

2009

Abnormal ECG search in long-term electrocardiographic recordings from an animal model of heart failure

Pisut Raphisak
West Virginia University

Follow this and additional works at: <https://researchrepository.wvu.edu/etd>

Recommended Citation

Raphisak, Pisut, "Abnormal ECG search in long-term electrocardiographic recordings from an animal model of heart failure" (2009). *Graduate Theses, Dissertations, and Problem Reports*. 4518.
<https://researchrepository.wvu.edu/etd/4518>

This Dissertation is protected by copyright and/or related rights. It has been brought to you by the The Research Repository @ WVU with permission from the rights-holder(s). You are free to use this Dissertation in any way that is permitted by the copyright and related rights legislation that applies to your use. For other uses you must obtain permission from the rights-holder(s) directly, unless additional rights are indicated by a Creative Commons license in the record and/ or on the work itself. This Dissertation has been accepted for inclusion in WVU Graduate Theses, Dissertations, and Problem Reports collection by an authorized administrator of The Research Repository @ WVU. For more information, please contact researchrepository@mail.wvu.edu.

Abnormal ECG Search in Long-term Electrocardiographic Recordings from An Animal Model of Heart Failure

Pisut Raphisak

Dissertation submitted to the
College of Engineering and Mineral Resources
at West Virginia University
in partial fulfillment of the requirements
for the degree of

Doctor of Philosophy
in
Electrical Engineering

Mark A. Jerabek, Ph.D., Chair
Matthew C. Valenti, Ph.D.
Natalia A. Schmid, Ph.D.
Bojan Cukic, Ph.D.
Stephanie C. Schuckers, Ph.D.
Sherman D. Riemenschneider, Ph.D.
Conard F. Failingner, M.D.

Lane Department of Computer Science and Electrical Engineering

Morgantown, West Virginia
2009

Keywords: Abnormality Search; Electrocardiography; Heart Failure;
Long-term Recordings; Massive Data
Copyright 2009 Pisut Raphisak

Abstract

Abnormal ECG Search in Long-term Electrocardiographic Recordings
from An Animal Model of Heart Failure

Pisut Raphisak

Doctor of Philosophy in Electrical Engineering

West Virginia University

Mark A. Jerabek, Ph.D., Chair

Heart failure is one of the leading causes of death in the United States. Five million Americans suffer from heart failure. Advances in portable electrocardiogram (ECG) monitoring systems and large data storage space allow the ECG to be recorded continuously for long periods. Long-term monitoring could potentially lead to better diagnosis and treatment if the progression of heart failure could be followed. The challenge is to analyze the sheer mass of data. Manual analysis using the classical methods is impossible. In this dissertation, a framework for analysis of long-term ECG recording and methods for searching an abnormal ECG are presented.

The data used in this research were collected from an animal model of heart failure. Chronic heart failure was gradually induced in rats by aldosterone infusion and a high Na and low Mg diet. The ECG was continuously recorded during the experimental period of 11-12 weeks through radiotelemetry. The ECG leads were placed subcutaneously in lead-II configuration. In the end, there were 80 GB of data from five animals. Besides the massive amount of data, noise and artifacts also caused problems in the analysis.

The framework includes data preparation, ECG beat detection, EMG noise detection, baseline fluctuation removal, ECG template generation, feature extraction, and abnormal ECG search. The raw data was converted from its original format and stored in a database for data retrieval. The beat detection technique was improved from the original algorithm so that it was less sensitive to signal baseline jump and more sensitive to beat size variation. A method for estimating a parameter required for baseline fluctuation removal is proposed. It provides a good result on test signals. A new algorithm for EMG noise detection was developed using morphological filters and moving variance. The resulting sensitivity and specificity are 94% and 100%, respectively. A procedure for ECG template generation was proposed to capture gradual change in ECG morphology and manage the matching process if numerous ECG templates are created. RR intervals and heart rate variability parameters are extracted and plotted to display progressive changes as heart failure develops. In the abnormal ECG search, premature ventricular complexes, elevated ST segment, and split-R-wave ECG are considered. New features are extracted from ECG morphology. The Fisher linear discriminant analysis is used to classify the normal and abnormal ECG. The results provide classification rate, sensitivity, and specificity of 97.35%, 96.02%, and 98.91%, respectively.

Acknowledgements

I would like to express my gratitude to my academic advisor, Dr. Stephanie Schuckers. Without her generous guidance, support, patient, and invaluable effort, this dissertation would not exist. My acknowledgement also goes to my committee chairperson, Dr. Mark Jerabek, and my committee members, Dr. Matthew Valenti, Dr. Natalia Schmid, Dr. Bojan Cukic, Dr. Sherman Riemenschneider, and Dr. Conard Failingner, for their valuable suggestions and insights.

I also would like to thank the graduate coordinator, Dr. Muhammad Choudhry, the associate dean for academic affairs, Dr. Warren Myers, and the department chairperson, Dr. Brian Woerner, for their understanding, generosity, and patience.

My appreciation also goes to Dr. Amy de Jongh Curry from the department of biomedical engineering, the University of Memphis, who provided the research data.

I would like to sincerely thank my friends, Dr. Dulpichet Rerkpreedapong, Miss Duangruthai Pokaratsiri, and Mrs. Simona Crihalmeanu, for their kindly helps and thoughtful suggestions. My gratefulness goes to Dr. Patama Kittidhaworn for her assistance in English. I would like to specially thank Miss Nattaya Tankratoke, who always encourages me in my work and supports me in all hard circumstances.

Lastly, but most importantly, I would like to thank my mother, my father, and my sister with all my heart for their endless love and warm support.

Contents

Acknowledgements	iii
List of Figures	vii
List of Tables	ix
1 Introduction	1
2 General background	5
2.1 The Heart	5
2.1.1 Electrophysiology of the heart	7
2.1.2 Electrocardiogram	8
2.2 Cardiac Arrhythmias	9
2.2.1 Premature ventricular contraction	9
2.2.2 Ventricular tachycardia	10
2.2.3 Ventricular fibrillation	11
2.2.4 Atrial fibrillation	12
2.2.5 Sinus bradycardia	13
2.3 Heart failure	13
2.4 Sudden cardiac death	15
2.5 Heart Rate Variability	15
2.5.1 Time domain analysis	16
2.5.2 Frequency domain analysis	18
2.6 Circadian rhythm	19
3 Experiment and data collection	20
3.1 Renin-Angiotensin-Aldosterone System	20
3.2 Cardiac Fibrosis	22
3.3 Animal Model of Heart Failure	23
3.4 Experiment	24
3.5 Data Collection	25

4	Literature review	27
4.1	Previous work	27
4.1.1	Experiment	27
4.1.2	Data collection	28
4.1.3	Review of the previous work	29
4.2	ECG analysis algorithms	33
4.2.1	ECG template	33
4.2.2	Fisher linear discriminant analysis	36
4.2.3	Similarity measure	37
4.2.4	Time series representation	39
4.2.5	ECG beat detection	41
4.2.6	Noise estimation	42
4.2.7	Abnormal ECG classification	43
4.3	Discussion	45
5	Data processing and feature extraction	48
5.1	Overview	48
5.2	Types of noise	52
5.3	Morphological filter	56
5.4	Beat detection	60
5.5	Baseline fluctuation removal	63
5.6	Electromyogram detection	66
5.6.1	Training	71
5.6.2	Testing	72
5.7	Data preprocessing	74
5.8	ECG template generation	75
5.8.1	ECG template	76
5.8.2	Matching procedure	77
5.8.3	Generating the ECG template	78
5.8.4	Searching for ECG template duplicates	82
5.9	Feature extraction	84
5.10	Heart rate variability visualization	90
5.11	Discussion	96
5.12	Summary	98
6	Abnormal ECG search	100
6.1	Problem statement	103
6.2	Dataset	103
6.3	ECG and noise classification	106
6.4	Normal and abnormal ECG classification	115
6.4.1	QRS classification	119
6.4.2	Split R-wave classification	125
6.5	Discussion	127
6.6	Summary	128

7 Conclusion and future work	130
7.1 Summary of work	131
7.2 Contributions of this dissertation	136
7.3 Future work	138
References	139

List of Figures

2.1	Anatomy of the heart	6
2.2	Sinus rhythm ECG and intervals in PQRST complex	8
2.3	Normal sinus rhythm	9
2.4	Premature ventricular contraction	10
2.5	Ventricular tachycardia	11
2.6	Ventricular fibrillation	12
2.7	Atrial fibrillation	12
4.1	Pairs of points used to calculate Minkowski metrics and DTW	38
4.2	Warping path calculated by using dynamic programming	39
4.3	Similarity measures for PLA	41
5.1	An overview of the data processing and analysis method	49
5.2	Example of a clean recorded ECG	53
5.3	Burst noise	54
5.4	Fullscale noise	54
5.5	Movement artifact	55
5.6	EMG noise	55
5.7	Baseline fluctuation	56
5.8	Block diagram of morphological filter	57
5.9	Plots of input signal and signal at different locations in morphological filter	58
5.10	Block diagram of beat detection algorithm	61
5.11	Flow chart of modified beat detection algorithm	62
5.12	Block diagram of baseline removal algorithm	64
5.13	ECG with additive baseline and the additive baseline	64
5.14	Illustrations for explaining the baseline removal algorithm	65
5.15	Baseline removed ECG	66
5.16	ECG from Fig. 5.7 after eliminating baseline	66
5.17	Diagram of EMG detection algorithm	67
5.18	ECG signal which contains EMG noise	67
5.19	Morphological filter for extracting EMG noise	68
5.20	Extracted EMG noise and by-product QRS complexes	69

5.21	Extracted EMG noise and normalized moving variance $\times 10^4$	70
5.22	Input ECG signal black and detected EMG sections	71
5.23	Flowchart of data preprocessing	75
5.24	ECG template	77
5.25	Flowchart of ECG template generation algorithm	79
5.26	Examples of ECG templates	81
5.27	Examples of features for subject 3	87
5.28	Histograms of features for subject 3	88
5.29	Two-dimensional histograms of features for subject 3	90
5.30	Image plots of HRV parameters for subject 3	95
6.1	Example of normal ECG using in abnormal ECG search	101
6.2	Examples of abnormal ECGs in abnormal ECG search	102
6.3	Examples of noise using in abnormal ECG search	103
6.4	Density histograms for RR interval features of ECG and noise	108
6.5	Density histograms for beat norm features of ECG and noise	109
6.6	Density histograms for beat level features of ECG and noise	110
6.7	Density histograms for EMG level of ECG and noise	110
6.8	An example of a normalized beat and its smoother version	112
6.9	Examples of \mathbf{b} and \mathbf{s} for ECG and noise classification	113
6.10	Density histograms of N_f and D_f for ECG and noise classification	113
6.11	Density histograms for RR interval features of normal and abnormal ECG	117
6.12	Density histograms for beat norm features of normal and abnormal ECG	118
6.13	Density histograms for beat level features of normal and abnormal ECG	119
6.14	Examples of extracted QRS complexes for QRS classification	121
6.15	Examples of \mathbf{c} and \mathbf{z} for QRS classification	122
6.16	Density histogram for h and v features for QRS classification	123

List of Tables

3.1	Information about experimental rats	26
5.1	Opening morphological filter algorithm	59
5.2	Closing morphological filter algorithm	60
5.3	Results from testing the EMG detection algorithm on 3 test sets	73
5.4	Bin configuration of 2D histograms of features	90
6.1	Information about dataset	104
6.2	Training set for ECG and noise classification	107
6.3	Test set for ECG and noise classification	107
6.4	Resulting confusion matrices for ECG and noise classification using D_f . .	114
6.5	Resulting confusion matrices for the ECG classification using D_f and ΔBm	115
6.6	Training set for QRS classification	120
6.7	Test set for QRS classification	120
6.8	Resulting confusion matrices for QRS classification using QRS morphology	121
6.9	Resulting confusion matrices for QRS classification using h and v features .	124
6.10	Resulting confusion matrices of QRS classification using h , v , and \mathbf{q}	125
6.11	Comparison of QRS classification using different features	125
6.12	Training and test sets for split R-wave classification	126

Chapter 1

Introduction

Approximately five million Americans suffer from heart failure. Approximately 550,000 new patients are diagnosed annually [1]. Heart failure can lead to heart attacks which is one of major causes of death in the United States. Early detection can save lives and preserve a normal lifestyle [2]. The electrocardiogram (ECG) is the electrical signal measured from the skin or directly from the heart. The ECG is generated by the heart and reflects heart function and characteristics. It is normally used for medical diagnosis of heart diseases. Advances in technology have allowed ECG to be collected continuously using advanced sensors and a portable data recorder [3]. Moreover, massive storage media and high-performance computers make possible long-term recording and massive data analysis. This enhanced capability encourages scientists and researchers to perform studies in long-term recording ECG on humans or an animal model [4, 5]. The better understanding of ECG signals and newly discovered ECG patterns may provide benefits in prediction or early detection of heart diseases and ultimately improve quality of life.

In a previous study, a long-term ECG recording was performed on an animal model. The research aimed to study the ECG of heart failure and sudden cardiac death in a rabbit model of heart failure. The ECG was collected 24 hours a day continuously throughout the

13 weeks of the experiment period to monitor progressive change and abnormalities in the heart rhythm. In the end, there was approximately 192 MB generated per animal subject every day. The total amount of data was 68 GB from three rabbits. With this massive amount of data, it was impossible to manually analyze by using classical methods. Methods to summarize and visualize massive data were introduced in [6] and [7] to view the changes in the entire range of data. RR intervals, which are intervals between ECG beat locations, were calculated and then followed by heart rate variability parameters. One, two, and three dimensional plots were invented to display the features in one graph. It should be noted that the third dimension is portrayed as colors. The plots showed the feature value changes, while the heart was deteriorating. However, an automatic procedure for finding abnormal ECG has not yet been established. In [8], parallel algorithms were presented for ECG templating. ECG beats which deviate from the master template (a template of normal ECG) were stored for viewing purposes. This work emphasized the development of parallel algorithms and no algorithm for identifying ECG types was presented. In the literature, methods for ECG analysis and techniques for surveying large time series data, called time series data mining, have been proposed. These algorithms are surveyed for use in this research.

Data used in this dissertation were collected in an experiment directed by Dr. Amy de Jongh Curry, Department of Biomedical Engineering, University of Memphis, Memphis TN. Chronic heart failure was induced in rats using the aldosterone-salt rat model [9]. The rats had long-term electrocardiograms from their normal condition to heart failure and eventually sudden cardiac death. The experiment lasted for 12 weeks. The ECG was sensed from subcutaneous (under skin) ECG leads and transmitted to a computer. It was continuously collected 24 hours a day throughout the experiment. There are 80 GB of data from five animals. The primary goal of this research project was to understand heart

electrophysiology changes which are associated with changes in the ECG in the animal model. This understanding will potentially contribute to prevention and prediction of heart failure or sudden cardiac death. In the big picture, electrophysiology patterns need to be drawn from the ECG and linked to progressive variations to the medical condition. Patterns can be established from irregular events such as arrhythmias (abnormalities in the heart rhythm) and ECG morphology deviations. Moreover, RR intervals and heart rate variability can also be used. In this dissertation, a framework for analyzing long-term ECG recording is established for data exploration, and a paradigm for locating abnormal ECG events is proposed. The techniques developed can be utilized as a part of knowledge discovery. The challenge is the sheer mass of data. High computation, memory usage, and time-consuming techniques are avoided. Some existing algorithms are improved for lower complexity. Moreover, the bigger challenge is ECG signal quality. The subcutaneous leads are prone to various kinds of artifacts including electromyogram (signal from muscle activities), ECG baseline fluctuation, and movement noise. Ambient noise and signal error from data recording are also involved.

Objectives of this dissertation

The main objective of this dissertation is to develop a framework for analysis of massive ECG data and methods to locate abnormal ECG. The framework includes the following:

1. Data preparation: The raw data is converted from its original format, cleaned, and stored in a database for data retrieval in the later processes. Some noise and error from recordings are first rejected in this step.
2. ECG beat detection and data preprocessing: A beat detection technique is modified to be less sensitive to signal baseline jump and more sensitive to beat to beat size vari-

ation, such a premature ventricular complex followed by a normal ECG beat. Beat locations are required in the further steps. A new algorithm for EMG noise detection is introduced and used to remove sections of ECG which are heavily interfered with by EMG. The fluctuated ECG baseline is adjusted by a technique extended from the original one.

3. ECG template generation: ECG beats with a similar shape are grouped and represented by ECG templates. A new procedure for ECG template generation is developed for large data.
4. Feature extraction: Features include RR intervals and heart rate variability parameters. Visualizations of the features are displayed to illustrate progressive changes in the data.
5. Abnormal ECG search: The ECG beats detected in the data are scanned and searched for abnormal ECG which includes premature ventricular complex, ST elevated, and split R wave beats. New features are developed and the Fisher linear discrimination analysis is applied for the abnormality detection.

Dissertation outline

This dissertation is organized as follows. Chapter 2 provides necessary medical background on heart electrophysiology, electrocardiogram, heart failure, and heart rate variability. In Chapter 3, data used in this dissertation is explained, including the animal model, experiment setup, and data collecting method. Chapter 4 addresses the related works and literature survey. Data processing and abnormal ECG search are described in Chapter 5 and 6, respectively. Finally, the contributions of this dissertation and future research possibilities are included in Chapter 7.

Chapter 2

General background

This chapter describes the basic function of the heart, including mechanical and electrical components. Abnormal electrical cardiac activity or arrhythmias are described. Lastly, heart rate variability, which is linked to the nervous system control of the heart, is explained.

2.1 The Heart

The heart functions as a pump propelling blood throughout the body and collecting blood circulating back from the body. The role of the circulation mechanism is to deliver oxygen and essential metabolites to the tissues of the body and eliminate waste products and carbon dioxide. The heart has four chambers divided into left and right, and upper and lower sides (Fig. 2.1). Two upper chambers are called atria (receiving chambers), while two lower chambers are called ventricles (pumping chambers). The blood can flow from atrium to ventricle via atrioventricular valves. Following the path of the blood, the right atrium receives unoxygenated blood from the lower and upper parts of the body and the heart itself via the inferior vena cava, superior vena cava, and coronary sinus, respectively.

After the collected blood fills the right atrium, the tricuspid valve (right atrioventricular valve) opens while the right atrium pumps blood to the right ventricle. After a short delay for filling, the right ventricle pumps the low oxygen blood passing through the pulmonary valve into the pulmonary arteries and to the lungs. The oxygenated blood circulates back to the left atrium via pulmonary veins, and passes to the left ventricle via the bicuspid valve (left atrioventricular valve) by pumping of the left atrium. Finally, the left ventricle contracts to eject oxygenated blood passing through the aortic valve to the aorta which distributes it to the peripheral tissues [10]. The right and left atria function simultaneously, as do the right and left ventricles.

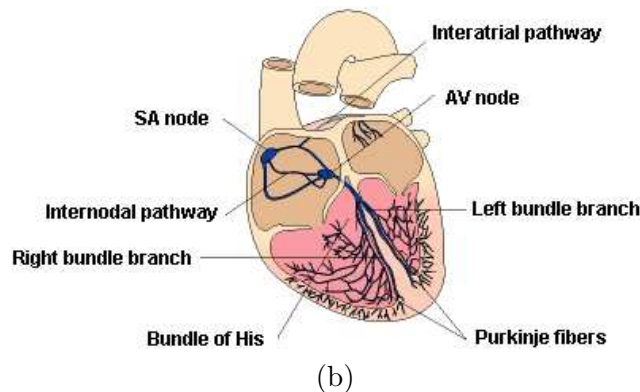
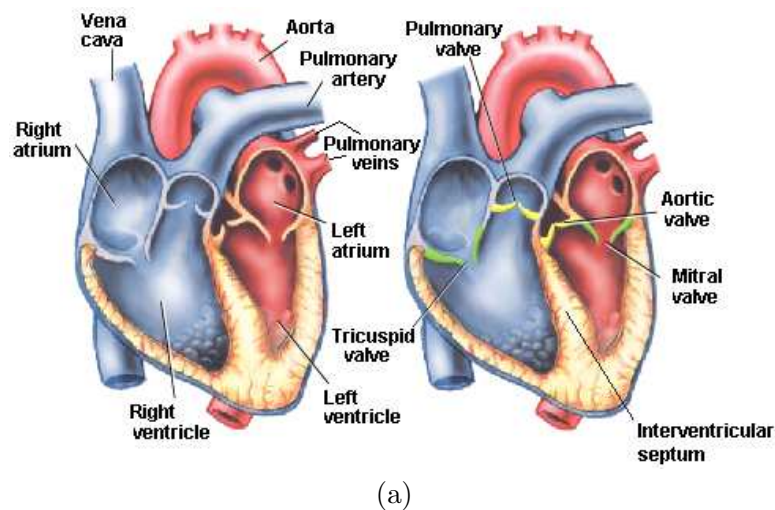


Figure 2.1: (a) Anatomy of the heart, (b) Specialized conduction system of the heart; The figures are taken from [11].

2.1.1 Electrophysiology of the heart

The electrical stimulation of the heart normally starts in the sinoatrial (SA) node – also called sinus node (Fig. 2.1). The electrical stimulation of each cardiac cell causes mechanical contraction. The location of the SA node is in the right atrium near the opening of the superior vena cava. The SA node functions as a pacemaker, and automatically generates electrical pulses at 60 to 100 cycles per minute. The electrical stimulus distributes through cardiac cells in the right atrium and then into the left atrium. As a result, right and left atria pump blood simultaneously to right and left ventricles. The spread of electrical stimulus stops at the junction between atrium and ventricle, except at small conduction tissues located at the end of the interatrial septum. This collection of conducting tissue is called the atrioventricular (AV) junction. The AV junction acts as an electrical bridge connecting the atria and ventricles. Besides conducting an electrical stimulus to the ventricles, the AV junction also delays the stimulus to ensure that the blood flows completely from atria to ventricles. The AV junction includes the AV node, which is the distal (upper) part of AV junction, and the bundle of His, which is the proximal (lower) part of the AV junction. The transmission of the stimulus is conducted to the left and right ventricular myocardium (ventricular muscle) via left and right branches of the bundle of His, respectively. The electrical stimulus spreads out broadly over the ventricular muscle by way of Purkinje fibers connected to the branches. The purpose of the bundle and Purkinje fibers is that they conduct faster than the AV node or regular cardiac muscle cells distributing electrical potential to many parts of the ventricle at once. In response, the ventricles pump blood into pulmonary arteries (for the right ventricle) and aorta (for the left ventricle) [12].

2.1.2 Electrocardiogram

An electrocardiogram (ECG) is a graphical recording of electrical voltages generated by the heart (atrium and ventricle muscles). In humans, the measurement may be performed by means of patch electrodes placed on the surface of the body or metal electrodes invasively attached to heart muscle fibers [12]. The noninvasive ECG measurement can be gathered from extremity leads and chest leads. The extremity leads record six voltage differences from electrodes on the limbs, and each lead has two subgroups, unipolar and bipolar. The chest leads attach to six positions on the chest. For invasive ECG measurement, the electrodes can be placed on specific areas of the heart in order to measure electrical activities occurring such as atrial or ventricular activity.

The waveform of one cycle of normal heart rhythm, called sinus rhythm, can be labeled corresponding to each deflection as shown in Fig. 2.2. The deflections of the ECG represent events from certain parts of the heart. The P wave represents atrial depolarization, and the QRS complex is ventricular depolarization. The interval from Q to S is the time required for an electrical stimulus to spread through the ventricles. The T wave represents ventricular repolarization [12]. Each electrical event corresponds to a mechanical event where the atrial contraction follows the P wave and ventricular the QRS. An example of ECG is illustrated in Fig. 2.3.

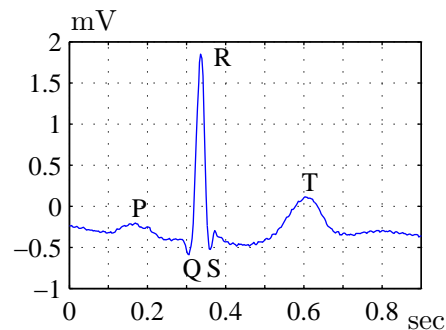


Figure 2.2: Sinus rhythm ECG and intervals in PQRST complex

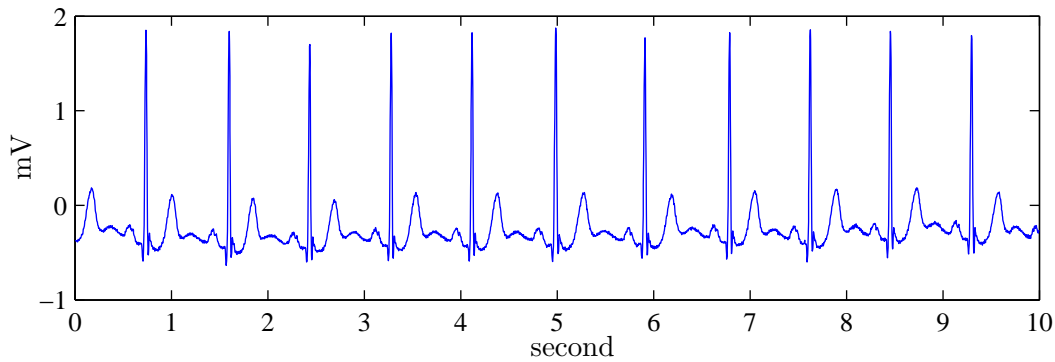


Figure 2.3: Normal sinus rhythm

2.2 Cardiac Arrhythmias

Cardiac arrhythmias are abnormalities in the heart rhythm. On the electrocardiogram, they appear as abnormal shapes of PQRST complexes and/or abnormal rates (too fast or too slow). There are several causes. For example, 1) the SA node develops an abnormal rate or rhythm, 2) the normal electrical pathway is interrupted due to heart tissue damage, and 3) another part of the heart tries to take over as the pacemaker. Some arrhythmias require immediate treatments, while no treatment is necessary for others [13]. However, arrhythmias in ventricles almost always require treatment immediately due to their control of pumping blood to the body. Some frequently occurring types of arrhythmias in patients with heart diseases are introduced below. Note that the arrhythmia figures are from the MIT-BIH database.

2.2.1 Premature ventricular contraction

A premature ventricular complex (PVC) is a premature depolarization arising in the ventricles (Fig. 2.4). PVCs occur before the next normal beat and have an abnormal

shape [14]. The QRS interval is abnormally wider than usual, > 0.12 sec. T wave and QRS complex usually are oriented in opposite directions, and have a fixed coupling interval between PVC and the preceding normal beat. An upright P wave may follow a PVC due to retrograde (reverse) conduction through the AV node [10]. Occasionally, PVCs may arise after a P wave, but before a QRS complex. PVCs may be combined in various fashions. Bigeminy is PVCs occurring every other beat, and every third beat is referred to as trigeminy. PVCs are only considered dangerous, if more than six PVCs appear in a row [12]. Moreover, PVCs can develop into more severe arrhythmias, such ventricular tachycardia or ventricular fibrillation.

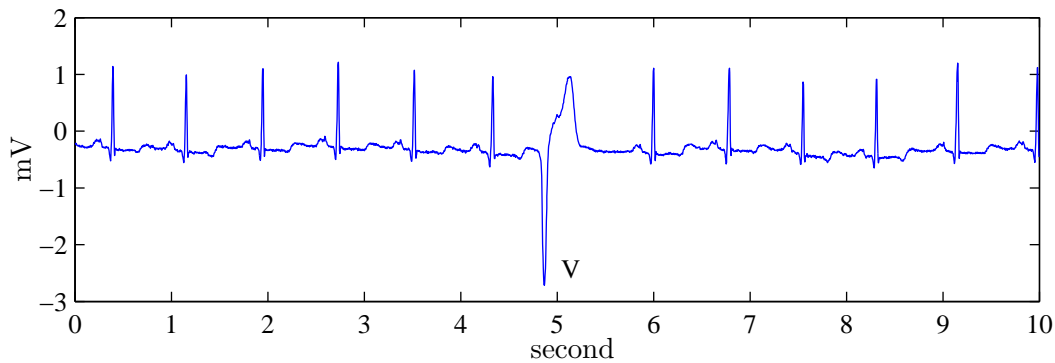


Figure 2.4: Premature ventricular contraction (denoted as “V”)

2.2.2 Ventricular tachycardia

Ventricular tachycardia (VT) is an arrhythmia consisting of three or more consecutive PVCs with uniform beat-to-beat QRS appearance (Fig. 2.5). The shape of the QRS complex is abnormal, and the duration of each complex is 0.12 seconds or greater (usually greater than 0.14 seconds) [13]. Occasionally, a P wave will occur following a QRS complex due to retrograde conduction through the AV node. The rate usually ranges from 150 to 250 beats per minute. In some cases, the heart rate may be as low as 120 beats per minute.

VT can be treated by DC voltage pacing or defibrillation [15].

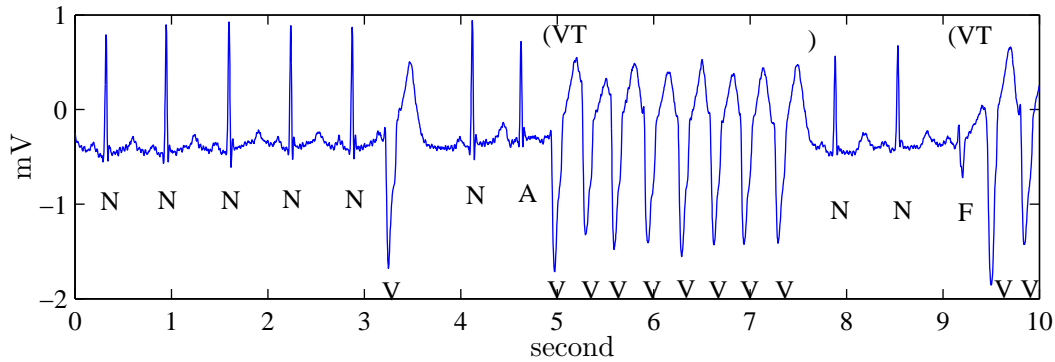


Figure 2.5: Ventricular tachycardia; “N” is denoted as a normal sinus rhythm, “V” is a premature ventricular contraction, “A” a is atrial premature beat, “F” is a fusion of ventricular and normal beat, “(VT” is the onset of VT passage, and “)” is the end of VT

2.2.3 Ventricular fibrillation

Ventricular fibrillation (VF) occurs when electrical activity in the ventricles is fractioned or chaotic (Fig. 2.6). The ventricular myocardial fibers do not contract in any coordinated fashion, but fibrillate or quiver ineffectively and asynchronously. Therefore, blood is not pumped to the body, and a patient will become unconscious and collapse within 10 to 20 seconds. Defibrillation is the only therapy and is required immediately, before any damage is done to the brain cells and the body. VF is one of three sources of cardiac arrest and the primary cause of sudden cardiac death. The ECG in VF is a chaotic, undulating pattern without discrete P waves or QRS complexes. The waveform may be either coarse or fine (Fine fibrillation waveform is below 0.2 mV) [16]. Usually, VF has a larger waveform at the onset. The rate varies from 150 to 500 beats per minute [13].

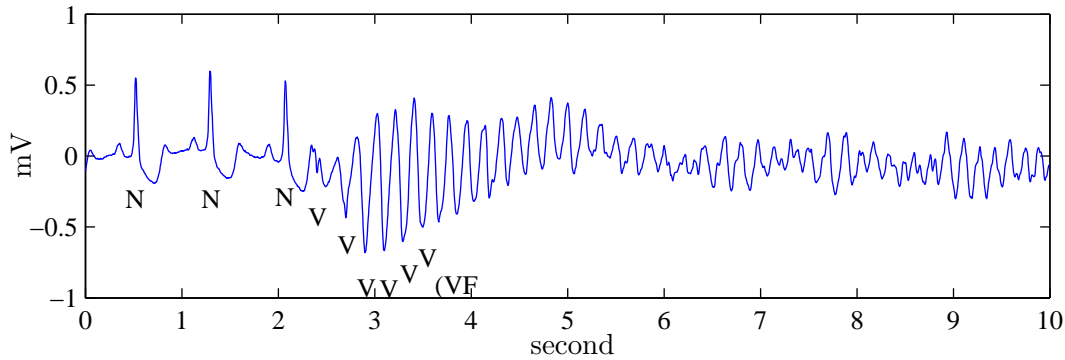


Figure 2.6: Ventricular fibrillation; “N” is denoted as a normal sinus rhythm, “V” is a premature ventricular contraction, and “(VF” is the onset of VF.

2.2.4 Atrial fibrillation

Atrial fibrillation is characterized by a notable absence of P waves in the rhythm, which are replaced with chaotic waves – fibrillation waves (Fig. 2.7). The QRS complexes are usually normal, but the ventricular response rate may be irregular, 140 - 200 beats/min. The ST and T waves are usually normal, unless myocardial diseases are present. Atrial fibrillation results from disorganized electrical activity in the atria. It does not directly cause death. However, the rapid ventricular response must be slowed down to within a range of 80 to 100 beats/min [17, 16].

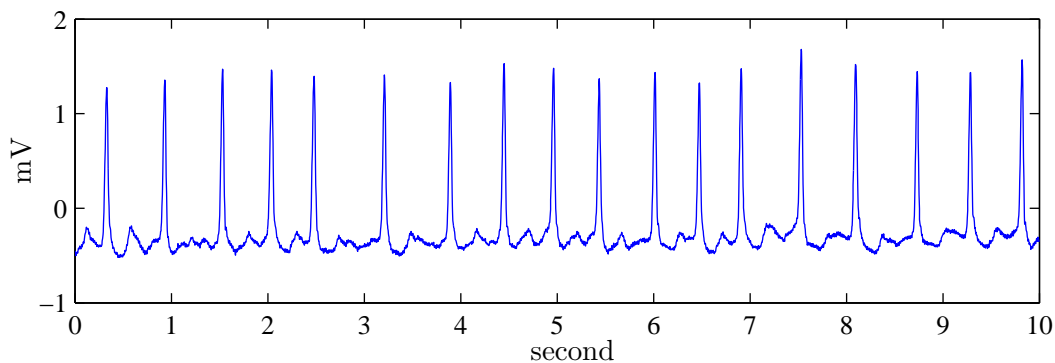


Figure 2.7: Atrial fibrillation

2.2.5 Sinus bradycardia

Sinus bradycardia has a very similar shape to normal sinus rhythm. Yet, the rate is below the lower bound of the normal rate, around 50 – 60 beats/min. When sinus bradycardia occurs, it causes a person to feel fatigued, dizzy, and lightheaded and may trigger fainting spells. However, sinus bradycardia is commonly seen in young, well-conditioned athletes, or in individuals who are sleeping [17].

2.3 Heart failure

Heart failure is a state which the heart loses the ability to maintain adequate blood circulation for the needs of the body [18]. The condition, when the heart is unable to pump all blood and, as a result, blood and fluid are congested in organs such lungs and livers, is called congestive heart failure. Such a situation requires the heart to work harder than normal. The causes of congestive heart failure include:

1. Ischemic heart disease: The usual cause is atherosclerosis which is the buildup of atheromas¹ on the walls of arteries, such as the coronary artery, that supply blood to the myocardium². Artery occlusion resulting from atheromas leads to an insufficient oxygen supply to the heart muscle, called ischemia, and, therefore, decreases the heart's ability to perform.
2. Heart attack: A heart attack causes acute damage to the myocardium due to oxygen deprivation. Scar tissue forms in replacement of injured and lost myocardium. The extent and location of scarring result in the development of heart failure.

¹fatty and plaque deposit

²heart muscle

3. Chronic high blood pressure: The heart needs to work harder to overcome greater resistance from high blood pressure. Chronic high blood pressure simulates enlargement of the heart muscle. This results in weakening of the heart muscle and, eventually, lowering the pumping ability of the heart.
4. Major cardiac arrhythmia: A severe, prolonged, and rapid rate arrhythmia affect the pattern of cardiac systole and diastole and reduce cardiac output to the point of heart failure.
5. Valvular heart disease: Defective heart valves, such as narrowed or leaking valves, cannot regulate the flow of blood through the heart chamber effectively. The heart's workload increases, as well as, a risk of heart failure [19].
6. Cardiomyopathy: Cardiomyopathy is a type of heart disease in which the heart muscle becomes inflamed and weakened. The heart loses its ability to pump blood effectively due to heart muscle damage. The condition normally begins in the walls of the ventricles. However, more severe cases may affect the walls of the atria as well. The damage commonly results in congestive heart failure and arrhythmias [20]. Its causes include infection, alcohol abuse, and cocaine abuse [19].

The pathophysiology following damage to the myocardium, for example in myocardial ischemia, includes primarily cardiomyocyte³ remodeling and hypertrophy to maintain cardiac function. Myocyte remodeling involves forming of fibrous and scar tissue following myocyte apoptosis to preserve the structure of the heart tissue. Note that fibrous and scar tissue formation is described in Chapter 3

³myocardium cell

2.4 Sudden cardiac death

The definitions of sudden cardiac death presented in [21] are:

Sudden cardiac death describes the unexpected natural death from a cardiac cause heralded by abrupt loss of consciousness within a short time period, generally less than one hour from the onset of symptoms.

Sudden cardiac death is any cardiac death occurring out of the hospital or taking place in the emergency room or dead on arrival in the emergency room.

Sudden cardiac death is primarily due to structural cardiac abnormality which alters myocardial electrophysiology enough to induce a potentially fatal tachyarrhythmia, bradyarrhythmias, or pump failure. From 157 records of ambulatory ECG from patients at their time of cardiac arrest⁴, 8% were VF, 62% were VT degenerating to VF, and 13% were torsades de pointes⁵. Among patients who have died suddenly while wearing an ambulatory ECG monitor, 17% had bradyarrhythmias as the initial rhythm [21].

2.5 Heart Rate Variability

Heart rate is principally governed by two mechanisms; 1) automatic firing of the pacemaker tissues, led by sinoatrial cells, at the intrinsic rate, and 2) the overriding control by the autonomic nervous system, consisting of the sympathetic and parasympathetic nervous systems [22]. Activation of the sympathetic results in increasing of the heart rate, while the parasympathetic results in a counteractive influence. Variability of the heart rate is a consequence of interchanging dominance of the sympathetic and parasympathetic. In many studies, analysis of heart rate and heart rate variability (HRV) has been long recog-

⁴absence of systole; failure of the ventricles of the heart to contract (usually caused by ventricular fibrillation) with consequent absence of the heart beat leading to oxygen lack and eventually to death

⁵a potentially deadly form of ventricular tachycardia.

nized as an important method for assessment of cardiac function. HRV represents one of the most promising markers for a propensity for lethal arrhythmias [23, 24, 25].

Measurement of HRV is based on normal-to-normal RR intervals (also called NN intervals), intervals between two adjacent R waves resulting from sinus node depolarization. There are two major approaches, time domain analysis and frequency domain analysis, to analyze HRV [26, 27].

2.5.1 Time domain analysis

Time domain analysis is performed on short RR-interval segments (lasting from 5 to 30 minutes). Fluctuation of heart rate is assessed by calculating parameters based on statistical operations on the RR-intervals segments. The commonly used parameters are listed as follows [28, 29].

1. Standard deviation of RR intervals (SDNN):

$$SDNN = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (RR_i - RR_{average})^2} \quad (2.1)$$

2. Interquartile range of RR intervals (IQR): The difference between the upper and lower 25 percent of the RR interval time series.

3. Normalized interquartile range of RR intervals (NIQR):

$$NIQR = \frac{IQR(RR)}{Median(RR)} \quad (2.2)$$

4. Coefficient of variation of RR (CV):

$$CV = 100 \times \frac{RMSM}{Mean(RR)} \quad (2.3)$$

where

$$RMSM = \sqrt{\frac{1}{N} \sum_{i=1}^N (RR_i - RR_{average})^2} \quad (2.4)$$

5. Mean of the standard deviation of successive 5-min RR intervals (SDNNIDX)
6. Standard deviation of the means of successive 5-min RR intervals (SDANN)
7. Standard deviation of successive differences between RR intervals (SDSD):

$$SDSD = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (\Delta RR_i - \Delta RR_{average})^2} \quad (2.5)$$

where ΔRR is the successive difference between RR intervals.

8. Interquartile range of ΔRR (DDIQR): The difference between the upper and lower 25 percentile of the ΔRR .
9. Normalized IQRSD (NDDIQR):

$$NDDIQR = \frac{IQR(\Delta RR)}{Median(RR)} \quad (2.6)$$

10. Root mean square of ΔRR (RMSSD):

$$RMSSD = \sqrt{\frac{1}{N-1} \sum_{i=1}^{N-1} (RR_i - RR_{i+1})^2} \quad (2.7)$$

11. Coefficient of variation of ΔRR (CVS):

$$CVS = 100 \times \frac{RMSSD}{Mean(RR)} \quad (2.8)$$

12. The percent of absolute ΔRR that are greater than 50ms (pNN50)
13. Counts: The absolute number of difference in successive RR intervals greater than 50ms.

