introduced on the transition from one block to another. This could be done by modifying the normalization step imposing the continuity of the DC level and amplitude.

REFERENCES

- [1] S.H. Jambukia, V.K. Dabhi, and H.B. Prajapati, “Classification of ECG signals using machine learning techniques: A survey,” in *2015 International Conference on Advances in Computer Engineering and Applications (ICACEA)*, March 2015, pp. 714 – 721.
- [2] M.M. Tantawi, K. Revett, M.F. Tolba, and A. Salem, “On the use of the electrocardiogram for biometric authentication,” in *Informatics and Systems (INFOS), 2012 8th International Conference on*, May 2012, pp. BIO– (48 – 54).
- [3] S. Sangurmath and N. Daimiwal, “Application of photoplethysmography in blood flow measurement,” in *International Conference on Industrial Instrumentation and Control (ICIC)*, May 2015, pp. 929 – 933.
- [4] Michael W. Wukitsch, “Pulse oximetry: Historical review and ohmeda functional analysis,” *International Journal of Clinical Monitoring and Computing*, vol. 4, no. 3, pp. 161 – 166, 1987.
- [5] L. Capdevila, J. Moreno, J. Movellan, E. Parrado, and J. Ramos-Castro, “HRV based health and sport markers using video from the face,” in *Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, August 2012, pp. 5646 – 5649.
- [6] Giuseppe Moretti, Richard A. Ellis, and Herbert Mescon, “Vascular patterns in the skin of the face,” *Journal of Investigative Dermatology*, vol. 33, pp. 103–112, September 1959.
- [7] Ufuk Bal, “Non-contact estimation of heart rate and oxygen saturation using ambient light,” *Biomedical Optics Express*, vol. 6, no. 1, pp. 86 – 97, January 2015.
- [8] J.-P. Couderc, S. Kyal, L.K. Mestha, B. Xu, D.R. Peterson, Xiaojuan Xia, and B. Hall, “Pulse harmonic strength of facial video signal for the detection of atrial fibrillation,” in *Computing in Cardiology Conference (CinC), 2014*, September 2014, pp. 661 – 664.
- [9] Xiaobai Li, Jie Chen, Guoying Zhao, and Matti Pietikainen, “Remote heart rate measurement from face videos under realistic situations,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [10] Shuchang Xu, Lingyun Sun, and Gustavo Kunde Rohde, “Robust efficient estimation of heart rate pulse from video,” *Biomedical Optics Express*, vol. 5, no. 4, pp. 1124 – 1135, April 2014.
- [11] Yong-Poh Yu, Ban-Hoe Kwan, Chern-Loon Lim, Siaw-Lang Wong, and P. Raveendran, “Video-based heart rate measurement using short-time fourier transform,” in *International Symposium on Intelligent Signal Processing and Communications Systems (ISPACS)*, November 2013, pp. 704 – 707.

