Technical University Of Crete
Department of Electronic and Computer Engineering

# QRS Complex Detection in long term Electrocardiograms (ECG) using Swarm Intelligence based Algorithms

## Thomas G. Ntafoulis

Thesis for the acquisition of the diploma of Electronic and Computer Engineer

**Comittee**

Konstantinos KALAITZAKIS (supervisor)
Professor

Michalis ZERVAKIS
Professor

Eftichios KOUTROULIS
P.D. 407/80

**Chania, October 2009**

# Acknowledgments

This thesis was submitted for the acquisition of the diploma of Electronic and Computer Engineering.

I would like to thank the supervisor Professor Mr. Kalaitzakis Konstantinos for his invaluable support, encouragement, supervision and for providing useful suggestions throughout this thesis. Additionally, many thanks towards Professor Mr. Zervakis Michalis and Doctor Mr. Eftichios Koutroulis for reading this thesis has to be spoken.

Moreover, MSc student mr. Papadakis Konstantinos has to be mentioned as a significant partner during the development and implementation of this thesis.

My family and friends deserve special thanks for the continuous support throughout my academic life in Chania.

**Abstract**

The subject of this thesis is the implementation of an efficient method to detect QRS complexes in real ECG recordings. Specifically, anatomy of the human heart along with a theoretical background of the interpretation of a typical ECG is presented. Two algorithms were selected for clustering. The first one is categorized in Ant Colony Optimization Algorithms and is called Ant Colony Optimization Clustering. The second one is the popular K-means clustering algorithm. An extensive description of both algorithms is presented along with the results from tests using synthetic data sets. A method of transforming an ECG in a more smooth and easy-to-read form is described. This method is applied to real ECG signals and the resulting data are clustered using both algorithms. Results of the time cost and QRS complex detection accuracy are given.

# Contents

# Chapter 1

# Introduction

In medical science there is always the need for better tools in the section of diagnosis. One of the most significant organs in the human body, the heart, is a subject of research. In order to diagnose a heart disease an electrocardiogram (ECG) recording the myocardium electrical activity on the body surface, is used. ECG is a periodic signal. Typically, the various features presented on an ECG are labeled using the letters P, Q, R, S, and T. A diagnosis is based on features extracted from the timing and morphology of such findings. Therefore, ECG detection is very important for the doctors as a guide to correct clinical diagnosis[1].

The detection of QRS complex on an ECG signal has been a subject of research for the past three decades. The most important pieces of information on an ECG signal can be found during the P wave, the QRS complex and the T wave[2][3]. These include the positions and/or magnitude of the PR, QRS, QT and ST intervals and the PR and the ST segments. The detection of QRS complexes is a difficult task because of various reasons, such as a

noise signal and power-line interference. To properly evaluate an ECG, the aforementioned problems must be overcome.

There have been several studies dealing with QRS complex detection for ECG signals. Pan and Tompkins [4] proposed an algorithm (the so-called PT method) to recognize QRS complexes. Also, the Wavelet Transforms (WT) has been proposed as method for detecting QRS complexes[5]. The Geometrical Matching Approach algorithm has been proposed to detect the ECG beat[6]. Based on the estimation of the first-order derivative, the SVW algorithm has also been proposed[7].

In this work, ECG is transformed with a method that keeps the morphological characteristics and additionally provides noise reduction. Candidate QRS complexes are evaluated through that process and the need for discriminating real complexes was created.

There is a large variety of clustering algorithms with different characteristics. In terms of this thesis two algorithms were used: Ant Colony Optimization Clustering and K-means clustering.

In **chapter 2** the anatomy of the human heart and its function are described. **Chapter 3** summarizes the particulars of ECG signal recording while in **chapter 4** two methods of clustering are proposed and compared. In **chapter 5** a method of transformation of ECG signals for QRS complex detection is proposed. The application of this method on real ECG signals is detailed in **Chapter 6**. Finally in **chapter 7** conclusions and future work are presented.

# Chapter 2

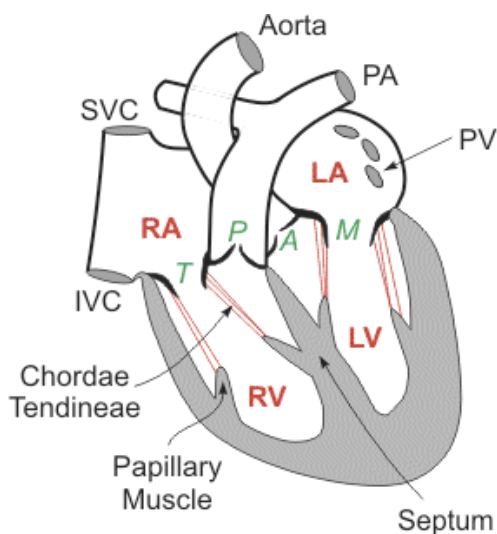# The human heart

## 2.1 Introduction

The human heart is a muscular organ responsible for pumping blood through the blood vessels by periodic contractions. The term cardiac (as in cardiology) means "related to the heart" and comes from the Greek word, *"kardia"*, for heart.

The heart of a vertebrate is composed of cardiac muscle, an involuntary striated muscle tissue which is found only within this organ. The average human heart, beating at 72 beats per minute, will beat approximately 2.5 billion times during a lifetime. It weighs on average 250 g to 300 g in women and 300 g to 350 g in men (Fig.2.1)[8].