2.5.2 Frequency domain analysis

The main advantage of frequency domain analysis is the possibility to study frequency-specific oscillations of RR intervals. Methods based on the fast Fourier transform are most commonly used to transform a series of RR intervals into the frequency domain. In humans, the spectrum is divided into four bands: ultra low frequency (< 0.0033 Hz), very low frequency (0.0033 - 0.04 Hz), low frequency (0.04 - 0.15 Hz), and high frequency (0.15 - 0.4 Hz). The low and high frequency bands are most commonly used where the low is associated with sympathetic and the high with parasympathetic influence on heart variability. In rats, the low frequency band is 0.04 to 1.0 Hz and the high frequency band is 1.0 to 3.0 Hz [30, 31]. The frequency domain parameters are written as follows [28, 22].

1. Total power: The total area under the power curve
2. Absolute high frequency power (HF): The area under the power spectral curve related to the high frequency band.
3. Absolute low frequency power (LF): The area under the power spectral curve related to the low frequency band.

4. Normalized HF (NHF):

$$\frac{HF}{HF + LF} \quad (2.9)$$

5. Normalized LF (NLF):

$$\frac{LF}{HF + LF} \quad (2.10)$$

6. Ratio of LF over HF (LF/HF):

$$\frac{LF}{HF} \quad (2.11)$$

2.6 Circadian rhythm

Humans and most organisms coordinate their activities with the day-light cycle, so their biological processes repeat at a cycle of 24 hours, or the so-called circadian rhythm. The circadian rhythm is controlled by an internal biological clock that runs approximately on a 24-hour cycle. A pattern of circadian rhythm also appears in HRV. For example, in young adults, HRV shows repeated patterns with maximum SDNN in the early morning before waking up [32]. However, these rhythms can be altered by actions of drugs and diseases, including heart disease. For example, a disturbed circadian rhythm of HRV appears in patients with coronary artery disease [33]. Circadian rhythm and its deviation are one of the important features to be examined in the study of heart failure and sudden cardiac death.

Chapter 3

Experiment and data collection

The purpose of this research is to develop ways to study the ECG in massive datasets. In this research, we are studying long-term changes over several weeks in the ECG as heart disease progresses. This research is performed using an animal model of heart failure. The animal is monitored using an implanted ECG monitor and the continuously stored ECG created a massive dataset. The next section describes the experimental steps in the animal model. The renin-angiotensin-aldosterone system and cardiac fibrosis are first explained as a background for an animal model of heart failure. Then, the experiment and data collection are described.

3.1 Renin-Angiotensin-Aldosterone System

Cells of living organisms function in fluid with a constant environment of essential composition. Therefore, it is essential that living organisms maintain their internal fluid balance in order to survive. Marine animals can easily maintain their internal fluid balance by exchanging salt and water with the surrounding sea water. In contrast, terrestrial animals, which are not surrounded by sea water, developed kidneys (or comparable organs)

which function to maintain a constant internal fluid environment. In mammals, the kidney functions include:

1. Maintaining H₂O balance
2. Regulating concentrations of Na⁺, Cl⁻, K⁺, HCO₃⁻, Ca²⁺, Mg²⁺, SO₄²⁻, PO₄³⁻, and H⁺
3. Maintaining proper plasma volume and osmolarity
4. Secreting hormones including renin¹ and erythropoietin²
5. Helping maintain the proper acid-base balance
6. Excreting the end products and foreign compounds
7. Converting vitamin D into its active form [11].

Filtration occurs inside a complex structure in the cortex and medulla of the kidney. Regulating the balance of salt and water is manipulated by various control mechanisms. The main control is the renin-angiotensin-aldosterone system [34]. Renin is secreted by the kidney in response of a reduction in Na⁺ concentration and water. It converts angiotensinogen, secreted by the liver, into angiotensin I. Angiotensin I is further changed to angiotensin II by the angiotensin-converting enzyme produced in the lungs. The adrenal cortex responds to increasing level of angiotensin II and produces aldosterone. Aldosterone promotes retention of Na⁺ by the kidneys. In consequence, internal osmotic pressure rises and more water is captured in the body. An increase in Na⁺ concentration and water volume inhibits renin secretion [11, 34]. Aldosterone also stimulates rising plasma K⁺ concentration and causes K⁺ secretion. Production of aldosterone decreases as the level of

¹Renin is an enzymatic hormone that triggers a chain reaction important in the process of salt conservation by the kidneys [11].

²Erythropoietin is a hormone that stimulates blood cell reproduction [11].

K^+ returns to normal. Aldosterone can be stimulated by two different stimuli. Reduction in Na^+ concentration increases the aldosterone level by using the renin-angiotensin-aldosterone pathway, while an elevation in K^+ concentration directly induces aldosterone secretion [34].

3.2 Cardiac Fibrosis

Adult cardiomyocytes are undividable and cannot be reproduced. Following myocyte necrosis or injury, such as from cardiac ischemia, inflammatory and tissue repair process response and reparative fibrosis, which is the development of excessive fibrous matter in connective tissue, and scar tissue appears to preserve the structure and formation of cardiac tissue which is crucial to heart function [35]. Angiotensin II and aldosterone take part in the tissue repair process by simulating the growth of fibroblasts³ and synthesis of type I and II collagen⁴ which are fibrillar collagens produced by fibroblasts and are involved in formation of scar tissue [34]. Microscopic scars which are initiated by a small number of necrotic cardiomyocytes can develop into macroscopic scars by continuing loss of heart tissue, as occurs in response to chronic administration of aldosterone and enhanced urinary potassium excretion [35]. Cardiac fibrosis changes myocardial structure and affects the heart function. The deposit of fibrillar collagen in cardiomyocytes and great vessels impairs tissue physical properties including increase in stiffness and resistance to deformation and decrease in contraction ability and coronary reserve. Furthermore, cardiac remodeling affects the electrical property of the heart as the fibrillar collagen obstructs the electrical link between adjacent cardiomyocytes. As a result, the ability to coordinate and synchronize contractions and blood pumping efficiency are decreased [36].

³This type of cell makes up connective tissues and the support matrix (stroma) of the skin and secretes connective tissue proteins such as collagen.

⁴A fibrous protein that makes up a major part of connective tissue

3.3 Animal Model of Heart Failure

The pathologic condition of chronic heart failure can be introduced in rats with uninephrectomy⁵ and continual administration of angiotensin II or aldosterone together with a high salt diet. In such a condition, cardiac and great vessel fibrosis develop and thus cardiac physiologic and electrophysiologic conditions are altered and result of heart failure [34, 35, 9]. In one study [9], rats were divided into four experimental groups including (1) control rats which were not operated on and were untreated, (2) intact rats with angiotensin II infusion at a rate of $9\mu\text{g}$ per hour, (3) uninephrectomized control rats receiving a high sodium diet for six weeks, and (4) uninephrectomized rats with aldosterone infusion ($0.75\mu\text{g/h}$) and high sodium diet for six weeks – called the chronic aldosterone-salt rat model. When comparing the results from control rats and uninephrectomized control rats, the histology of rats with either aldosterone or angiotensin II infusion showed diffusely distributed microscopic scarring in both left and right atria and perivascular fibrosis of the pulmonary artery and aorta. As in [35], the chronic aldosterone-salt rat model also shows perivascular fibrosis and, furthermore, accumulation of interstitial collagen in both ventricles. The mechanism is that excessive aldosterone promotes urinary potassium excretion above the normal level which can cause cardiac myocyte necrosis and injury due to reduction of intracellular potassium. In consequence, inflammatory and tissue repair process respond and stimulate production of transforming growth factor B1 (TGF-B1) which produces collagens to form fibrous and scar tissue [35]. Another suggestive pathway is that aldosterone produces fibrosis through an increase of angiotensin-converting enzyme binding and angiotensin II receptors in the myocardium which promotes myocardial fibrogenetic action. In addition, local bradykinin regulated by angiotensin converting enzyme found in fibrosis tissue stimulates production of collagens in fibroblasts. In [9], there was evidence

⁵One kidney removed

of a significant increase in angiotensin-converting enzyme binding at sites of fibrous tissue in both atria and great vessels and bradykinin receptor binding at atrial and perivascular fibrosis sites.

The level of salt intake is crucial to cardiac fibrosis in aldosterone-infused rats. It was shown that aldosterone-infused rats with a restricted salt diet have no statistical difference in the amount of fibrosis tissue when compared to control rats [35]. Aldosterone with a high salt diet in uninephrectomized rat can stimulate chronic heart failure [9]. The structural remodeling presented in the aldosterone-salt rat model is morphologically indistinguishable from the failing human heart [34].

3.4 Experiment

The experiment was directed by Dr. Amy de Jongh Curry, Department of Biomedical Engineering, University of Memphis, Memphis TN. Chronic heart failure was induced in rats using the aldosterone-salt rat model in [9] which is described as follows. Ten eight-week-old male Sprague-Dawley rats were obtained from Harlan Sprague-Dawley (Indianapolis, IN). They were divided into two groups: control ($n = 2$) and treatment ($n = 8$). The control rats were untreated and were not operated on. Treatment rats received a uninephrectomy and subcutaneous osmotic mini-pump implant (Alzet model 2004). These pumps function to deliver aldosterone at a rate of $0.75 \mu\text{g}/\text{h}$ and last for 28 days. During the experiment, after the pump expired, a new pump was installed to replace the old pump and lasted for another 14 days. In total, treatment rats received aldosterone infusion for 42 days. In addition, treatment rats also received a high Na (1% NaCl in drinking water) and low Mg diet (1.7 mmol/kg, Harlan Teklad, TD 81088). After the aldosterone treatment stopped, treatment rats continued to receive the 1% NaCl and low Mg diet and

were under observation for an additional six weeks. A 12:12-hour light-dark cycle was controlled by turning on and off light at approximately 6 AM and 6 PM, respectively.

In [37], the same experiment was performed. Myocardial fibrosis was identified after day 75 in all treatment rats. Coronal sections of the ventricles were stained with the collagen-specific stain picosirius red. Interstitial and perivascular collagen volume fraction was determined using videodensitometry.

3.5 Data Collection

A commercial biotelemetry system (Data Science International, St. Paul, MN.) was used. The system is composed of two separate units, biosensor/transmitter and receiver. The biosensor and transmitter are contained in one device which is implanted under the animal skin. The two biopotential leads of the biosensor are placed at two different places subcutaneously to measure ECG. The data are transmitted to the receiver through a modulated radio frequency. The receiver connects to a computer which is where the signal is unmodulated and sampled at 1000 Hz and uniformly digitized with a resolution of 16 bits per sample. The setting ECG amplitude was typically in range of -10 mV to 10 mV. However, it might be adjusted in some rats to a proper range.

The rats underwent surgery to implant a biosensor/transmitter between the shoulder blades. The biopotential electrodes were tunneled under the skin and placed subcutaneously in specific positions. Lead-II and modified lead-II configurations were used. In lead-II, the electrodes were placed just below the pectoral muscles against the rib cage near the first interspace and within the muscle tissue in the lower left abdominal wall. For modified lead-II, the subcutaneous leads were located above the xiphoid process and anterior mediastinum. The recording started after the rats were recovered from surgery

(to implant the mini-pump and ECG sensor). In the first week, rats received no treatment and recorded ECG was used as the baseline for comparing ECG during and after treatment. Following the first week, treatment rats received the treatment for 42 days and spent another 42 days under observation. The ECG was recorded 24 hours a day continuously throughout the 12 weeks of the experiment. There were approximately 192 MB generated for each rat every day, with a total of around 16 GB for the entire period of the experiment. Data from five rats, which include one control rat and four treatment rats, were selected for the study in this dissertation. Table 3.1 provides an information about the selected rats.

Table 3.1: Information about experimental rats: treatment group and ECG lead configuration

Subject ID	Group	ECG lead
1	Treatment	modified lead II
2	Treatment	lead II
3	Treatment	modified lead II
4	Treatment	lead II
5	Control	lead II

Chapter 4

Literature review

4.1 Previous work

This section describes previous research on analyzing long-term ECG recordings. The work was performed on data from an experiment conducted by Dr. Stephanie Schuckers, Lane Department of Computer Science and Electrical Engineering, West Virginia University (a different experiment from the experiment in Chapter 3). The strategies and techniques applied in data analysis are explained below.

4.1.1 Experiment

To provide an understanding about the data, the experiment is first explained. The experiment was designed to monitor progressive changes in the ECG in the doxorubicin model of heart failure. Three white New Zealand rabbits were selected. The experiments lasted 13 weeks, three control weeks without any treatment followed by 10 weeks with treatment. The rabbits were divided into two groups, one control rabbit and two treatment rabbits. Adriamycin was administered in the treatment rabbits once weekly via ear veins in

the amount of 2mg/kg. The control rabbit was injected with 3.75 ml of saline. To assess the progress of cardiomyopathy, blood was collected and echocardiography was also performed weekly. ECG was recorded for 24 hours for the entire duration of the experiment. Details about the doxorubicin model of heart failure and adriamycin are provided below.

Doxorubicin is an anthracycline antibiotic. It is an effective antitumor and chemotherapeutic agent for cancer treatments. Doxorubicin prepared for administration, doxorubicin hydrochloride, is called Adriamycin [38]. The type of doxorubicin-induced cardiomyopathy is mostly a chronic dilated cardiomyopathy which is usually manifested by congestive heart failure [39]. In some cases, an acute form of cardiomyopathy occurs at the very first injections. This is shown in the form of ST-T wave alterations and arrhythmias [38].

In the doxorubicin model of heart failure, chronic heart failure was induced in rabbits by administering adriamycin [40, 41, 42]. In a previous study, doxorubicin was injected in the amount of 1 mg/kg twice a week in weeks six to 10 of the study. The blood test and autopsy results showed that cardiomyopathy had developed in the animal. At the end, the hemodynamic and humoral responses were shown to be similar in the animal model and in human cardiomyopathy [40].

4.1.2 Data collection

ECGs were recorded through a commercial biotelemetry system from Data Sciences International, St. Paul, MN. Implantable biopotential transmitters were placed subcutaneously in a pocket located on the scapulae (in the back) and wired ECG leads were tunneled beneath the skin to specific positions. There were three ECG lead configurations used which include: (1) Lead II measures the ECG voltage between the right upper chest and the upper part of the left thigh. (2) MX configuration is improved from lead II configuration in order to reduce noise due to movement. It is located at the manubrium sternum

(cathode) and xiphoid process (anode) [6]. (3) An intracardiac lead is a bipolar electrode which is placed inside the right ventricle. The electrode was introduced at the jugular vein and passed through the superior vena cava to the right atrium and then the right ventricle. The intracardiac lead has benefits over subcutaneous leads because the movement does not introduce interference to the sensing ECG. The control rabbit had a lead II configuration. One of the treatment rabbits had an MX configuration while the other had two ECG leads recorded which included an MX configuration and an intracardiac lead.

The ECG was transmitted over a modulated radio frequency to a receiver which was connected to a local computer. At the local computer, the signal was sampled at 1000 Hz. Then, it was uniformly digitized with a resolution of 16 bits per sample. The acceptable amplitude was in the range of -10 mV to 10 mV, but was varied for each rabbit to maximize the input range. The ECGs were recorded 24 hours a day continuously throughout the experiment. There were approximately 192 MB generated per subject per ECG lead every day – around 68 GB totally for the entire experiment.

4.1.3 Review of the previous work

The experiment generated a massive amount of data, approximately 17 GB per animal (from one ECG lead). Therefore, it is impossible to manually analyze by using classical methods. Techniques which were used by the previous works to analyze the data are explained below.

In [6], S. Crihalmeanu suggested ways to analyze and visualize the project data. The heart rate variability of the ECGs was computed and used in the analysis. In the computation process, R-R intervals were derived from the ECGs by using a beat detection program [43]. Then, an artifact rejection algorithm was applied to remove artifacts, which included missing beats, extra beats, misplaced triggers, abnormal beats, noise where the

QRS was discernable, and noise where the QRS was undetectable [44]. HRV was computed for every usable 5-minute segment of R-R intervals, including the following parameters: MEAN, MEDIAN, SDNN, CV, IQR, SDD, DDIQR, RMSSD, LF, HF, and LHF. Note that, in rats, the low frequency band is 0.0625 to 0.1875 Hz and the high frequency band is 0.4373 to 0.5625 Hz [45]. There were in total 288 values for each parameter for 24 hours of data. To visualize the results, six different plots were introduced to either display changes in the parameters by time on a single subject or compare the parameters between study subjects. The plots are explained below.

1. Individual plots from hourly and weekly averaged data: Each parameter was averaged hourly and again averaged for each hour over the week. The results from all weeks were plotted on a graph with a 24-hour horizontal axis.
2. Individual plots where each five-minute segment was averaged weekly: This plot was the same as the previous plot but each five-minute segment was averaged over the week only.
3. Mixed plots per week from hourly and weekly averaged data: The hourly and weekly averaged data from all rabbits were compared in one graph per week.
4. Individual and mixed plots from averaged data per day and night: Daytime data from 6:00 AM to 6:00 PM was averaged weekly, as is the nighttime data. The averaged daytime and nighttime data for all rabbits were plotted on the same graph where the horizontal axis is week number. This plot compares day and night circadian rhythms.
5. Individual plots of one parameter versus another parameter: All values for any two parameters from every single five-minute segment were plotted on a graph where the horizontal and vertical axes are scales of one and the other parameter.

6. Individual plots of all data for five weeks plotted versus hour: A bar graph plotted the hourly averaged data for five weeks. It is a 24-hour plot summarizing the data for the five weeks.

By observation, MEAN, MEDIAN, SDNN, CV, IQR, SDSD, RMSSD, LF, and HF began to decrease after week 5 in the treatment group, while they remained unchanged in the control group. In addition, the circadian rhythm in the treatment group became less and less varied.

S. Crihalmeanu also manually scanned through the ECGs with the help of Dr. M. Finkel, a cardiologist from the Department of Internal Medicine, School of Medicine, WVU, for abnormal events and arrhythmias. ECGs from one of the treatment rabbits show progressive morphology changes after the first injection, and bradycardia and ventricular bigeminy in week 11. The other treatment rabbit also had progressive morphology changes after the first injection, bradycardia, and signs of myocardial infarction in week 9.

Extending on the work of S. Crihalmeanu, T. Yan presented another visualization technique in [46]. Instead of using an x-y coordinate plot and averaging across each week, the parameters were plotted by using color-coded palettes. Every 24-hour segment was plotted as a horizontal strip beginning with the first day on the top, followed by the other days. By using this visualization, additional information could be observed. For example, this tool includes features, such as zoom, rolling up data, and creating contours. In addition, the parameters can be plotted in 3D. (the x-axis is hour from 0 to 24, the y-axis is days or weeks, and the z-axis is parameter values.)

T. Yan also performed other work in this area in [7]. Each 24-hour segment was compared to the 24-hour segment from the first day to measure changes in circadian rhythm. Euclidean and dynamic time warping distances were used to measure similarities in the segments. There was a significant change in circadian rhythm during the few days after

injection in all rabbits. The difference in circadian rhythm rose in the treatment rabbits, while the circadian rhythm of the control rabbit remained stable. In addition, the fast and single-pass quantile estimation algorithm proposed in [47] was also studied by T. Yan. The algorithm was implemented and tested on generated random sequences. It provides a very small margin of error and requires little time to implement. As of yet, the algorithm has not been used to estimate quantiles from the project data.

S. Kratsas presented parallel algorithms for ECG template-based abnormality detection in [8]. The purpose is to apply fast parallel computing on abnormal ECG detection in massive data. Abnormal ECGs were detected by using the concept of template matching where the master template generated from the normal ECG cycles was used to compare with other ECG cycles in the sequence. Types of abnormal ECGs, however, need to be identified by having medical experts view the templates. Locations of unmatched ECG cycles were identified and templates of unmatched ECG cycles were created for viewing purpose. Details of the ECG template generation are given in Section 4.2. Four parallel computing algorithms, which are root-only template formation, node data server, multiple process correlation computation, and thread data server, were implemented. The root-only template formation provided the best result. The amount of time used to process 17 GB of ECG from one rabbit on a 500 MHz Pentium III computer decreased from 11 days to 26-28 hours when processing on a six-node Beowulf cluster by using the root-only template formation approach.

A database of the project was designed by Dr. S. Schuckers, et al. It was first implemented by C. Dong on the HP 9000 K450 Server by using Oracle 8, and then modified by L. Guo [48, 49]. This design can store only the raw ECGs and ECG events, such as arrhythmias.

In one of my studies, I extended T. Yan's work in visualization [5]. In this, we made the tool more flexible to account for two-variable visualization and outliers (out-of-considerable-range data). I also developed new approaches for scanning the data automatically, instead of manually as performed by S. Crihalmeanu. The database was reconfigured so that extracted features and calculation results could be stored and interesting events could be queried.

4.2 ECG analysis algorithms

4.2.1 ECG template

The ECG template matching technique groups ECG cycles with a similar shape represented by an ECG template and stamps the time where each beat of a particular shape is located in the data. These generated templates summarize the data and provide an overview as to what have occurred in the data, such as arrhythmias [15].

Euclidean distance

The matching score between a template, $\{t_1, t_2, \dots, t_N\}$, and a waveform, $\{s_1, s_2, \dots, s_N\}$, can be simply defined by the Euclidean distance, d , which is formulated as below. Note that N is a waveform dimension.

$$d = \sqrt{\sum_{i=1}^N (t_i - s_i)^2}$$

Correlation waveform analysis

The similarity of two waveforms is determined by their correlation coefficient. The correlation coefficient has a value within the range, $-1 \leq \rho \leq 1$. A perfect match gives

a correlation coefficient of 1. The correlation coefficient is defined by $\{t_1, t_2, \dots, t_N\}$ and $\{s_1, s_2, \dots, s_N\}$ which are two waveforms with a sample size of N and average values of \bar{t} and \bar{s} , respectively. The correlation coefficient, ρ , which is used as a matching score is computed as:

$$\rho = \frac{\sum_{i=1}^{i=N} (t_i - \bar{t})(s_i - \bar{s})}{\sqrt{\sum_{i=1}^{i=N} (t_i - \bar{t})^2 \sum_{i=1}^{i=N} (s_i - \bar{s})^2}}.$$

Bin area method

Bin area method is a method to approximate correlation coefficient. Compared to the above method, it requires lower computational power with a comparative result. Two comparing waveforms are segmented to bins where a bin is a summation of three consecutive samples. Thus, the number of compared points decreases from N to M , where $N = 3M$. The bins can be represented as $\{S_1, S_2, \dots, S_M\}$, where $S_j = \sum_{i=3j-2}^{i=3j} s_i$ and $\{T_1, T_2, \dots, T_M\}$, where $T_j = \sum_{i=3j-2}^{i=3j} t_i$. The matching score is the correlation coefficient which is calculated as the following equation.

$$\rho = 1 - \sum_{j=1}^{j=M} \left| \frac{T_j - \bar{T}}{\sum_{k=1}^{k=M} |T_k - \bar{T}|} - \frac{S_j - \bar{S}}{\sum_{k=1}^{k=M} |S_k - \bar{S}|} \right|,$$

where \bar{S} is computed as $1/M \sum_{j=1}^{j=M} S_j$.

Comparing the three matching scores, correlation waveform analysis requires highest computation power followed by bin area method and the Euclidean distance. The Euclidean distance is used for exact match, while correlation waveform analysis and bin area method allow some degree of variation.

Method

ECG cycles are extracted by using the trigger program in [43] to locate beat position and applying a window with the size of N samples. The first template is created from 20 beats of sinus rhythm from the control week. The selected cycles must occur near in time and closely resemble each other. Moreover, the matching score between the generated template and each beat used to create the template must exceed a specific value or that beat is excluded from the templates. For the Euclidean distance and the correlation, a template is created by averaging the beats. For the bin area method, the bins of the selected cycles are normalized and baseline shifts are removed by subtracting its average value before creating the template. The template is formed by averaging the corresponding bins, where the bins of the selected cycles must be properly aligned to the others. Once the first template is created, it is used to compare other ECG cycles in the data. If the matching score is above a set threshold, that cycle is added to the template and its count is increased by one. Note that the count of the first template is set to 20 at the beginning. The new cycle is added to the template by summing its bins to the template bins which are multiplied by the previous count, and then dividing the summation by the current count. In the case where the matching score is below the threshold, that cycle is assigned as a new template with count of 1. A new template is generated whenever a cycle does not match other templates. If new templates are generated over 50 consecutive times, this section is marked. This happens when ventricular fibrillation occurs. Correct alignment is necessary when matching templates to a cycle. Thus, a window with a size larger than N samples is employed and the windowed cycle is shifted to find the best alignment for the template [50, 15, 8].

4.2.2 Fisher linear discriminant analysis

Fisher linear discriminant analysis is a method used in statistics and machine learning. Best linear combination of features is determined to separate classes of data. The linear combination is further used as a classifier (linear classifier) [51]. The Fisher linear discriminant for two data classes in [52] and [53] is explained as follows.

Define a training set vector as $\mathbf{x}_i \in \mathbf{R}^n$ and a set of class labels as $y \in \{1, 2\}$. Indices of training vectors belonging to class 1 and 2 are stored in sets $\mathcal{I}_y = \{i : y_i = y\}$, $y \in \{1, 2\}$. The linear discrimination function is written as $f(\mathbf{x}) = \langle \mathbf{w} \cdot \mathbf{x} \rangle + b$, and the class separability in a direction is

$$F(\mathbf{w}) = \frac{\langle \mathbf{w} \cdot \mathbf{S}_B \mathbf{w} \rangle}{\langle \mathbf{w} \cdot \mathbf{S}_W \mathbf{w} \rangle},$$

where \mathbf{S}_B is the between-class scatter matrix,

$$\mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T, \quad \boldsymbol{\mu}_y = \frac{1}{|\mathcal{I}_y|} \sum_{i \in \mathcal{I}_y} \mathbf{x}_i, \quad y \in \{1, 2\},$$

and \mathbf{S}_W is the within class scatter matrix defined as

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2, \quad \mathbf{S}_y = \sum_{i \in \mathcal{I}_y} (\mathbf{x}_i - \boldsymbol{\mu}_y)(\mathbf{x}_i - \boldsymbol{\mu}_y)^T, \quad y \in \{1, 2\}.$$

\mathbf{w} of the linear discrimination function is determined by maximizing the class separability.

$$\mathbf{w} = \operatorname{argmax}_{\mathbf{w}'} F(\mathbf{w}') = \operatorname{argmax}_{\mathbf{w}'} \frac{\langle \mathbf{w}' \cdot \mathbf{S}_B \mathbf{w}' \rangle}{\langle \mathbf{w}' \cdot \mathbf{S}_W \mathbf{w}' \rangle}$$

The result \mathbf{w} is

$$\mathbf{w} = \mathbf{S}_W^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

The bias of linear discrimination function, b , is determined such that it meets the following condition:

$$\langle \mathbf{w} \cdot \boldsymbol{\mu}_1 \rangle + b = -(\langle \mathbf{w} \cdot \boldsymbol{\mu}_2 \rangle + b).$$

From the above approach, \mathbf{w} and b of the linear discrimination function were derived such that class $y = 1$ and $y = 2$ can be separated by a hyperplane $\mathcal{H} = \{\mathbf{x} : \langle \mathbf{w} \cdot \mathbf{x} \rangle + b = 0\}$. The classes, therefore, can be determined by using the following rule:

$$\langle \mathbf{w} \cdot \mathbf{x}_i \rangle + b \geq 0, \quad y_i = 1,$$

$$\langle \mathbf{w} \cdot \mathbf{x}_i \rangle + b < 0, \quad y_i = 2.$$

4.2.3 Similarity measure

A simple way to determine similarity between two signals is by subtracting one signal from the other signal on point-by-point basis, and summarizing absolute values of the differences. The result is called Manhattan distance. It measures dissimilarity between the two signals. In other words, two similar signals give a small distance, and vice versa. Similar types of methods which determine similarity by using the distance between two signals are called distance measures. These are required properties of the distance measures: 1) symmetry ($D(T_1, T_2) = D(T_2, T_1)$), 2) constancy of self-similar ($D(T_1, T_1) = 0$), positivity ($D(T_1, T_2) = 0$ iff $T_1 = T_2$), and 3) triangular inequality ($D(T_1, T_2) \leq D(T_1, T_3) + D(T_2, T_3)$) [54]. Some useful distance measures are discussed below.

Minkowski metrics

Given two time series $Q = \{q_1, \dots, q_n\}$ and $C = \{c_1, \dots, c_n\}$. Their Minkowski metrics are defined as:

$$D(Q, C) = \sqrt[p]{\sum_{i=1}^n |q_i - c_i|^p},$$

where $p = 1$ is Manhattan distance, $p = 2$ is the Euclidean distance, and $p = \infty$ is the Supremum or "sup". Pairs of points used to calculate Minkowski metrics are displayed in Fig. 4.1(a).

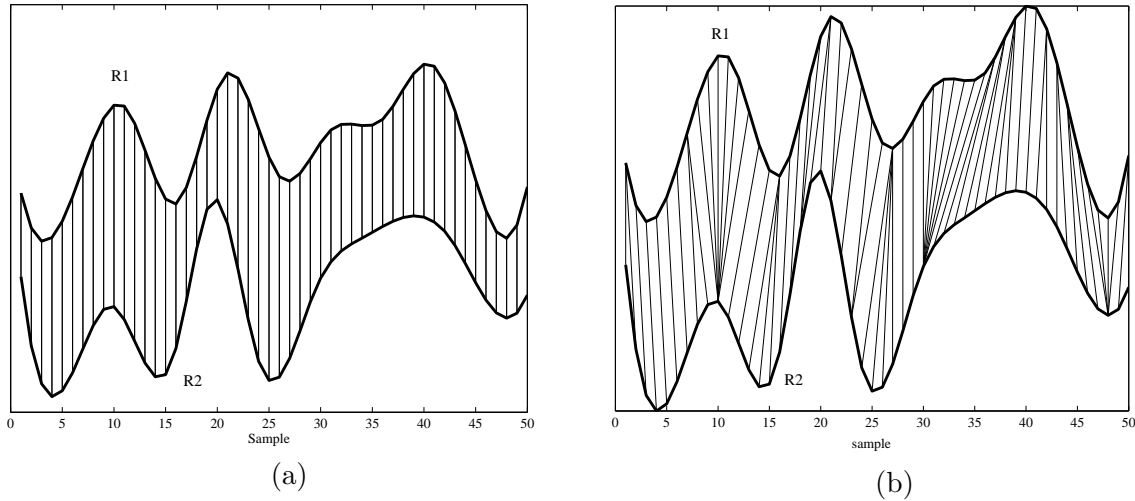


Figure 4.1: Pairs of points used to calculate distances. R1 and R2 are two passages of signal to measure distance. (a) Pairs of points for Minkowski metrics, (b) Pairs of points matching by dynamic algorithm which are used to calculate DTW distance. The pictures are both from [55].

Dynamic time warping

Dynamic time warping (DTW) is another distance measure which incorporates the capability of imprecise matching. This extends to comparing two uneven-length signals. The algorithm defines optimal matching points by using dynamic programming, as illustrated in Figs. 4.1(b) and 4.2 [55, 56, 57]. Then, the resulting distance is the Euclidean distance between the optimal matching points. The drawback is that the process time-consuming due to dynamic programming computation. However, this can be speeded up by using low-dimension representations which will be described below.

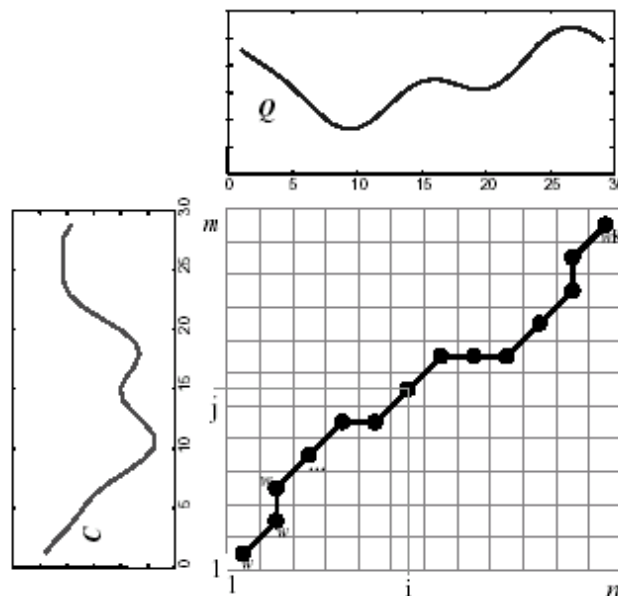


Figure 4.2: Warping path calculated by using dynamic programming. The picture is from [55].

4.2.4 Time series representation

The major hurdle in similarity search is due to the inherent high dimensionality of data. This high dimensionality causes difficulties with high computation including high memory demand and high time consumption. To solve this problem, several papers suggested various approaches to represent time series in lower-dimension forms [58, 59, 60, 61, 62].

Time series is represented by K linear segments, called piecewise linear approximation (PLA), as shown in Fig. 4.3. The number of segments, K , is selected such that it is not too small so that important features are lost, and it is not too large so that unnecessary details are captured [63]. By using the bottom-up algorithm, optimum segmentation can be achieved [64]. The algorithm starts by combining two adjacent points to create the finest segments. The cost of merging each pair of adjacent segments is computed with the representation error. Iteratively, a pair of adjacent segments which produced the lowest cost is merged to create a longer segment. The algorithm stops when the representation

error of all segments meets a certain criteria. PLA is denoted by a 4-tuple vector for each segment, $[A_{XL_i}, A_{XR_i}, A_{YL_i}, A_{YR_i}]$, where X is time, Y is magnitude, L is left, and R is right. Another time series representation is piecewise aggregate approximation [59]. A time series X of length n is equally divided into N segments and each segment is represented by its mean. Therefore, the representation is a vector $\bar{\mathbf{X}} = [\bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_N]$. It is computed by using the following equation:

$$\bar{x}_i = \frac{N}{n} \sum_{j=\frac{n}{N}(i-1)+1}^{\frac{n}{N}i} x_j.$$

Adaptive piecewise constant approximation (APCA) is one of time series representations [59]. It represents a signal by a series of the means of uneven-length segments. With equal dimensions, the number of segments of APCA is half that of PAA because the endpoint of each segment is stored along with the mean value. The APCA is denoted as $\{[cv_1, cr_1], \dots, [cv_M, cr_M]\}$, where cv_i is the mean value of the data points in the i^{th} segment and cr_i is the right endpoint of the i^{th} segment. Even though the APCA has fewer segments than the PAA, it produces a lower representation error than the PAA.

In [55], two methods for matching PLA are proposed. In the first method, as illustrated in Fig. 4.3(a), every endpoint on each sequence is projected onto the other sequence. The variance of the length of the projected lines is measured and used as the distance measure. Two similar sequences will produce a small variance, and vice versa. The second method utilizes DTW. The comparison is performed on the endpoints of two sequences as displayed in Fig. 4.3(b). From a comparison of the two methods in [55], the second method gives much higher accuracy than the first method, while the computational time for the second method is comparable to that of the first method.

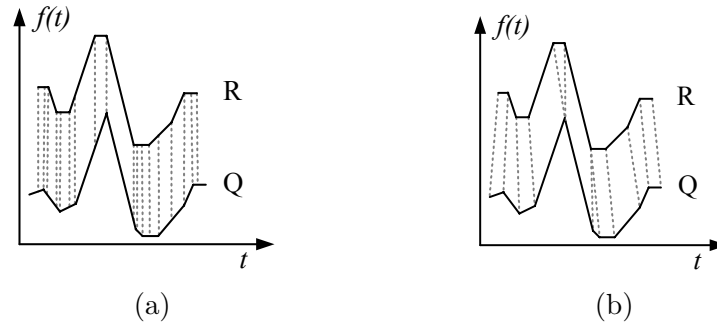


Figure 4.3: Similarity measures for PLA; R is a signal and Q is a query signal. (a) using Minkowski metrics, (b) using segmented dynamic time warping. The picture (a) and (b) are from [63] and [55], respectively.

The basic idea of similarity search is to slide the query waveform along the data signal while comparing the windowed signal to the query waveform. The window size may differ from the length of query waveform to allow unequal-length matching, yet the distance measure used must permit unequal-length comparison [63, 55].

4.2.5 ECG beat detection

In [65] and [43], low computation and high speed algorithms were presented for ECG beat detection. For the algorithm in [65], the slope of ECG signal is measured using a digital filter. QRS complexes are detected when the signal slope exceeds the adaptive threshold level. The algorithm can discriminate abnormal and distorted QRS complex from noise. The algorithm in [43] uses a filter to remove noise in the ECG and an adaptive threshold level to detect QRS complexes. The threshold level changes according to ECG baseline wandering due to breathing. Therefore, the algorithm can detect ECG beats when the ECG baseline fluctuates. The algorithm in [43] was used in the previous work explained in Section 4.1 and showed good performance in ECG beat detection. Therefore, it will be used in this research. However, its problems will be addressed and a solution will be presented.

4.2.6 Noise estimation

Surface ECG leads contain various kinds of noise, such as interference from the environment, muscle noise, movement artifacts, EMG noise. Noise is the major source of error for ECG analysis. Moody and Mark proposed a method for noise estimation using the Karhunen-Loeve transformation [66]. QRS complexes are represented using the five-term Karhunen-Loeve expansion. Noise contaminated in the ECG signal can be estimated from the residual error which exceeds the detecting level.

Gritzali, *et al.* estimated noise in ECG signals based on the distribution of the frequencies of the slopes in an ECG waveform [67]. The existence of noise magnitude in a test ECG was measured by applying the Chi-squared test. The test compared the distribution from the test ECG to the distribution from the ECG without noise. The algorithm required two cardiac cycles for noise detection. No QRS complex detection was needed. The method was evaluated using real ECGs.

Brouse, *et al.* developed an algorithm for detecting electrocautery noise in the ECG using a wavelet decomposition of the signal [68]. Electrocautery is an instrument for directing a high-frequency current to a specific area of the body in order to stop bleeding from small vessels or cut through soft tissue. An ECG block (segment) was decomposed to a single level using a Daubechies-10 analyzing wavelet. The coefficients with amplitude smaller than the threshold were set to zero. Groups of nonzero coefficients were joined together recursively, if the space between adjacent groups was within a user-specified range. Any group which was longer than a user-specified number of coefficients is categorized as noise. The approach was tested on ECG data from 15 operations spanning 38.5 hours. The result false positive and false negative rates were 0.71% and 0.33%, respectively.

4.2.7 Abnormal ECG classification

Many methods for automatic detection and classification of abnormal ECG have been described in the literature as follows. Moody and Mark presented a method for QRS morphology representation using the Karhunen-Loeve transformation [66]. QRS complexes from normal ECG and PVC are represented using the five-term Karhunen-Loeve expansion. By testing on the MIT-BIH and AHA databases, the method achieved a PVC sensitivity of 97.06% and positive predictivity of 97.56%.

Lin and Chang developed an algorithm for PVC detection [69]. The feature for recognition was derived from the residual error signal after processing the ECG sequence by Durbin's linear prediction. The optimal prediction order is 2 for signals with a sampling rate less than 1000 Hz. A nonlinear transformation was applied to transform the residual error signal to a three-state pulse-code train. If the cross-correlation coefficient obtained from a test pulse-code train and a template pulse-code train representing a normal ECG exceeds the threshold, a PVC is found. The test set contained 2857 PVC beats from the MIT/BIH arrhythmia database. The algorithm yielded a 95.3% sensitivity for detection, and 8.01% rate for incorrect detection.

Novel features, which were extracted using a filter bank, for PVC detection were introduced by Wieben, *et al.* [70]. The filter bank contained seven sub-bands. The feature set was derived from the energies in the sub-bands. Two distinguishing classifiers were implemented, a decision tree and a fuzzy rule-based system. The decision tree gave a sensitivity of 85.3% and a positive predictivity of 85.2%, while the fuzzy rule-based system resulted in 81.3% sensitivity and 80.6% positive predictivity.

Cheng, *et al.* presented an automatic ECG classification (normal ECG, PVC, and fusion of ventricular and normal beats) using hidden Markov models (HMM) [71]. ECG signals were approximated by line segments. The features were the differential and normalized

intervals of each line segment. HMMs were separately used for training for each type of ECG. The resulting classification rate is 93% for normal ECG, 65.55% for PVC, and 56.38% for fusion of ventricular and normal beats.

In [72], a neural-network-based classification of PVCs using wavelet transformation and timing interval features was described. Forty ECG records were selected from the MIT/BIH database, focusing on modified-lead II signals. The selected signals consisted of normal, left bundle branch block, right bundle branch block, PVC, atrial flutter, and paced beats. Six feature sets were constructed based on the original ECG data and its dyadic wavelet transformation at the first five scales, in conjunction with timing information (RR intervals). A feed-forward multi-layer perceptron neural network with a single hidden layer was implemented for classification. Inan, *et al.* found that the feature set, which contains the fourth scale of a dyadic wavelet transformation with a quadratic spline wavelet and the pre/post RR-interval ratio, is very effective in discriminating normal and PVC from other beats. Accuracy of 95.16% was achieved over 93,281 beats from all 40 records.