- [12] T Pursche, J Krajewski, and R Moeller, “Video-based heart rate measurement from human faces,” *Electronics (ICCE)*, pp. 544 – 545, 2012.
- [13] J.B. Bolkhovskiy, C.G. Scully, and K.H. Chon, “Statistical analysis of heart rate and heart rate variability monitoring through the use of smart phone cameras,” in *Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, August 2012, pp. 1610 – 1613.
- [14] Ming-Zher Poh, Daniel J McDuff, and Rosalind W Picard, “Non-contact, automated cardiac pulse measurements using video imaging and blind source separation.,” *Optics Express*, vol. 18, no. 10, pp. 10762 – 10774, 2010.
- [15] Ming-Zher Poh, D.J. McDuff, and R.W. Picard, “Advancements in noncontact, multiparameter physiological measurements using a webcam,” *IEEE Transactions on Biomedical Engineering*, vol. 58, no. 1, pp. 7–11, January 2011.
- [16] Wim Verkruyse, Lars O. Svaasand, and J. S. Nelson, “Remote plethysmographic imaging using ambient light,” *Biomedical Optics Express*, vol. 16, no. 26, pp. 21434 – 21445, December 2008.
- [17] Sungjun Kwon, Hyunseok Kim, and Kwang Suk Park, “Validation of heart rate extraction using video imaging on a built-in camera system of a smartphone,” *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS*, pp. 2174 – 2177, 2012.
- [18] Aymen A. Alian and Kirk H. Shelley, “Photoplethysmography,” *Best Practice & Research Clinical Anaesthesiology*, vol. 28, no. 4, pp. 395 – 406, 2014, Hemodynamic Monitoring Devices.
- [19] B. M. Jayadevappa, G. H. Kiran Kumar, L. H. Anjabeya, and S. Holi Mallikarjun, “Design and development of electro-optical system for acquisition of ppg signals for the assessment of cardiovascular system,” *International Journal of Research in Engineering and Technology*, vol. 3, pp. 520 – 525, June 2014.
- [20] W. G. Zijlstra, A. Buursma, and W. P. Meeuwssen-van der Roest, “Absorption spectra of human fetal and adult oxyhemoglobin, de-oxyhemoglobin, carboxyhemoglobin, and methemoglobin,” *Clinical Chemistry*, vol. 37(9), pp. 1633 – 1638, 1991.
- [21] J.R. Estep, E.B. Blackford, and C.M. Meier, “Recovering pulse rate during motion artifact with a multi-imager array for non-contact imaging photoplethysmography,” in *IEEE International Conference on Systems, Man and Cybernetics (SMC)*, October 2014, pp. 1462 – 1469.
- [22] J. Allen, “Photoplethysmography and its application in clinical physiological measurement,” *Physiological Measurement*, vol. 28, pp. R1 – R39, 2007.

- [23] Alrick B. Hertzman, “Photoelectric plethysmography of the fingers and toes in man,” *Proceedings of the Society for Experimental Biology and Medicine*, vol. 37, pp. 529 – 534, 1937.
- [24] Alrick B. Hertzman and C. Spielman, “Observations on the finger volume pulse recorded photoelectrically,” *American Journal of Physiology*, vol. 119, pp. 334 – 335, 1937.
- [25] J. de Trafford and K. Lafferty, “What does photoplethysmography measure?,” *Medical & Biological Engineering & Computing*, vol. 22, pp. 479 – 480, 1984.
- [26] P. Pelegris, K. Banitsas, T. Orbach, and K. Marias, “A novel method to detect heart beat rate using a mobile phone,” in *Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, Aug 2010, pp. 5488 – 5491.
- [27] Arpan Pal, Aniruddha Sinha, Anirban Dutta Choudhury, Tanushyam Chattopadhyay, and Aishwarya Visvanathan, “A robust heart rate detection using smart-phone video,” in *Proceedings of the 3rd ACM MobiHoc Workshop on Pervasive Wireless Healthcare*, New York, NY, USA, 2013, MobileHealth 13, pp. 43 – 48, ACM.
- [28] Hao-Yu Wu, Michael Rubinstein, Eugene Shih, John Guttag, Frédo Durand, and William T. Freeman, “Eulerian video magnification for revealing subtle changes in the world,” *ACM Trans. Graph. (Proceedings SIGGRAPH 2012)*, vol. 31, no. 4, pp. 651 – 658, 2012.
- [29] Richard Szeliski, *Computer Vision: Algorithms and Applications*, Springer, 2011.
- [30] Neal Wadhwa, Michael Rubinstein, Frédo Durand, and William T. Freeman, “Phase-based video motion processing,” *ACM Trans. Graph. (Proceedings SIGGRAPH 2013)*, vol. 32, no. 4, 2013.
- [31] J. Rajeshwari, K. Karibasappa, and M.T. GopalKrishna, “Survey on skin based face detection on different illumination, poses and occlusion,” in *International Conference on Contemporary Computing and Informatics (IC3I)*, November 2014, pp. 728 – 733.
- [32] Ying Liu, Dengsheng Zhang, Guojun Lu, and Wei-Ying Ma, “A survey of content-based image retrieval with high-level semantics,” *Pattern Recognition*, vol. 40, no. 1, pp. 262 – 282, 2007.
- [33] C. Santos, E. Souto, and E.M. dos Santos, “Andimage: An adaptive architecture for nude detection in image,” in *Information Systems and Technologies (CISTI), 2015 10th Iberian Conference on*, June 2015, pp. 1 – 6.
- [34] Seong G. Kong, Jingu Heo, Bisma R. Abidi, Joonki Paik, and Mongi A. Abidi, “Recent advances in visual and infrared face recognition—a review,” *Computer Vision and Image Understanding*, vol. 97, no. 1, pp. 103 – 135, 2005.
- [35] Diego A. Socolinsky, Andrea Selinger, and Joshua D. Neuheisel, “Face recognition with visible and thermal infrared imagery,” *Computer Vision and Image Understanding*, vol. 91, no. 1–2, pp. 72 – 114, 2003, Special Issue on Face Recognition.