## 2.2 Anatomy of the human heart

Venous blood enters the **right atrium** (RA) of the heart through the **superior vena cava** (SVC) and **inferior vena cava** (IVC)(Fig.2.1). The right atrium has a relatively thin muscular wall and easily expands with blood as it fills. Because of its high compliance, the RA pressure is normally very low $(0 - 3$ mmHg). It also undergoes spontaneous contractions to aid in the filling of the **right ventricle** (RV). Blood passes from the RA to the RV through the tricuspid valve. The free wall of the right ventricle is not as thick as the left ventricle, and anatomically it wraps itself around part of the larger, and thicker, left ventricle. The RV wall, however, is thicker and more muscular than the RA, so that when it contracts, it can develop considerably more pressure ($\sim 25$ mmHg) than the RA. As the RV contracts and generates pressure, blood leaves the RV, flows across an open semilunar pulmonic valve, and enters the pulmonary artery that distributes the output of the right ventricle to the lungs where exchange of oxygen and carbon dioxide occur.

The pulmonic valve, like all healthy heart valves, permits blood to flow in only one direction. Blood returns to the heart from the lungs through four pulmonary veins that enter the **left atrium** (LA). This chamber is similar to the RA in that it is very distensible, although the blood pressure within the LA is several mmHg higher than the RA ($6 - 10$ mmHg in the LA compared to $0 - 3$ mmHg in the RA).

Figure 2.1: **The human heart.**

Blood flows from the LA, across the mitral valve, and into the **left ventricle** (LV). The LV wall is very thick so that it can generate high pressures when it contracts (normally $\sim$ 120 mmHg at rest. ). When the LV contracts, blood is expelled through the semilunar aortic valve and into the aorta, which then distributes blood to the arterial system[9]

The tricuspid and mitral valves (also called *atrioventricular*, or *AV valves*) have fibrous strands (*chordae tendineae*) on their leaflets that attach to papillary muscles located on the respective ventricular walls. The papillary muscles contract during ventricular contraction and generate tension on the valve leaflets via the chordae tendineae to prevent the AV valves from bulging back into the atria and becoming incompetent. The semilunar valves (*pulmonic* and *aortic*) do not have analogous attachments[10].
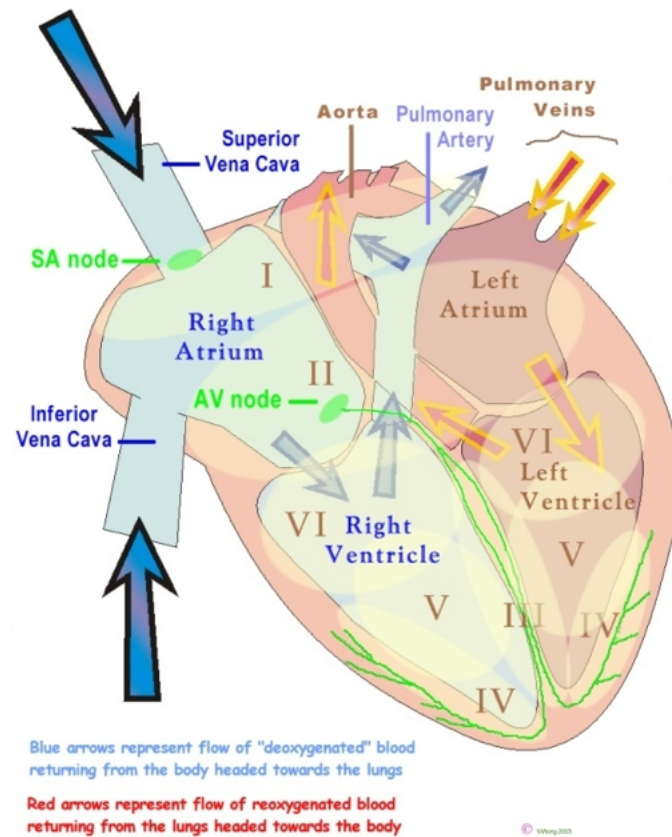
## 2.3   The Cardiac Cycle

A single cycle of cardiac activity can be divided into two basic phases. The first phase is diastole, which represents ventricular filling and a brief period just prior to filling at which time the ventricles are relaxing. The second phase is systole, which represents cardiac contraction and ejection of blood from the ventricles[10].

To analyze these two phases in more detail, the cardiac cycle is usually divided into *seven* stages[11].

**The *seven* stages of the cardiac cycle in brief:**

- Stage 1 - Atrial Contraction

- Stage 2 - Isovolumetric Contraction

- Stage 3 - Rapid Ejection

- Stage 4 - Reduced Ejection

- Stage 5 - Isovolumetric Relaxation

- Stage 6 - Rapid Filling

- Stage 7 - Reduced Filling

Figure 2.2: **The Cardiac Cycle.**

## The Seven Stages

### Atrial Contraction (Stage 1)

This is the first stage of the cardiac cycle which represents electrical depolarization of the atria. Atrial depolarization then causes contraction of the atrial musculature. As the atria contract, the pressure within the atrial chambers increases, which forces more blood flow across the open atrioventricular (AV) valves, leading to a rapid flow of blood into the ventricles. Blood does not

flow back into the vena cava because of inertial effects of the venous return and because the wave of contraction through the atria moves toward the AV valve thereby having a "milking effect" (Fig.2.3)[12].
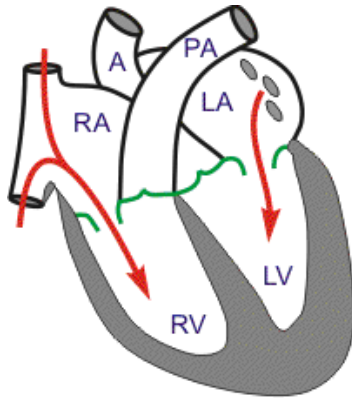


Figure 2.3: ***Atrial contraction.***

Atrial contraction normally accounts for about 10% of left ventricular filling when a person is at rest because most of ventricular filling occurs prior to atrial contraction as blood passively flows from the pulmonary veins, into the left atrium, then into the left ventricle through the open mitral valve.