An automated system for detection of transient ischemic and heart rate-related ST-segment episodes was investigated by Smrdel, *et al.* [73]. A sequence of average heart-beats was constructed and used for searching for the position of the iso-electric level in the P-Q interval and for the position of the J point. An ST level function was derived from the ST-segment amplitude measured on average with respect to its iso-electric level. A ST-segment reference function was constructed using the first order Mahalanobis functions and the Karhunen-Loeve transformation. A ST deviation function, which was derived from the ST level function and the ST-segment reference function, was used to detect ST-segment episodes according to the detection threshold. The test result gave a detection sensitivity of 81.3% and positive predictivity of 89.2%. The algorithm consists of many steps and requires many parameter settings.

Stamkopoulos, *et al.* performed ischemia classification on the European ST-T database using a nonlinear-principal-component-analysis (NLPCA) neural model [74]. Features were extracted by nonlinear principal component analysis (two principal components). The neural model was used for training using the radial-basis function network, and only normal ECGs were used in the training process. A threshold was determined for ischemia detection from the output of the neural model. The method accomplished correct classification rates higher than 90%.

Jeong and Yu developed an algorithm for detecting the ST level change and classifying the ST shape type using the polynomial approximation [75]. The algorithm divided ECG data into small groups, where the duration of each group was 0.04 second. The Q and S waves were searched for within these groups. Then, the ST segment was approximated into a polynomial formula. Either one of the two following methods was chosen according to the magnitude of noise. The first is a polynomial formula of the 9th order over the whole ST. The second method is three polynomial formulae of under the 5th order for the three-segmented ST. The slope values of four selected points between the S and T waves were measured. The algorithm determines ST type by comparing the slope values to the slopes of the reference ST type.

4.3 Discussion

In previous studies on analyzing massive ECG data, RR intervals, HRV parameters, and ECG templates were extracted from ECG data. Methods to summarize and visualize massive data were introduced in [6] and [7] to view the changes in the entire range of data. One, two, and three dimension plots were proposed to display the features in one graph. However, the algorithms for data processing were developed separately and an integrated

procedure was not established. In additions, the results needed to be analyzed manually and an automatic procedure for finding abnormal ECG has not yet been implemented. In [8], the work emphasized the development of parallel algorithms and no algorithm for identified ECG types was presented.

The time series representations reduce the dimension of data. However, distance measures for the representations are quite complex. The time series representations do not gain any benefit if it decreases only a small number of dimensions of data.

Several approaches for noise estimation and detection include algorithms based on the Karhunen-Loeve transformation [66], the distribution of the frequencies of the slopes in an ECG waveform [67], and wavelet decomposition [68]. The algorithm in [66] measures noise only in the QRS section. The method in [68] can detect noise within a single ECG beat, while the algorithm in [67] detects noise in an ECG passage and requires two ECG cycles in the detection process. [66] and [67] do not provide a classification rate for noise detection. [68] achieved very high true positive and true negative rates.

Automated algorithms for normal ECG and PVC classification in the literature review use approaches based on the Karhunen-Loeve transformation [66], linear prediction [69], filter banks [70], hidden Markov models [71], wavelet transformation, and neural network techniques [72]. The results from the Karhunen-Loeve transformation, linear prediction, wavelet transformation, and neural network techniques show high performance in the classification. The linear prediction approach needs to determine the optimal linear predictive coefficients using methods which are complex [76, 77]. Algorithms for detecting ST level change previously utilized techniques as follows: the Karhunen-Loeve transformation [73], a nonlinear-principal-component-analysis (NLPCA) neural model [74], and polynomial approximation [75]. The polynomial approximation approach contains many steps to find the slopes of ST segment, the least squares curve for samples between the S wave and

T wave. Nevertheless, the algorithm can recognize seven types of ST level change. The Karhunen-Loeve expansion is used to express any observation as a linear combination of the basis functions. The basis functions are determined by using the singular value decomposition (SVD) technique. The SVD technique requires extremely high memory space for determining the basis functions, especially when many training signals are included [78]. Wavelet transformation is an intensive computational technique [79]. It is not suitable for analyzing large amounts of data recordings.

Chapter 5

Data processing and feature extraction

5.1 Overview

The experimental method explained in Chapter 3 describes a model of chronic heart failure in rats. During the progression of the disease, the ECG was collected nonstop for approximately 12 weeks. From this experiment, approximately 16 GB of ECG data were obtained for each rat. The sheer mass of data requires techniques to automatically process the data. This chapter describes techniques developed to analyze the data collected from the rats as heart failure developed. Abnormalities such as arrhythmias and abnormal heart rate variability are expected to be seen in the deteriorating heart. This section explains an overview of the entire process of data analysis including data preprocessing and processing and abnormality search methods. Fig. 5.1 depicts the relation between sections of data analysis. Each block is described briefly below and the details are explained in the following sections. This research involved data preparation, data preprocessing, feature extraction template generation, HRV analysis and visualization, and abnormality search.

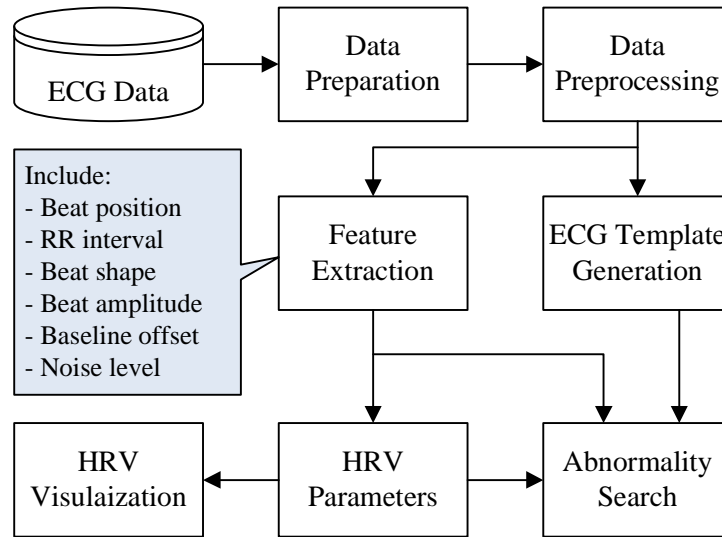


Figure 5.1: An overview of the data processing and analysis method

1. Data preparation: During the experiment, the data is collected and recorded by a software from Data Science International (DSI) in the DSI format [80]. The data which consist of a continuous stream of ECG are recorded into separate files. Each file contains approximately one hour of ECG data and information about the data (called *metadata*) such as subject ID and time of recording. In data preparation, metadata of the recorded files are scanned and structurally stored in a database for data allocation. A section of ECG at a particular time can be retrieved by querying its location, which includes DSI filename and position in the file, from the database and then extracting from the DSI file. Besides creating data allocation tables, this section also scans through all of the recording data and divides them into usable and unusable data. To provide a better understanding, necessary technical information of DSI format is given. Recording data is coded by using a two-byte signed integer where $(2^{15} - 1)$ represents the highest amplitude and $-(2^{15} - 1)$ is used for the lowest amplitude. $\frac{V_{max}}{2^{15}-1}$ is a coefficient used to convert coded value to actual value

in millivolts where V_{max} is the setting value of maximum voltage magnitude. 2^{15} stands for out of recording range. A recording section of 2^{15} is called *fullscale signal*. The DSI data is categorized into three groups including raw ECG signal, no signal which is caused by either unsuccessful communication between the transmitter and the receiver or temporary turn-off of the recording unit, and fullscale signal. The raw ECG signal is usable data and sent to the data preprocessing section, while no signal and fullscale signal are unusable data.

2. Data preprocessing: Some sections of the raw ECG signal are corrupted by noise from the recording such as environmental noise and transmission signal disruption and from the animal such as electromyogram (EMG) noise and baseline fluctuations. The types of noise are further explained in Section 5.2. The noise is detected and minimized in this part by using the algorithm described in Sections 5.5 and 5.6. At this point, the raw ECG is further categorized into EMG, cleaned ECG, and short ECG. ECG with an acceptable noise level (by observation) is considered as cleaned ECG. Short ECG is defined as passages of cleaned ECG which are too short to use in analysis – i.e. shorter than seven seconds. The database for signal allocation is then updated according to the new categories of ECG. The categories include (1) cleaned ECG, (2) short ECG, (3) EMG, (4) fullscale signal, and (5) no signal. The cleaned ECG will be called only “ECG” after the data preprocessing section.
3. ECG template generation: An ECG template is a representation of ECG morphology in the data. Its size is fixed and covers considering ECG waves such as QRS complexes. By sequentially scanning through the data, an ECG template is created for any distinct ECG beat shape. The templates are used to find the deviation in the ECG morphology and arrhythmias. The details are explained in Section 5.8

4. Feature extraction: The features are extracted from the cleaned ECG and include beat position, RR interval, beat shape, beat norm, baseline offset, and noise level. The beat position is a series of time stamps specifying the locations of the R-wave of ECG beats in the data. It is computed by using an algorithm explained in Section 5.4. The RR intervals are derived from the beat positions as a distance between consecutive R-waves. The beat shape is the morphology of cropped ECG. The beat norm is the Euclidean norm of windowed ECG. The baseline offset is the mean or median amplitude of beat samples. The noise level is measured in the EMG detection part. The RR interval is used in HRV parameter calculation, and all of features are submitted for abnormality search.
5. HRV parameter calculation: The following HRV parameters are calculated for every specified interval (90 seconds) of RR intervals: mean, median, interquartile range, coefficient of variance, standard deviation, high frequency power, low frequency power, ratio of low over high frequency power. The parameters computed from successive differences between RR intervals include mean, standard deviation, and interquartile range. The details and formulae are described in Section 2.5.
6. HRV visualization: The HRV parameters are displayed by an HRV visualization tool to observe the change. The HRV visualization is in Section 5.10.
7. Abnormality search: The calculated features and ECG templates are used in searching for ECG signals which are deviated from normal ECG. This approach will be used to find abnormal ECG and arrhythmias, such as premature beats, ventricular tachycardia, and ventricular fibrillation. In Chapter 6, algorithms are developed to search for premature beats, elevated ST segments, and split R-waves.

The following sections are ordered as follows: Types of noise in the recording data are explained in Section 5.2 to provide a background for baseline fluctuation removal and electromyogram detection. Morphological filter, which is described in Section 5.3, will be applied in the beat detection and baseline fluctuation removal sections. Beat detection, baseline fluctuation removal, and electromyogram detection (Sections 5.4, 5.5, and 5.6, respectively) are algorithms used in data preprocessing (Section 5.7). ECG template generation is in Section 5.8. Sections 5.9 and 5.10 explain results from feature extraction and heart rate variability visualization, respectively.

5.2 Types of noise

In the recording process, noise and artifact are introduced in the data. The noise comes from several sources including the environment, the radio transmitter-receiver recording device, and the animal itself. This section explains the types of noise that corrupt the data and need to be eliminated. To provide a comparison, a clean recorded ECG is illustrated in Fig. 5.2.

1. Burst noise: Fig. 5.3 depicts an example of burst noise. The source of burst noise is unknown. However, the cause of burst noise is suspected to be bursts of electromagnetic interference from the environment.
2. Fullscale noise: Fig. 5.4 depicts an example of fullscale noise. The DSI software cannot record the monitored signal over a specific range of amplitude. The default setting range for ECG is -10mV to 10mV . Nevertheless, this range can be configured to suit the best recording. The signal with an amplitude outside the setting amplitude range is stored as the upper-limit amplitude or fullscale. For example, in the default setting, an out-of-range signal is stored as 10mV or 32767 , the maximum positive

number of the 16-bit signed integer. The origin of the fullscale noise includes high-amplitude ECG, sustained burst noise, and rapid movement such as jumping.

3. Motion artifact: Motion artifact occurs mostly during position change and limb movements. Movement of muscle(s) in contact with the electrode(s) may create a shift up or down in the baseline and return to the previous baseline level as shown in Fig 5.5. Other shapes of motion artifact include a rapid upstroke of signal and random-fashion baseline wandering. Motion artifact is usually the most troublesome type of noise for arrhythmias.
4. EMG noise: An electromyogram is a composite signal of muscle fiber action potentials occurring in the muscle tissue. It is characterized as a rapid fluctuation which contains a wide range of frequencies, as shown in Fig. 5.6.
5. Baseline fluctuation: The ECG baseline fluctuates due to the animal's breathing and slow movement. It is composed of low frequency components. A figure for baseline fluctuation is shown in Fig. 5.7

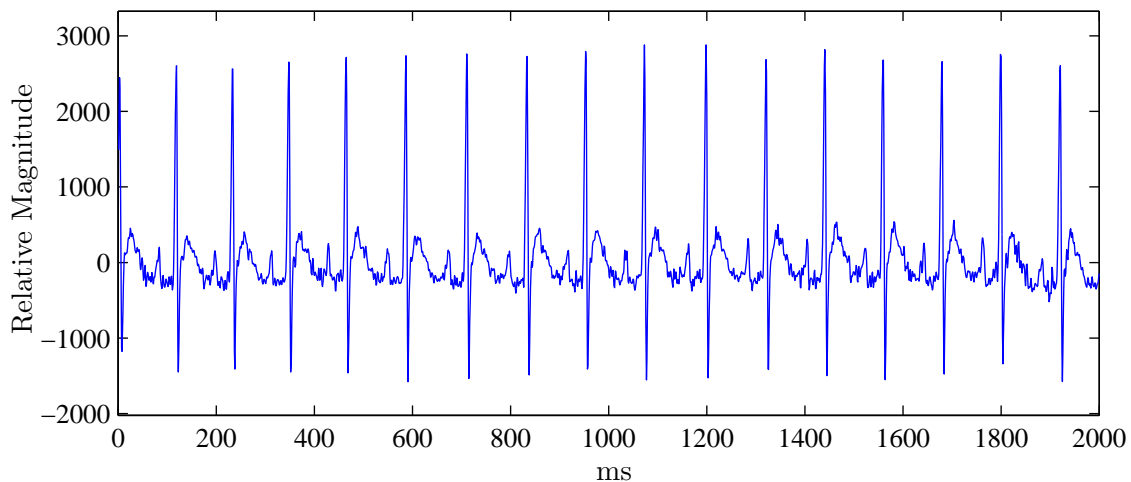


Figure 5.2: Example of a clean recorded ECG

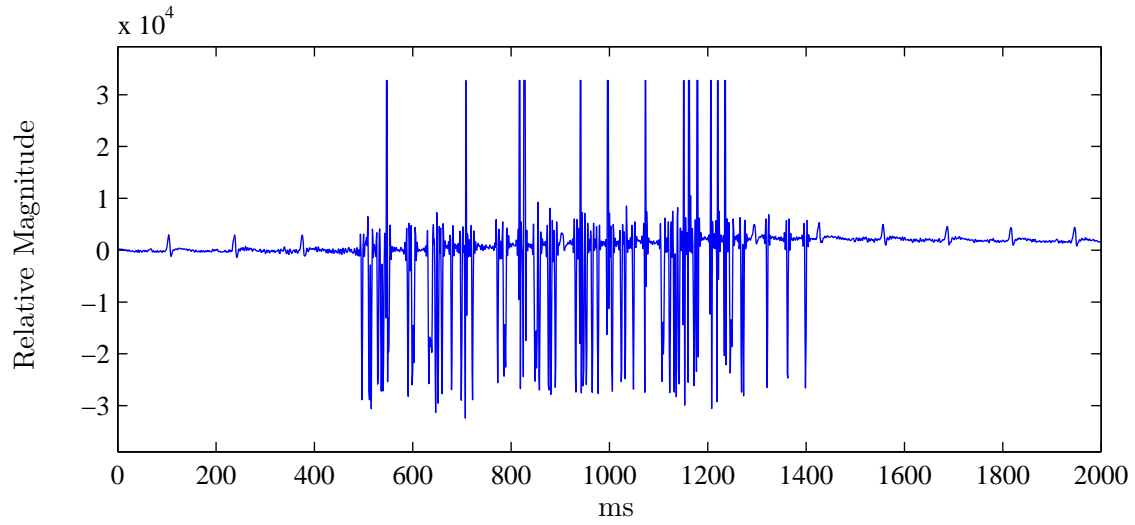


Figure 5.3: Burst noise

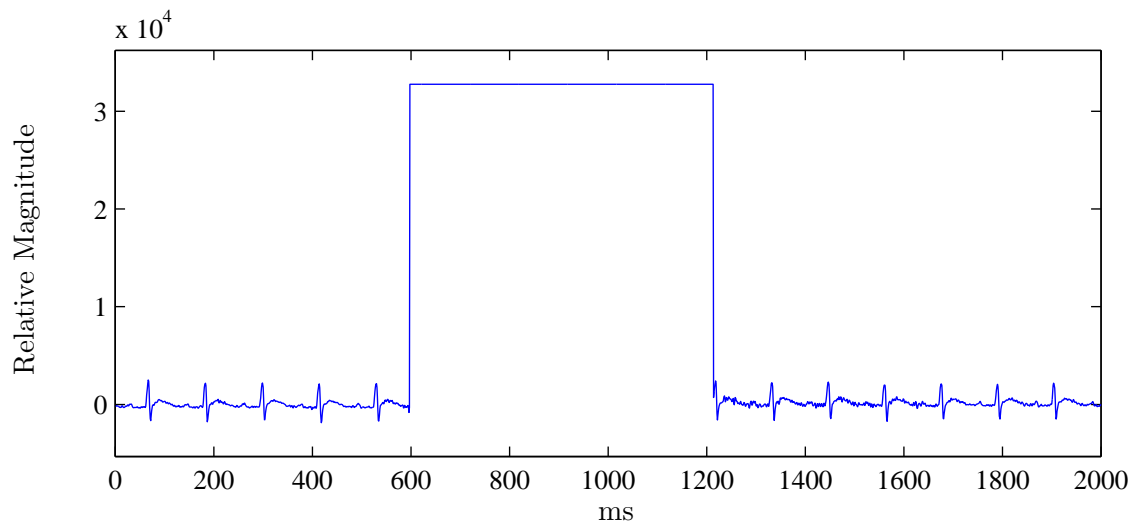


Figure 5.4: Fullscale noise

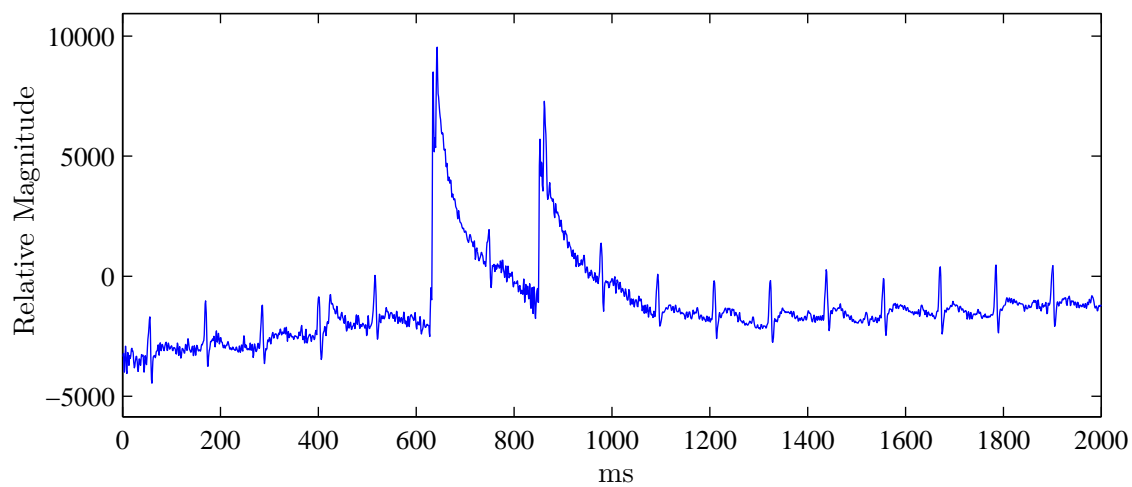


Figure 5.5: Movement artifact

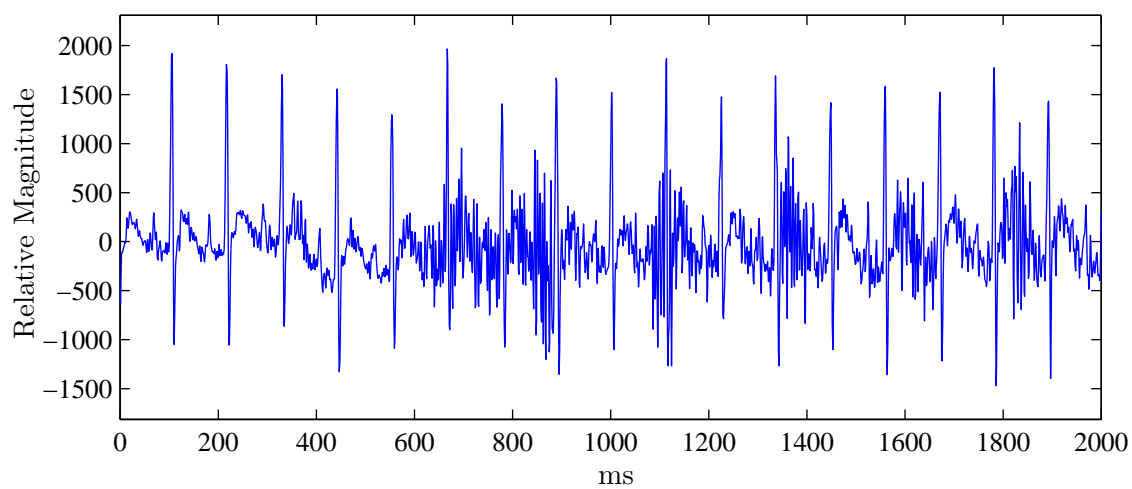


Figure 5.6: EMG noise

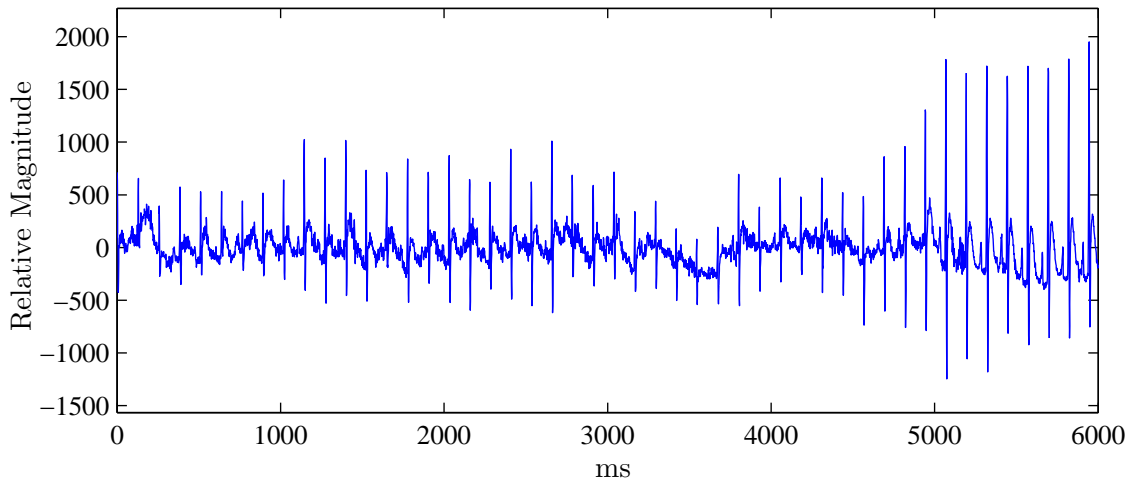


Figure 5.7: Baseline fluctuation

5.3 Morphological filter

Morphological filters have been widely used in the field of image processing. They have applications in filtering unwanted shapes of signals while not changing the other parts of the signal. The fundamental operators of the morphological filter are erosion and dilation which are applied in tandem to form opening and closing operators of the morphological filter. The shape information of the signal extracted is defined as a *structuring element*. It can be a square, triangular wave with arbitrary base width, or any arbitrary shape of signal. In one dimension, erosion and dilation operators are formulated as in (5.1) and (5.2) [81]. x and s denote signal and structuring elements and are defined on $X = \{1, 2, \dots, L\}$ and $S = \{1, 2, \dots, B\}$ where $L > B$, respectively. Erosion is a “shrinking” operator. The resultant signal, $x \ominus s$, has a magnitude of minimum difference between the input signal and the translating structuring element. Therefore, $x \ominus s$ always has a smaller magnitude than x . In contrast, dilation is an “expansion” operator and is denoted by $x \oplus s$ and the output is greater than the input [81]. The complexity of erosion and dilation is $O(LB)$.

$$(x \ominus s)(i) = \min_{j=1, \dots, B} x(i+j) - s(i) \quad \text{for } i = 1, \dots, L - B + 1 \quad (5.1)$$

$$(x \oplus s)(i) = \max_{j=i-B+1, \dots, i} x(j) + s(i-j) \quad \text{for } i = B, B+1, \dots, L \quad (5.2)$$

Opening and closing are two operators defined based on erosion and dilation operators. $x \circ s$ and $x \bullet s$ denote opening and closing and are described in (5.3) and (5.4), respectively. A morphological filter is formed by two cascades of opening and closing, *i.e.* opening followed by closing and closing followed by opening. The output is an average of the terminal signals of these cascades. Its mathematical expression and diagram are displayed in (5.5) and Fig. 5.8, respectively.

$$x \circ s = (x \ominus s) \oplus s \quad (5.3)$$

$$x \bullet s = (x \oplus s) \ominus s \quad (5.4)$$

$$y = (x \circ s \bullet s + x \bullet s \circ s) / 2 \quad (5.5)$$

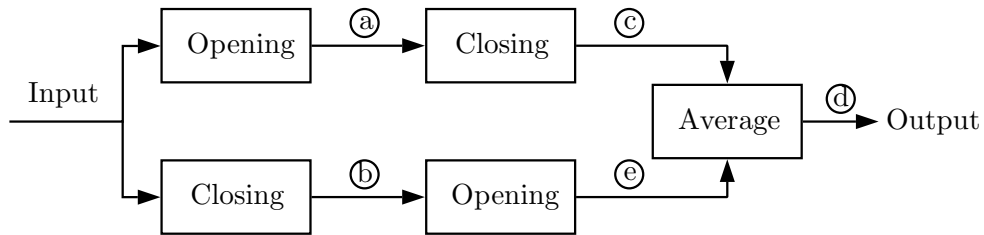


Figure 5.8: Block diagram of morphological filter

Fig. 5.9 depicts an ECG signal as an input signal and signals from different locations in morphological filter. The structuring element is a square wave with a width equal to a base width of R and S waves, seven samples. Therefore, the morphological filter eliminates

QRS wave and small fluctuations which have a base width equal to or smaller than seven samples (0.007 second). The opening operator removes upright spikes while the closing operator purges downward deflections. Fig. 5.9 (d) shows the output signal where QRS waves and small fluctuations are removed.

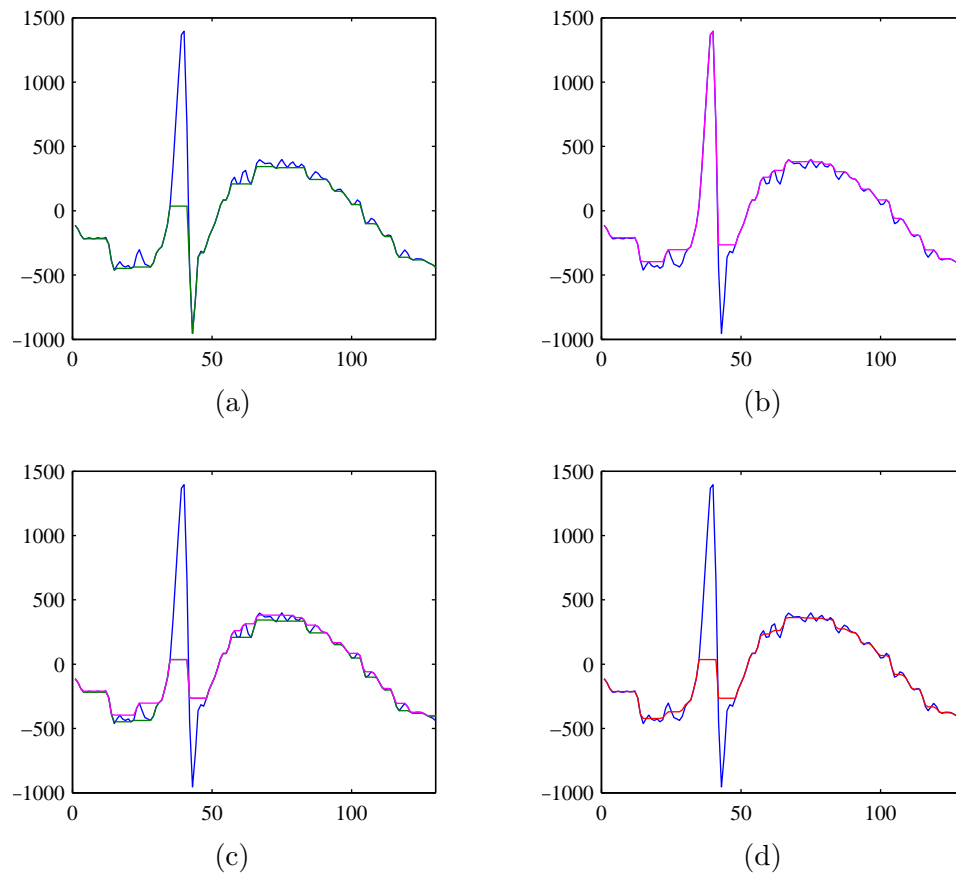


Figure 5.9: Plots of input signal (blue) and signal at different locations (green, magenta, and red) in morphological filter; (a) signal at point *a* in Fig. 5.8; (b) signal at point *b*; (c) Green is signal at point *c* and magenta is point *d*.; (d) output signal

A morphological filter is composed of eight erosion and dilation operators. Therefore, the algorithm complexity is high. In [82], a fast algorithm for the morphological filter with a unit square-wave structuring element is presented. The algorithm is constructed based on the properties expressed below. Its pseudo code for opening and closing is presented in

Tables. 5.1 and 5.2. Note that the closing algorithm is similar to the opening except for some few lines and Table. 5.2 displays only those different lines.

Property 1: $x \circ s(j) = x(j)$ if and only if there exists a point $k \in \bar{W}_s + j$ such that $x(j) = \min\{x(z) \mid z \in W_s + k\}$, where W_s is $\{1, 2, \dots, B\}$, $W_s + j$ is the translation of W_s by point j , and \bar{W}_s is the 180° rotation of W_s about the origin.

Property 2: $x \bullet s(j) = x(j)$ if and only if there exists a point $k \in W_s + j$ such that $x(j) = \max\{x(z) \mid z \in \bar{W}_s + k\}$

Property 3: If $x(m_1)$ and $x(m_2)$ with $m_1 < m_2$ are two successive invariant samples of $x \circ s(j)$, then $x \circ s(j) = \max\{x(m_1), x(m_2)\}$ for $m_1 < j < m_2$, whereas if $x(n_1)$ and $x(n_2)$ with $n_1 < n_2$ are two successive invariant samples of $x \bullet s(j)$, then $x \bullet s(j) = \min\{x(n_1), x(n_2)\}$ for $n_1 < j < n_2$.

Table 5.1: Opening morphological filter algorithm

1: % Initialization	17: $i_{newsample} = \max(\mathbf{i}_{window});$
2: $[x_{min}, \mathbf{i}_{invariant}] = \min(x(1 : B));$	18: if $x(i_{newsample}) \leq x(i_{rightmost})$
3: $i_{rightmost} = \max(\mathbf{i}_{invariant});$	19: $\mathbf{i}_{invariant} = [\mathbf{i}_{invariant}, i_{newsample}];$
4: $i_{leftmost} = \min(\mathbf{i}_{invariant});$	20: $N_{new} = 1;$
5: $y(1 : i_{rightmost}) = x_{min};$	21: end
	22: end
6: % Slide window to the right one step	23: if $N_{new} > 0$
% at a time	24: $i_{rightmost} = \max(\mathbf{i}_{invariant});$
7: for $k = 2 : L - B$	25: $i_{leftmost} = \min(\mathbf{i}_{invariant});$
8: $N_{new} = 0;$	26: $y(i_{leftmost} : i_{rightmost}) = \dots$
9: $\mathbf{i}_{window} = (k - 1) + (1 : B);$	$x(i_{leftmost});$
10: $i_{k-1} = \max(\mathbf{i}_{invariant} < k)$	27: $y((i_{k-1} + 1) : (i_{leftmost} - 1)) = \dots$
11: $\mathbf{i}_{invariant} = (\mathbf{i}_{invariant} \geq k)$	$\max(x(i_{k-1}), x(i_{leftmost}));$
12: if $(i_{rightmost} < k)$	28: end
13: $[x_{min}, \mathbf{i}_{invariant}] = \dots$	29: end
$\min(x(\mathbf{i}_{window}));$	30: $i_{rightmost} = \max(\mathbf{i}_{invariant});$
14: $N_{new} = \text{length}(\mathbf{i}_{invariant});$	31: $y(i_{rightmost} : L) = x(i_{rightmost});$
15: else	
16: $i_{rightmost} = \max(\mathbf{i}_{invariant});$	

Table 5.2: Closing morphological filter algorithm

2: $[x_{max}, \mathbf{i}_{invariant}] = \max(x(1 : B));$ 5: $y(1 : i_{rightmost}) = x_{max};$ 13: $[x_{max}, \mathbf{i}_{invariant}] = \max(x(\mathbf{i}_{window}));$ 18: if $x(i_{newsample}) \geq x(i_{rightmost})$ 27: $y((i_{k-1} + 1) : (i_{leftmost} - 1)) = \min(x(i_{k-1}), x(i_{leftmost}));$
--

5.4 Beat detection

Determining ECG beat location is crucial in preprocessing and ECG analysis. It is involved in and is required as a first step in several processes including ECG baseline removal, EMG noise elimination, template matching, heart rate variability analysis, and abnormality detection. [43] presented an ECG beat detection program or trigger program which performs online detection and requires low memory and computation. The algorithm diagram is shown in Fig. 5.10. The key of this algorithm is an exponentially decayed threshold, v_i , which decays at rate of c per sample. i locates a sample position. This threshold is compared with the absolute value of the input signal filtered by a 20-60 Hz bandpass filter, $|x_i|$. Once the threshold exceeds $|x_i|$, the algorithm labels that position as a beat. The lower bound of the threshold is defined as a beat amplitude multiplied by a coefficient b . This lower bound allows the algorithm to detect the next beats. b is a positive number and determines sensitivity of beat detection. The lower b is, the higher the beat detection sensitivity is. To avoid detecting a beat more than once, an additional condition is defined as detected beats must be after their previous beat by at least a certain time interval, w . Default w is 50ms.

The algorithm has drawbacks including that it is sensitive to severe baseline fluctuation which is present in surface ECGs and misses beats following an abrupt increase in beat

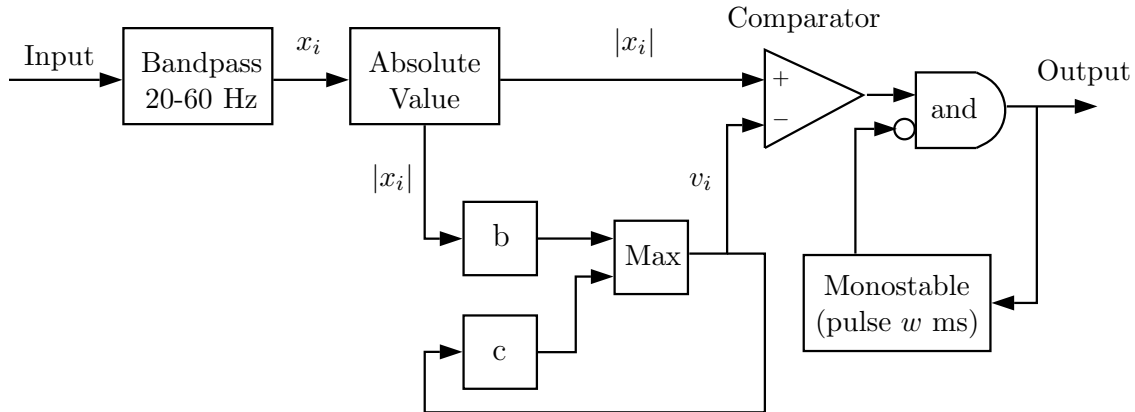


Figure 5.10: Block diagram of beat detection algorithm

magnitude such as premature beats or spurious noise. The location of beat is preferred at the R or S-wave for consistency of detection. Nevertheless, the algorithm confuses the T-wave with an R or S-wave in the case that the T-wave is as high as the R or S-wave. Therefore, the algorithm in [43] is adjusted for the above problems. The modified trigger program is presented in Fig. 5.11. Note that Fig. 5.11 is a flowchart version of Fig. 5.10 with some additive modified features. The modification includes as follows: QRS complexes are extracted by a morphological filter and used in beat detection instead of the filtered ECG. The fast morphological filter described in Section 5.3 is applied. The width of the structuring element is 15 ms to capture both normal and abnormal QRS complexes. In the flowchart, the QRS signal is denoted as y_i where i denotes sample position. The exponentially decayed threshold is v_i . The upper bound of the exponentially decayed threshold is added in the algorithm. It is equal to the amplitude of the previous beat multiplied by $(1 + g)$ where g is a growth rate and must be a positive number. The upper bound prevents the threshold from being abnormally high when detecting a beat with unusually high amplitude and the algorithm misses the following normal beats. b , c , and w are defined in the previous paragraph. The algorithm detects beats when $|y_i|$ exceeds v_i

and beats must depart from each other by at least w ms. The parameters are initiated to the following values before the algorithm starts: $v_{i-1} = 0$, $w = 50$, $b = 0.5$, $c = 0.9992$, and $g = 0.1$. The problems stated above are solved in this algorithm. The QRS signal from the morphological filter does not contain high T-waves (the top parts of T-waves may contain in the signal) or baseline fluctuations which can cause beat misdetection. The threshold cannot be too high after an abrupt increase in signal magnitude. As a result, the algorithm does not miss beats following such an increasing amplitude.

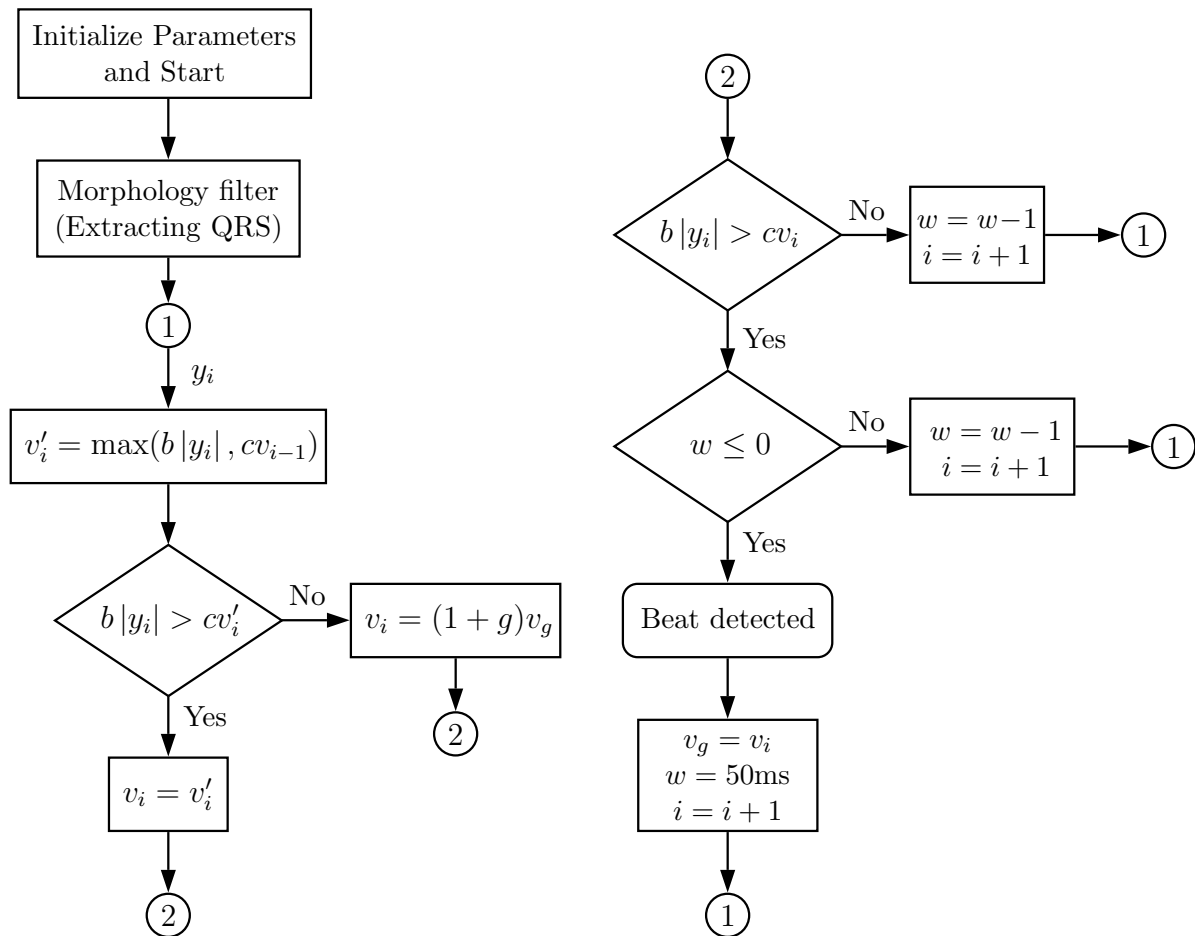


Figure 5.11: Flow chart of modified beat detection algorithm

The detected beat positions may be not located at the peak of the R-wave. Thus, the peak of the R-wave is determined by searching for the maximum in the vicinity of the

detected location and used as the beat reference point. In the case that the R-wave is very small and cannot be detected, the S-wave will be detected instead – i.e. the minimum point in the surrounding neighborhood is located and used as the beat reference point. The second case usually happens when abnormal beats are detected.