- [36] P. Kakumanu, S. Makrogiannis, and N. Bourbakis, “A survey of skin-color modeling and detection methods,” *Pattern Recognition*, vol. 40, no. 3, pp. 1106 – 1122, 2007.
- [37] Vladimir Vezhnevets, Vassili Sazonov, and Alla Andreeva, “A survey on pixel-based skin color detection techniques,” in *IN PROC. GRAPHICON-2003*, 2003, pp. 85 – 92.
- [38] B.D. Zarit, B.J. Super, and F.K.H. Quek, “Comparison of five color models in skin pixel classification,” in *International Workshop on Recognition, Analysis, and Tracking of Faces and Gestures in Real-Time Systems*, 1999, pp. 58 – 63.
- [39] J.-C. Terrillon, M.N. Shirazi, H. Fukamachi, and S. Akamatsu, “Comparative performance of different skin chrominance models and chrominance spaces for the automatic detection of human faces in color images,” in *Fourth IEEE International Conference on Automatic Face and Gesture Recognition*, 2000, pp. 54 – 61.
- [40] G. Gomez, “On selecting colour components for skin detection,” in *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, 2002, vol. 2, pp. 961 – 964.
- [41] H. Stern and B. Efron, “Adaptive color space switching for face tracking in multi-colored lighting environments,” in *Fifth IEEE International Conference on Automatic Face and Gesture Recognition*, May 2002, pp. 249 – 254.
- [42] Jie Yang, Weier Lu, and Alex Waibel, “Skin-color modeling and adaptation,” in *Proceedings of the Third Asian Conference on Computer Vision-Volume II*, London, UK, UK, 1997, ACCV 98, pp. 687 – 694, Springer-Verlag.
- [43] Rafael C. Gonzalez and Richard E. Woods, *Digital Image Processing*, Prentice Hall, third edition edition, 2007.
- [44] J. Brand and J.S. Mason, “A comparative assessment of three approaches to pixel-level human skin-detection,” in *Pattern Recognition, 2000. Proceedings. 15th International Conference on*, 2000, vol. 1, pp. 1056 – 1059.
- [45] M.J. Jones and J.M. Rehg, “Statistical color models with application to skin detection,” in *Computer Vision and Pattern Recognition, 1999. IEEE Computer Society Conference on.*, 1999, vol. 1, p. 280.
- [46] T.S. Caetano and D.A.C. Barone, “A probabilistic model for the human skin color,” in *International Conference on Image Analysis and Processing*, September 2001, pp. 279 – 283.
- [47] T. Wark and S. Sridharan, “A syntactic approach to automatic lip feature extraction for speaker identification,” in *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech and Signal Processing*, May 1998, vol. 6, pp. 3693 – 3696.
- [48] David A. Forsyth and Jean Ponce, *Computer Vision: A Modern Approach*, Prentice Hall, second edition edition, 2011.