At high heart rates, however, the atrial contraction may account for up to 40% of ventricular filling. This is sometimes referred to as the "atrial kick." The atrial contribution to ventricular filling varies inversely with duration of ventricular diastole and directly with atrial contractility[13].

**Isovolumetric Contraction (Stage** 2)

**All Valves Closed**

This stage of the cardiac cycle begins with the triggering of excitation-contraction coupling, myocyte contraction and a rapid increase in intraventricular pressure. Early in this stage, the rate of pressure development becomes maximal. This is referred to as **maximal dP/dt**[12].

During the time period between the closure of the AV valves and the opening of the aortic and pulmonic valves, ventricular pressure rises rapidly without a change in ventricular volume (i.e., no ejection occurs).
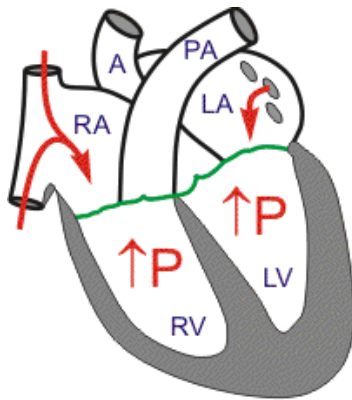


Figure 2.4: **Isovolumetric contraction.**

Ventricular volume does not change because all valves are closed during this stage. Contraction, therefore, is said to be "isovolumic" or "isovolumetric". Individual myocyte contraction, however, is not necessarily isometric because individual myocyte are undergoing length changes.

Therefore, ventricular chamber geometry changes considerably as the heart becomes more spheroid in shape; circumference increases and atrial base-to-apex length decreases (Fig.2.4)[13].

**Rapid Ejection (Stage 3)**

This stage represents the initial and rapid ejection of blood into the aorta and pulmonary arteries from the left and right ventricles, respectively. Ejection begins when the intraventricular pressures exceed the pressures within the aorta and pulmonary artery, which causes the aortic and pulmonic valves to open. Blood is ejected because the total energy of the blood within the ventricle exceeds the total energy of blood within the aorta. In other words,

there is an energy gradient to propel blood into the aorta and pulmonary artery from their respective ventricles. During this stage, ventricular pressure normally exceeds outflow tract pressure by a few mmHg. This pressure gradient across the valve is ordinarily low because of the relatively large valve opening (i.e., low resistance). Maximal outflow velocity is reached early in the ejection stage, and maximal (systolic) aortic and pulmonary artery pressures are achieved[12].
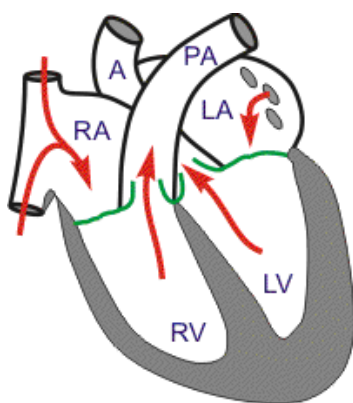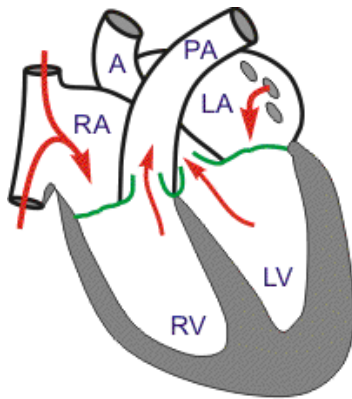


Figure 2.5: ***Rapid ejection.***

Left atrial pressure initially decreases as the atrial base is pulled downward, expanding the atrial chamber. Blood continues to flow into the atria from their respective venous inflow tracts and the atrial pressures begin to rise, and continue to rise until the AV valves open at the end of stage 5 (Fig.2.5)[13]

**Reduced Ejection (Stage** 4**)**

Approximately 200 msec after the beginning of ventricular contraction, ventricular repolarization occurs. Repolarization leads to a decline in ventricular active tension and therefore the rate of ejection (ventricular emptying) falls[12].

Figure 2.6: **Reduced ejection.**

Ventricular pressure falls slightly below outflow tract pressure; however, outward flow still occurs due to kinetic (or inertial) energy of the blood. Left atrial and right atrial pressures gradually rise due to continued venous return from the lungs and from the systemic circulation, respectively (Fig.2.6)[13]

**Isovolumetric Relaxation (Stage 5)**

**All Valves Closed**

When the intraventricular pressures fall sufficiently at the end of stage 4, the aortic and pulmonic valves abruptly close (aortic precedes pulmonic). Valve closure is associated with a small backflow of blood into the ventricles and a characteristic notch n the aortic and pulmonary artery pressure tracings[12].

After valve closure, the aortic and pulmonary artery pressures rise slightly (dicrotic wave) following by a slow decline in pressure.
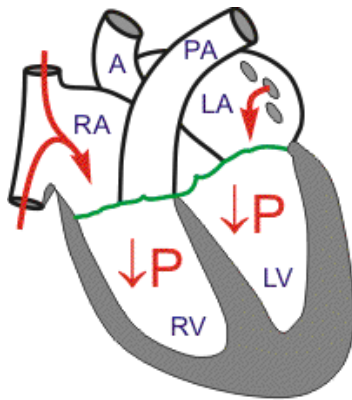
Figure 2.7: **Isovolumentric relaxation.**

The rate of pressure decline in the ventricles is determined by the rate of relaxation of the muscle fibers, which is termed lusitropy. This relaxation is reg-. ulated largely by the sarcoplasmic reticulum that are responsible for rapidly re-sequestering calcium following contraction (Fig.2.7)[13]