5.5 Baseline fluctuation removal

A low-complexity algorithm in [83] was applied for removing baseline fluctuation. The diagram of the algorithm is shown in 5.12. Two morphological filters are connected in sequence. The first removes QRS complexes while the second eliminates residual QRS complexes and P/T waves. The fast morphological filter algorithm explained in Section 5.3 is used. The structuring element is a unit square wave. The widths of the structuring element of the first (for R-wave) and second (for T-wave) morphological filters needed to be determined by this dissertation. The width of first structuring element is set to 7 ms. The output signal contains base parts of QRS complexes. However, it will be eliminated by the second morphological filter. Fig. 5.13 displays a test signal which is created by adding a generated baseline to a passage of ECG. The generated baseline has highest frequency component of 5 Hz. The output of the first morphological filter is shown in Figs. 5.14 (a) and (b). The fundamental frequency for a signal overriding the baseline is determined by its power spectral density. The peak frequency (approximately 10 Hz), shown in subfigure (c), is converted into a period and set as the width of the second structuring element. This width is approximately the length of the T-wave base. This is the method that this dissertation uses to determine the structuring element width of the second morphological filter which varies according to T-wave base width. The last step applies a moving average to smooth the step-like output from the morphological filters. The estimated baseline and

the actual baseline are compared in subfigure (d). Fig. 5.15 displays the result of the test ECG after having the baseline removed. The algorithm was also tested on the ECG in Fig. 5.7. The result is shown in Fig.5.16. The first morphological filter also acts as a noise filter, such as EMG, which helps increase the accuracy in calculating the width of the second structuring element. The limitations of this algorithm are that (1) the estimated baseline at the beginning and the end of the signal passage is inaccurate, and (2) baseline fluctuation with frequency components higher than the T-wave cannot be removed.

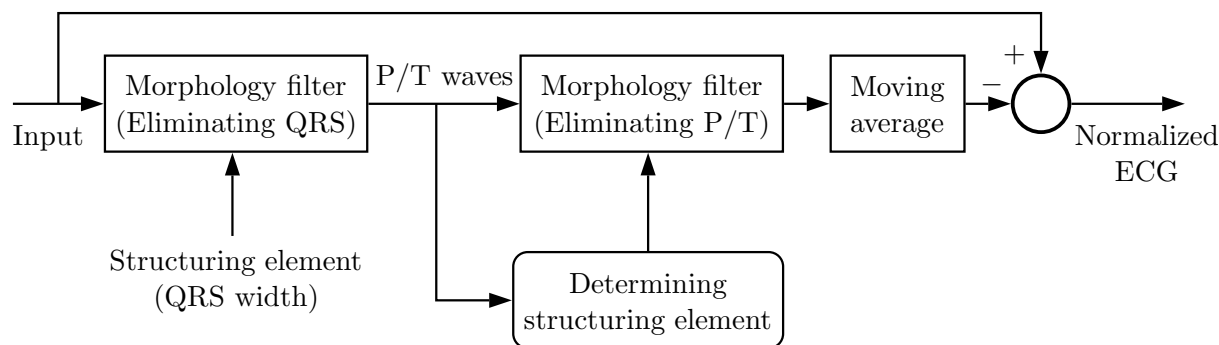


Figure 5.12: Block diagram of baseline removal algorithm

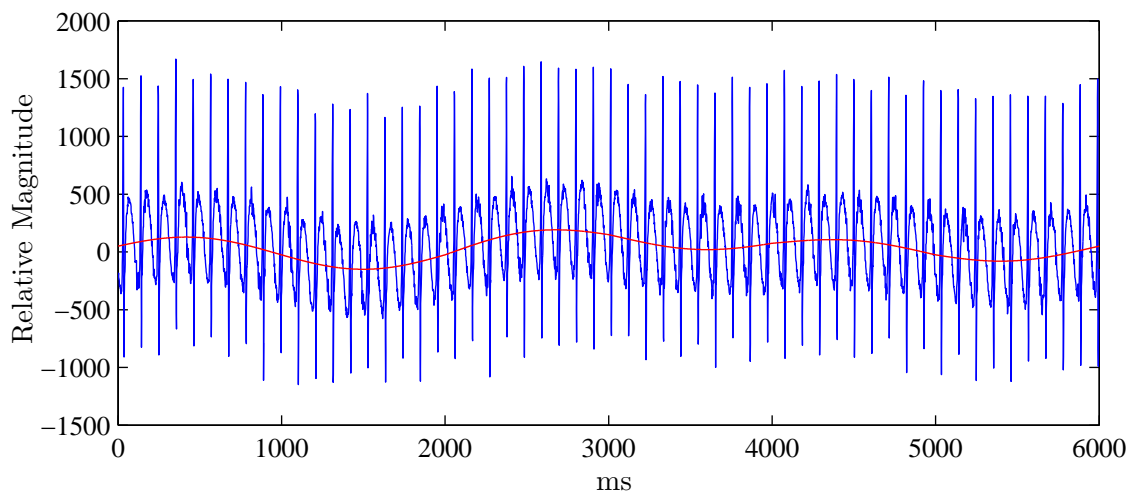


Figure 5.13: ECG with additive baseline (blue) and the additive baseline (red)

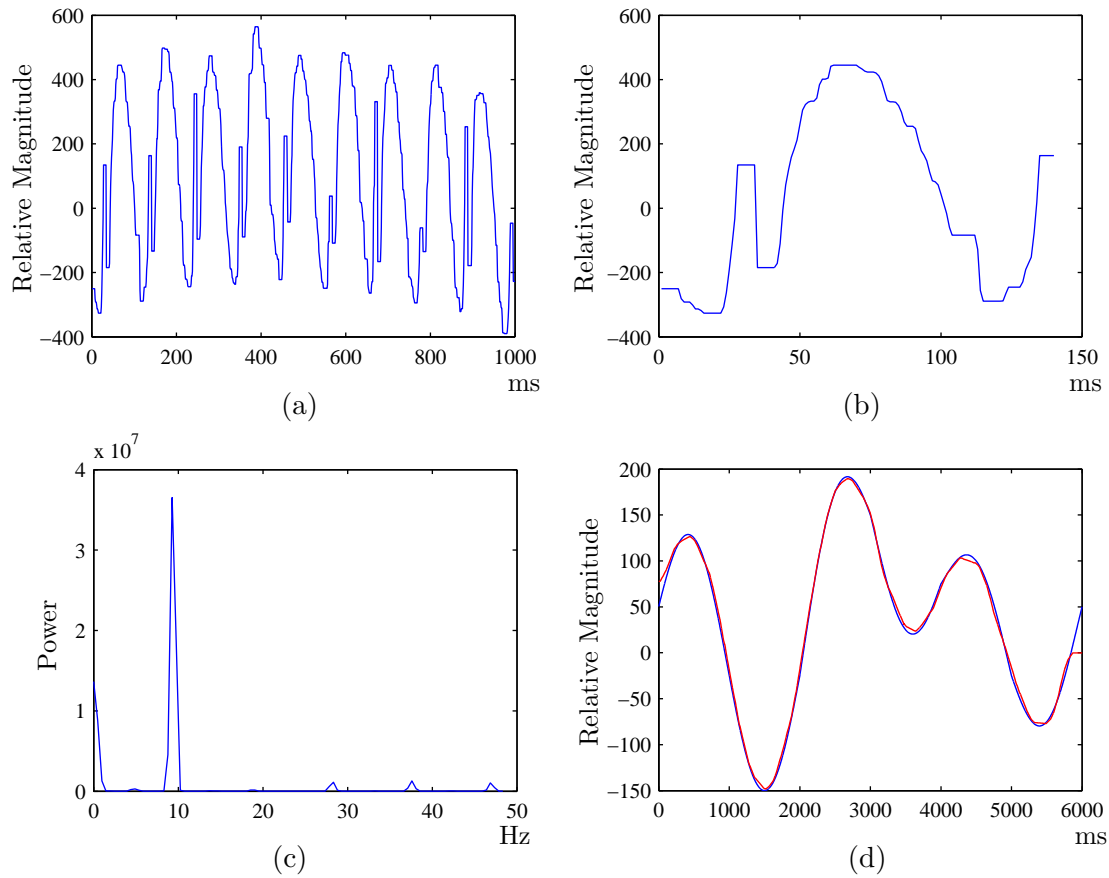


Figure 5.14: (a) Signal after filtering QRS complexes consists of partial QRS complexes, P waves and T waves. (b) Zoom-in of subfigure (a) (c) Power spectral density of signal in subfigure (a); The period ($1/\text{frequency}$) of the peak power is selected as a structuring element width for the morphological filter for eliminating residual QRS complexes and P/T waves (d) actual baseline (blue) V.S. calculated baseline (red)

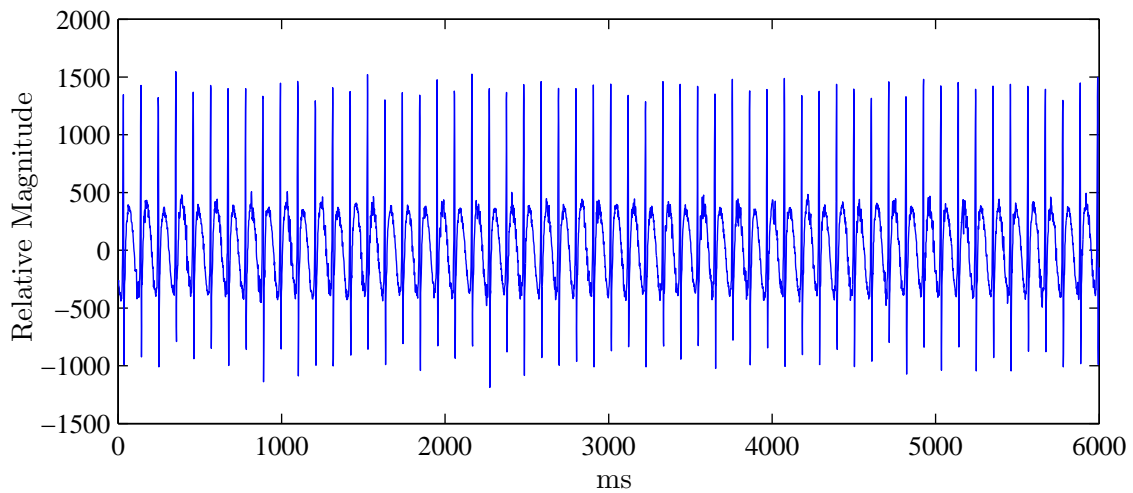


Figure 5.15: Baseline removed ECG

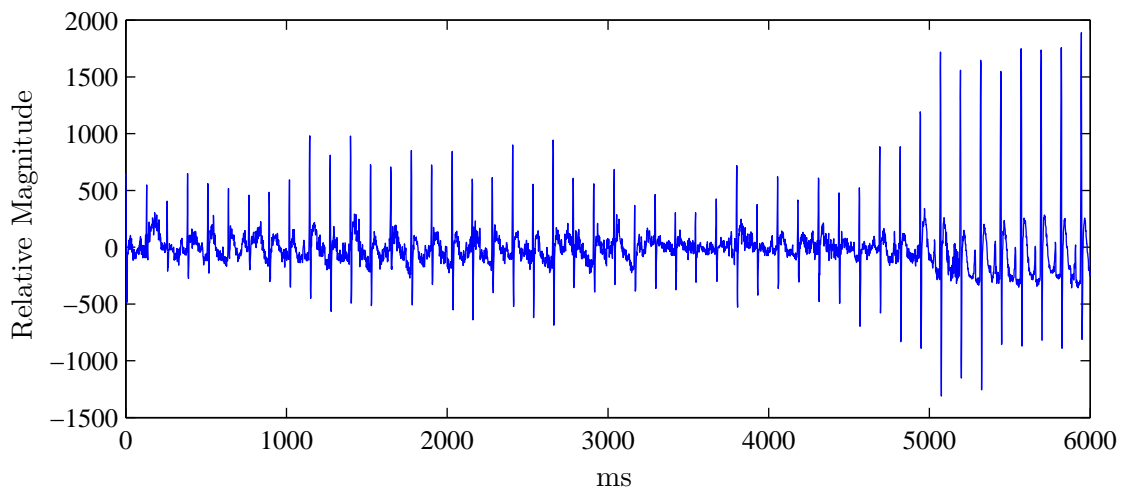


Figure 5.16: ECG from Fig. 5.7 after eliminating baseline

5.6 Electromyogram detection

The electromyogram (EMG) is generated from the electrical activity of the muscle. Sections of ECG may be interfered with and corrupted by surface EMG which causes difficulties in data processing and analysis. Motion artifact may also be present in the signal. Thus, EMG noise needs to be detected and filtered. In a recorded ECG, EMG

interference appears as rapid fluctuations which vary faster than ECG waves. An example of EMG in recorded ECG is shown in Fig. 5.6. EMG noise in the ECG can be detected by measuring the degree of signal fluctuation excluding the fluctuations of QRS complexes. Fig. 5.17 depicts a diagram of the EMG detection algorithm developed in this research. Fig. 5.18 depicts an ECG containing EMG noise used in the algorithm explanation.

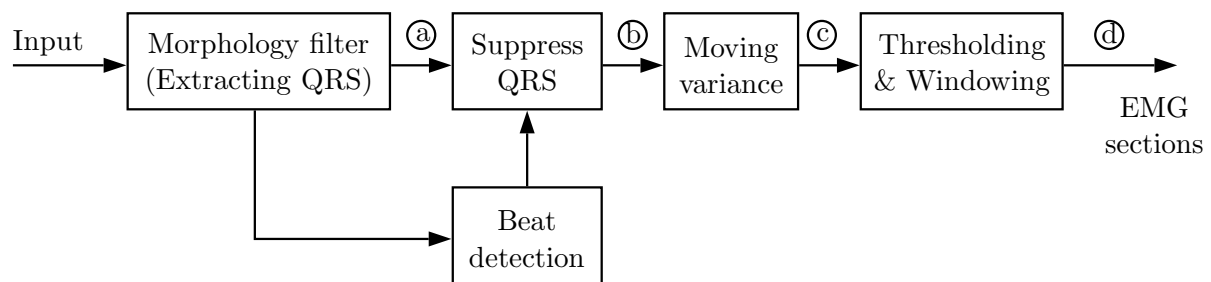


Figure 5.17: Diagram of EMG detection algorithm

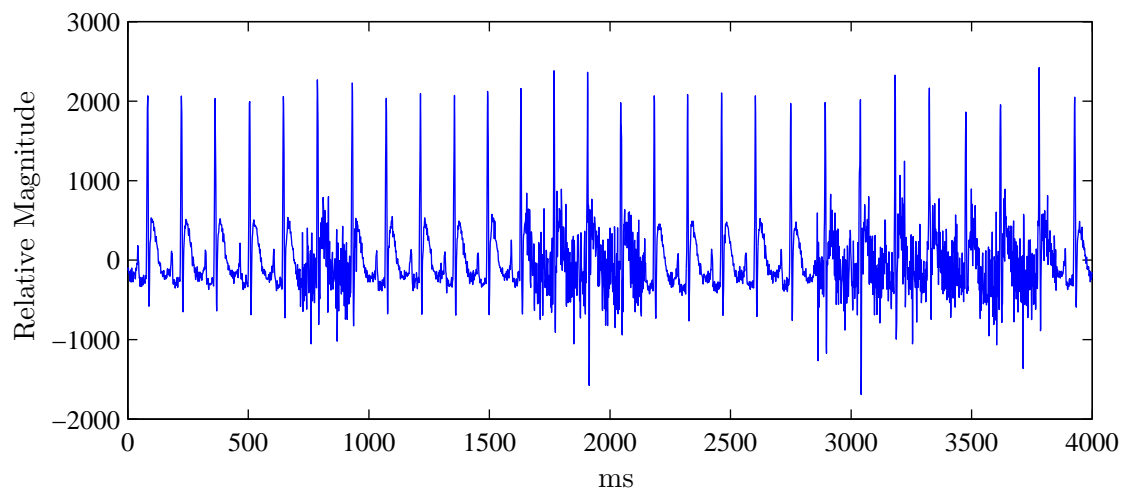


Figure 5.18: ECG signal which contains EMG noise

The algorithm is described as follows:

1. *Extracting EMG noise*: Impulsive noise, such as EMG, can be separated from ECG by using a morphological filter with a dome-like structuring element which is smaller than the ECG waves [81]. In Fig. 5.19, the output of the morphological filter results in an ECG with EMG noise suppressed. The EMG noise can be retrieved by subtracting the output of the morphological filter from the input ECG. Because we emphasize extracting EMG and do not need to preserve the quality of ECG waves, a square-wave structuring element produces a similar result. Thus, a fast morphological filter for a unit square-wave structuring in [82] can be applied and helps speed up the computation process. By experimentation, a square-wave structuring element with a width of 7 ms provides the best result. However, a portion of the QRS complex is also present in the output. Suppressing the QRS complexes is the next step.

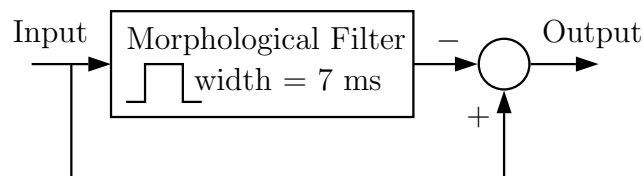


Figure 5.19: Morphological filter for extracting EMG noise

2. *Suppressing QRS*: The output from the EMG-extraction morphological filter contains EMG noise and partial QRS complexes which need to be suppressed (Fig. 5.20). By using the beat detection algorithm in Section 5.4, the QRS complexes can be located. The samples around the beat reference point (left 8 ms and right 11 ms) are reduced in their magnitude to one-tenth of their original size.
3. *Calculating moving variance*: Moving variance is denoted as signal variance within a sliding window. Degree of EMG noise can be measured from the level of signal

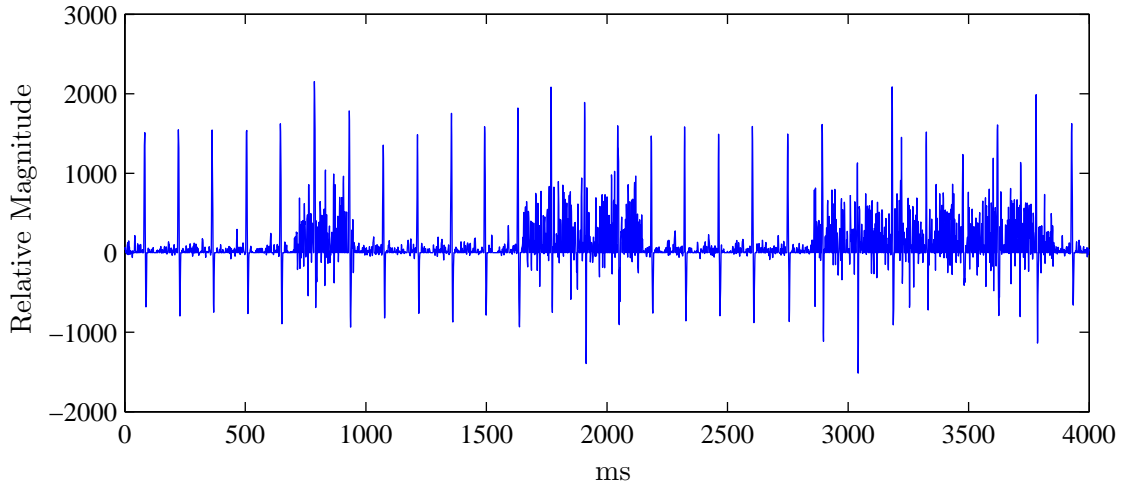


Figure 5.20: Extracted EMG noise and by-product QRS complexes

fluctuation which is associated with the moving variance of signal. A low computation moving variance is introduced and expressed in (5.6) to (5.8).

For $i = W_2 + 2, \dots, L - W_2$

$$m_1(i) = m_1(i - 1) - \frac{x(i - W_2 - 1) - x(i + W_2)}{W} \quad (5.6)$$

$$m_2(i) = m_2(i - 1) - \frac{x(i - W_2 - 1)^2 - x(i + W_2)^2}{W} \quad (5.7)$$

$$v(i) = m_2(i) - m_1(i)^2 \quad (5.8)$$

$v(i)$ is denoted as a moving variance of a window centered at sample i . The sliding window is stretched to the left and right by W_2 samples. Its total length is $W = 2W_2 + 1$. Input signal, x , has a length of L samples. $m_1(i)$ and $m_2(i)$ are the first and second moments, respectively, of a window centered at sample i . The exception is that (5.6) to (5.8) cannot calculate moving variance, first and second moments of first and last $W_2 + 1$ samples of signal passage. They need to be calculated by using the standard formulae. The size of the sliding window will be described in Section 5.6.1.

4. *Thresholding and windowing*: Sections of EMG noise can be identified by setting a threshold on the calculated moving variance. However, the magnitude of recorded signals is subjective and depends on recording settings and configurations. Therefore, moving variance is normalized by the square of the average of R-wave amplitudes of the ECG passage. Fig. 5.21 displays EMG noise and its normalized moving variance. Sections where the normalized moving variance exceeds the threshold are labeled as EMG noise. Nevertheless, moving variance may rise after the onset of EMG and fall prior to the end of EMG because it needs a sufficient amount of EMG in the sliding window in order to increase the variance. Thus, the detected EMG section is expanded to the left and right for 0.05 seconds. Fig. 5.22 shows EMG sections detected by the algorithm. A process to determine the threshold value will be explained in the following section.

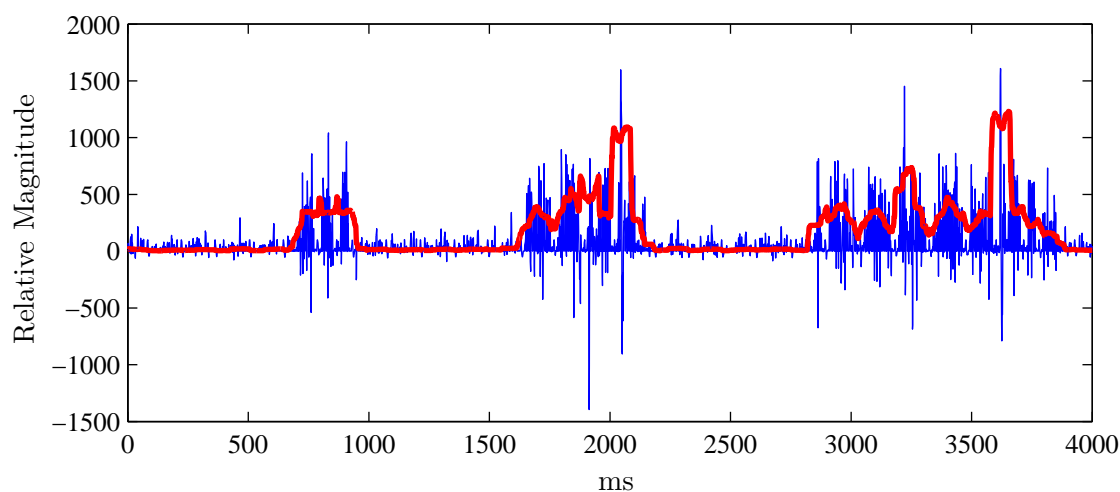


Figure 5.21: Extracted EMG noise (black) and normalized moving variance $\times 10^4$ (red)

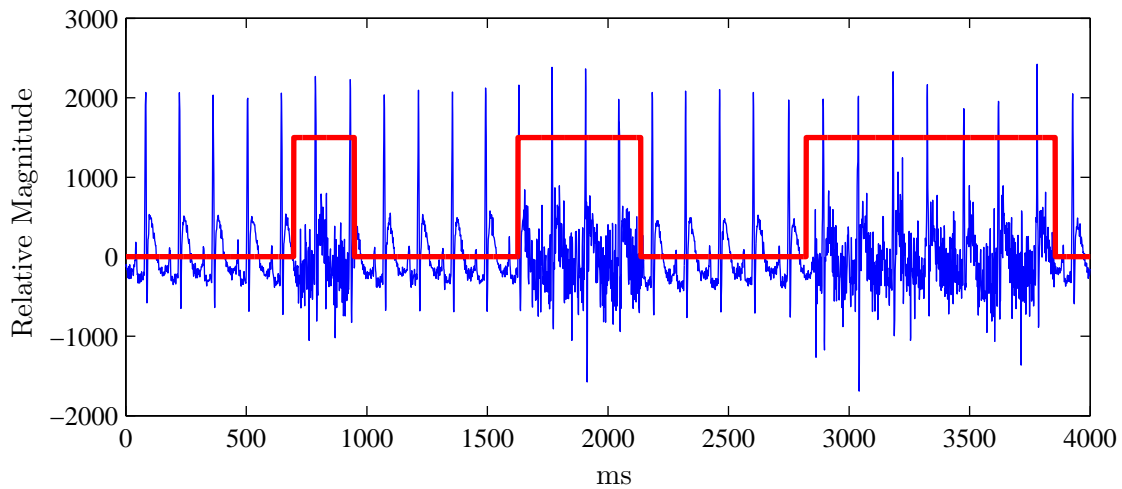


Figure 5.22: Input ECG signal black and detected EMG sections (red)

5.6.1 Training

The algorithm parameters, size of the moving variance window and the moving variance threshold were determined by training on a set of data and tuned to achieve the best detection rate.

Training set

Training signals were generated by adding EMG bursts at different lengths and strengths to EMG-free ECGs. Twenty EMG-free ECGs were selected from all rats. Each ECG has a length of 30 seconds. Six EMG bursts were added to each ECG passage at different lengths of 0.4, 0.5, 0.7, 1, 3, and 5 seconds at 2 seconds intervals. By reusing the same ECG, the amplitude of the added EMG bursts was adjusted to 0.4, 0.5, 0.7, and 0.8 of the average of the R-wave amplitudes. In total, eighty training signals were created from twenty EMG-free ECGs. Note that all additive EMG bursts were selected from different EMG signals.

Procedure

The criterion for correct detection is defined as follows. The algorithm is considered to have a correct detection when the algorithm detects the entire EMG section and the detected section is not longer than 0.08 second from the onset and the end of the actual EMG interval. The process for determining the size of the moving variance window and threshold value for EMG detection is divided into two stages. First, trial and error was applied to find a reasonable range of window size (0.04 to 0.2 second) and threshold value (0.008 to 0.02). Then, the threshold value was selected at 0.015 and window size was adjusted for the best detection rate, achieved at $W = 0.081$ second. Then, at the best window size, the threshold value was varied. The algorithm achieved a 100% detection rate at a threshold equal to 0.01.

5.6.2 Testing

Test set

Three sets of signals were used to test the algorithm. Each set contains fifty signals from all rats. There are in total 150 test signals.

Set 1: Test signals were created by using the same method for the training set except that the signals had a length of ten seconds and contained one burst of EMG noise. The bursts began after two seconds of ECG. Their lengths and amplitudes were randomly chosen from lengths and amplitudes used for the training set.

Set 2: Signals were manually selected from the recorded signals which have spontaneous EMG noise. Each signal had one section of EMG noise.

Set 3: To test the performance of the algorithm in the real data, set 3 contained signals which were randomly selected from the recorded signals. Each signal had a length of 30 seconds. In this set, EMG noise may not be present in every signal.

Results

After using the correct detection criterion for training, the results are presented in Table 5.3. In this case, the sensitivity was computed as number of EMG sections which were correctly detected divided by total number of EMG sections and multiplied by 100%. The specificity was calculated in the same way as the sensitivity but for non-EMG sections. Sets 1 and 3 have a sensitivity of 100%, while set 2 has a sensitivity of 94%. Three out of fifty signals in set 2 have detection intervals of approximately 0.06 second shorter than the actual EMG intervals. Specificities are 100% for all test sets. Therefore, none of the ECG sections are detected as EMG noise

Table 5.3: Results from testing the EMG detection algorithm on 3 test sets

	Sensitivity (%)	Specificity (%)
Set 1	100	100
Set 2	94	100
Set 3	100	100

By using the fast implementation of the morphological filter in [82], each opening and closing operation requires fewer than four comparisons on average per sample. Therefore, the fast morphology filter has a complexity of $O(N)$. The other parts of the algorithm also have a complexity of $O(N)$. As a result, the overall computational complexity is $O(N)$.

The QRS detection is very crucial to the algorithm. Missing the QRS causes an increase in moving variance and a normal ECG will be detected as EMG noise. The possibilities of missing the QRS are when there is spurious noise which may cause the EMG to be detected when none is present.

5.7 Data preprocessing

The purpose of data preprocessing is to select ECG with acceptable quality for data processing. The flowchart of data preprocessing is depicted in Fig. 5.23. This uses algorithms previously described including beat detection, baseline fluctuation removal, and EMG detection. The recorded data, after passing data preparation, were originally categorized as raw ECG signal, no signal, and fullscale signal. In data preprocessing, the raw ECG is conditioned and trimmed off as unusable signal. As mentioned in Section 5.2, the raw ECG is contaminated by burst noise, motion artifact, EMG, and baseline fluctuation. These noises are detected and removed or minimized in this section. As shown in Fig. 5.23, the first step is to detect beat position because it is required for baseline fluctuation removal and EMG detection. The raw ECG and beat positions are fed into the baseline fluctuation removal program which outputs ECG without baseline fluctuation. The resulting ECG is then inputted to the EMG detection block which locates sections of raw ECG that have EMG or EMG-like noise above the setting threshold. The high EMG sections are labeled as EMG. ECG sections between any non-ECG sections (including no signal, fullscale, and EMG) are defined as clean ECG if their length is longer than 7 seconds. Otherwise, they are short ECG. These short ECG section are not submitted to data processing because they are too short to allow any performance analysis. It should be noted that the clean ECG is not perfectly clean, but contains some noise at a level that can be managed by the following procedures.

Five types of noise mentioned above are removed or minimized in data preparation and data preprocessing. Fullscale noise is cut out of the recording ECG in the data preparation. Baseline fluctuation and motion artifact are minimized by the baseline fluctuation removal process. Burst noise and EMG are detected by EMG detection and removed in the data categorization.

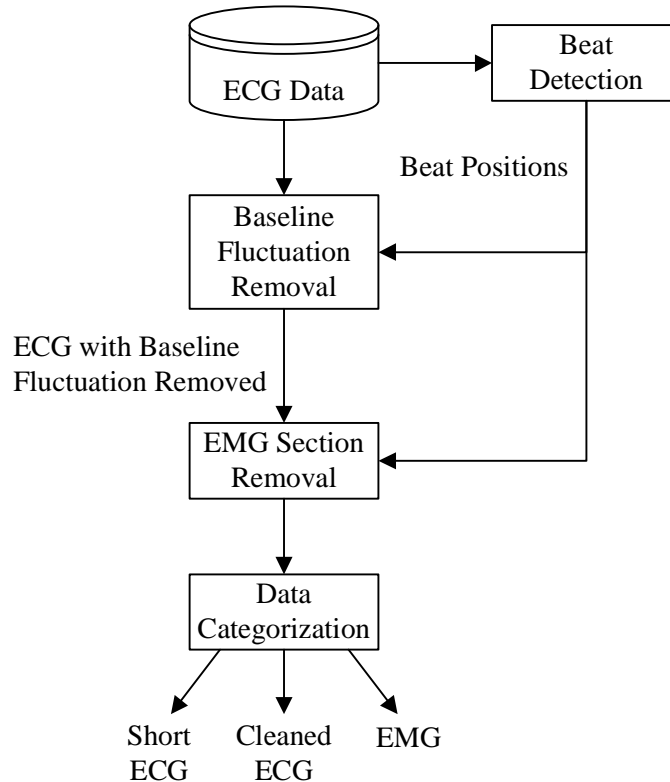


Figure 5.23: Flowchart of data preprocessing

It can be summarized that, after data preparation and data preprocessing, the recording signal is divided into five categories: clean ECG, short ECG, EMG, fullscale signal, and no signal. Only the clean ECG is submitted for data analysis. For ease in writing, clean ECG will be called as only “ECG” below.

5.8 ECG template generation

There are several types of ECG in the data, such as normal sinus rhythm and arrhythmias, including ectopic beats, ventricular tachycardia, and ventricular fibrillation. These ECG types need to be located and their positions used in searching for patterns. ECG template generation is a procedure which scans through ECG data and creates a set of ECG

templates. Each template represents a group of ECG beats that share the same shape. The set of ECG templates contains distinct ECG templates where any two templates have a distance measure exceeding a setting value. In the procedure, each beat is stamped with its ECG template along with its position in the data stream. The following explains the procedure of ECG template generation.

5.8.1 ECG template

An ECG template is a window of ECG beat. The window expands to the left w_l and right w_r milliseconds from the beat reference point. The sample mean is subtracted from the windowed ECG beat to adjust the baseline shift, which is then normalized by its Euclidean norm. The ECG template is designed to capture the QRS complex of ECG beat. Unlike P and T waves, the QRS complex is much less disturbed by noise and baseline fluctuation due to its high amplitude and high frequency, respectively. In addition, sometimes, P and T waves have a small amplitude and are buried by noises. The suitable w_l and w_r are 20 and 26 ms, respectively. There are, in total, 47 samples in the window (sampling rate is 1000 Hz). An example of an ECG template is given in Fig. 5.24. There are some cases when the R-wave is very small and undetectable. In such cases, the S-wave is used as the reference point instead. w_l and w_r equaling 20 and 26, respectively, can also be used when the S-wave is the reference point. It should be noted that the reference points are determined by the beat detection algorithm. For abnormal ECG beats, the highest amplitude peak or valley in the vicinity of the detected location is set as the reference point.

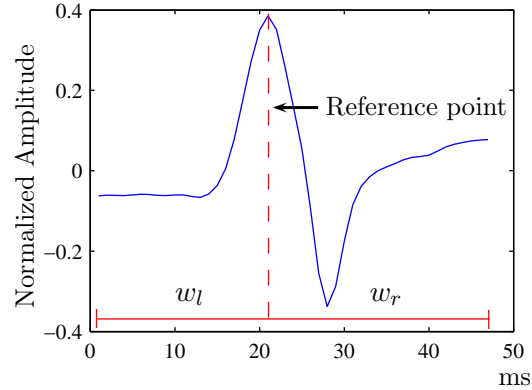


Figure 5.24: ECG template

5.8.2 Matching procedure

ECG template matching is a process for measuring the similarities between an ECG beat and an ECG template. The comparison requires that both the ECG beat and ECG template are in the same scale. Therefore, the ECG beat needs to be converted to a template beforehand. This template will be called a *beat template*. Among the matching scores in Section 4.2.1, the Euclidean distance was selected for large database, since it has the lowest complexity. However, other matching scores can also be applied. The matching score between the ECG template and beat template is calculated as shown in (5.9). d_β is the Euclidean distance. $\boldsymbol{\tau} = \{\tau_1, \tau_2, \dots, \tau_n\}$ and $\boldsymbol{\beta} = \{\beta_1, \beta_2, \dots, \beta_n\}$ are an ECG template and beat template, respectively. n is the number of samples in the template. A vertical alignment between the beat template and ECG template, b , is adjusted to achieve the lowest distance. By solving the equation, an offset b is the mean of the sample difference between $\boldsymbol{\tau}$ and $\boldsymbol{\beta}$ as written in (5.10). It is considered as a match when d_β is lower than the setting threshold, $d_{\beta th}$. Since the reference point can be either the R-wave or S-wave, the ECG beat and ECG template can be compared only when they have the same type of reference point.

$$d_{\beta} = \min_b \frac{\sqrt{\sum_{i=0}^n (\tau_i - \beta_i + b)^2}}{n} \quad (5.9)$$

$$b = \frac{\sum_{i=0}^n (\tau_i - \beta_i)}{n} \quad (5.10)$$

5.8.3 Generating the ECG template

ECG template generation is a procedure for creating a set of ECG templates which represent ECG beats in the data. The algorithm sequentially scans ECG from the first to the last beat. Distinct beat morphologies are recorded as ECG templates. The flowchart of the ECG template generation algorithm is displayed in Fig. 5.25. The algorithm begins with an empty set of ECG templates. Successively, ECG beat positions are located by ECG beat detection as explained in Section 5.4, ECG beats are cropped, and their beat templates are computed. The first beat template is automatically assigned as an ECG template. Then, each of the following beat templates is sequentially compared to the templates in the ECG-template set. If the distance between a considered beat template and a compared ECG template is less than the matching threshold ($d_{\beta} < d_{\beta th}$), a matched ECG template is found. The first matched ECG template is used immediately to speed up the procedure. Then, the matched ECG template is updated by averaging the beat template to itself as expressed in (5.11), where k is denoted as a k^{th} version of the ECG template.

$$\tau_i^{(k+1)} = \frac{\tau_i^{(k)}k + \beta_i}{k + 1} \quad (5.11)$$

The ECG template update step allows ECG templates to become representatives of their groups of matched ECG beats. In addition, white noise in the ECG template is more

attenuated as the number of templates averaged increases. However, the shape of ECG beats may gradually change over time. To capture this morphological change, the number of updates has to be limited to a certain number, k_{limit} . Otherwise, the ECG template would follow the gradual change in ECG beat morphology. For the case where no ECG template matches a considered beat template, the beat template is assigned as a new ECG template. In this research, the matching threshold, $d_{\beta th}$, used is 0.014 and k_{limit} is 100.

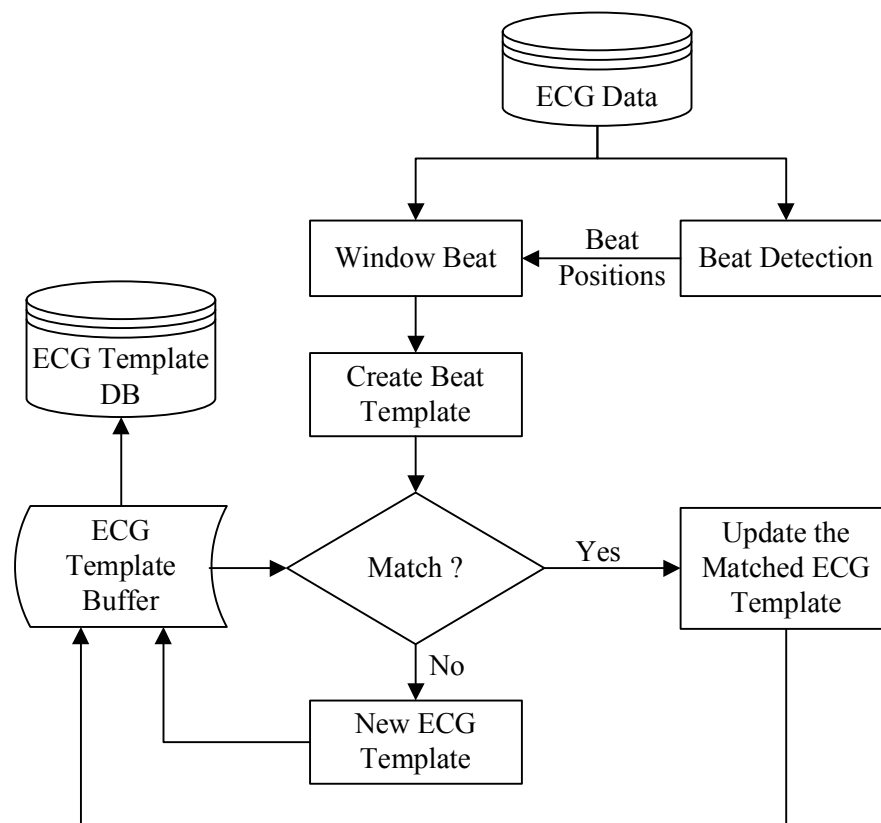


Figure 5.25: Flowchart of ECG template generation algorithm

Due to a large number of ECG beats in the data, a large number of ECG templates could be generated and result in a relatively large number of comparisons per ECG beat. Therefore, the number of ECG templates used in the matching procedure is limited to approximately 600 templates. These ECG templates are stored in a memory buffer and

sorted by time of the last match. If the number of ECG templates grows over the limit, a group (40%) of the oldest templates is transferred to the ECG template database.

The following information is recorded for individual ECG templates.