- [49] Charles A. Poynton, “Frequently asked questions about color,” March 1997.
- [50] J. Kovac, P. Peer, and F. Solina, “Human skin color clustering for face detection,” in *EUROCON 2003. Computer as a Tool. The IEEE Region 8*, September 2003, vol. 2, pp. 144 – 148.
- [51] Sofia Tsekeridou and Ioannis Pitas, “Facial feature extraction in frontal views using biometric analogies,” in *European Signal Processing Conference*, 1998, pp. 315 – 318.
- [52] D. Chai and K.N. Ngan, “Face segmentation using skin-color map in videophone applications,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 9, no. 4, pp. 551 – 564, June 1999.
- [53] Rehanullah Khan, Zeeshan Khan, Muhammad Aamir, and Syed Qasim Sattar, “Static filtered skin detection,” *International Journal of Computer Science Issues*, vol. 9, no. 3, pp. 257 – 261, March 2012.
- [54] M. H. Yang and N. Ahuja, “Gaussian mixture model for human skin color and its application in image and video databases,” *Proceedings of SPIE: Conference on Storage and Retrieval for Image and Video Databases*, vol. 12, pp. 987 – 1003, 2000.
- [55] Stuart Russell and Peter Norvig, *Artificial Intelligence: A Modern Approach*, Pearson Education, second edition edition, 2003.
- [56] Hani K. Al-Mohair, Junita Mohamad-Saleh, and Shahrel Azmin Suandi, “Human skin color detection: A review on neural network perspective,” *International Journal of Innovative Computing, Information and Control*, vol. 8, no. 12, pp. 8115 – 8131, December 2012.
- [57] G. Balakrishnan, F. Durand, and J. Guttag, “Detecting pulse from head motions in video,” in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, June 2013, pp. 3430 – 3437.
- [58] A. Asthana, S. Zafeiriou, Shiyang Cheng, and M. Pantic, “Robust discriminative response map fitting with constrained local models,” in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, June 2013, pp. 3444 – 3451.
- [59] Peter D. Welch, “The use of fast Fourier transform for the estimation of power spectra: A method based on time averaging over short, modified periodograms,” *IEEE Transactions on Audio and Electroacoustics*, vol. 15, no. 2, pp. 70 – 73, June 1967.
- [60] P. Viola and M. Jones, “Rapid object detection using a boosted cascade of simple features,” in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2001, vol. 1, pp. I-(511 – 518).
- [61] R. Lienhart and J. Maydt, “An extended set of Haar-like features for rapid object detection,” in *International Conference on Image Processing*, 2002, vol. 1, pp. I-(900–903).

- [62] Jean-François Cardoso, “High-order contrasts for independent component analysis,” *Neural Comput.*, vol. 11, no. 1, pp. 157 – 192, january 1999.
- [63] K.A. Hausman and E.B. Merrick, “Pulse oximeter,” november 1989, US Patent 4,883,353.
- [64] F. P. Wieringa, F. Mastik, and Steen, “Contactless multiple wavelength photoplethysmographic imaging: A first step toward spo2 camera technology,” *Annals of Biomedical Engineering*, vol. 33, no. 8, pp. 1034 – 1041, 2005.
- [65] Kenneth Humphreys, Tomas Ward, and Charles Markham, “Noncontact simultaneous dual wavelength photoplethysmography: A further step toward noncontact pulse oximetry,” *Review of Scientific Instruments*, vol. 78, no. 4, pp. 044304 – (1–6), 2007.
- [66] Eric Weisstein, *The CRC Concise Encyclopedia of Mathematics*, CRC Press LLC, New York, 2002.
- [67] Serge Beucher and Christian Lantuéjoul, “Use of watersheds in contour detection,” workshop published, september 1979.
- [68] C. Tomasi and R. Manduchi, “Bilateral filtering for gray and color images,” in *Sixth International Conference on Computer Vision*, January 1998, pp. 839 – 846.
- [69] J. Shi and C. Tomasi, “Good features to track,” in *Proceedings CVPR 94., 1994 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, June 1994, pp. 593 – 600.
- [70] Bouguet Jean Yves, “Pyramidal implementation of the Lucas–Kanade feature tracker,” *Microsoft Research Labs, Tech. Rep*, 1999.
- [71] Charles L. Lawson and Richard J. Hanson, *Solving Least Squares Problems*, Series in Automatic Computation. Prentice-Hall, Englewood Cliffs, NJ 07632, USA, 1974.
- [72] Khalid Sayood, *Introduction to Data Compression*, Morgan Kaufmann, fourth edition edition, 2012.
- [73] B.E. Bayer, “Color imaging array,” july 1976, US Patent 3,971,065.