**Rapid Filling (Stage** 6**)**

**A-V Valves Open**

As the ventricles continue to relax at the end of stage 5, the intraventricular pressures will at some point fall below their respective atrial pressures.
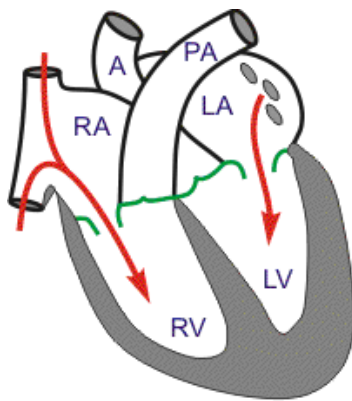


Figure 2.8: **Rapid filling.**

When this occurs, the AV valves rapidly open and ventricular filling begins[12].

Despite the inflow of blood from the atria, intraventricular pressure continues to briefly fall because the ventricles are still undergoing relaxation.

Once the ventricles are completely relaxed, their pressures will slowly rise as they fill with blood from the atria (Fig.2.8)[13].

**Reduced Filling (Stage 7)**

**A-V Valves Open**

As the ventricles continue to fill with blood and expand, they become less compliant and the intraventricular pressures rise. This reduces the pressure gradient across the AV valves so that the rate of filling falls[12].
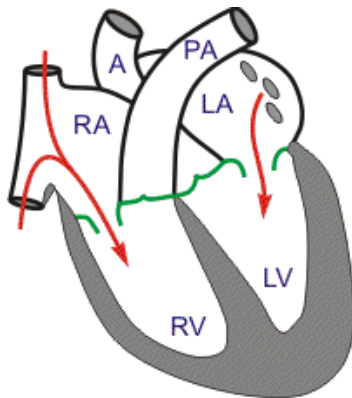


Figure 2.9: ***Reduced filling (diastasis).***

In normal, resting hearts, the ventricle is about 90% filled by the end of this stage. In other words, about 90% of ventricular filling occurs before atrial contraction (stage 1)[13]. Aortic pressure and pulmonary arterial pressures continue to fall during this period (Fig.2.9).

# Chapter 3

# The Electrocardiogram (ECG)

## 3.1 Introduction

An Electrocardiogram (ECG) is composed of a series of waves ordered into some repeatable pattern. The height of the tracing represents millivolts while the width of the ECG macs a time interval. An ECG is composed of a series of waves, intervals and segments.
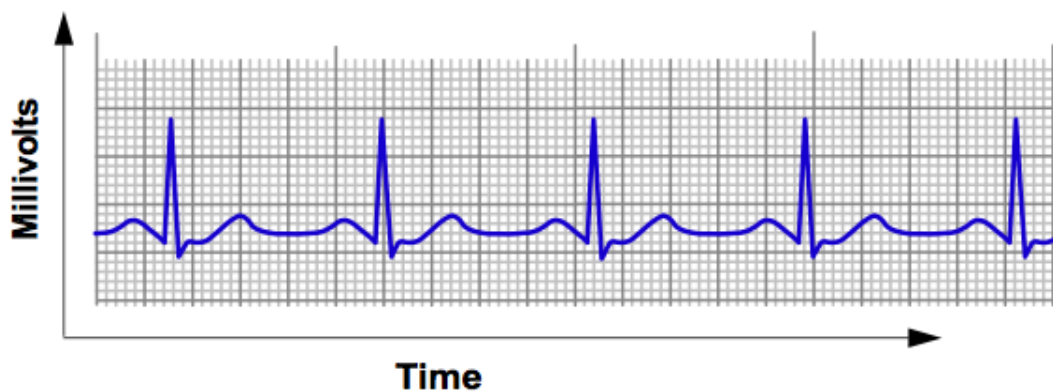


Figure 3.1: *A typical ECG*

- A wave is every deflection on the ECG.

- A segment is the region between two waves.

- An interval includes one segment and one or more waves.

While waves are fairly self-explanatory, intervals measure time from the start of one wave to the start of another wave (an interval includes at least one wave) and segments measure time between waves (waves are not included in a segment)[14].

Timed interpretation of an ECG was once incumbent to a stylus and paper speed. Computational analysis now allows considerable study of the heart rate variability. A typical electrocardiograph runs at a paper speed of 25 mm/s, although faster paper speeds are occasionally used. Each small block of ECG paper is 1 mm$^2$. At a paper speed of 25 mm/s, one small block of ECG paper translates into 0.04 s (or 40 ms). Five small blocks make up one large block, which translates into 0.20 s (or 200 ms). Hence, there are 5 large blocks per second. A diagnostic quality 12 lead ECG is calibrated at 10 mm/mV, so 1 mm translates into 0.1 mV. A calibration signal should be included with every record. A standard signal of 1 mV must move the stylus vertically 1 cm, that is two large squares on ECG paper[15].