- *Template ID (T)*: All ECG templates are assigned distinct ID numbers.
- *Type of reference wave*: Either the R-wave or S-wave can be the reference wave for the ECG template.
- *Number of beat templates averaged to ECG template (N_t)*: An ECG template is an average of matched beat templates as explained above. The maximum N_t possible is k_{limit} . With a higher number, the impact by the noise in the ECG template is less (a smoother the ECG template).
- *Number of ECG beats that share a particular ECG template (N_b)*: A particular ECG template represents a number of ECG beats. In the other words, a number of ECG beats has the beat templates matched to a particular ECG template.

Examples of the ECG templates are displayed in Fig. 5.26. Subfigure (a) shows an ECG template of a normal ECG rhythm. It has the R-wave as the reference wave, $N_t = 100$, and $N_b = 103761$. This means that the ECG template is an average of the first 100 matched beat templates and represents a total of 103,761 ECG beats in the data. Subfigure (b) displays a noise which is mistakenly detected as an ECG beat and created as an ECG template. It has N_t and N_b both equal to 1. Such a template should be rejected and not used in analysis. Subfigure (c) depicts the ECG template of an ECG beat with an abnormal shape. For subfigures (d) to (f), ECG templates of premature beats are displayed.

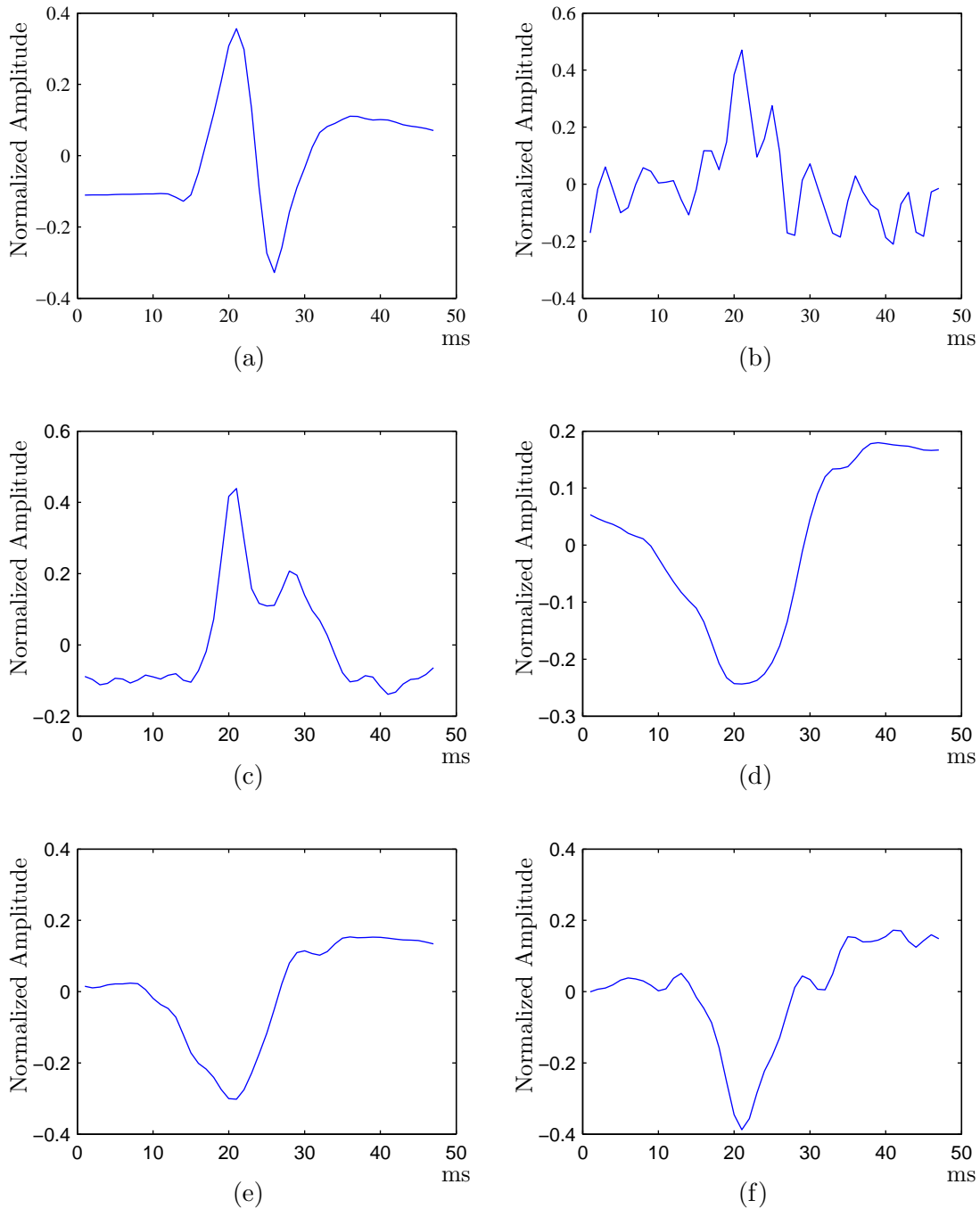


Figure 5.26: Examples of ECG templates – (a) depicts the ECG template of an ECG beat. (b) displays a mistaken ECG template which the algorithm extracts from noise. (c) to (d) are ECG templates of abnormal ECG beats.

5.8.4 Searching for ECG template duplicates

As mentioned in Section 5.8, a portion of the ECG templates is used in the template matching procedure to reduce computation time. As a result, there may be more than one ECG template which represents the same ECG morphology. In this section, duplicates for an ECG template are searched for to reduce the number of templates used in analysis.

By selecting an ECG template as a query template, duplicates of the query template are defined as templates which are located within a specified Euclidean distance, r_{dup} . It begins with the ECG template with the highest number of matching (N_b) as the first query template. The nearest neighbors search in [84, 85] is applied with a setting r_{dup} to determine duplicated templates. The resulting templates are purged from the database and their matching beats are relabeled to the ID of the query template. The searching process is repeated by setting an ECG template with the next highest N_b as the query template and so on. The drawback is that the nearest-neighbors-search algorithm is relatively more computationally intensive because of high dimensionality and the high number of templates. Therefore, ECG templates are downsampled to half and, alternatively, ECG templates with a very low N_b can be eliminated from the search for duplicates. The algorithm of the nearest neighbors search is briefly explained as follows:

Nearest neighbors search

The following description is the nearest neighbor search algorithm explained in [84]. The algorithm is divided into two stages, preprocessing and search. The preprocessing stage constructs a cluster tree to use in the search phase. Once the cluster tree is generated, it can be used repetitively for the nearest neighbor search.

Preprocessing phase: In this phase, the data points are divided into small clusters containing approximately 30 to 200 members – called terminal nodes (clusters). Recursively, data are divided into two child clusters and child clusters become parent clusters which are further separated into two clusters until the number of members meets the requirement. At the end of process, the distances from the clusters center to all points belonging to this cluster are computed and stored. The farthest is recorded as the cluster radius.

Search phase: The algorithm utilizes the cluster tree to determine the searching space. The terminal nodes which are distant from the query point, q , and cannot possibly contain nearest neighbors are not considered in the search phase. To determine k nearest neighbors, the algorithm initially searches for k nearest points to q . The k^{th} nearest point is denoted as p_k . $d(\cdot, \cdot)$ denotes the distance between two points. The distance from the query point to the k^{th} nearest point, $d(p_k, q)$, provides an upper limit or radius for the searching area. $d(p_k, q)$ is set to r_{dup} for the searching for ECG template duplicates. The triangular inequality is applied to prune off distant terminal nodes and distant points in a considered terminal node as follows:

Exclude cluster i , if

$$d(p_k, q) < d(c_i, q) - r_i, \quad (5.12)$$

where c_i and r_i are the center and radius of cluster i , respectively.

Any point, x , inside a considered terminal node is excluded, if

$$d(p_k, q) < |d(c_i, q) - d(c_i, x)|. \quad (5.13)$$

5.9 Feature extraction

ECG beats are detected using the algorithm in Section 5.4. The DSI time of the beat reference point is assigned as the beat ID and used to refer to a specific beat. The features explained below are computed for every beat.

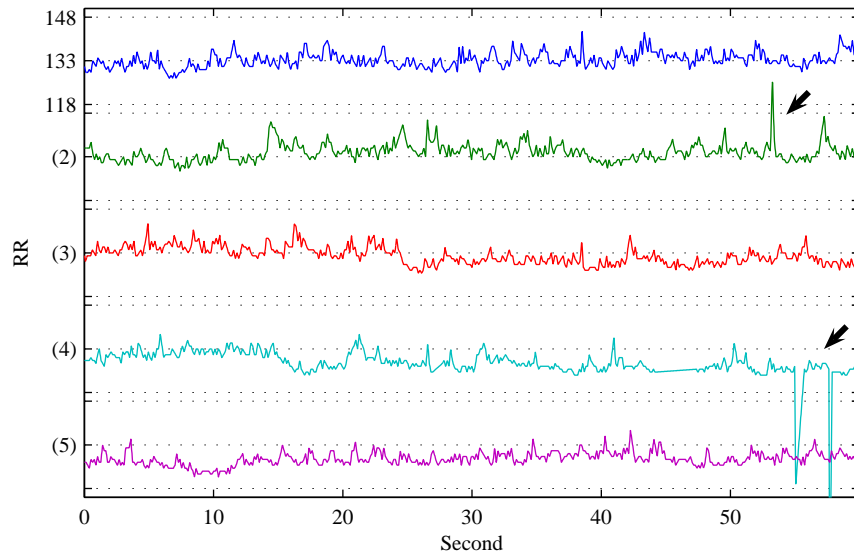
1. *RR interval*: This is a time interval in milliseconds from the previous beat to the current beat.
2. *Beat norm*: Each ECG beat is cropped by a window expanding from the reference points to the left and right w_l and w_r , respectively, and calculating the Euclidean norm as its beat norm. w_l and w_r are 20 and 26 ms, respectively. In other words, beat norm is the Euclidean norm used to calculate the beat template.
3. ΔRR : It is the difference between RR intervals of the current and previous beats.
4. *Beat EMG level*: EMG noise level within a beat window calculated by (5.8) is accumulated and divided by the window size squared.
5. *ID of the matched ECG template*: ID of the matched ECG template from section 5.8 is stored as one of the beat features.

Example plots of features are displayed in Fig. 5.27 where RR interval, beat norm, ΔRR , and beat EMG level are plotted in subfigures (a) to (d), respectively. Each plot starts from the top row continuing to the row below and so on. Magnitude levels are shown in the first row and also applied to the other rows. The plots contain five minutes of data. Missing beat detection causes artifacts in features depicted as arrows in the plots. The artifacts include abnormally high RR intervals and ΔRR as shown in the second row plot of subfigures (a) and (b). High amplitude noise displays as high beat norm and high beat

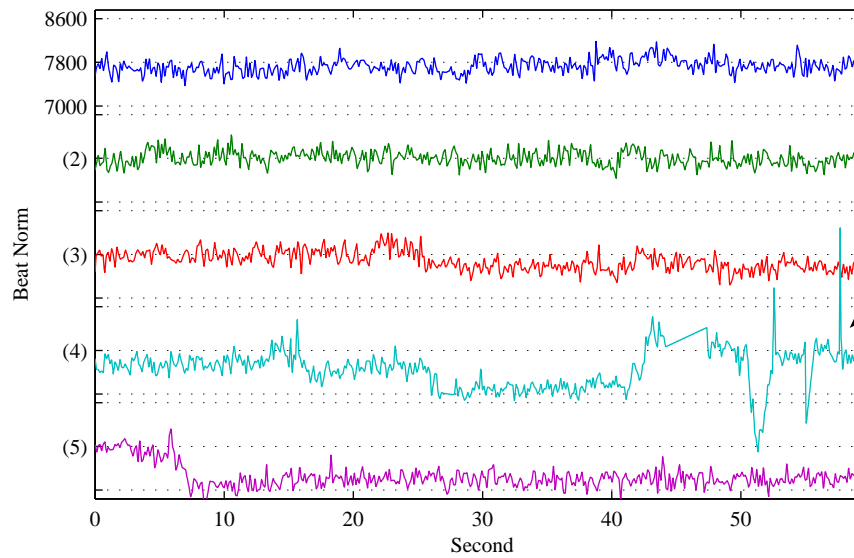
EMG level in the fourth row plot in subfigures (c) and (d). It results in a high RR interval and ΔRR in the fourth row plot in subfigures (a) and (b). To avoid missing beat detection, the beat detection algorithm is set to be sensitive (low parameter c in the beat detection algorithm) but not too sensitive. It has the drawback that more artifacts and noise are detected as a beat. Therefore, a later analysis algorithm must cope with these artifacts.

Distributions of RR interval, ΔRR , beat norm, and beat EMG level from subject 3 are illustrated as histograms in Fig. 5.28. Total number of beats from one animal is 27,450,194. Histogram ranges (min,max) for RR interval, ΔRR , beat norm, and beat EMG level are (1,299), (1,9999), (-25,25), and $(0, 2 \times 10^{-4})$, respectively. Their bin sizes are 1 except for beat EMG level which uses 10^{-6} . Values smaller than the minimum value are collected in the first bin, and the last bin also counts values larger than the maximum value as indicated by arrows in subfigures (c) and (d). In subfigure (a), the arrow shows counts of outliers whose values are out of typical ranges. These can be noise or abnormal events. Note that the typical ranges are determined by experts.

A two-dimensional histogram depicts spatial distribution and correlation of two features. Fig. 5.29 illustrates 2D histograms for subject 3. Table 5.4 lists details for each histogram. Values smaller than the minimum value and larger than the maximum value are put in the first and last bin, respectively, as shown by the arrows in Figs. 5.29 (b) and (c). Counts are coded according to the color bar at the right of each histogram and white indicates zero count. The arrows in all figures point at clusters which are suspected of being noise or abnormal events.



(a)



(b)

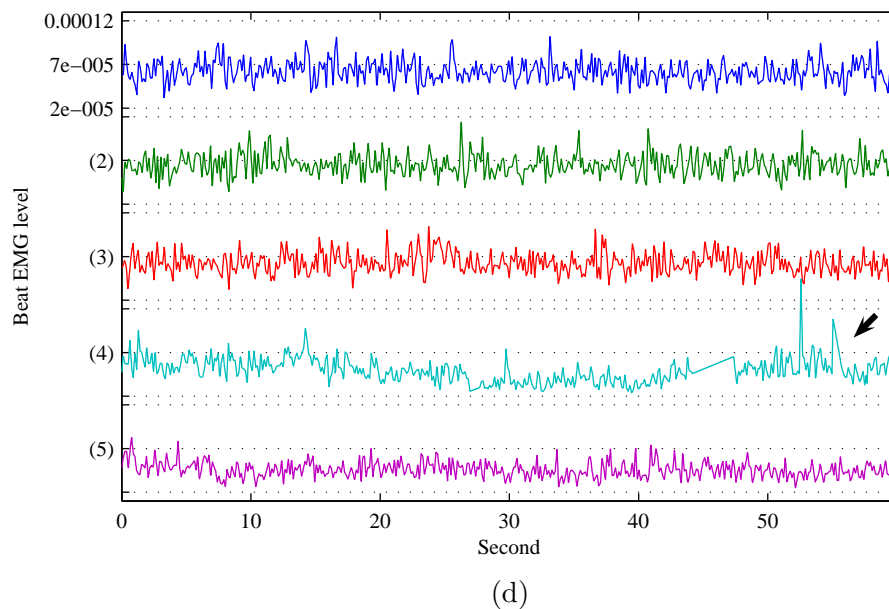
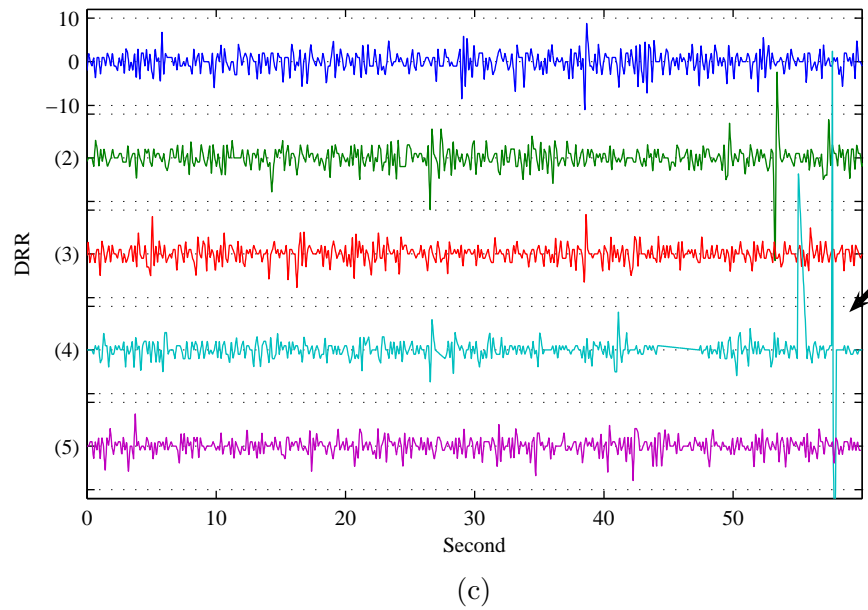


Figure 5.27: Examples of features for subject 3; The horizontal axis is time in seconds. The vertical axes are (a) RR interval in milliseconds, (b) beat norm displayed as relative magnitude, (c) Δ RR in milliseconds, and (d) beat EMG level. Each plot starts from the top row continuing to the row below and so on. Each row has 3 dotted lines which depict magnitude levels where the values written on the top row are applied to the rows below as well. The numbers in parentheses show row number. Feature artifacts are indicated by arrows.

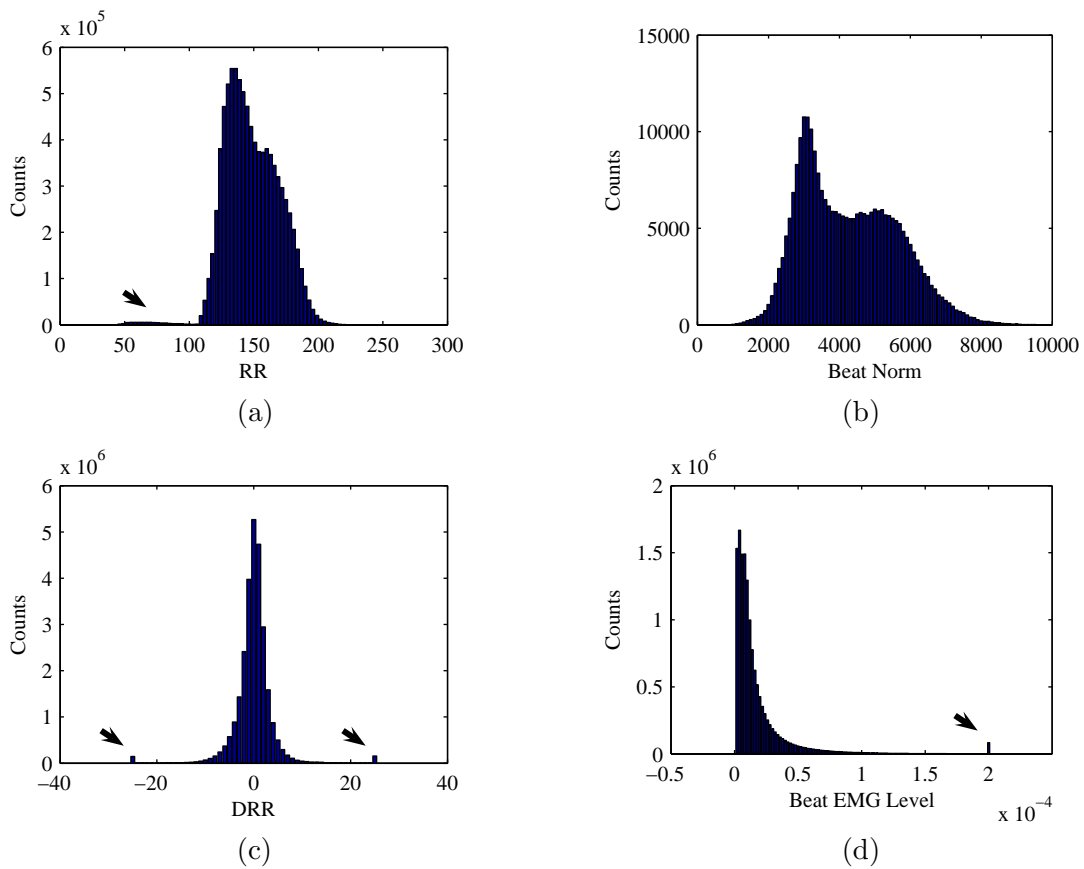
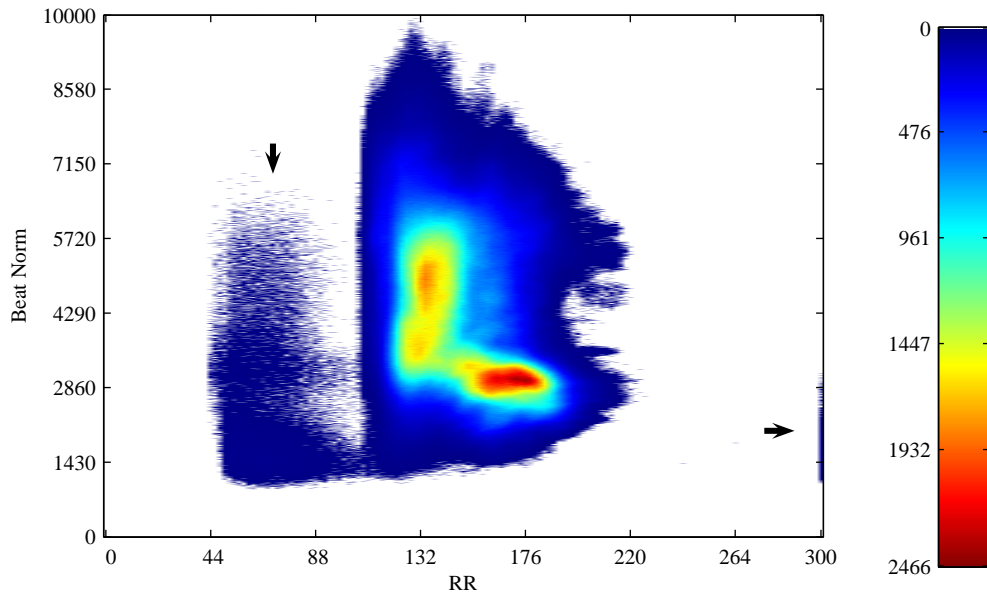
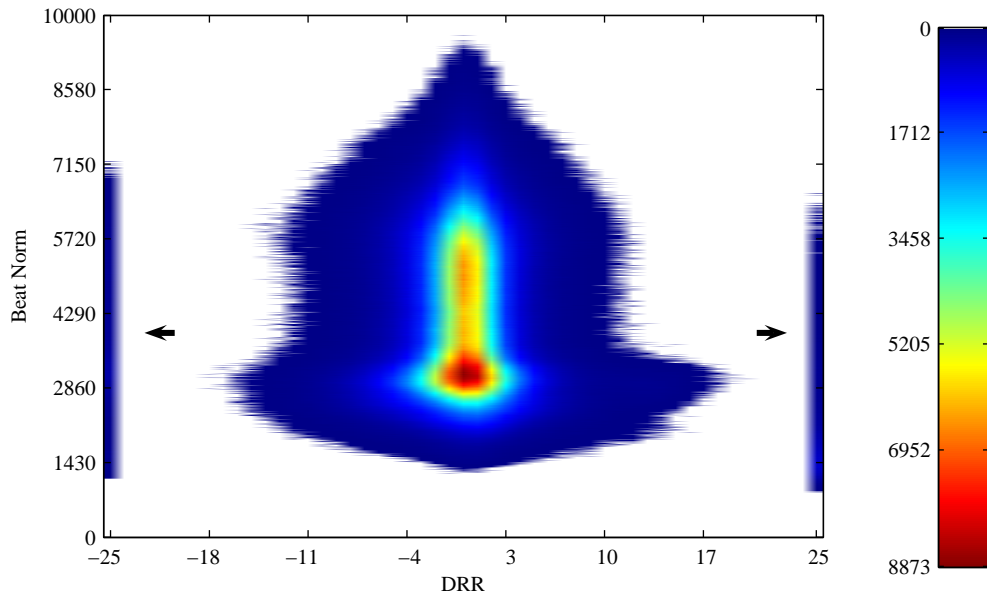


Figure 5.28: Histograms of features for subject 3; The vertical axis is count. The horizontal axes are (a) RR interval, (b) beat norm, (c) ΔRR , and (d) beat EMG level. The arrows point to counts of noise in the features.



(a)



(b)

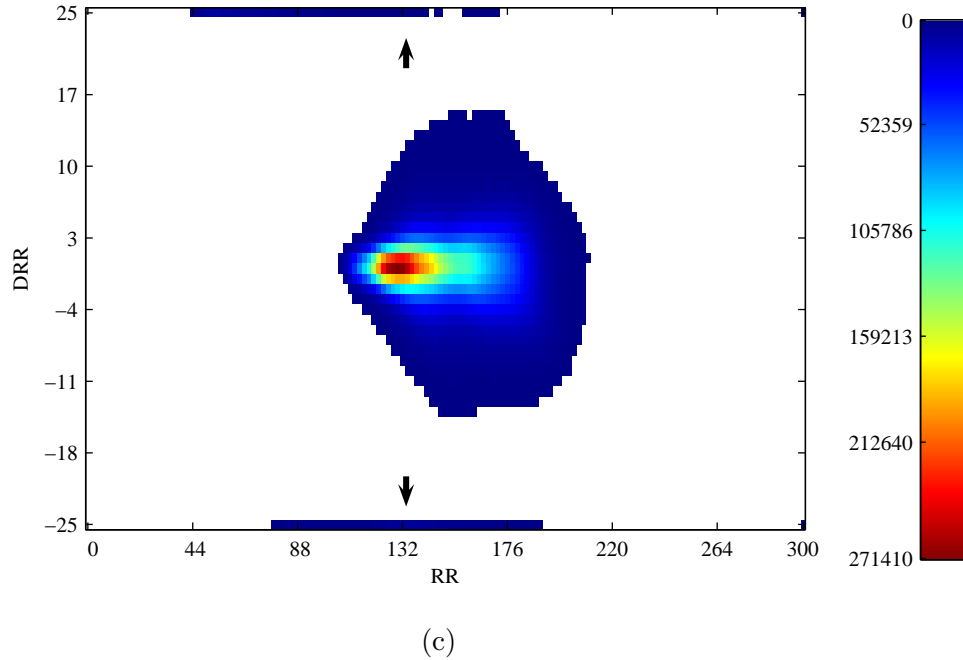


Figure 5.29: Two-dimensional histograms of features for subject 3 (horizontal axis v.s. vertical axis) – (a) RR interval v.s. Beat norm, (b) ΔRR v.s. Beat norm, and (c) RR interval v.s. ΔRR : Counts are coded according to the color bar at the right of each histogram and white indicates zero count. The arrows point to groups of feature artifacts.

Table 5.4: Bin configuration of 2D histograms in Fig. 5.29; Values smaller than the minimum value and larger than the maximum value are put in the first and last bin.

Subfigure	Feature	Min	Max	Bin size
(a)	RR	0	300	2
	Beat norm	0	10^4	5
(b)	ΔRR	-25	25	1
	Beat norm	0	10^4	5
(c)	RR	0	300	2
	ΔRR	-25	25	1

5.10 Heart rate variability visualization

In data preprocessing, unanalyzable portions of the signals are removed, ECG beats are detected and RR intervals are computed. The artifact rejection program in [44] is applied to RR intervals to remove the ECG beat detection artifact. The algorithm compares each

RR interval with the median of 25 surrounding RR intervals and the last non-artifact RR interval. It is determined as an artifact RR interval, if both differences exceed 20%. The following time-domain and frequency-domain heart rate variability (HRV) parameters are calculated for every usable 90-second segment of RR intervals: mean (MEAN), median (MEDIAN), standard deviation (SDNN), coefficient of variance (CV), interquartile range (IQR), mean of ΔRR (MUDRR), standard deviation of ΔRR (SDSD), interquartile range of ΔRR (DDIQR), root mean square of ΔRR (RMSSD), low frequency (0.04 - 1.0 Hz) power (LF), high frequency (1.0 - 3.0 Hz) power (HF), and ratio LF/HF (LHF) [30, 31]. A usable segment is defined as a segment containing non-artifact RR intervals in more than 85% of its time interval. The definitions and formulae for calculated HRV parameters were given in Section 2.5.

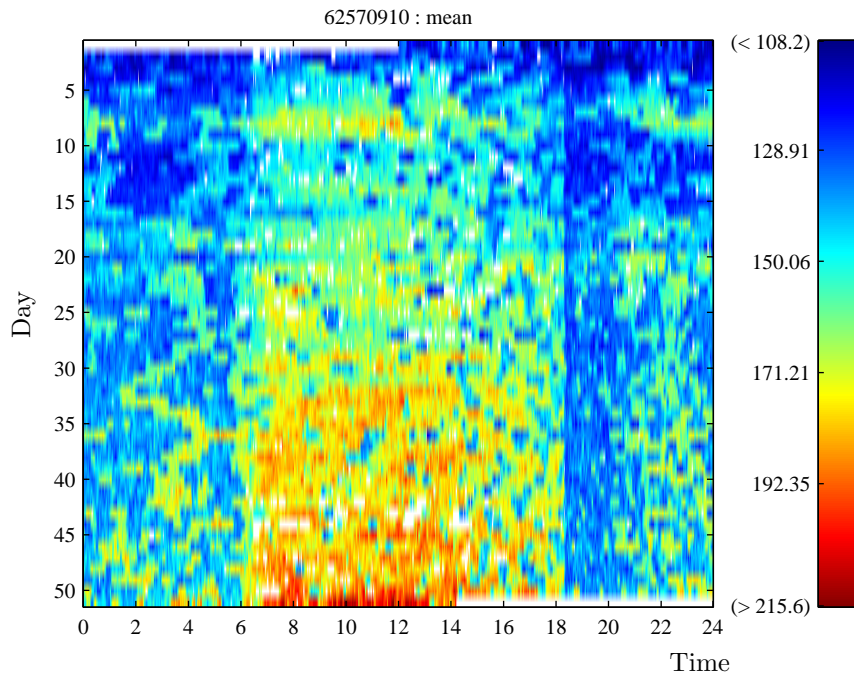
The results are visualized by using an image plot as displayed in Fig. 5.30 [7, 4]. The horizontal axis is time of day starting from midnight (0) to midnight of the next day (24 hours). The vertical axis is day of experiment. Each pixel is a 90-second interval corresponding to time of day and day of experiment. HRV parameter values for each 90-second interval are represented by colors coded according to the color indicator bar at the right of the graph. There are in total 255 colors. Missing data or unusable segments are indicated by white pixels or stripes. Selected HRV parameters for subject 3 are plotted in Fig. 5.30 to show changes in sample mean, deviation and frequency components of RR interval and ΔRR . For this treatment rat, day of experiment begins from the day after implanting an aldosterone pump and transmitter (day 2 in the graphs). The aldosterone treatment with 1% NaCl and low Mg diet starts at day 2 and ends at day 43. After day 43, the 1% NaCl and low Mg diet are continued.

A 12:12-h light-dark cycle is controlled by turning on and off lights at approximately 6 AM and 6 PM, respectively. The results show several interesting observations as follows.

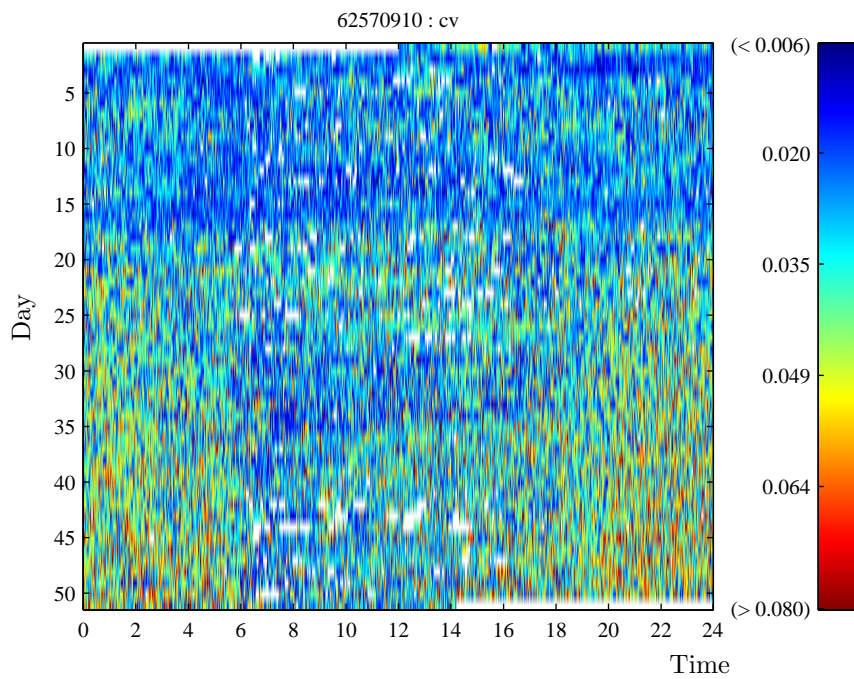
There is a day/night difference in HRV as seen by the changes associated with lights at 6 AM and 6 PM – as shown in graphs (a) - (f) which establish a 24-hour circadian rhythm. MEAN is high during daytime and low during nighttime – i.e. the heart rate is low during daytime and high at night. As the study progresses, there is a decrease in heart rate and increase in variability over the time course. It can be seen that the day-to-day MEAN continuously increases from the beginning to the end of the experiment – especially during daytime where MEAN turns from blue in the first day to red in day 51. Note that increase in MEAN means decrease in heart rate. Besides the day-to-day change, the heart period (RR interval) variation across 24 hours also changes. It is shown as a pixel-color variation in each row of graph MEAN. Pixel-color fluctuation increases as number of days increase. The RR interval variation can also be observed in MUDRR, DDIQR, and RMSSD. High value of these parameters means high variability. MUDRR, DDIQR, and RMSSD also show an increase in RR interval variation as the experiment progresses. MUDRR determines the shift of RR intervals. It can be positive (shift up) or negative (shift down) values. Swing in heart rate shows as change in color shade of horizontally consecutive pixels. CV measures the variation of RR intervals. It is normalized by the mean of RR intervals. Therefore, CV is low when MEAN is high, as displayed in the CV graph. LHF reflects sympathetic and parasympathetic influence on heart variability. As the days of the experiment pass, MUDRR and CV increase and LHF decreases. The differences in the patterns can be noticed in DDIQR, RMSSD, and LHF at approximately day 30 where the subcutaneous aldosterone pump is replaced.

The visualization tool is a very useful tool that provides an insight into data characteristics and short/long term changes as disease develops. Feature values from the entire experiment can be summarized in one graph. Differences early in the study versus later in the study can be observed in the data. This visualization tool can serve as a prelimi-

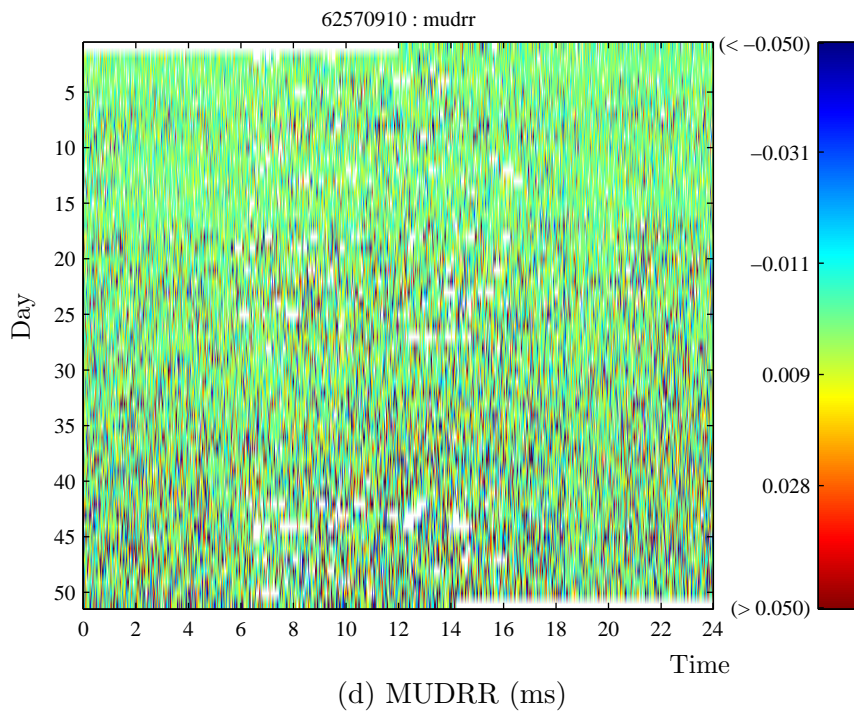
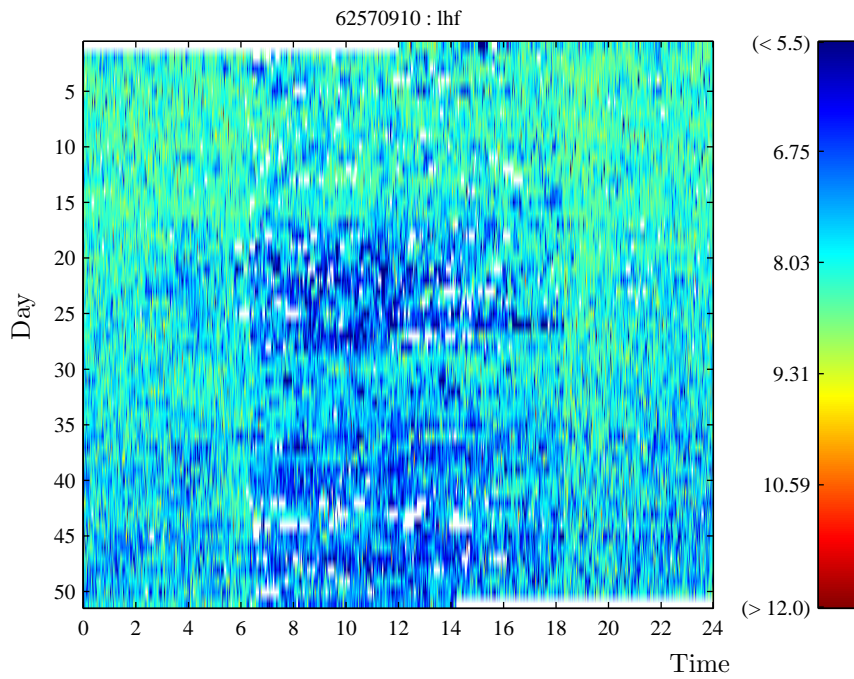
nary tool for locating interesting events and searching for/locating abnormalities by human experimenters.



(a) MEAN (ms)



(b) CV (ms)



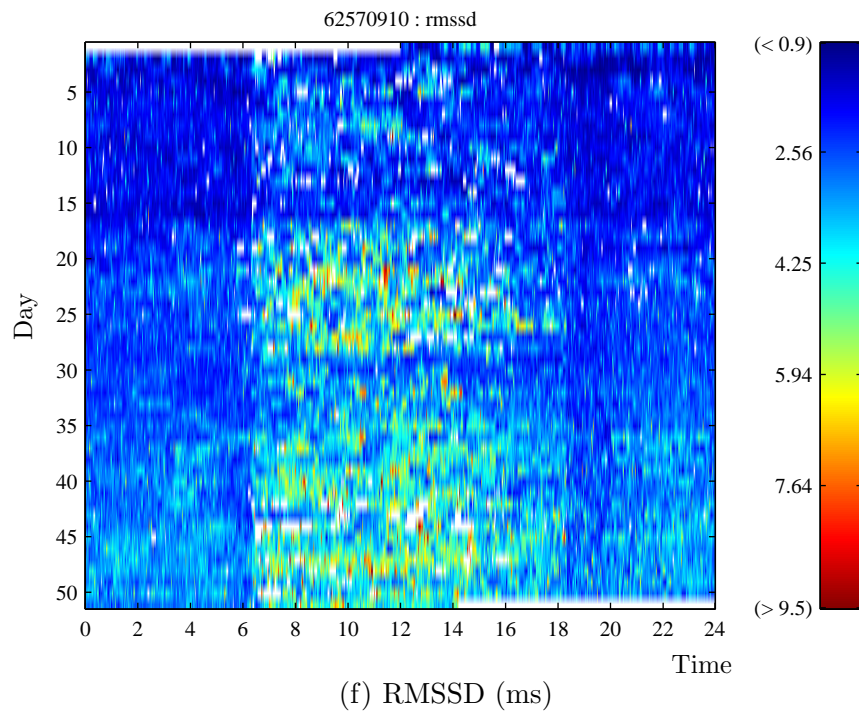
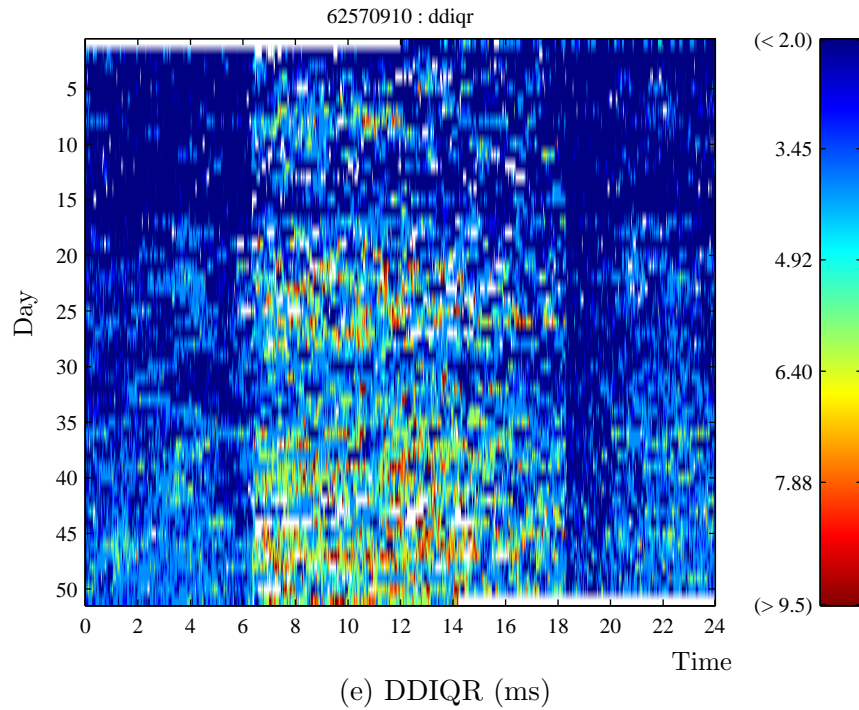


Figure 5.30: Image plots of HRV parameters for subject 3: Each pixel has the HRV parameter value of a 90-second interval, corresponding to time of day on the horizontal axis and day of experiment on the vertical axis. The color of each palette is coded according to the color indicator bar at the right of the figure.