APPENDIX

I. AFFINE TRANSFORMATION MATRIX

The affine transformation matrix can be given by a combination of simple transformations, as

$$R = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \alpha \begin{bmatrix} (-1)^m & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} 1 & \alpha_h \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \alpha_v & 1 \end{bmatrix}$$

vertical shear
horizontal shear
rotation
mirroring
scaling

(I.1)

We assume that $\det(R) \neq 0$. Therefore, we can find the values of m , α , θ , α_h and α_v that represent this transformation.

Scaling and mirroring can be determined directly from the matrix determinant:

$$m = \begin{cases} 1, & \text{if } \det(R) < 0 \\ 0, & \text{otherwise} \end{cases} \quad (\text{I.2})$$

$$\alpha = \sqrt{|\det(R)|} \quad (\text{I.3})$$

Let us now modify matrix R in such a way that the determinant of the modified matrix is equal to 1, resulting in

$$R' = \begin{bmatrix} a'_{11} & a'_{12} \\ a'_{21} & a'_{22} \end{bmatrix} = \frac{1}{\alpha} \begin{bmatrix} (-1)^m & 0 \\ 0 & 1 \end{bmatrix} R \quad (\text{I.4})$$

From this modified matrix, we can compute the other parameters as:

$$\alpha_h = \pm \sqrt{a'^2_{12} + a'^2_{22} - 1} \quad (\text{I.5})$$

$$\alpha_v = \frac{a'_{11}a'_{12} + a'_{21}a'_{22} - \alpha_h}{a'^2_{12} + a'^2_{22}} \quad (\text{I.6})$$

$$\cos(\theta) = \frac{a'_{12}\alpha_h + a'_{22}}{1 + \alpha_h^2} \quad (\text{I.7})$$

$$\sin(\theta) = \frac{a'_{22}\alpha_h - a'_{12}}{1 + \alpha_h^2} \quad (\text{I.8})$$

There are two possible solution for α_h , α_v and θ depending on the sign of α_h , but both solutions lead to the same matrix R .

II. RANDOM ESTIMATION OF THE HEART RATE

Consider two independent random variables I and E which represent the input and estimated frequency, respectively. Assume that I is uniform between 60 and 200 and E is uniform between 30 and 240. Their joint probability distribution, $f(i, e)$, is depicted in Figure II.1. As they are independent, $f(i, e)$ is given by the product of the probability distribution function of I and E and is, therefore, constant inside the drawn rectangle and zero outside.

The joint probability distribution is restricted to the constraint:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(i, e) di de = \int_{60}^{200} \int_{30}^{240} K di de = 1, \quad (\text{II.1})$$

where K is the constant level of $f(i, e)$ inside the rectangle and is found to be equal to the inverse of the rectangle area given by $(200 - 60)(240 - 30)$.

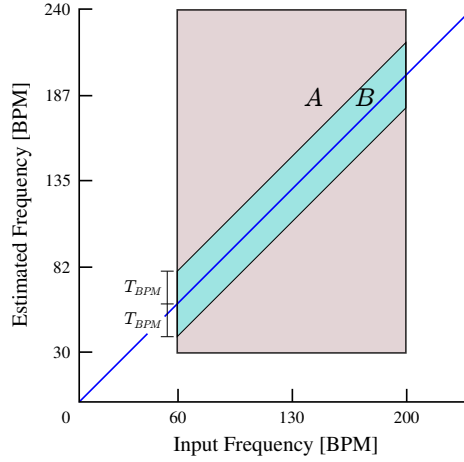


Figure II.1: Joint probability distribution.

The blue line corresponds to the points where $I = E$ and the region marked in blue correspond to all points where $|I - E| \leq T_{BPM}$. Thus, $P\{|I - E| \leq T_{BPM}\}$ is given by the integral of $f(i, e)$ inside the blue region and is equivalent to the ratio of the blue region area by the total area of the rectangle.

When $T_{BPM} \leq 30$, the upper and bottom frontiers of $|I - E| \leq T_{BPM}$ do not cross the upper and bottom lines of the rectangle and the area of the blue region is given by $2T_{BPM}(200 - 60)$. Therefore,

$$P\{|I - E| \leq T_{BPM}\} = \frac{2T_{BPM}(200 - 60)}{(200 - 60)(240 - 30)} = \frac{T_{BPM}}{105}, \quad (\text{II.2})$$

which, even though it may be small, it is greater than zero.