## 3.2 Origin of electrical current in the heart

### Flow of Electrical Current

Typically, the heart is located in the middle of the chest to the left of the mediastinum. The sinoatrial (SA) node is located in the top of the right

atrium, the atrioventricular (AV) node is located in the bottom of the atrium, and the bundle branches conduct through the septum and ventricles. Because of this normal flow, the direction of the electrical flow (vector) is mainly downward, from right to left (Fig.3.2)[14].
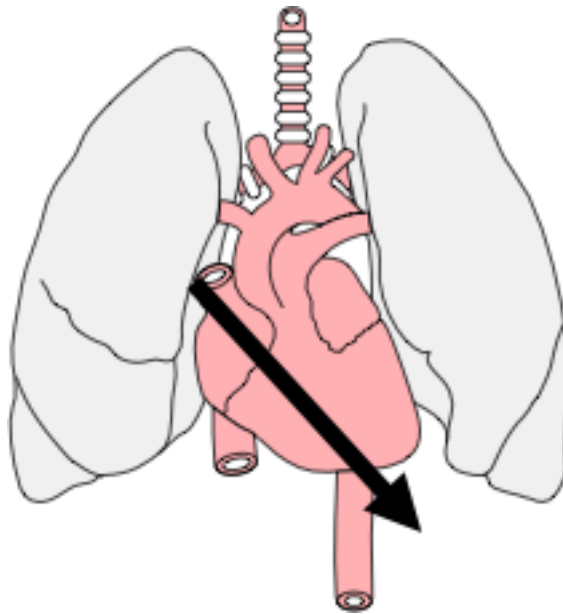


Figure 3.2: **Flow of electrical current.**

## Impulse origin and atrial depolarization

When the SA node, a pacemaker cell, fires off an impulse, the impulse travels down and toward the right and left atria. The direction – or vector – of this flow looks like Fig.3.3. The electrical flow is translated to the ECG as the P wave[14].
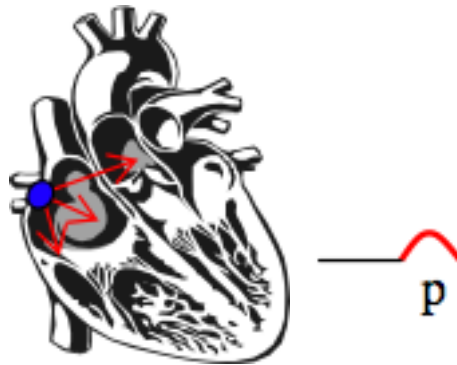
Figure 3.3: *Atrial depolarization.*

## Septal depolarization

The electrical flow stops briefly at the AV node, and then travels quickly down the common bundle (Bundle of His) and through the right and left bundle branches to the interventricular septum. The depolarization of the septum causes a small negative deflection a q wave in some leads; and a small positive deflection or "r" wave in others (Fig.3.4)[14].
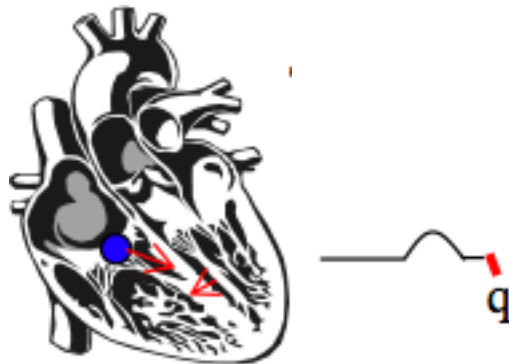


Figure 3.4: *Septal depolarization.*

## Apical and early ventricular depolarization

After depolarizing the septum, the impulse moves downward and to the left. This results in a large waveform  either an R wave or an S wave (Fig.3.5)[14].
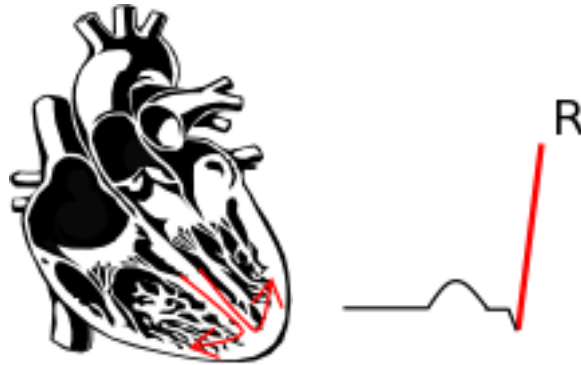


Figure 3.5: **Apical and early ventricular depolarization**

## Late ventricular depolarization

The final stage of depolarization takes place in the furthest stretches of the ventricle. The electrical stimulus moves upward, resulting in either a taller R wave or a smaller S wave (Fig.3.6)[14].
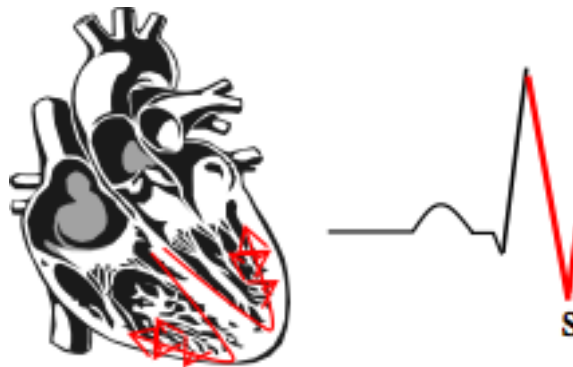


Figure 3.6: **Late ventricular depolarization.**

## Ventricular repolarization

Finally, the electrical stimulus is completed, ending depolarization. The ions in the cells move back into their normal resting positions, from top to bottom, causing the T wave (Fig.3.7)[14].



Figure 3.7: *Ventricular repolarization.*

## The whole cardiac cycle

1. Atrial depolarization (P wave).

2. Septal depolarization (Q wave).

3. Early ventricular depolarization (R/S wave).

4. Late ventricular depolarization (S/R wave).
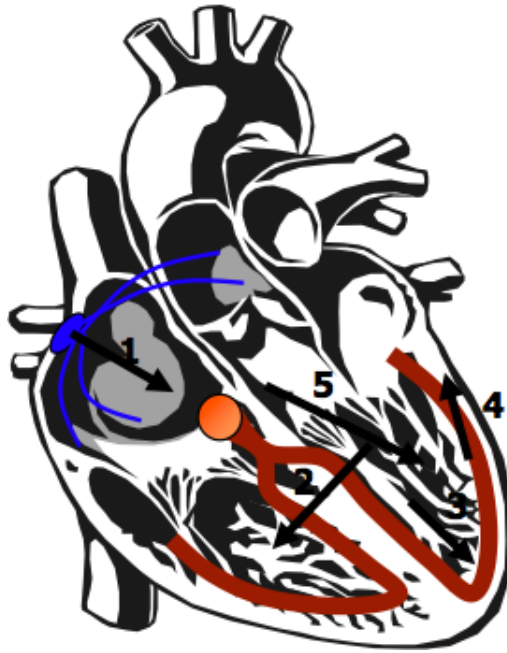
5. Ventricular repolarization (T wave).