5.11 Discussion

The framework for analyzing long-term ECG recordings consists of data preparation, ECG beat detection, baseline fluctuation removal, EMG section removal, ECG template generation, feature extraction, HRV parameter calculation, and abnormal ECG search. Beat detection is the most important step in the entire procedure. Locations of beat are used in other steps of the framework. The ECG beat detection program in [43] was modified to eliminate its drawbacks which include misdetection of severe baseline shift as an ECG beat and missing beats following an increase in ECG high amplitude. A morphological filter was added to eliminate baseline shift before detecting the ECG beat. The width of the structuring element must be set to the base width of abnormal ECG such as the base width of S-wave of PVC. The morphological filter allows deflections which have a base width less than the width of the structure element to pass through it for beat detection. If the width is too low, abnormal ECG will not be detected. And, if the width is too high, abrupt baseline shift will be detected as a beat. Other parameters, which include b , c , and g , need to be set manually until the beat detection algorithm can detect ECG beats properly. g is used to limit the growth of the threshold level so that ECG beats which follow a high amplitude are not missed. Low g increases beat detection sensitivity, and vice versa.

In the baseline-fluctuation-removal step, an algorithm in [83] was implemented. The algorithm had two morphological filters. The width of the structuring element of one morphological filter was fixed, while that of the other was determined by an algorithm developed for this research. The algorithm could not estimate the baseline at the beginning and the end of the signal passage accurately. Therefore, a program for the baseline-fluctuation-removal algorithm was written to process the ECG signal as a continuous stream of data, not segments of data. In addition, the algorithm could not remove a baseline which swings

faster than a T-wave. Therefore, moving artifacts were still presented in the output signal.

Sections of EMG noise in the data were detected and cropped off using the algorithm explained in Section 5.6. Correct beat detection was very crucial to the EMG detection algorithm. Consecutive missing beats, especially normal ECG beats, caused fault acceptance. Another cause of fault acceptance was spurious noise. The limitation of the algorithm is that it cannot detect short EMG sections, as it needs a sufficient amount of EMG to raise the moving variance above a set threshold level. In Section 5.6, sections of EMG noise, which are longer than 0.4 seconds, were tested and could be detected by the algorithm. EMG detection sensitivity can be increased by decreasing the threshold level.

The ECG-template-generation, which is explained in Section 5.8, is a suggestive procedure and can be adjusted as follows. The template window sizes, w_l and w_r , can be changed to capture interesting ECG waves, such as the P-wave and T-wave. Other matching scores, which are correlation waveform analysis and bin area method, can be used to measure the similarity between an ECG template and beat template instead of the matching distance (d_β) in (5.9). Correlation waveform analysis and bin area method are explained in Section 4.2.1. They require higher computation power than d_β . But, they can match two templates which horizontally shift away from each other. It should be noted that d_β is calculated regardless of the vertical shift in templates. d_β and k_{limit} can be reduced to capture a slight change in ECG morphology. However, the number of generated ECG templates increases as d_β and k_{limit} decrease. Numerous ECG templates cause difficulties in memory management and searching for ECG template duplicates. Therefore, this tradeoff needs to be considered when selecting d_β and k_{limit} values. The ECG-template-buffer size in Fig. 5.25 can be varied, depending on available memory space. The number of duplicate ECG templates reduces as the buffer size increases. As a result, the computation time for ECG generation increases while searching for ECG template duplicates requires less time

and memory usage. The ECG-template-buffer size should be increased, if a large number of ECG templates is generated.

Features extracted in Section 5.9 include RR interval, beat norm, ΔRR , and beat EMG level. RR interval and ΔRR were used to calculate HRV parameters, while all of them were prepared for abnormal ECG search in the next chapter. Feature values for suspect abnormal ECG and noise could be estimated using one and two-dimension histograms in Figs. 5.28, and 5.29. Other features can be extracted for specific analyses. HRV parameters were computed and visualized using the image plot in [46]. The plot shows a potential for manual pattern discovery. Besides HRV parameters, other features can also be displayed using this image plot. Before plotting, feature artifacts should be rejected.

5.12 Summary

A procedure for analyzing long-term ECG recordings was introduced in this chapter. ECG data were collected from rats in a model of chronic heart failure for 12 weeks. There are in total 80 GB of data. The entire process of data analysis includes data preparation, ECG beat detection, baseline fluctuation removal, EMG section removal, ECG template generation, feature extraction, HRV parameter calculation, and abnormal ECG search. In data preparation, ECG data were decoded from the recorded files (DSI files) and stored in a database for data retrieval. Recording errors and noise were discarded. ECG beats were located using a modified algorithm which improves beat-detection sensitivity. The detected beat positions were used in the rest processes. ECG-baseline fluctuation was estimated using two morphological filters. A method to determine proper widths of the structuring elements was implemented. ECG sections which are highly corrupted by EMG and cannot be used in analysis were detected and windowed off by an algorithm invented for

this research. The algorithm has a sensitivity of 94% and specificity of 100%. Beat EMG levels were computed by the EMG detection algorithm and stored in the database. The usable ECG sections were submitted for ECG template generation. A suggestive algorithm was introduced for ECG template generation for use with large amounts of ECG data. The algorithm can capture gradual change in ECG morphology. The following features were computed: RR intervals, Δ RR, and beat norm. HRV parameters were extracted and visualized. Progressive changes in the HRV parameters could be observed as the number of days of experiment increase. Abnormal ECG search will be focused on and explained in the next chapter.

Chapter 6

Abnormal ECG search

The heart abnormalities depicted in the ECG signal represent irregularities such as irregular RR intervals and defective morphologies. In this chapter, abnormalities are searched for based on ECG morphology. Examples of ECG passages in the data are provided below. Fig. 6.1 (a) is a passage of normal ECG. Various abnormal ECGs are shown in Fig. 6.2. They include premature ventricular complexes in (a) and (b), elevated ST segment in (c), and split R-wave in (d). As explained in the previous chapter, the ECG beat positions are located by the beat detection algorithm. The recorded ECG is obtained by surface leads which are sensitive to subject movement and noise. The detected beats are not only ECG beats but also noise. Fig. 6.3 (a) displays movement artifact misdetected as a beat. Subfigures (b) and (c) show ECG beats corrupted by noise and burst noise – note that the detected locations are shown in black. The problem is not only distinguishing abnormal from normal ECG, but also the noise. Therefore, it is necessary to divide the larger problem into two subproblems. The first is to discriminate between noise (or misdetected beat) and ECG. This is explained in Section 6.3. The second subproblem is searching for abnormal ECG. This is a very difficult problem since there are many types of abnormalities. Medical experts detect an abnormal ECG from deviations in the ECG

components such as P, Q, R, S, T waves, and P-Q, S-T, Q-T segments. In a noisy environment, P, T waves and segment lengths can be easily corrupted and unmeasurable. Therefore, abnormalities which affect the QRS complex are focused on in this chapter. The types of abnormal ECG considered are PVC, elevated ST segment, and split R-wave. These are selected because they are found in sufficient numbers to develop algorithms and they are related to important heart problems. The search algorithm will utilize standard classifiers to find the specified abnormal ECG. It is explained in Section 6.4.

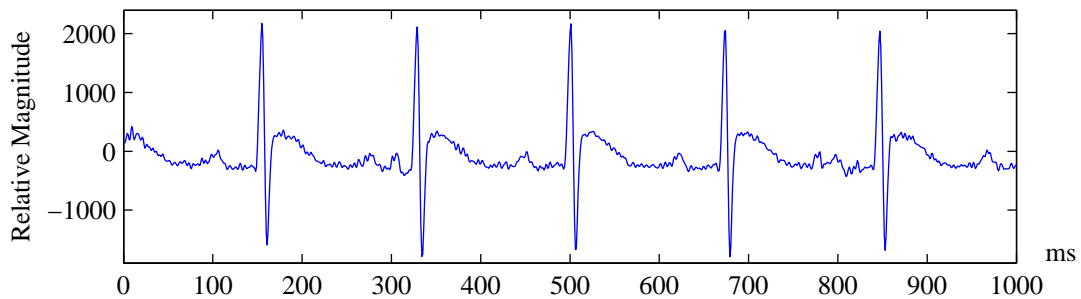
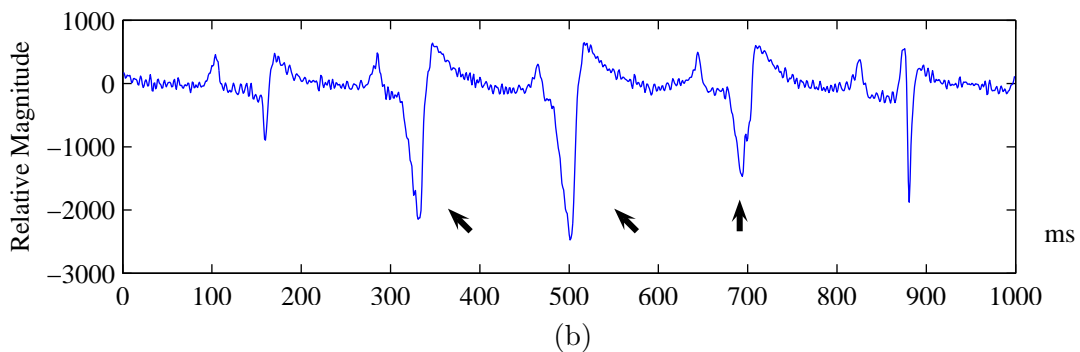
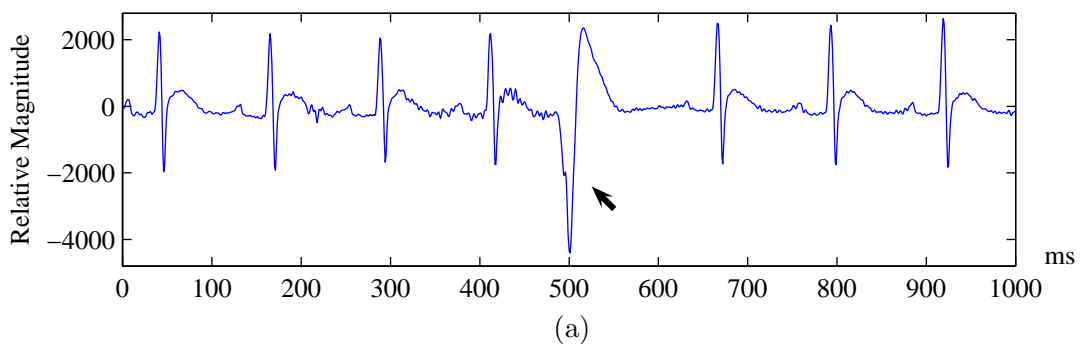


Figure 6.1: Examples of normal ECG



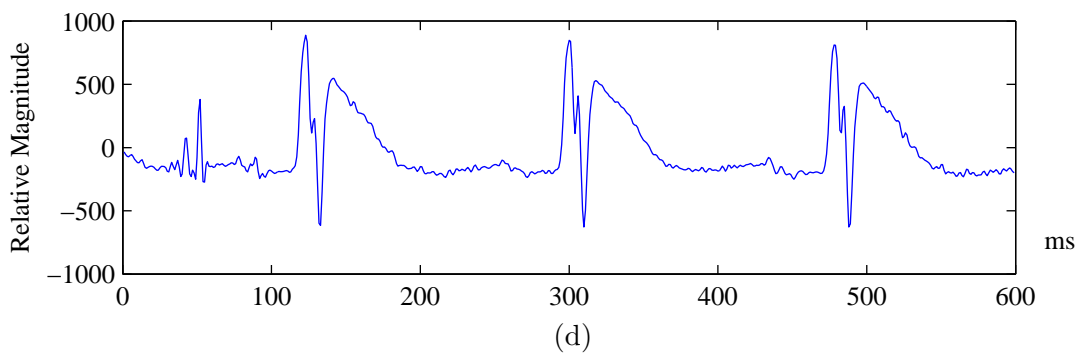
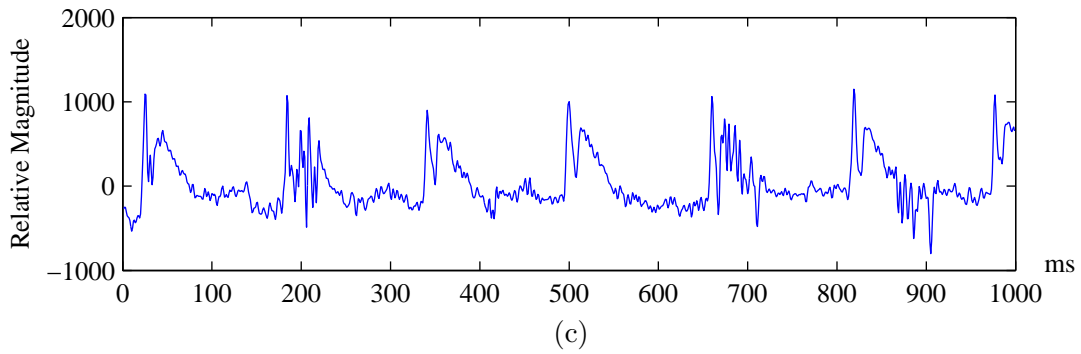
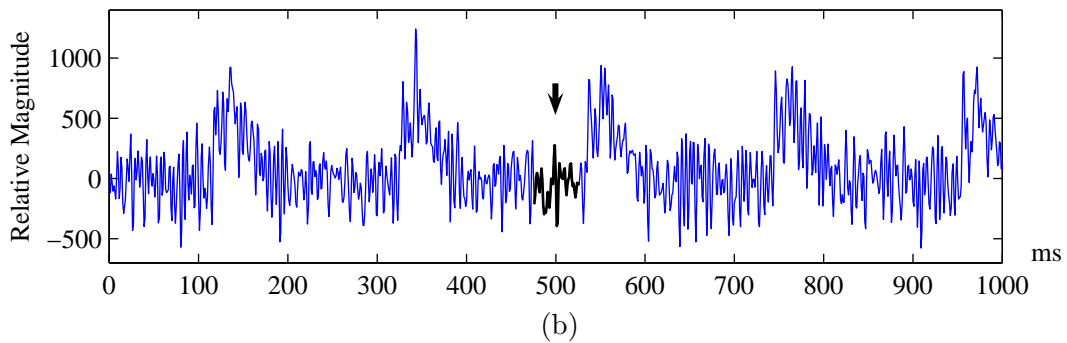
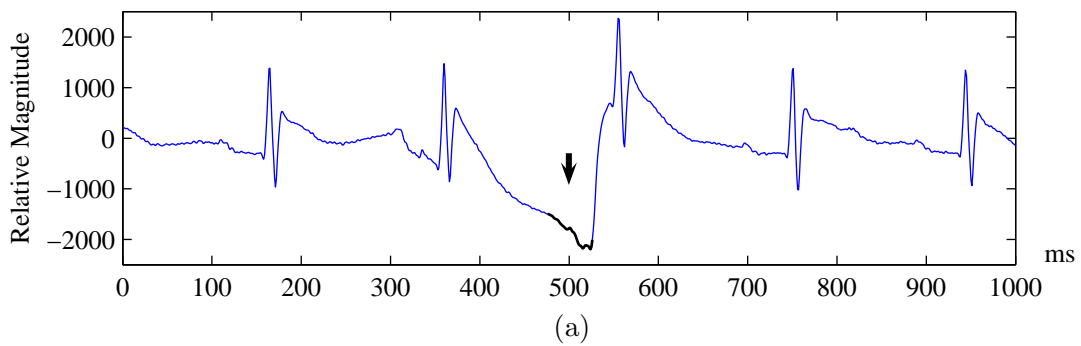


Figure 6.2: Examples of abnormal ECGs; (a) and (b) display premature ventricular complexes (They are indicated by arrows). (c) is a passage of elevated ST segment beats. (d) demonstrates split R-wave beats.



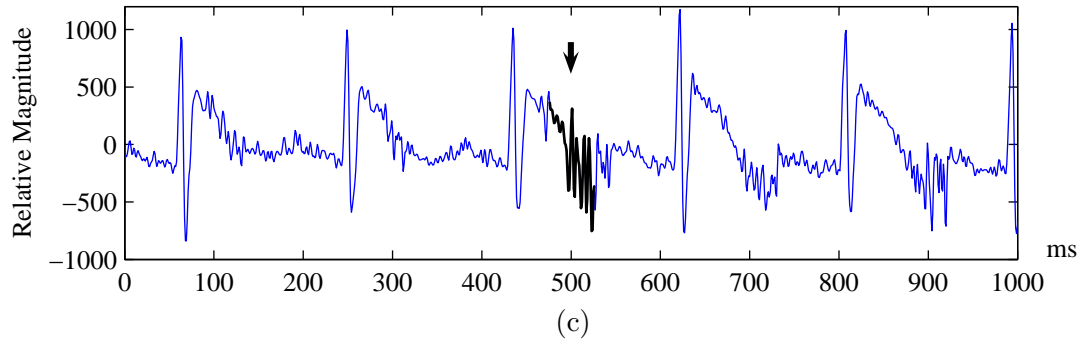


Figure 6.3: Examples of noise, artifact, and misdetections; Noise is indicated by arrows and shown in black. (a) is moving artifact detected as a beat. (b) displays an ECG passage highly corrupted by noise. (c) depicts EMG noise.

6.1 Problem statement

As introduced in the previous paragraph, the problems are formally stated in this section. Beat means a beat detected by the beat detection algorithm from the previous chapter. Therefore, beats can be ECG or noise. Due to various types of abnormal ECG, the abnormal ECG in this chapter refers to PVC, elevated ST segment, and split R-wave. Noise includes noise and ECG which is highly corrupted by noise. This chapter aims to distinguish abnormal ECG from detected beats. The problem is divided into two stages.

1. Discriminating normal and abnormal ECG beats from noise beats.
2. Isolating abnormal ECG beats from normal and abnormal ECG beats.

6.2 Dataset

The recorded ECG for all subjects have been pre-processed which includes beat detection, baseline removal, and EMG noise detection. The detected beats were manually

scanned for normal ECG, abnormal ECG, and noise. Table 6.1 shows the number of beats for each type and subject.

Table 6.1: Information about dataset: number of beats for each type and subject

Type	Subject					Total
	1	2	3	4	5	
Normal ECG	93	197	438	334	86	2560
Abnormal ECG	81	3	269	575	67	2663
Noise	1624	1633	1517	1328	811	7689
Total	1798	1833	2224	2237	964	12912

The following features are extracted from the selected beats:

1. RR interval (RR): The time interval in milliseconds from the previous beat to the current beat.
2. Successive difference between RR intervals (ΔRR): The difference between RR intervals for the current and the previous beats.
3. Percentage difference between the current RR interval and the mean of its surrounding RR intervals (ΔRRu): It is defined as:

$$\Delta RRu_i = \frac{RR_i - \mu s_i}{\mu s_i} \times 100\%, \quad (6.1)$$

$$\mu s_i = \frac{\sum_{\substack{j=-N, \\ j \neq i}}^N RR_{i+k}}{2N - 1},$$

where ΔRRu_i and RR_i are the current ΔRRu and RR , respectively, and N is the number of the surrounding RR intervals and is set to 2.

4. Percentage difference between the current RR interval and the median of its surrounding RR intervals (ΔRRm): It is given as:

$$\Delta RRm_i = \frac{RR_i - m s_i}{m s_i} \times 100\%, \quad (6.2)$$

$$ms_i = \text{Median}(RR_{i-N}, RR_{i-N+1}, \dots, RR_{i-1}, RR_{i+1}, RR_{i+2}, \dots, RR_{i+N}),$$

where the number of the surrounding RR intervals, N , is 2.

5. Beat norm (B): An Euclidean norm of a beat cropped by a window expanding to the left and right 20 and 26 ms, respectively, from the beat reference point. The reference point is a peak or valley in the vicinity of the location of the detected beat.
6. Successive difference between beat norms (ΔB): It is defined as $\Delta B_i = B_i - B_{i-1}$ where i is for the current beat.
7. Percentage difference between the current beat norm and the mean of its surrounding beat norms (ΔBu): It can be computed by using (6.1) and replacing RR interval, RR , with beat norm, B .
8. Percentage difference between the current beat norm and the median of its surrounding beat norms (ΔBm): It is calculated by using (6.2) and replacing RR interval, RR , with beat norm, B .
9. Beat level (L): The mean amplitude of the windowed beat. It indicates the jumping up or down of the signal.
10. Successive difference between beat levels (ΔL): It is calculated as $\Delta L_i = L_i - L_{i-1}$.
11. Percentage difference between the current beat level and the mean of its surrounding beat levels (ΔLu): The calculation repeats (6.1), replacing RR interval, RR , with beat level, L .
12. Percentage difference between the current beat level and the median of its surrounding beat levels (ΔLm): It is the same as (6.2), with RR interval, RR , replacing with beat level, L .

13. Beat EMG level: Accumulated EMG noise level (explained in Section 5.6) within the beat window.

The percentage difference features are normalized by the mean or median of the surrounding value so that they are less subjective compared to unnormalized features such as RR , ΔRR , B , ΔB , L , and ΔL . RR , ΔRR , ΔRRu , and ΔRRm are called in short RR interval features, so as beat norm features (B , ΔB , ΔBu , and ΔBm) and beat level features (L , ΔL , ΔLu , and ΔLm).

It should be noted that the ECG templates were investigated in the preliminary stage. One template can match more than one signal type, such as noise or abnormal ECG. Therefore, the features extracted from the ECG templates, such as beat morphology, size, and time intervals, are not taken into consideration.

6.3 ECG and noise classification

The recorded data contains noise and artifact. The beat detection algorithm sometimes detects noise and artifact as ECG beats. Therefore, noise and artifact needs to be recognized and eliminated from the dataset of ECG beats. In the literature review, noise can be estimated and detected using algorithms based on the Karhunen-Loeve transformation [66], the distribution of the frequencies of the slopes in an ECG waveform [67], and wavelet decomposition [68]. The Karhunen-Loeve-transformation method estimates noise only in the QRS complex portion. In some cases, noise is very similar to the QRS complex. The algorithm in [67] requires two consecutive ECG cycles to detect noise in an ECG passage. The computation of the wavelet approach is very intensive and not suitable for analyzing large amounts of data recordings. Low complexity algorithms for classifying ECG and noise beats are developed in this section.

Data for developing the algorithm were taken from the dataset in Section 6.2. The data were divided into training and test sets. Each set contains 3500 ECG beats and the same number for noise. The ECG beats include 2425 normal beats and 1075 abnormal beats in each set. The details are in Tables 6.2 and 6.3.

Table 6.2: Information about training set using for ECG and noise classification: number of beats from each type and subject

Type	Subject					Total
	1	2	3	4	5	
ECG (Normal)	482	447	518	449	529	3500
ECG (Abnormal)	318	253	182	251	71	
Noise	700	700	700	700	700	3500
Total	1500	1400	1400	1400	1300	7000

Table 6.3: Information about test set using for ECG and noise classification: number of beats from each type and subject

Type	Subject					Total
	1	2	3	4	5	
ECG (Normal)	483	447	517	448	530	3500
ECG (Abnormal)	317	253	183	252	70	
Noise	700	700	700	700	700	3500
Total	1500	1400	1400	1400	1300	7000

Figs. 6.4 to 6.7 display density histograms of the features to illustrate characteristics of the data. It should be noted that a density histogram is an ordinary histogram where the count is divided by the total area under the curve. In the density histogram, blue and red lines represent ECG and noise, respectively. RR intervals of noise are mainly lower than those of ECG. The same is true of beat norm density histograms. It shows that misdetected beats or noise are in between beats and have a lower amplitude than the surrounding detected beats. This is also reflected in the density histograms for ΔRR , ΔRRu , ΔRRm , ΔB , ΔBu , and ΔBm where noise has a higher density than ECG below a value of zero.

There are overlaps between noise and ECG density histograms for these features. This also occurs in density histograms for beat level features and EMG levels. EMG level is used to detect signal sections containing EMG. A significant distinction between noise and ECG in the EMG level is not evident in the density histogram, because noise beats can also be abrupt baseline changes, interference noise, short burst noise, etc. and EMG level is not sensitive to a very short section of signal, for instance a one-beat-length section. The feature values from ECG and noise beats are significantly differentiated. New features are investigated below.

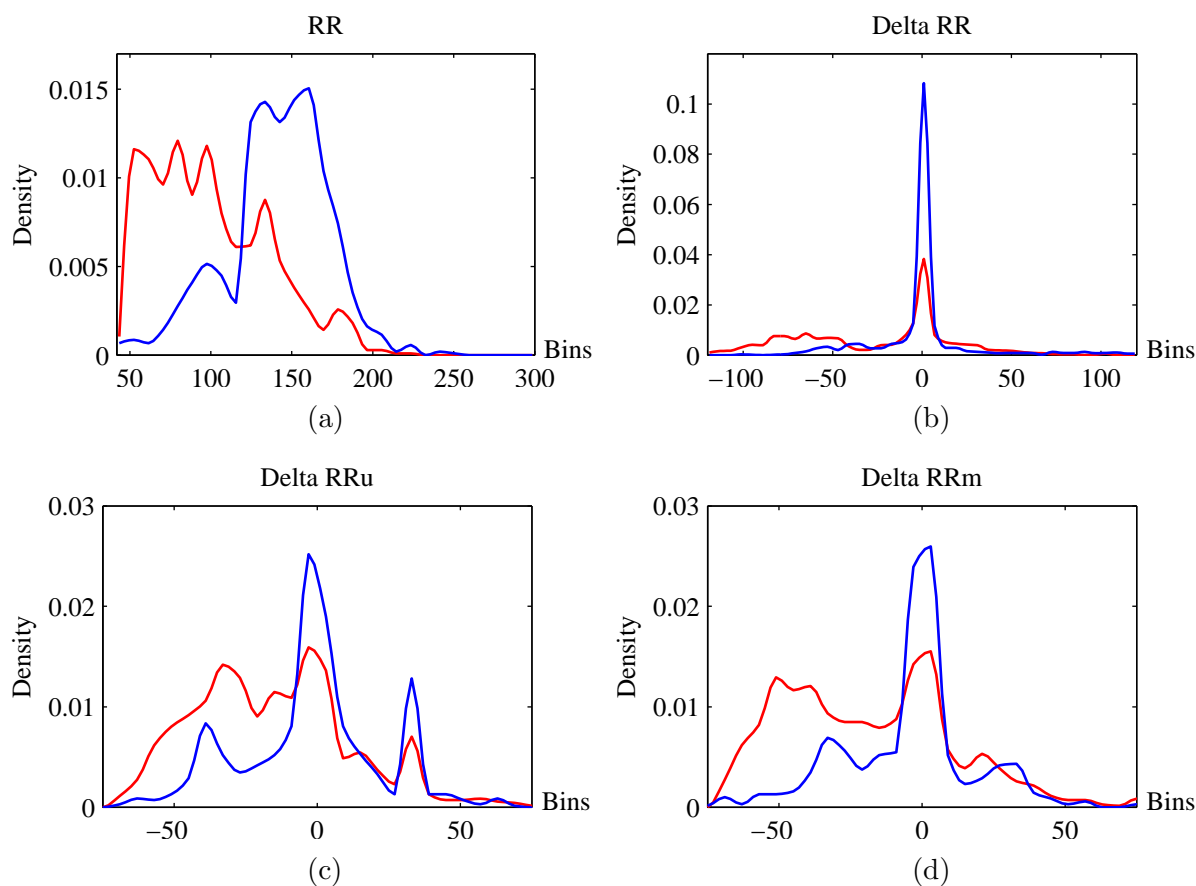


Figure 6.4: Density histograms for RR interval features of ECG and noise: Blue and red lines represent ECG and noise, respectively. (a) to (d) are RR , ΔRR , ΔRRu , and ΔRRm , respectively.

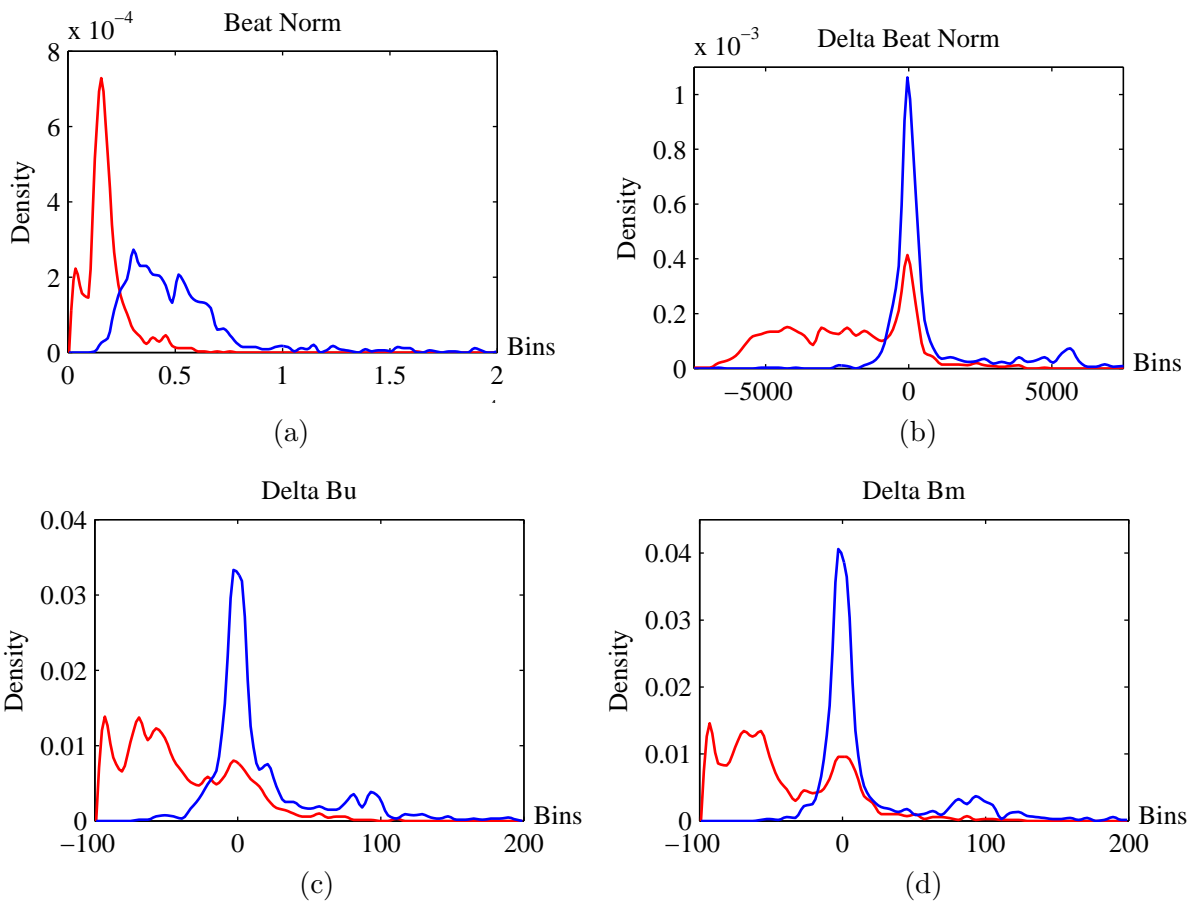
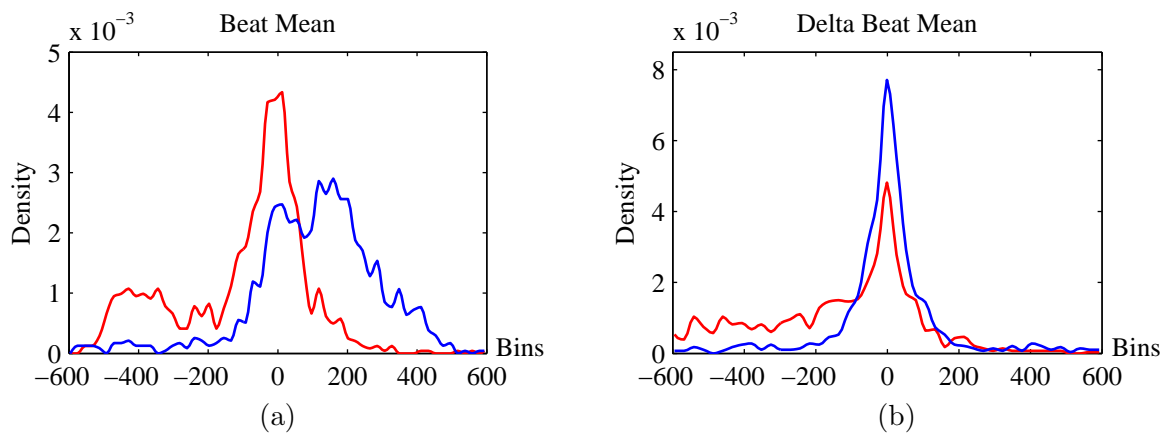


Figure 6.5: Density histograms for beat norm features of ECG and noise. Blue and red lines represent ECG and noise, respectively. (a) to (d) are B , ΔB , ΔBu , and ΔBm , respectively.



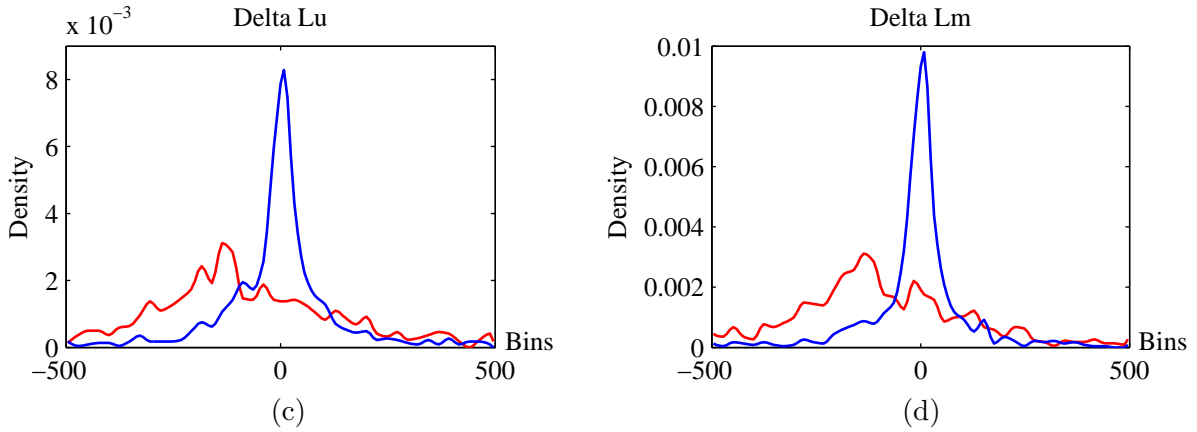


Figure 6.6: Density histograms for beat level features of ECG and noise. Blue and red lines represent ECG and noise, respectively. (a) to (d) are L , ΔL , ΔLu , and ΔLm , respectively.

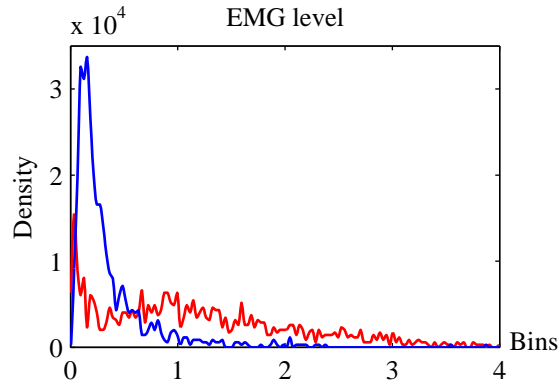


Figure 6.7: Density histograms for beat level features of ECG and noise. Blue and red lines represent ECG and noise, respectively.

New features are extracted for classifying the ECG and noise. They are the number of fluctuations (N_f) and the accumulated fluctuation distance (D_f) for the difference between a normalized beat and its smoother version, respectively. Normalized beat computation is described as follows: a beat cropped by a window expanding to the left and right 20 and 26 ms, respectively, from the beat reference point. Then, the windowed beat is normalized by its beat norm. The smoother version of a normalized beat is calculated by applying the moving average with a specific window size twice on the signal passage containing the beat, windowing at the location of the beat, normalized by the beat norm. The normalized

beat and its smoother version can be formulated as follows. Define $\mathbf{x} = \{\dots, x_i, \dots\}$ as a signal passage where x_i is a sample point. Normalized beat, $\mathbf{b} = \{b_{k-wl}, \dots, b_k, \dots, b_{k+wr}\}$, is expressed in (6.3). k is the location of the beat reference point. wl and wr are left and right window expansions and are set to 20 and 26 ms, respectively. B is the Euclidean norm of the windowed signal. Define $m(\mathbf{x})$ as the moving average of \mathbf{x} which is formulated in (6.4). xm_i is the average value at sample i . ml and mr are left and right expansions of the moving average window, respectively. The smooth version of the normalized beat is denoted as $\mathbf{s} = \{s_{k-wl}, \dots, s_k, \dots, s_{k+wr}\}$ and expressed in (6.5) where ml and mr of $m(\cdot)$ are both 3 ms.

$$\mathbf{b} = \frac{\{x_{k-wl}, \dots, x_k, \dots, x_{k+wr}\}}{B} \quad (6.3)$$

$$m(\mathbf{x}) = \{\dots, xm_i, \dots\}, \text{ where } xm_i = \frac{\sum_{j=i-ml}^{i+mr} x_j}{ml + mr + 1} \quad (6.4)$$

$$\mathbf{s} = \frac{\{y_{k-wl}, \dots, y_k, \dots, y_{k+wr}\}}{B}, \text{ where } \mathbf{y} = \{\dots, y_i, \dots\} = m(m(\mathbf{x})). \quad (6.5)$$

An example of \mathbf{b} (blue line) and \mathbf{s} (red line) are given in Fig. 6.8 (a). Fluctuations of \mathbf{b} over/under \mathbf{s} are measured. Noise beats should have a higher fluctuation than ECG beats do. Normal ECG has QRS complex which produces a high fluctuation. Therefore, the section of QRS complex is omitted. Define $\mathbf{f} = \{b_{k-wl} - s_{k-wl}, \dots, b_{k-ql} - s_{k-ql}, b_{k+qr} - s_{k+qr}, \dots, b_{k+wr} - s_{k+wr}\}$ where ql and qr determine the omitted section and are set to 8 and 11 ms, respectively. Fig. 6.8 (b) illustrates intervals used to determine \mathbf{f} (red-line intervals). The red-dot interval is omitted. The green-line intervals are shown to display the complete beat but are not used in the calculation. Examples of \mathbf{b} and \mathbf{s} are shown in Fig. 6.9. Subfigures (a) to (d) are ECG beats which display normal ECG, R-wave

split, and two premature beats, respectively. Subfigures (e) to (f) display noise beats. \mathbf{f} contains sections of positive and negative amplitudes. Sections whose highest absolute amplitudes are lower than 0.03 are not considered in the calculation, because they contain only insignificant swings of \mathbf{b} over/under \mathbf{s} . N_f is defined as the total number of sections, whereas D_f denotes a summation of the sections' highest absolute amplitudes. Density histograms of N_f and D_f are plotted in Figs. 6.10 (a) and (b), respectively.

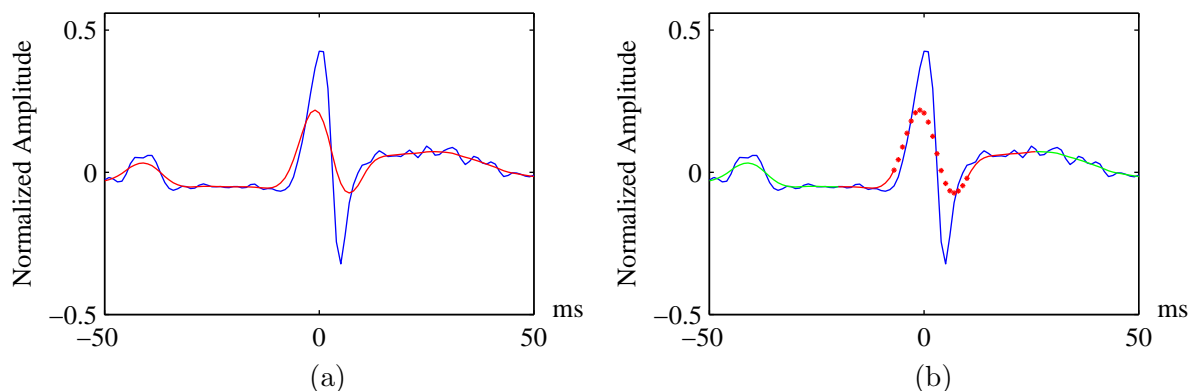
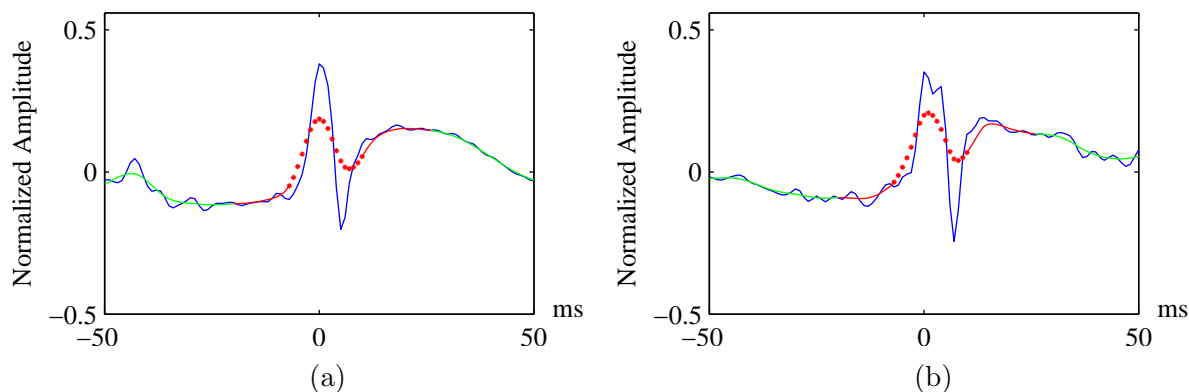


Figure 6.8: An example of a normalized beat and its smoother version. (a) The normalized beat and its smoother version are shown using blue and red lines, respectively. (b) Intervals of the signals used to calculate N_f and D_f are displayed using red lines. The red-dot interval is omitted in the calculation. The intervals represented by the green lines are shown to display the entire beat signal.



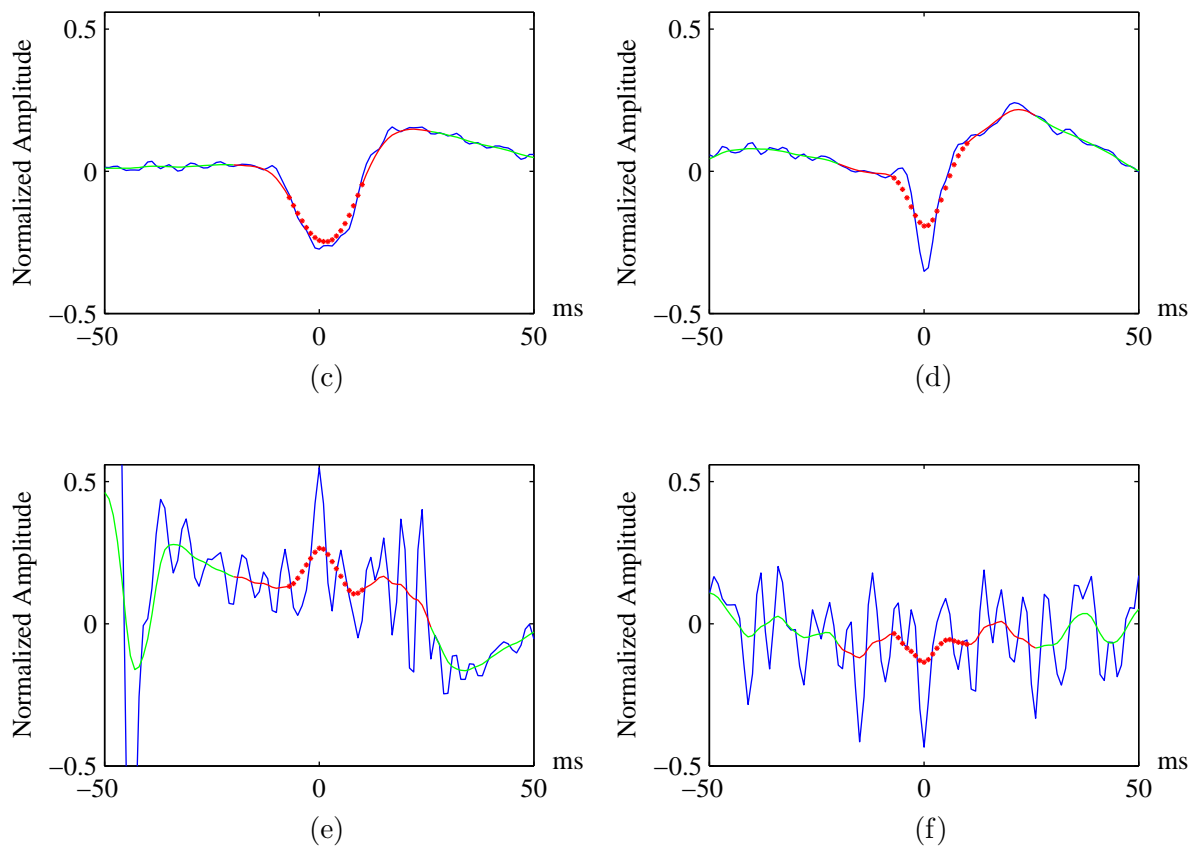


Figure 6.9: Examples of b and s . Blue lines represent b . s is displayed in red lines, red dotted lines, and green lines. (a) shows a normal ECG beat. (b) is an R-wave-split ECG beat. (c) and (d) are premature beats. (e) and (f) are noise beats.

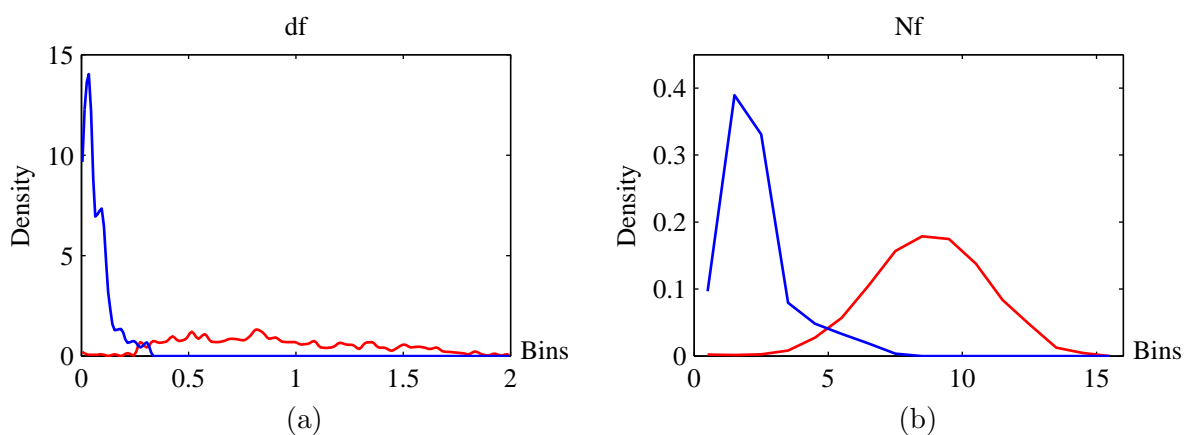


Figure 6.10: Density histograms of N_f and D_f . Blue and red lines represent ECG and noise, respectively. (a) is N_f and (b) is D_f .

D_f shows a significant separation between ECG and noise, while N_f does not. The classification rule is determined from the training set as follows:

If $D_f \leq 0.3124$, beats are classified as ECG.

Otherwise, they are noise.

In the training set, the rule provides a classification rate of 97.83%. Sensitivity for ECG detection is 100% and the specificity is 95.66%. In the test set, the classification rate is 97.87%. The sensitivity and specificity are 100% and 95.74%, respectively. Confusion matrices for the training and test sets are shown in Table 6.4.

Table 6.4: Confusion matrices for training and test sets in ECG and noise classification using D_f .

Training set			Test set		
Actual type	Classified as		Actual type	Classified as	
	ECG	Noise		ECG	Noise
ECG	3500	0	ECG	3500	0
Noise	152	3348	Noise	149	3351

The rule does not misclassify ECG as noise, but there is some noise misclassified as ECG. An additional rule is applied to the training set to achieve a higher classification rate. The new rule is as follows:

If $D_f \leq 0.3124$ and $\Delta Bm > -50.74$, beats are classified as ECG.

Otherwise, they are noise.

In the training set, the classification rate, sensitivity and specificity are 99.04%, 99.89%, and 98.2%, respectively. In the test set, the classification rate, sensitivity and specificity are 98.74%, 99.71%, and 97.77%, respectively. Confusion matrices for the training and test sets are shown in Table 6.5.

Table 6.5: Confusion matrices for training and test sets in ECG and noise classification using D_f and ΔBm .

Training set			Test set		
Actual type	Classified as		Actual type	Classified as	
	ECG	Noise		ECG	Noise
ECG	3496	4	ECG	3490	10
Noise	63	3437	Noise	78	3422

A program for calculating D_f was written. \mathbf{s} can be calculated in a single pass using the moving average in (6.6), although (6.5) expresses \mathbf{s} as two passes of the moving average. In (6.6), x , y , mr , and ml are variables and parameters in (6.3) to (6.5). y_i at the first $2ml+1$ and last $2mr+1$ samples need to be calculated using the standard formula. Sections of \mathbf{f} are found by assigning 1 to points whose amplitudes are higher than 0.03 and -1 to points whose amplitudes are lower than -0.03 . The points with an amplitude between -0.03 and 0.03 are neglected. Section edges are the positions where the points change their value from 1 to -1 or -1 to 1. D_f is the summation of the maximum absolute magnitudes from the sections. The computation time is in proportion to the number of sections, N_f , of \mathbf{f} . From Fig. 6.10, the mode of N_f of the ECG beats is two sections, while the mode of N_f of the noise beats is nine sections.

$$y_i = y_{i-1} + \frac{x_{i+2mr} - x_{i-2ml-1}}{2(mr + ml) + 1} \quad (6.6)$$

6.4 Normal and abnormal ECG classification

Beats detected by the beat detection algorithm are classified as ECG beats and noise beats in the previous section. The ECG beats are further classified into normal and abnormal ECG beats as stated in Section 6.1. In the literature review, automated algorithms for abnormal ECG classification use approaches based on the Karhunen-Loeve transfor-

mation [66, 73], linear prediction [69], filter banks [70], polynomial approximation [75], nonlinear-principal-component-analysis [74], hidden Markov models [71], wavelet transform [72], and neural network [74, 72]. The filter-bank and hidden-Markov-model algorithms give a lower classification rate than other methods. There are in total about 25 million beats per one animal. Therefore, algorithms, which need a high computation effort, such as linear prediction, polynomial approximation, nonlinear-principal-component analysis, and wavelet transformation, are avoided. Note that the neural network approaches in [74, 72] use nonlinear-principal-component analysis and wavelet transformation for feature extraction. The Karhunen-Loeve-transformation approach in [66] detects only abnormal QRS complexes, and the algorithm in [73] has many steps for feature extraction. In this section, low computational algorithms are developed for discriminating abnormal ECG (PVC, elevated ST segment, and split R-wave) from normal ECG. Searching for a group of abnormal ECG helps to decrease the number of ECG beats for consideration. Therefore, complex methods can be applied for further analyses.

The characteristics of the ECG beats are illustrated using density histograms of their features which are shown in Figs. 6.11 to 6.13. Normal ECG features have values more concentrated on one value than those of abnormal ECG, which are spread out. The density histograms for both ECG types are on top of each other and do not show significant separation. The features were evaluated in a preliminary study and did not show potential for classification. For example, feature matrices were formed from all features and classification models were generated based on the Fisher linear discriminant analysis. The resulting classification rates did not exceed 70%. Therefore, the morphology of a beat can be considered for use in the classification. The time series representations, which are explained in the literature review chapter, can reduce the dimension of an ECG beat. However, distance measures for the representations are not simple. The time series repre-

sentations do not gain any benefit if only a small number of dimensions is reduced. The algorithms described below are developed based on the beat morphology without applying the time-series-representation techniques. The Euclidean distance is used for the measure of similarity because of its simplicity.

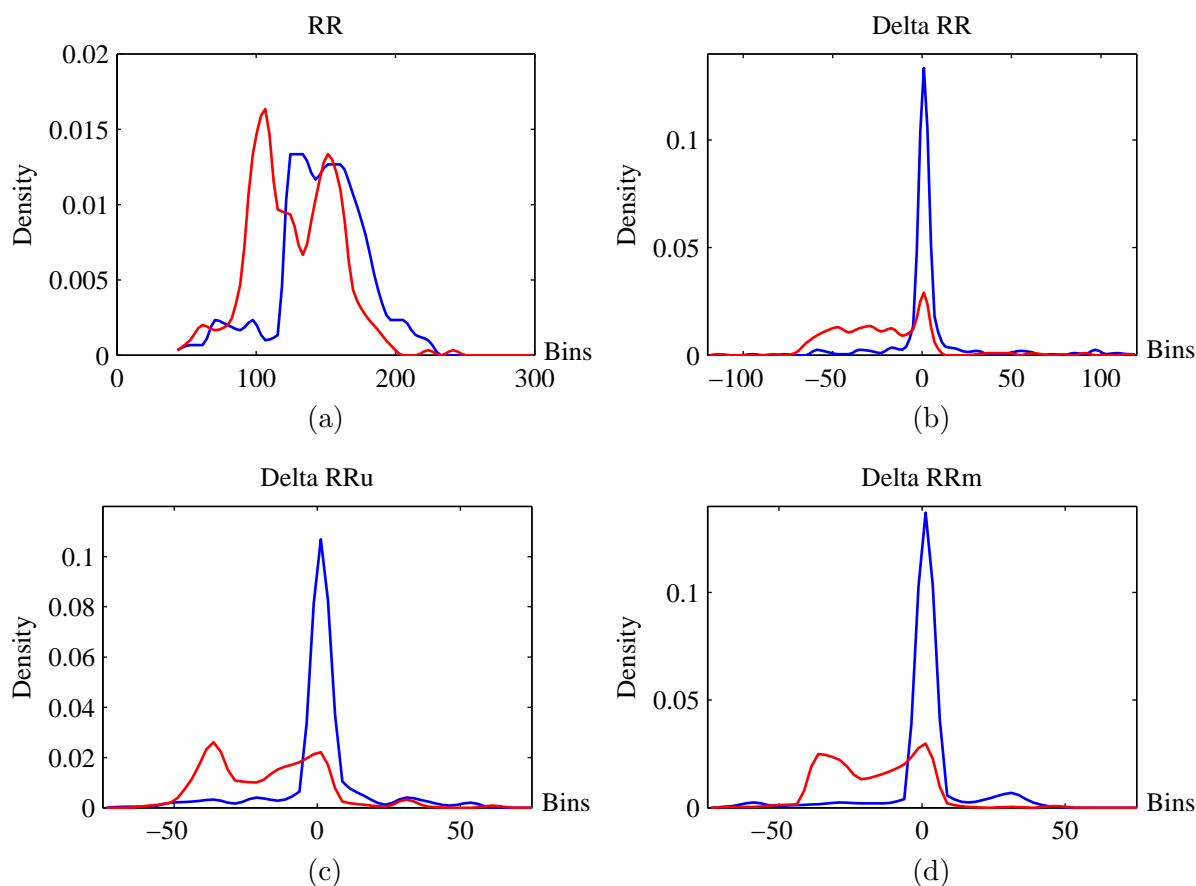


Figure 6.11: Density histograms for RR interval features of normal and abnormal ECG. Blue and red lines represent normal and abnormal ECG, respectively. (a) to (d) are RR , ΔRR , ΔRRu , and ΔRRm , respectively.

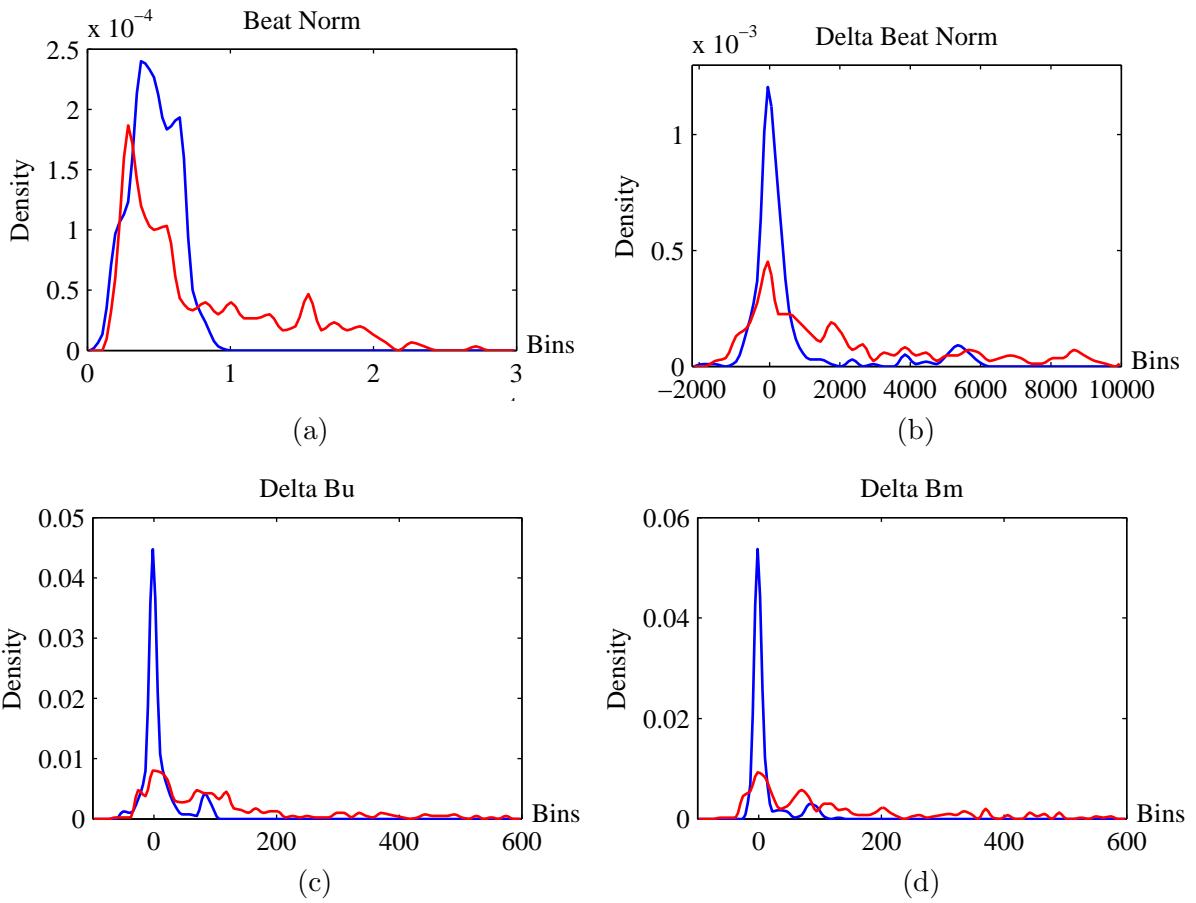
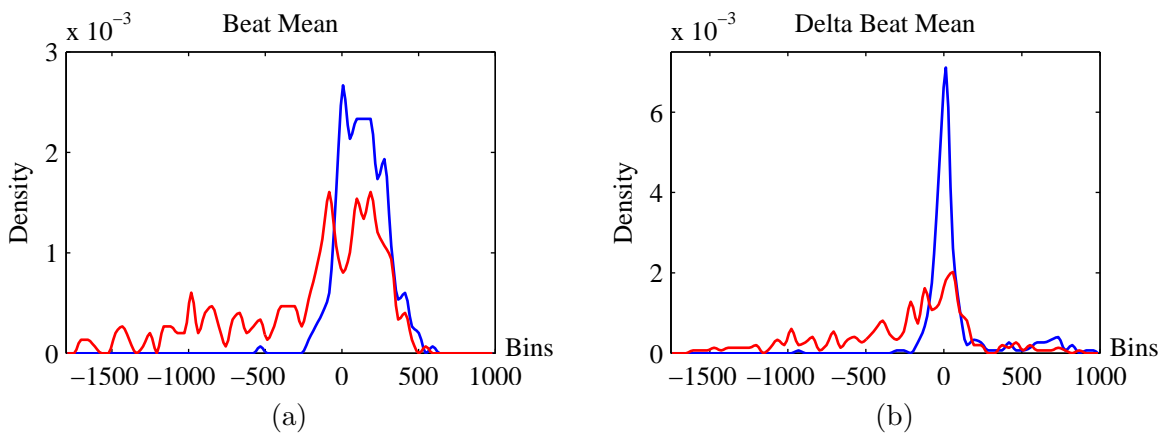


Figure 6.12: Density histograms for beat norm features of normal and abnormal ECG. Blue and red lines represent normal and abnormal ECG, respectively. (a) to (d) are B , ΔB , ΔBu , and ΔBm , respectively.



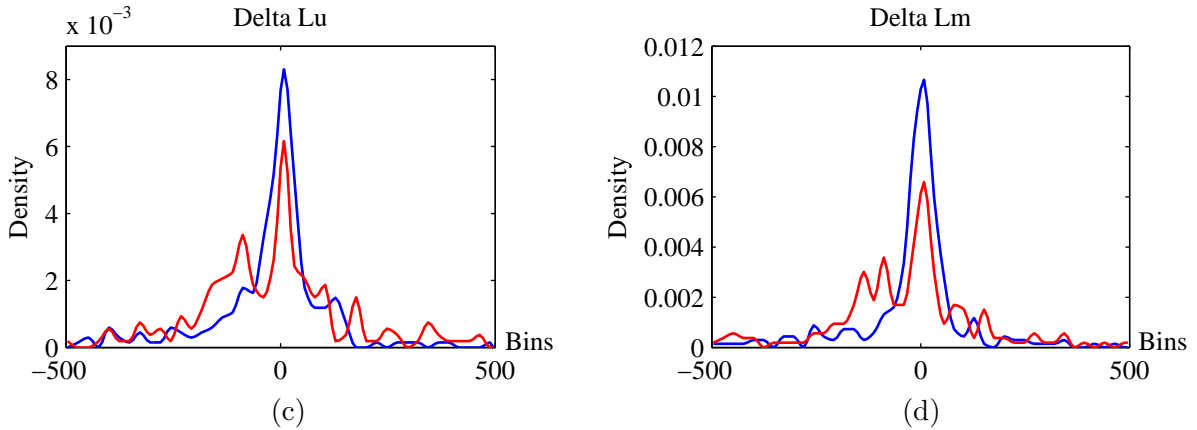


Figure 6.13: Density histograms for beat level features of normal and abnormal ECG. Blue and red lines represent normal and abnormal ECG, respectively. (a) to (d) are L , ΔL , ΔLu , and ΔLm , respectively.

The method is divided into two phases. The first phase is to classify QRS complexes. Normal ECG and ECG with normal QRS width and no elevated ST segment is categorized in class Qn . Besides the normal ECG, split R-wave beats with normal QRS width are included in this class. Other abnormal types (PVC, elevated ST segment, and split-R-wave with wide QRS complexes) are assigned as class Qa . The second phase is to determine R-wave split in class Qn . The class of normal R-wave is named Rn , while Ra is for split-R-wave class.

6.4.1 QRS classification

In this subsection, two feature types are extracted from beat shape. They are individually applied to the Fisher linear discriminant analysis. The resulting classifiers are evaluated. Then, the two sets of features are utilized to generate a model. Lastly, the performances are compared.

Data were selected from the dataset in Section 6.2. Both training and test sets have class Qn and Qa in an equal number of 1000 beats. The total number of beats in each set is 2000. Tables 6.6 and 6.7 provide the details for the datasets.

Table 6.6: Information about the training set used for QRS classification: number of beats for each type and subject

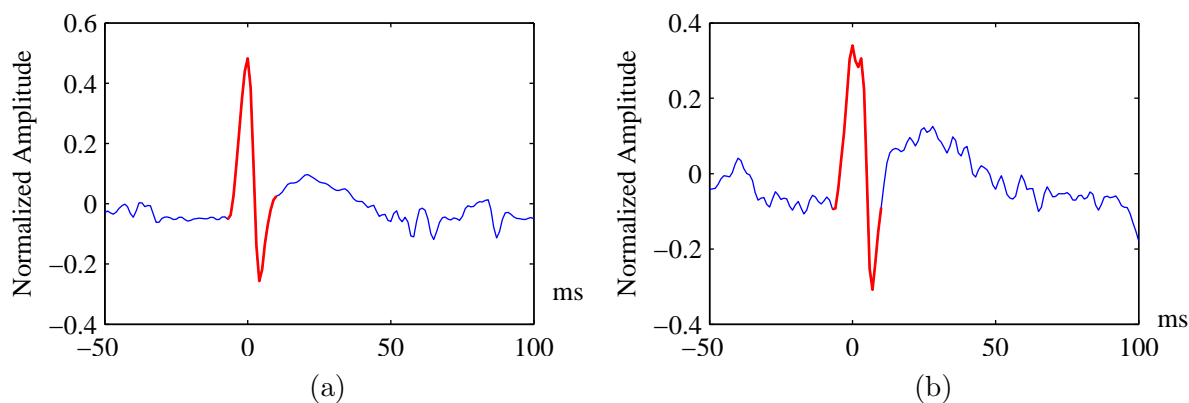
Type	Subject					Total
	1	2	3	4	5	
Qn (Normal ECG)	180	180	180	180	200	1000
Qn (Split R-wave)	20	20	20	20	0	
Qa	359	138	141	291	71	1000
Total	559	338	341	491	271	2000

Table 6.7: Information about the test set used for QRS classification: number of beats for each type and subject

Type	Subject					Total
	1	2	3	4	5	
Qn (Normal ECG)	180	180	180	180	200	1000
Qn (Split R-wave)	20	20	20	20	0	
Qa	359	139	141	291	70	1000
Total	559	339	341	491	270	2000

The beat QRS complex is extracted as a feature by windowing and then normalizing to make it less subjective. The mathematical expression of the feature, \mathbf{q} , is written in (6.7) where wl and wr are 7 and 10 ms, respectively. B_q and L_q are the Euclidean norm and mean of the windowed QRS complex. Fig. 6.14 shows examples of the cropped QRS complex (in red) for various beat types.

$$\mathbf{q} = \frac{\{x_{k-wl} - L_q, \dots, x_k - L_q, \dots, x_{k+wr} - L_q\}}{B_q} \quad (6.7)$$



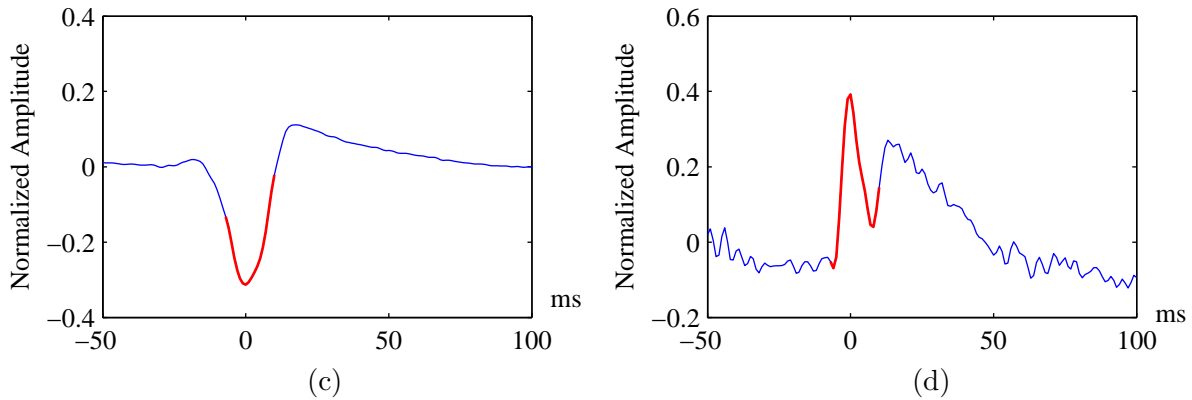


Figure 6.14: Examples of extracted QRS complexes for QRS classification. The blue and red lines represent ECG beats and their extracted QRS complexes, respectively. (a) shows a normal ECG beat. (b) is a split-R-wave beat. (c) is a PVC beat. And, (d) displays a ST elevated beat.

A set of \mathbf{q} was extracted from the training set and the Fisher linear discriminant analysis was applied. The model was evaluated. The results provide a classification rate of 89.75%. Sensitivity and specificity for Qa detection are 85.2% and 88.3%. For the test set, the classification rate, sensitivity, and specificity are 87.65%, 86%, and 89.3%, respectively. The confusion matrices for the training and test sets are shown in Table 6.8. By observation, the majority of cases of misclassified Qa (abnormal classified normal) are ST elevated beats. This may be because ST elevated ECG normally has a QRS complex similar to normal ECG expect for elevated S and T waves. Using only the QRS shape may not be enough to attain a better classification rate.

Table 6.8: Confusion matrices for training and test sets in QRS classification using QRS morphology

Training set			Test set		
Actual QRS class	Classified as		Actual QRS class	Classified as	
	Qa	Qn		Qa	Qn
Qa	852	148	Qa	860	140
Qn	117	883	Qn	107	893

As mentioned in the previous paragraph, new features were devised. They were obtained from the difference between the normalized beat signal, \mathbf{c} , and its smoother version, \mathbf{z} . \mathbf{c} is assigned as in (6.8). w_l and w_r are equally set to 50 ms. B is the beat norm. \mathbf{z} is described in (6.9). m_l and m_r of $m(\cdot)$ are both 30 ms. Examples of \mathbf{c} (blue line) and \mathbf{z} (red line) are displayed in Fig. 6.15. Subfigures (a) to (d) show normal ECG, split R-wave, PVC, and ST elevated beats, respectively.

$$\mathbf{c} = \frac{\{x_{k-w_l}, \dots, x_k, \dots, x_{k+w_r}\}}{B} \quad (6.8)$$

$$\mathbf{z} = \frac{\{y_{k-w_l}, \dots, y_k, \dots, y_{k+w_r}\}}{B}, \text{ where } \mathbf{y} = \{\dots, y_i, \dots\} = m(m(\mathbf{x})). \quad (6.9)$$

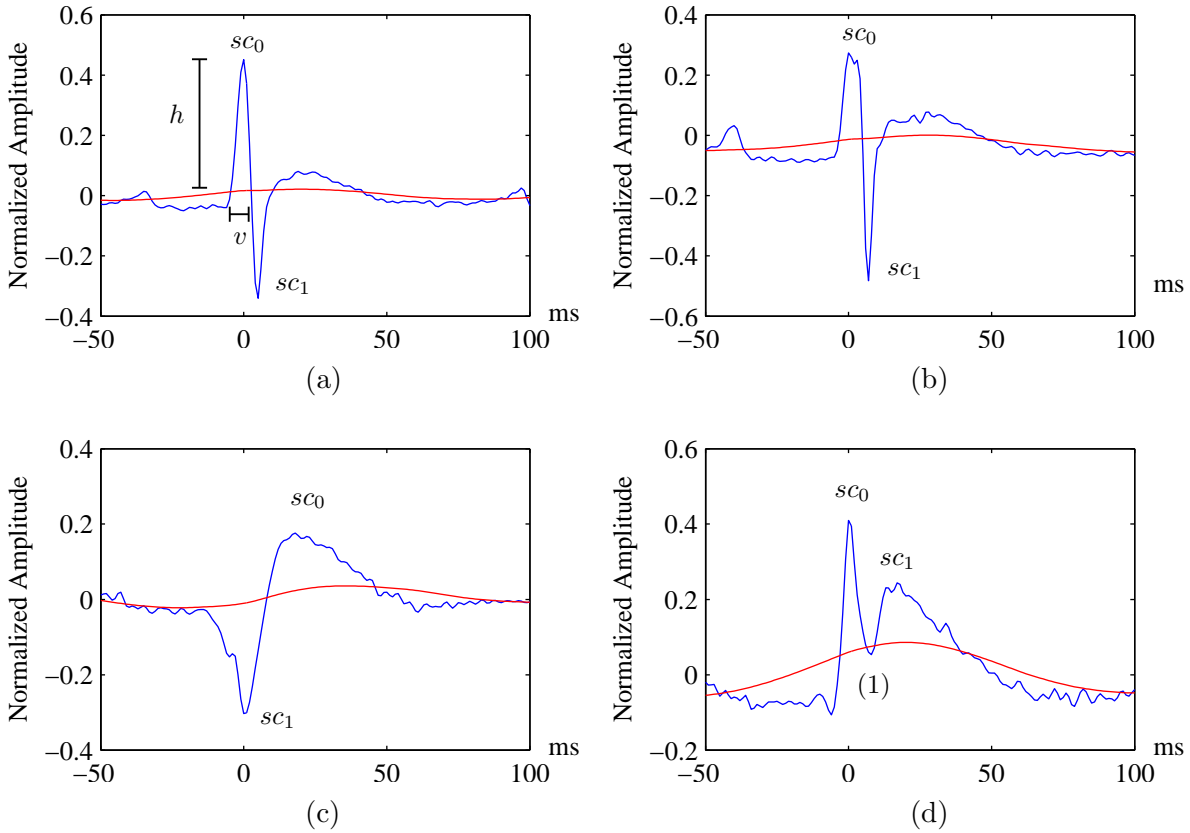


Figure 6.15: Examples of \mathbf{c} and \mathbf{z} for QRS classification. The blue and red lines represent \mathbf{c} and \mathbf{z} , respectively. (a) to (d) show normal ECG, split-R-wave, PVC, ST elevated beats, respectively. Section sc_0 and sc_1 are shown in the subfigures. In (c), (1) is marked as an unconsidered section. Measurement of h and v is illustrated in (a).

Define $\mathbf{g} = \{c_{k-wl} - z_{k-wl}, \dots, c_{k+wr} - z_{k+wr}\}$. \mathbf{g} includes positive and negative amplitude sections. In the other words, swings of \mathbf{c} over or under \mathbf{z} are considered as sections. Consecutive sections with swing amplitudes of less than 0.03 are joined as one section. Section height is defined as h and section base width is referred to as v . Fig. 6.15 illustrates h and v . sc_0 is named for the section where the beat reference point is. The right section of sc_0 is called sc_1 . Sections with absolute values of h less than 0.03 and v less than 4 ms are linked to the previous sections. In subfigure (c), (1) is joined to sc_0 . sc_0 has parameters of h_0 and v_0 , so as sc_1 . Density histograms for h_0 , h_1 , v_0 , and v_1 are plotted in Fig. 6.16.

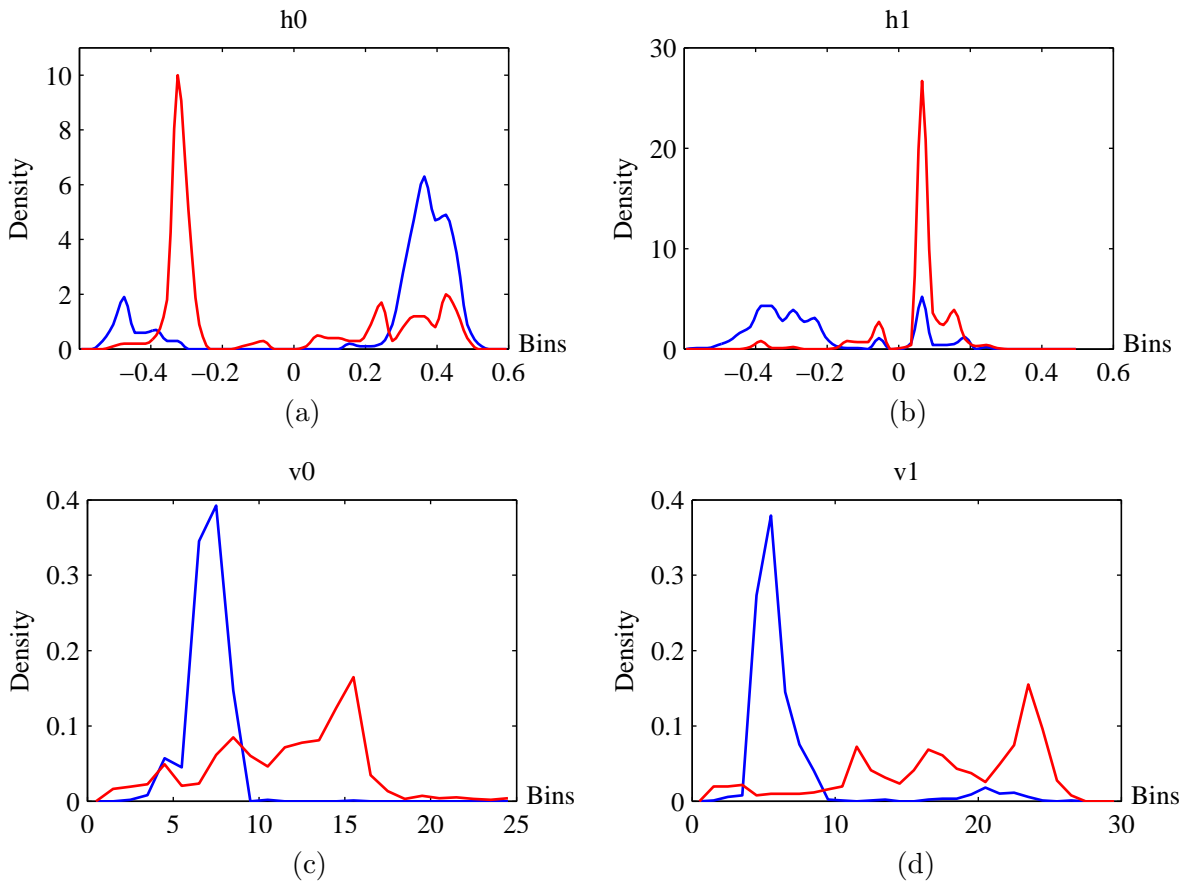


Figure 6.16: Density histogram for h_0 , h_1 , v_0 , and v_1 for QRS classification. The blue and red lines represent class Qn and Qa , respectively. (a) to (d) are h_0 , h_1 , v_0 , and v_1 , respectively.

The Fisher linear discriminant analysis was applied to h_0 , h_1 , v_0 , and v_1 of the training set to generate a classifier model. The model was evaluated on the training and test sets. Table 6.9 shows the resulting confusion matrices. For the training set, the classification rate is 96.75%. Sensitivity and specificity for detecting class Qa are 95.3% and 98.2%, respectively. The test set has a 97.25% classification rate, 95.8 % sensitivity, and 98.7% specificity. Compared to \mathbf{q} , h and v features render a better classification and have a lower dimension.

Table 6.9: Confusion matrices for training and test sets in QRS classification using h_0 , h_1 , v_0 , and v_1 .

Training set			Test set		
Actual QRS class	Classified as		Actual QRS class	Classified as	
	Qa	Qn		Qa	Qn
Qa	953	47	Qa	958	42
Qn	18	982	Qn	13	987

To get a better classification rate, the h and v features were concatenated with \mathbf{q} . The classifier created by the Fisher linear discriminant analysis produces results as shown in Table 6.10. For the training set, the classification rate, sensitivity, and specificity are 97.5%, 96.5%, and 98.5%. They are 98.25%, 97.5%, and 99%, respectively, for the test set. The concatenated feature has a higher dimension than the h and v features. Nevertheless, it does not provide much greater performance. Table 6.11 compares classification rates, sensitivity, and specificity for the three models. It should be noted that the values are taken from the test set. h and v features give a good performance. However, they require more computation when compared to \mathbf{q} .