Figure 3.8: *The whole cardiac cycle*

## 3.3   The ECG components

Each portion of a heartbeat produces a different deflection on an ECG. These deflections are recorded as a series of positive and negative waves.

### The Isoelectric Line

There is a part of the normal ECG rhythm that is electrically neutral - there is nothing electrically happening in the heart during that period. This is called the "isoelectric" line. This is a straight line passing from the end of the T wave and the beginning of the next P wave (Fig.3.9).

Figure 3.9: **Baseline or Isoelectric Line**

## The P Wave

The P wave begins with the first deviation from baseline and finishes when the wave meets the baseline once again. While the P wave is an electrical representation and not mechanical, a P wave strongly suggests that the atria have followed through with a contraction (Fig.3.10)[16].

The P wave indicates atrial depolarization. Its shape is round and smooth and the width of the normal P wave is less than 2.5 mm (0.11 seconds). The height of the normal P wave is less than 3 mm (0.3 mV).



Figure 3.10: **The P Wave.**

## The PR Interval
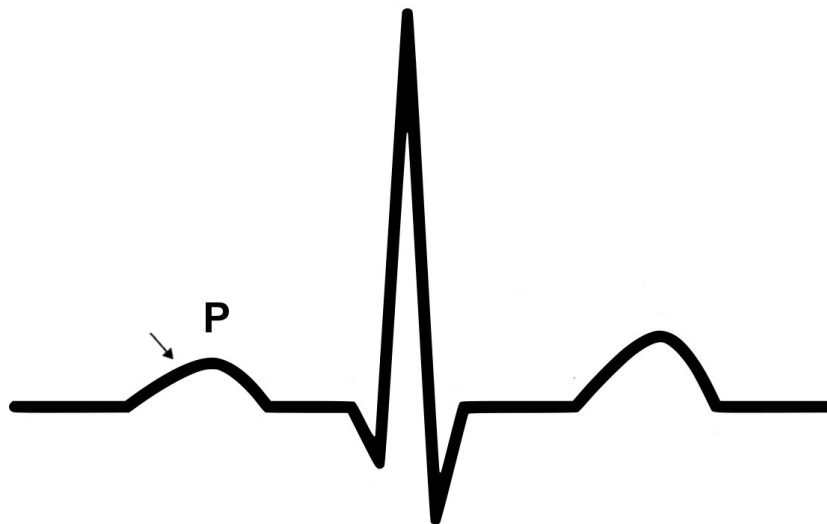
The PR segment is the part of the ECG between the end of the P wave and the beginning of the QRS complex. The PR segment signifies the time taken to stimulate the slow AV junction. This delay allows for atrial kick. The PR segment also serves as a guide for the isoelectric line (Fig.3.11).
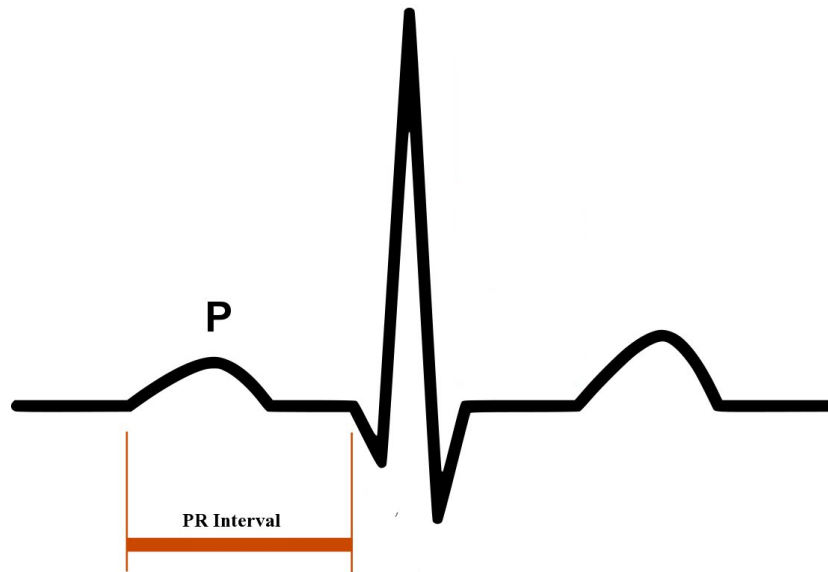


Figure 3.11: *The PR Interval.*

## The PR Segment

The PR interval is measured from the start of the P wave to the start of the QRS complex. While it might appear obvious that this is indeed a PQ interval, a Q wave is not always present on an ECG tracing. For consistency, the term "PR interval" is used whether a Q wave exists or not (Fig.3.12).

If the P wave is consistently followed by a QRS complex across a consistent PR interval, this is strong evidence that the originating impulse is supraventricular. A consistent PR interval is often sufficient to declare that this is a supraventricular rhythm[16].

The time from the beginning of atrial depolarization to the beginning of ventricular depolarization. Normal duration of the PR segment is 3 to 5 mm (0.12 to 0.20 seconds).



Figure 3.12: **The PR Segment.**

## The QRS Complex

ECG interpretation relies heavily on the QRS complex.  The QRS complex represents the depolarization of the ventricles.  The repolarization of the atria is lost within the QRS complex (Fig.3.13).



Figure 3.13: *The QRS Complex.*

Three distinct waveforms are often present in a normal QRS complex representing ventricular depolarization.  Depolarization of the ventricular septum begins first from left part of the heart to the right.  This early depolarization causes a small downward deflection called a Q wave. A Q wave is the first negative deflection of the QRS complex that is not preceded by a R

wave. A normal Q wave is narrow and small in amplitude. Note that a wide and/or deep Q wave may signify a previous myocardial infarction (MI).

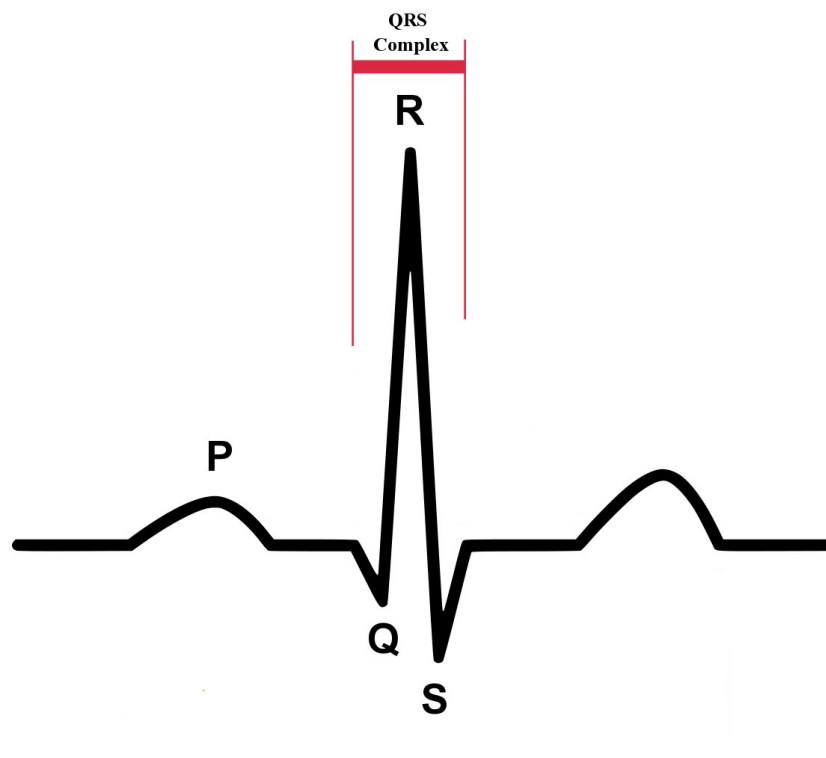Following the depolarization of the interventricular septum, ventricular depolarization then progresses from the endocardium through to the epicardium across both ventricles producing an R wave and an S wave. An R wave is the first positive deflection of the QRS complex. An S wave is the first wave after the R wave that dips below the baseline (isoelectric line). The end of the S wave occurs where the S wave begins to flatten out. This is called the J point (Fig.3.13)[16]. Various QRS Complex morphologies are represented in Figure 3.14. As a convention all these different morphologies are defined as QRS complexes.
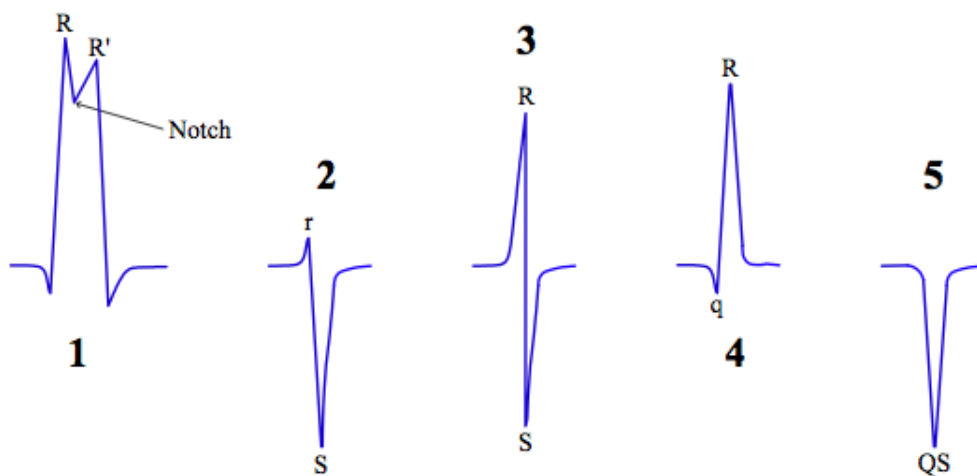


Figure 3.14: *Various QRS Complex Morphologies.*

The QRS Complex represents ventricular depolarization. Normal width is less than 3 mm and rarely less than 1.5 mm (0.12 seconds and rarely less than 0.06 seconds).

## The Q Wave and The QT Interval

A normal Q wave represents a depolarization of the ventricular septum, which usually travels from left to right. When present, a Q wave is the first downward. While ST segment deviation is a sign of present events, a prominent Q wave points to an MI that has already occurred, recently to some time ago (Fig.3.15).

A long QT interval wider than 1/2 the R-R interval is a significant risk factor for developing hemodynamically unstable dysrhythmias such as ventricular tachycardia and torsades de pointes. A prolonged QT interval is also associated with a higher incidence of sudden death.

The concern around a longer QT interval centers around the possibility of the next QRS coming at the tail end of the T wave, called an R-on-T phenomenon. This phenomenon can potentially cause dangerous dysrhythmias. Causes of prolonged QT intervals include long QT syndrome, antiarrythmics such as quinidine and procainamide, tricyclic antidepressants, and hypokalemia[16].