An implementation of a program for extracting h and v features is similar to the program for calculating D_f . The two-pass moving average for computing \mathbf{z} is determined using (6.6). Points in \mathbf{z} are set to 1 if their amplitudes are higher than 0.03. -1 is assigned

to points with an amplitude lower than -0.03 . For the points with the amplitude between -0.03 and 0.03 , zero is assigned. Section edges are located at the positions where the point value changes. The sections in the vicinity of the beat reference point are considered. The zero sections are linked to the section before them, and the point values are changed according to the point value in the section before them. Sections sc_0 and sc_1 are determined from the beat reference point. Then, h_0 , v_0 , h_1 , and v_1 are measured. A program for the Fisher linear discriminant classifier is simple. The class identifier is a dot product of a feature vector and the vector of the classifier coefficients. If the class identifier is higher than or equal to zero, one class is assigned. Otherwise, the other class is assigned.

Table 6.10: Confusion matrices for training and test sets in QRS classification using h_0 , h_1 , v_0 , v_1 and QRS morphology

Training set			Test set		
Actual QRS class	Classified as		Actual QRS class	Classified as	
	Qa	Qn		Qa	Qn
Qa	965	35	Qa	975	25
Qn	15	985	Qn	10	990

Table 6.11: Comparison of performance of QRS classification using different features. The table shows the classification rate, sensitivity and specificity for Qa detection.

	Features		
	\mathbf{q}	h_0, h_1, v_0, v_1	h_0, h_1, v_0, v_1 and \mathbf{q}
Classification rate (%)	87.65	97.25	98.25
Sensitivity (%)	86	95.8	97.5
Specificity (%)	89.3	98.7	99

6.4.2 Split R-wave classification

This is phase two of the classification where beats in Qn are further categorized as normal or split R-wave. Training and test sets have equal numbers of beats for each type and subject as illustrated in Table 6.12.

Table 6.12: Information about training and test sets used for split R-wave classification: number of beats for each type and subject

Type	Subject					Total
	1	2	3	4	5	
Normal ECG	16	16	16	16	16	80
Split R-wave	20	20	20	20	0	80
Total	36	36	36	36	16	160

Sliding windows, $\{1, -1, 0\}$ and $\{0, -1, 1\}$, are simultaneously applied to the section sc_0 of beat. Down swing of the R-wave vertex is detected, if the dot products of the sliding windows and the beat samples are both greater than the set threshold. The setting threshold is used to avoid a small deflection due to noise, and set to 0.008. The algorithm achieves a 100% classification rate, sensitivity, and specificity on both training and test sets.

Considering the whole process of normal and abnormal ECG classification, the QRS classification using h_0 , h_1 , v_0 , and v_1 is selected as phase one and the split-R-wave classification as phase two. The results from phase one were applied to phase two. The outcome has a classification rate, sensitivity, and specificity of 96.75%, 95.64%, and 98.04%, respectively, for the training set and 97.35%, 96.02%, and 98.91%, respectively, for the test set. The algorithms described above need to be applied in tandem, i.e. ECG and noise classification, and then normal and abnormal ECG classification. For the ECG and noise classification using feature D_f , there are some noise beats classified as ECG, 152 and 149 beats in training and test sets, respectively. These beats are tested using the normal and abnormal ECG classifier described previously. 95.4% and 94% are classified as abnormal ECG in the training and test sets, respectively.

6.5 Discussion

The ECG and noise classification uses D_f to distinguish between ECG and noise. The beat's QRS complex, which fluctuates in a fashion similarly to noise, especially in normal ECG, is excluded from D_f calculation. The position of the beat reference point, which is used to determine the location of the QRS complex, is crucial for classifying ECG and noise. Normal ECG can be incorrectly categorized as noise, if the beat reference point is not placed at R-wave or S-wave. In this case, the beat's QRS complex is included in D_f calculation. As a result, D_f is higher than the classification threshold and beats are recognized as noise. Correct position of the beat reference point is also important in the normal and abnormal ECG classification, since \mathbf{q} , h , and v features are extracted according to the beat reference point. Therefore, the beat detection algorithm and its parameters must be carefully tuned and tried on a portion of data. For ECG and noise classification, the classification threshold for D_f can be increased to gain higher specificity. Less noise will be recognized as abnormal ECG. However, more ECG will be categorized as noise and discarded from consideration. For normal and abnormal ECG classification, the values of h and v features are sensitive to the ECG baseline fluctuation. Therefore, the baseline fluctuation must be removed before calculating h and v features to get correct detection. The baseline-fluctuation-removal algorithm explained in Section 5.5 cannot remove the baseline which swings faster than the T wave. Therefore, h and v features may provide unreliable classification when a fast-swing baseline exists. \mathbf{q} can be used if only PVC needs to be isolated from normal ECG. In Table 6.8, such a misclassification occurred, because shapes of normal ECG and elevated ST segment beats are similar. \mathbf{q} provided a good feature for classification of normal ECG and PVC. In addition, it requires less computation than h and v features.

6.6 Summary

Abnormal ECG search is a part of the framework, which is explained in the previous chapter, for analyzing long-term ECG recordings. The detected beats from the previous chapter were selected for developing algorithms for distinguishing normal from abnormal ECGs. In the study, the abnormal ECG includes PVC, elevated ST segment, and split R-wave. The problem was divided into two stages in sequence. The first is to discriminate normal and abnormal ECG beats from noise beats. Noise beats are movement artifacts, EMG, and burst noise which are misdetected as ECG beats. ECG beats which are highly corrupted by noise were also included in noise beats. New features were extracted from each detected beat. They are the number of fluctuations (N_f) and the accumulated fluctuation distance (D_f) for the difference between a normalized beat and its smoother version. Training and test sets contained 3500 ECG beats and the same number of noise beats. D_f showed a significant distinction between ECG and noise, while N_f does not. The classification rule was determined from the training set as follows: *If $D_f \leq 0.3124$, beats are ECG. Otherwise, they are noise.* In the training set, the classification rate, sensitivity, and specificity for ECG detection were 97.83%, 100%, and 95.66%, respectively. In the test set, the classification rate, sensitivity, and specificity were 97.87%, 100%, and 95.74%, respectively. The percentage difference between the current beat and the median of its surrounding beat norms (ΔBm) was added to the rule as follows: *If $D_f \leq 0.3124$ and $\Delta Bm > -50.74$, beats are ECG. Otherwise, they are noise.* This new rule gave a classification rate, sensitivity, and specificity of 99.04%, 99.89%, and 98.2% in the training set. In the test set, the classification rate, sensitivity, and specificity were 98.74%, 99.71%, and 97.77%, respectively. The second stage is to isolate abnormal ECG beats from normal and abnormal ECG beats. At first, normal ECG and ECG with normal QRS width or no elevated ST segment (class Qn) were distinguished from PVC, elevated ST segment, and split

R-wave with wide QRS complexes (class Qa). It should be noted that class Qn contains normal ECG and split R-wave with normal QRS-complex width. Three sets of features were used in classifying class Qn and Qa . They included a normalized beat-QRS-complex (\mathbf{q}), height (h) and width (v) of sections of the difference between the normalized beat and its smoother version, and a concatenate of \mathbf{q} , h and v . The dimensions of these three types of features are 18, 4, and 22, respectively. Both training and test sets had class Qn and Qa in an equal number of 1000 beats. The Fisher linear discriminant analysis was used to determine a classifier for each feature set. h and v features provided the best performance in the aspect of low dimension and high classification rate. For the training set, the classification rate, sensitivity, and specificity for detecting class Qa were 96.75%, 95.3%, and 98.2%, respectively. The test set had a 97.25% classification rate, 95.8% sensitivity, and 98.7% specificity. After classifying class Qn and Qa , beats in class Qn were further categorized as normal R-wave (class Rn) and split R-wave (class Ra). Training and test sets contained class Rn and Ra in an equal number of 80 beats. Class Ra was detected, if the dot products of sliding windows, $\{1, -1, 0\}$ and $\{0, -1, 1\}$, and the extracted R-wave section are both greater than 0.008. The results provided a 100% classification rate, sensitivity, and specificity on both training and test sets. The entire procedure of stage two was evaluated by using h and v features to categorize class Qn and Qa . The results had a classification rate, sensitivity, and specificity of 96.75%, 95.64%, and 98.04%, respectively, for the training set and 97.35%, 96.02%, and 98.91%, respectively, for the test set.

Chapter 7

Conclusion and future work

Advanced radiotelemetry devices and large storage media allow researchers to study ECG in the long term. Newly discovered ECG patterns may provide benefits in prediction or early detection of heart diseases. The challenge is to analyze hundreds gigabytes of recorded data. Many techniques for ECG analysis were presented in literatures – for example, ECG beat detection, ECG template generation, and ECG classification. Some of the techniques require high computation and memory usage and are not suitable if applied to large ECG data. Therefore, ECG analysis algorithms with low complexity were studied to develop a procedure for analyzing massive ECG data in this dissertation. In previous studies on analyzing massive ECG data, RR intervals, HRV parameters, and ECG templates were extracted from ECG data. However, the algorithms were developed separately and an integrated procedure for data analysis was not established. In addition, the results needed to be analyzed manually using visualization tools. This dissertation presents a framework, which integrates processes from data preparation to data visualization, and a method to locate abnormal ECG. The algorithms were developed on real data. Problems from noise were addressed when developing the framework.

The ECG used in this research was collected in an experiment directed by Dr. Amy de Jongh Curry, Department of Biomedical Engineering, University of Memphis, Memphis TN. Chronic heart failure was gradually induced by aldosterone infusion and a high Na and low Mg diet. The ECG was continuously recorded during the experimental period of 12 weeks through radiotelemetry. ECG leads were placed subcutaneously in a lead-II configuration. Signals were recorded at a sampling rate of 1000 Hz and data resolution of 16 bits. This resulted in 80 GB of data for five animals.

By manually scanning the data, abnormal ECG present in the data two to three weeks after the beginning of the experiments. Many abnormal ECGs were found – including PVC, bigeminy, trigeminy, ventricular tachycardia, elevated ST segment, split R-wave, and abnormal QRS complexes. There was also a lot of noise and artifacts in the data. Abnormal ECG and noise can cause irregularity in extracted features. Therefore, noise must be isolated from abnormal ECG before searching for interesting events in data. The framework developed is a very useful tool for processing massive ECG data. ECG, some abnormal ECG, and noise can be identified automatically. Once ECG can be distinguished from noise and abnormal ECG can be located, ECG patterns can be further studied.

7.1 Summary of work

The entire framework includes data preparation, ECG beat detection, baseline fluctuation removal, EMG noise detection, ECG template generation, feature extraction, and abnormal ECG search. In data preparation, the recorded files are read from the beginning to the end and structurally stored in the database for data retrieval. The ECG signal can be queried by using signal time. Sections of data which have no signal or fullscale signal are marked as unusable data.

An ECG beat detection program with low memory usage and computation presented in [43] was utilized. It has drawbacks including misdetecting severe baseline fluctuations as beats and missing beats following abrupt increase in beat magnitude such as premature beats or spurious noise. The fast morphological filter in [82] was inserted to improve misdetection due to baseline fluctuation. The width of the structuring element is 15 ms. Beats are detected by using an exponentially decayed threshold. An upper bound of the threshold was added to prevent the threshold from being abnormally high when detecting a beat with unusually high amplitude. Therefore, the following normal beats were not missed. Beats were detected at a location in the vicinity of R or S waves. The peak of the R-wave was determined and used as the beat reference point. The valley of the S-wave was searched for and used instead, if the R-wave was very small (the peak was very low). The second case usually happened when abnormal beats were detected. The beat locations were required for several of the following processes including ECG baseline removal, EMG noise elimination, ECG template matching, feature extraction, and abnormal ECG detection.

A low-complexity algorithm in [83] was modified and applied for removing baseline fluctuation. The algorithm uses two morphological filters connected in sequence. The first removes QRS complexes while the second eliminates residual QRS complexes and P/T waves. The width of the first structuring element is set to 7 ms, while the width of the second structuring element is set to the base width of the T-wave. A technique was devised to estimate the width. The power spectral density of the output from the first morphology filter was determined. The frequency with the highest power was converted to time and set as the width of the second structuring element. The last step applied a moving average and resulted in an estimated baseline. The result shows a good baseline estimation in test signals. The limitation is that a baseline which has a higher fluctuation than the T-wave cannot be removed.

Sections of ECG may be interfered with and corrupted by surface EMG which causes difficulties in data processing and analysis. A new algorithm for EMG detection is proposed. EMG is separated from ECG by using the morphological filter with a structuring element width of 7 ms. The output also contains partial QRS complexes. Then, samples 8 ms to the left and 11 ms to the right of the beat reference point are reduced in their magnitude multiplied by 0.1. A sliding window with a size of 81 ms is applied to the filtered signal. Sample variances are calculated as the window slides – called moving variance. Because the magnitude of the recorded ECG signals is subjective, the moving variance is then normalized by the square of average of the R-wave amplitudes. EMG sections are found if the normalized moving variance is higher than 0.01. These sections need to be expanded to 50 ms from the beginning and the end of the sections to capture the entire EMG passage. The algorithm achieves a sensitivity and specificity of 100% and 100% on the training set and 94% and 100% on the test set. The detected EMG sections are eliminated from further consideration.

A procedure for ECG template generation for long-term recording data is introduced. ECG template generation is a procedure which creates a set of ECG templates from the detected beats. Each template represents a group of ECG beats that share the same morphology. Any two ECG templates have a distance measure exceeding the set value. Each beat is stamped with its ECG template and its signal time. ECG beats are transformed to a template – called a beat template. The beat template is generated by windowing a beat, with the average value of the sample amplitudes subtracted, and normalizing by its Euclidean norm. The window expands 20 ms to the left and 26 ms to the right of the beat reference point. The first beat template in the data is automatically assigned as an ECG template and first put in an empty set of ECG templates. Then, each of the following beat templates is sequentially matched to the templates in the ECG-template set. A matching

score is computed as the minimum Euclidean distance regardless of the shift in templates. If the score is lower than the set threshold, 0.014, the ID of the matched ECG template is stamped to the considered ECG beat. If not, the beat template is inserted into the ECG-template set and becomes a new ECG template. A new ID is assigned to the new ECG template. Every time a beat template is matched to an ECG template, it is updated to the ECG template as formulated in (5.11). If a number of matchings in the ECG template reaches a certain number, 100, the update stops. Therefore, this prevents the ECG template from following a gradual change in ECG beat morphology. If the ECG template set contains more than 600 templates, 240 of them are transferred to the database to limit the number of comparisons per beat template. Once all beat templates are processed, duplicate ECG templates in the database are searched by the nearest neighbor search in [84]. The duplicates are deleted.

The following features were extracted: RR intervals, ΔRR , beat EMG level, and beat norm (window extending to the left 20 ms and the right 26ms from the beat reference point). One and two dimension histograms are displayed in Section 5.9. An artifact rejection program in [44] was applied to RR intervals to remove ECG beat detection artifact. The following heart rate variability parameters are calculated for every usable 90-second segment of RR intervals: MEAN, MEDIAN, SDNN, CV, IQR, MUDRR, SDSD, DDIQR, RMSSD, LF, HF, and LHF. A usable segment is defined as a segment containing non-artifact RR intervals for more than 85% of its time interval. The visualization tool in [7, 4] is used to display the results in Section 5.10. The plots establish 24-hour circadian rhythm and progressive changes during the experiment.

Algorithms for abnormal ECG search are devised. The algorithms aim to find abnormal ECG in the set of detected beats. The considered abnormal ECG include PVC, elevated ST segment, and split R-wave. However, the detected beats contain noise besides normal

and abnormal beats. Noise in this case refers to noise or artifact misdetected as beat and ECG beats heavily corrupted by noise. The algorithm is divided into two steps. The first step is discriminating between normal and abnormal ECG beats from noise beats. The second step is distinguishing between abnormal ECG beats and normal ECG beats. In the first step, new features were invented. They are the number of fluctuation (N_f) and accumulated fluctuation distance (D_f) of difference between normalized beat and its smoother version, respectively. Normalized beats are beats which are windowed 20 ms and 26 ms to the left and right, respectively, from the beat reference point and normalized by their Euclidean norm. Their smoother versions are calculated by applying a moving average with a window size of 6 ms twice on the signal. ECG beats were selected from the data to create training and test sets. Both sets contained 3500 noise beats and the same number of ECG beats. A decision rule was generated from the training set as follows. *If $D_f \leq 0.3124$, beats are classified as ECG. Otherwise, they are noise.* Note that N_f did not provide a good classification. It achieved a 97.83% and 97.87% classification rate on the training and test sets, respectively. A feature, ΔBm , was added to the rule. ΔBm is the percentage difference between norm of the current beat and median of beat norms from the four surrounding beats. The new rule is *If $D_f \leq 0.3124$ and $\Delta Bm > -50.74$, beats are classified as ECG. Otherwise, they are noise.* This rule achieved a 99.04% and 98.74% classification rate on the training and test sets, respectively. In the second step, the classified ECG beats are categorized in normal and abnormal ECG. This step is divided into two phases. The first phase is detecting PVC and elevated ST beats. Split-R-wave beats are classified in the second phase. In the first phase, three sets of features were used in the Fisher linear discriminant analysis and the results were compared. The first set is morphology of the QRS complex portion, \mathbf{q} . The feature is a window expanding left 7 ms and right 10 ms from the beat reference point. New features were created in the second

set. They are derived from the difference between the normalized beat and its smoother version. Both of them are computed in the way described above except as follows. The beat window is 50 ms to the left and right of the beat reference point. The previous beat norm is used. The moving average window has a length of 60 ms. Sections of the difference between the normalized beat and its smoother version were found as described in Section 6.4. h and v were defined as section height and base width, respectively. A set of h and v was extracted and used as the features. In the last feature set, \mathbf{q} , h and v features were concatenated. The results showed the features in set two provided the best performance in aspect of low dimension and high classification rate. The classification rates are 96.75% and 97.25% in the training (1000 normal beats and 1000 abnormal beats) and test sets (same as the training set), respectively. In the second phase of the abnormal ECG classification, split R-wave was determined by applying two windows and detecting the down flip of R wave. The training and test sets have 80 normal beats and the same number of split-R-wave beats. The results had classification rates of 100% on both data sets.

7.2 Contributions of this dissertation

The specific contributions that have been made in this dissertation are listed below:

- A technique for abnormal ECG search is introduced. PVC, elevated ST segment, and split R-wave can be found in a long-term recording by applying the Fisher linear discriminant analysis on new features, h and v . The algorithm requires low computational effort and can be used with massive ECG recordings. h and v features provide a high detection rate and have a low dimension. Moreover, the features reflect width and height of QRS complexes which are used in diagnosis by medical doctors.

- Decision rules for determining misdetections were formed based on a new feature, D_f , which is simple to compute. D_f is not subject to the recording amplitude setting. In addition, it has a potential for use in determining the noise corrupting level of ECG beats.
- An EMG detection algorithm in the ECG signal was invented for a long-term recording ECG. Signal sections which are corrupted by EMG can be eliminated before processing the signal.
- The beat detection in [43] was improved to address the drawbacks which include misdetecting severe baseline fluctuations as beats and missing beats following an abrupt increase in beat magnitude such as premature beats or spurious noise.
- An algorithm was devised to estimate the width of the structuring element which is required in the baseline removal algorithm in [83].
- A suggested procedure for ECG template generation in large ECG data was presented. Unlike the previous algorithms in [15, 8], the procedure includes a method to generate ECG templates so that gradual change in ECG morphology can be captured. In addition, a solution is proposed to manage numerous generated ECG templates.
- By using the above developments, a framework for analysis of massive ECG data is proposed.
- Noise and problems from data recording were explored and stated.

7.3 Future work

There is a much room for research in analysis of long-term ECG recording. Possible topics are listed below:

- There are other types of abnormal ECG that are yet not considered or found in this data. The developed algorithms can be used as a preliminary tool to find those abnormalities. h and v features can be further investigated for their utility and limitations. On the other hand, new techniques and new features can also be researched.
- h and v features can be studied for their potential as parameters for querying beats.
- Level of contaminated noise in ECG should be studied and quantified, because noise effects measurements from ECG signals. It will be useful if the relation between the reliability of measurement and noise level is determined. D_f and N_f can be investigated on this topic.
- The ECG template generation procedure can be integrated with the parallel algorithms in [8] to increase the performance.
- Abnormalities are not only expressed in ECG morphology. Abnormality search for RR intervals and heart rate variability parameters is a very interesting research.
- The ultimate goal of analyzing ECG of an animal model of heart failure is to discover patterns that link to heart disease or life-threatening arrhythmias. Pattern discovery in massive ECG data and associating the patterns to heart physiological measurements (such as blood test and echocardiogram) are also possible research.

References

- [1] D. Lloyd-Jones and et al., “Heart disease and stroke statistics–2009 update: A report from the american heart association statistics committee and stroke statistics subcommittee,” *Circulation*, vol. 119, pp. e21–e181, Jan. 2008. [Online]. Available: <http://circ.ahajournals.org/cgi/reprint/CIRCULATIONAHA.108.191261> 1
- [2] Centers for Disease Control and Prevention, “Heart disease,” March 2009. [Online]. Available: <http://www.cdc.gov/heartdisease/statistics.htm> [Accessed: March 10, 2009] 1
- [3] J. Bai, Y. Zhang, D. Shen, L. Wen, C. Ding, Z. Cui, F. Tian, B. Yu, B. Dai, and J. Zhang, “A portable ECG and blood pressure telemonitoring system,” *IEEE Engineering in Medicine and Biology Magazine*, vol. 18, no. 4, pp. 63–70, July-Aug. 1999. 1
- [4] S. Schuckers, P. Raphisak, T. Yan, and M. Schuckers, “Development of an experimental method for long-term electrocardiographic recordings in a heart failure rabbit model,” in *Computers in Cardiology*, Sept. 2002, pp. 333–336. 1, 5.10, 7.1
- [5] P. Raphisak, A. D. Curry, R. A. Malkin, and S. C. Schuckers, “Heart rate variability in rats with aldosterone-induced chronic heart failure,” in *Engineering in Medicine and Biology Society, 2003. Proceedings of the 25th Annual International Conference of the IEEE*, vol. 1, Sept. 2003, pp. 228–231. 1, 4.1.3
- [6] S. G. Crihalmeanu, “Representative ways to analyze and survey changes in long-term electrocardiographic recordings,” Master’s thesis, West Virginia University, WV, USA, 2000. 1, 4.1.2, 4.1.3, 4.3
- [7] T. Yan, “Analysis of long-term electrocardiographic data,” Problem Report, West Virginia University, WV, USA, 2002. 1, 4.1.3, 4.3, 5.10, 7.1
- [8] S. Kratsas, “Parallelization of ECG template-based abnormality detection,” Master’s thesis, West Virginia University, WV, USA, 2000. 1, 4.1.3, 4.2.1, 4.3, 7.2, 7.3
- [9] Y. Sun, F. J. Ramires, and K. T. Weber, “Fibrosis of atria and great vessels in response to angiotensin II or aldosterone infusion,” *Cardiovascular Research*, vol. 35, no. 1, pp. 138–47, July 1997. 1, 3.3, 3.4

- [10] K. M. Jones, *Interpretation of the Electrocardiogram : A Review for Health Professionals*. New Jersey: Appleton & Lange, 1990. 2.1, 2.2.1
- [11] L. Sherwood, *Human physiology : from cells to systems*, 4th ed. California: Brooks/Cole, 2001. 2.1, 7, 3.1, 1, 2
- [12] A. L. Goldberger, *Clinical Electrocardiography, A Simplified Approach*. Missouri: Mosby, 1999, ch. 3, pp. 21–31. 2.1.1, 2.1.2, 2.2.1
- [13] D. H. Bennett, *Cardiac Arrhythmias*, 4th ed. Oxford: Butterworth-Heinemann, 1993, ch. 11, pp. 107–114. 2.2, 2.2.2, 2.2.3
- [14] M. B. Canover, *Understanding Electrocardiography, Arrhythmias and the 12-lead ECG*. Missouri: Mosby, 1992, ch. 9, pp. 132–148. 2.2.1
- [15] J. M. Jenkins and S. A. Caswell, “Detection algorithm in implantable cardioverter defibrillators,” *Proceeding of the IEEE*, vol. 84, no. 3, pp. 428–445, March 1996. 2.2.2, 4.2.1, 7.2
- [16] J. A. Kastor, *Arrhythmias*. Pennsylvania: Saunders, 1993. 2.2.3, 2.2.4
- [17] B. M. Beasley, *Understanding EKGs: A Practical Approach*. New Jersey: Prentice Hall, 1999. 2.2.4, 2.2.5
- [18] D. B. Barnett, H. Pouleur, and G. S. Francis, *Congestive Cardiac Failure : Pathophysiology and Treatment*. New York: Marcel Dekker, Inc., 1993. 2.3
- [19] R. Soufer, *Heart Failure*. New York: Hearst Books, 1992, ch. 14, pp. 177–183. 5, 6
- [20] W. Alexander, V. Fuster, and E. Sonnenblick, *Hurst’s the Heart, Arteries and Vein*, 9th ed. USA: McGraw-Hill, 1998. 6
- [21] E. S. Antezano and M. Hong, “Sudden cardiac death,” *Journal of Intensive Care Medicine*, vol. 18, no. 6, pp. 313–329, 2003. 2.4
- [22] Z. Ori, G. Monir, J. Weiss, X. Sayhouni, and D. H. Singer, “Heart rate variability frequency domain analysis,” *Cardiology Clinics*, vol. 10, no. 3, pp. 499–533, 1992. 2.5, 2.5.2
- [23] J. J. T. Bigger, *Spectral Analysis of RR Variability to Evaluate Autonomic Physiology and Pharmacology and to Predict Cardiovascular Outcomes in Humans*. Pennsylvania: WB Saunders, 1994, ch. 101, pp. 1151–1170. 2.5
- [24] J. T. Bigger, P. Albrecht, R. C. Steinman, L. M. Rolnitzky, J. L. Fleiss, and R. J. Cohen, “Comparison of time- and frequency domain-based measures of cardiac parasympathetic activity in holter recordings after myocardial infarction,” *American Journal of Cardiology*, vol. 64, pp. 536–538, 1989. 2.5

- [25] G. G. Berntson, T. Bigger, D. L. Eckberg, P. Grossman, P. G. Kaufmann, M. Malik, H. N. Nagaraja, S. W. Porges, J. P. Saul, P. H. Stone, and M. W. V. D. Molen, "Heart rate variability: Origins, methods, and interpretive caveats," *Psychophysiology*, vol. 34, pp. 623–648, 1997. 2.5
- [26] J. P. Spiers, B. Silke, U. McDermott, R. G. Shanks, and D. W. G. Harron, "Time and frequency domain assessment of heart rate variability: A theoretical and clinical appreciation," *Clinical Autonomic Research*, vol. 3, pp. 145–158, 1993. 2.5
- [27] D. A. Litvack, T. F. Oberlander, L. H. Carney, , and J. P. Saul, "Time and frequency domain methods for heart rate variability analysis: A methodological comparison," *Psychophysiology*, vol. 32, pp. 492–504, 1995. 2.5
- [28] X. Xu, "Prediction of life-threatening events in infants using heart rate variability measurements," Ph.D. dissertation, West Virginia University, WV, USA, 2002. 2.5.1, 2.5.2
- [29] R. E. Kleiger, P. K. Stein, M. S. Bosner, and J. N. Rottman, "Time domain measurements of heart rate variability," *Cardiology Clinics*, vol. 10, no. 3, pp. 487–498, 1992. 2.5.1
- [30] M. Malik and A. J. Camm, *Heart Rate Variability*. New York: Futura Publishing Company, Inc., 1995. 2.5.2, 5.10
- [31] M. Kuwahara, K.-I. Yayou, K. Ishii, S.-I. Hashimoto, H. Tsubone, and S. Sugano, "Power spectral analysis of heart rate variability as a new method for assessing autonomic activity in the rat," *Journal Electrocardiology*, vol. 27, no. 4, pp. 333–7, Oct. 1994. 2.5.2, 5.10
- [32] H. V. Huikuri, K. M. Kessler, E. Terracall, A. Castellanos, and M. K. Linnaluoto, "Reproducibility and circadian rhythm of heart rate variability in healthy subjects," *The American Journal of Cardiology*, vol. 65, pp. 391–393, 1990. 2.6
- [33] H. V. Huikuri, M. J. Niemela, S. Ojala, A. Rantala, M. J. Ikaheimo, and K. E. J. Airaksinen, "Circadian rhythms of frequency domain measures of heart rate variability in healthy subjects and patients with coronary artery disease," *Circulation*, vol. 90, no. 1, pp. 121–126, 1994. 2.6
- [34] K. T. Weber, "Aldosterone in congestive heart failure," *New England Journal of Medicine*, vol. 345, no. 23, pp. 1689–97, Dec. 2001. 3.1, 3.2, 3.3
- [35] P. Lijnen and V. Petrov, "Induction of cardiac fibrosis by aldosterone," *Journal of Molecular and Cellular Cardiology*, vol. 32, no. 6, pp. 865–79, June 2000. 3.2, 3.3
- [36] A. Gonzalez, B. Lopez, R. Querejeta, and J. Diez, "Regulation of myocardial fibrillar collagen by angiotensin II. a role in hypertensive heart disease?" *Journal of Molecular and Cellular Cardiology*, vol. 34, no. 12, pp. 1585–93, Dec. 2002. 3.2

- [37] C. G. Brilla, R. Pick, L. B. Tan, J. S. Janicki, and K. T. Weber, “Remodeling of the rat right and left ventricles in experimental hypertension,” *Circulation Research*, vol. 67, no. 6, pp. 1355–64, Dec. 1990. 3.4
- [38] J. G. Hardman, L. E. Limbird, and A. G. Gilman, *Goodman & Gilman’s the Pharmacological Basis of Therapeutics*, 10th ed. New York: McGraw-Hill Medical, 2001. 4.1.1
- [39] W. I. Ganz, K. S. Sridhar, S. S. Ganz, R. Gonzalez, S. Chakko, and A. Serafini, “Review of tests for monitoring doxorubicin-induced cardiomyopathy,” *Oncology*, vol. 53, pp. 461–470, 1996. 4.1.1
- [40] D. Langton, B. Jover, B. P. McGrath, and J. Ludbrook, “Cardiovascular responses to graded treadmill exercise during the development of doxoribicin induced heart failure in rabbits,” *Cardiovascular Research*, vol. 24, pp. 959–968, 1990. 4.1.1
- [41] S. A. Bocherens-Gadient, U. Quast, J. Nussberger, H. R. Brunner, and R. P. Hof, “Chronic adriamycin treatment and its effect on the cardiac beta-adrenergic system in the rabbit,” *Journal of Cardiovascular Pharmacology*, vol. 19, pp. 770–778, 1992. 4.1.1
- [42] R. B. Wanless, I. S. Anand, P. A. Poole-Wilson, and P. Harris, “An experimental model of chronic cardiac failure using adriamycin in the rabbit: Central haemodynamics and regional blood flow,” *Cardiovascular Research*, vol. 58, pp. 18–24, 1996. 4.1.1
- [43] R. MacDonald, J. Jenkins, R. Arzbaecher, and R. Throne, “A software trigger for intracardiac waveform detection with automatic threshold adjustment,” in *Computers in Cardiology*, 1989, pp. 167–170. 4.1.3, 4.2.1, 4.2.5, 5.4, 5.4, 5.11, 7.1, 7.2
- [44] X. Xu and S. Schuckers, “Automatic detection of artifacts in heart period data,” *Journal of Electrocardiology*, vol. 34 Supplement, pp. 205–210, 2001. 4.1.3, 5.10, 7.1
- [45] A. V. Viatcheslev, L. Shiel, J. Oliver, and B. P. McGrath, “Power spectral analysis of heart rate variability reflects the level of cardiac autonomic activity in rabbits,” *Journal of Autonomic Nervous System*, vol. 21, pp. 7–13, 1996. 4.1.3
- [46] T. Yan, “Visualization of long-term electrocardiographic data,” Problem Report, West Virginia University, WV, USA, 2002. 4.1.3, 5.11
- [47] J. C. Liechty, D. K. J. Lin, and J. P. McDermott, “Single-pass low-storage arbitrary quantile estimation for massive datasets,” Pennsylvania State University, PA, USA, Tech. Rep. 01-04, July 2001. 4.1.3
- [48] C. Dong, “Building an electrocardiogram signal database with Oracle 8,” Problem Report, West Virginia University, WV, USA, 1999. 4.1.3

- [49] L. Guo, "Development of a web-based, multimedia database for collection, organization and analysis of biomedical signals," Problem Report, West Virginia University, WV, USA, 2000. 4.1.3
- [50] S. A. Caswell, "Reliable signal detection of ventricular fibrillation in intracardiac electrograms for precision of therapeutic choice," Ph.D. dissertation, University of Michigan, MI, USA, 1997. 4.2.1
- [51] Wikipedia, "Linear discriminant analysis," March 2009. [Online]. Available: http://en.wikipedia.org/wiki/Linear_discriminant_analysis [Accessed: March 10, 2009] 4.2.2
- [52] V. France and V. Hlavac, "Statistical pattern recognition toolbox for matlab users guide," Czech Technical University, Czech Republic, Tech. Rep., June 2004. 4.2.2
- [53] C. M. Bishop, *Neural Networks for Pattern Recognition*, 1st ed. USA: Oxford University Press, 1996. 4.2.2
- [54] E. Keogh, "A fast and robust method for pattern matching in time series database," in *Proceedings of the 9th International Conference on Tools with Artificial Intelligence (ICTAI-97)*, 1997, pp. 578–584. 4.2.3
- [55] E. J. Keogh and M. J. Pazzani, "Scaling up dynamic time warping to massive dataset," in *PKDD '99: Proceedings of the Third European Conference on Principles of Data Mining and Knowledge Discovery*, 1999, pp. 1–11. 4.1, 4.2.3, 4.2, 4.2.4, 4.3, 4.2.4
- [56] E. Keogh and M. Pazzani, "Derivative dynamic time warping," in *Proceedings of the First SIAM International Conference on Data Mining, Chicago*, 2001. 4.2.3
- [57] S. Chu, E. Keogh, D. Hart, and M. Pazzani, "Iterative deepening dynamic time warping," in *Proceedings of the Second SIAM International Conference on Data Mining*, 2002. 4.2.3
- [58] E. Keogh and M. Pazzani, "An enhanced representation of time series which allows fast and accurate classification, clustering and relevance feedback," in *Proceedings of the Fourth International Conference of Knowledge Discovery and Data Mining*, 1998, pp. 239–241. 4.2.4
- [59] E. Keogh, K. Chakrabarti, S. Mehrotra, and M. Pazzani, "Locally adaptive dimensionality reduction for indexing large time series databases," in *Proceeding of ACM SIGMOD Conference, Santa Barbara*, 2001, pp. 151–162. 4.2.4
- [60] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra, "Dimensionality reduction for fast similarity search in large time series databases," *Knowledge and Information Systems*, vol. 3, no. 3, pp. 263–286, 2000. 4.2.4
- [61] K. Chan and A. W. Fu, "Efficient time series matching by wavelets," in *Proceedings of the 15th International Conference on Data Engineering*, 1999, pp. 126–133. 4.2.4

- [62] K. Kanth, D. Agrawal, and A. Singh, "Dimensionality reduction for similarity searching in dynamic databases," in *Proceedings of SIGMOD Conference*, 1998, pp. 166–176. 4.2.4
- [63] E. Keogh, "Fast similarity search in the presence of longitudinal scaling in time series databases," in *Proceedings of the Ninth IEEE International Conference on Tools with Artificial Intelligence*, 1997, pp. 578–584. 4.2.4, 4.3, 4.2.4
- [64] E. Keogh, S. Chu, D. Hart, and M. Pazzani, "An online algorithm for segmenting time series," in *Proceedings IEEE International Conference on Data Mining*, 2001, pp. 289–296. 4.2.4
- [65] M. L. Kejariwal, "A QRS detection algorithm for discriminating artifacts in ECG records," in *Bioengineering Conference, 1989., Proceedings of the 1989 Fifteenth Annual Northeast*, March 1989, pp. 227–228. 4.2.5
- [66] G. B. Moody and R. G. Mark, "QRS morphology representation and noise estimation using the karhunen-loeve transform," in *Computers in Cardiology*, Sept. 1989, pp. 269–272. 4.2.6, 4.2.7, 4.3, 6.3, 6.4
- [67] F. Gritzali, G. Frangakis, and G. Papakonstantinou, "Noise estimation in ECG signals," in *Engineering in Medicine and Biology Society, 1988. Proceedings of the Annual International Conference of the IEEE*, vol. 1, Nov. 1988, pp. 152–153. 4.2.6, 4.3, 6.3
- [68] C. Brouse, G. Dumont, F. Herrmann, and J. M. Ansermino, "A wavelet approach to detecting electrocautery noise in the ECG," in *Engineering in Medicine and Biology Society, 2005. Proceedings of the 27th Annual International Conference of the IEEE*, Sept. 2005, pp. 3876–3880. 4.2.6, 4.3, 6.3
- [69] K. P. Lin and W. H. Chang, "QRS feature extraction using linear prediction," *IEEE Transactions on Biomedical Engineering*, vol. 36, no. 10, pp. 1050–1055, 1989. 4.2.7, 4.3, 6.4
- [70] O. Wieben, V. X. Afonso, and W. J. Tompkins, "Classification of premature ventricular complexes using filter bank features, induction of decision trees and a fuzzy rule-based system," *Medical and Biological Engineering and Computing*, vol. 37, no. 1, pp. 560–565, Jan. 1999. 4.2.7, 4.3, 6.4
- [71] W. T. Cheng and K. L. Chan, "Classification of electrocardiogram using hidden markov models," in *Engineering in Medicine and Biology Society, 1998. Proceedings of the 20th Annual International Conference of the IEEE*, vol. 20, no. 1, 1998, pp. 143–146. 4.2.7, 4.3, 6.4
- [72] O. T. Inan, L. Giovangrandi, and G. T. A. Kovacs, "Robust neural-network-based classification of premature ventricular contractions using wavelet transform and timing interval features," *IEEE Transactions on Biomedical Engineering*, vol. 53, no. 12, pp. 2507–2515, 2006. 4.2.7, 4.3, 6.4

- [73] A. Smrdel and F. Jager, "Automated detection of transient ST-segment episodes in 24h electrocardiograms," *Medical and Biological Engineering and Computing*, vol. 42, no. 3, pp. 303–311, May 2004. 4.2.7, 4.3, 6.4
- [74] T. Stamkopoulos, N. Maglaveras, K. Diamantaras, and M. Strintzis, "Ischemic classification techniques using an advanced neural network algorithm," in *Computers in Cardiology*, Sept. 1997, pp. 351–354. 4.2.7, 4.3, 6.4
- [75] G. Y. Jeong and K. H. Yu, "Morphological classification of ST segment using reference STs set," in *Engineering in Medicine and Biology Society, 2007. Proceedings of the 29th Annual International Conference of the IEEE*, Aug. 2007, pp. 636–639. 4.2.7, 4.3, 6.4
- [76] J. Marhoul, "Linear prediction: A tutorial review," *Proceedings of the IEEE*, vol. 63, pp. 561–580, 1975. 4.3
- [77] —, "Correction to "linear prediction: A tutorial review"," *Proceedings of the IEEE*, vol. 64, p. 285, 1976. 4.3
- [78] J. A. Cadzow, B. Baseghi, and T. Hsu, "Singular-value decomposition approach to time series modelling," *IEE Proceedings F Communications, Radar and Signal Processing*, vol. 130, pp. 202–210, 1983. 4.3
- [79] S. Mallat, *A Wavelet Tour of Signal Processing (Wavelet Analysis and Its Applications)*, 2nd ed. USA: Academic Press, 1999. 4.3
- [80] Data Sciences International, *Dataquest - User Manual Guide*, Data Sciences International, MN, USA, 1994. 1
- [81] C. H. Chu and E. J. Delp, "Impulsive noise suppression and background normalization of electrocardiogram signals using morphological operators," *IEEE Transactions on Biomedical Engineering*, vol. 2, no. 2, pp. 262–273, Feb. 1989. 5.3, 1
- [82] D. Wang and D.-C. He, "A fast implementation of 1-d grayscale morphological filters," *IEEE Transactions on Circuits and Systems II*, vol. 41, no. 9, pp. 634–636, Sept. 1994. 5.3, 1, 5.6.2, 7.1
- [83] P. Sun, Q. H. Wu, A. M. Weindling, A. Finkelstein, and K. Ibrahim, "An improved morphological approach to background normalization of ECG signals," *IEEE Transactions on Biomedical Engineering*, vol. 50, no. 1, pp. 117–121, Jan. 2003. 5.5, 5.11, 7.1, 7.2
- [84] C. Merkwirth, U. Parlitz, and W. Lauterborn, "Fast exact and approximate nearest neighbor searching for nonlinear signal processing," *Physical Review*, vol. E 62, no. 2, pp. 2089–2097, 2000. 5.8.4, 7.1

- [85] —, “TSTOOL - a software package for nonlinear time series analysis,” in *International Workshop on Advanced Black-Box Techniques for Nonlinear Modeling*, Katholieke Universiteit, Leuven, Belgium, July 1998, pp. 144–146. 5.8.4