Figure 3.15: *The Normal Q Wave and QT Interval.*

## The ST Segment

Between the QRS complex and the T wave, lies the ST segment. The ST segment represents early repolarization of the ventricles. Early repolarization includes a plateau phase where the cardiac cell membrane potential does not change. ECG leads do not record any electrical activity during early repolarization. The ST segment is usually aligned with the isoelectric line.

Determining where the ST segment begins is determined by the J point. The J point, the juncture of the QRS and the ST segment, defines the starting point of the ST segment. The J point marks where the QRS complex changes

direction, forming a notch or bump in the ECG tracing. The ST segment is evaluated for any deviation from the ECG baseline 0.04 seconds after the J point (Fig.3.16).



Figure 3.16: **The ST Segment.**

While ST deviations may be normal for some individuals, usually they are a sign of either myocardial ischemia, myocardial infarction and/or cardiac disease[16], making them an important finding in ECG interpretation.

ST depression of 1 mm or more in two contiguous leads (neighboring leads) is suggestive of myocardial ischemia, injury or infarction. ST elevation of 1 mm or more in two contiguous leads is highly suggestive of a myocardial

injury or infarction. Note that ST changes (elevation or depression) are highly suggestive of the acute coronary events that are happening at the time an ECG is taken.

The presence of ST elevation in several leads of a 12 lead ECG suggests pericarditis. Ventricular rhythms and supraventricular rhythms with left bundle branch block have wide and bizarre QRS complexes, making the detection of ST changes very difficult[16].



Figure 3.17: *Various types of ST Segment Deviations.*

## The T Wave

A T wave usually follows every QRS complex. The T wave corresponds to the repolarization of the ventricle. The T wave is typically about 0.10 to 0.25 seconds wide with an amplitude less than 5 mm. While ventricular depolarization occurs rapidly producing a tall QRS complex, ventricular repolarization is spread over a longer interval, resulting in a shorter and broader T wave (Fig.3.18).

The T wave is, under normal conditions, slightly asymmetrical and usu-

ally larger than the P wave. The T wave is normally upright in lead II. Note
that as heart rates increase, the P wave and the T wave begin to share the
same space on an ECG. The larger T wave often covers the P wave. Note
that the T wave is rarely notched. A notched T wave may also contain a P
wave trying to show itself.



Figure 3.18: **The T Wave.**

Abnormally shaped T waves can indicate acute cardiac ischemia, elec-
trolyte imbalances, and cardiac disease related medication. For example,
peaked T waves can occur early during periods of myocardial ischemia and
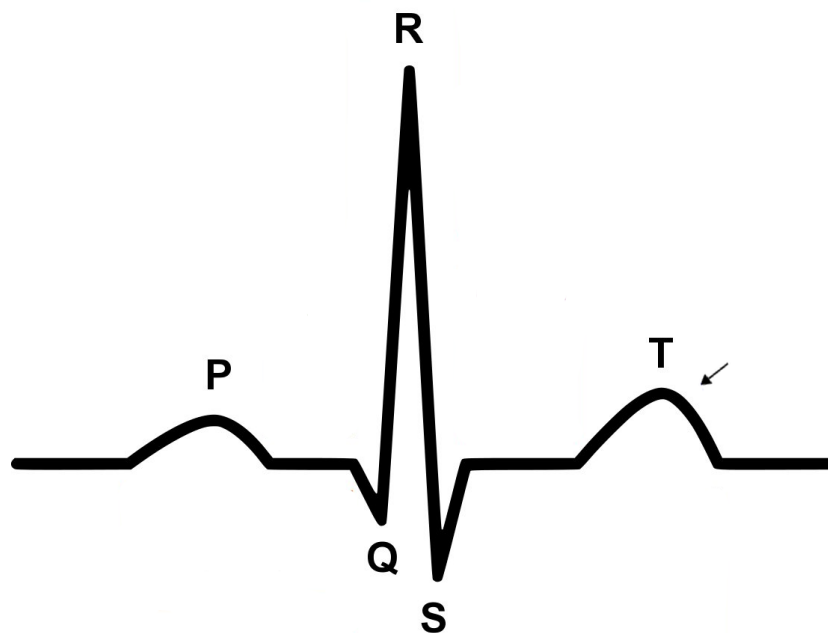infarction. Cardiac ischemia may cause the T wave to invert. Electrolyte im-

balances can also affect the T wave. Hyperkalemia is often associated with peaked T waves and can flatten the T wave. Quinidine can widen the T wave while digitalis can flatten the T wave.



Figure 3.19: *Normal and Abnormal T Waves.*

All morphologies of T waves, from normal to peaked to inverted can be found in healthy individuals without any evidence of disease, cardiac or otherwise. This makes the T wave a weak sign for any diagnosis. The T wave must be placed along side other clinical evidence. Rarely would treatment be based solely on the shape of the T wave[16].

## The U Wave

Occasionally, another wave, the U wave, is recorded immediately following the T wave and before the P wave. The U wave has yet to be fully explained but current studies suggest it represents a final stage of repolarization of certain ventricular cells in the mid-myocardium. The U wave will most often be oriented in the same direction as the T wave with an amplitude less than

2 mm. An abnormal U wave is inverted or tall with an amplitude of 2 mm or more[16].

An abnormally tall U wave is associated with conditions such as hypokalemia, diabetes, ventricular hypertrophy, and cardiomyopathy. Cardiac medications such as digoxin and quinidine can also cause a tall U wave[17].



Figure 3.20: *The whole ECG waveform with all waves, intervals and segments.*

Normal, abnormal parameters for every ECG component with possible causes of them are shown below (Table 3.1).
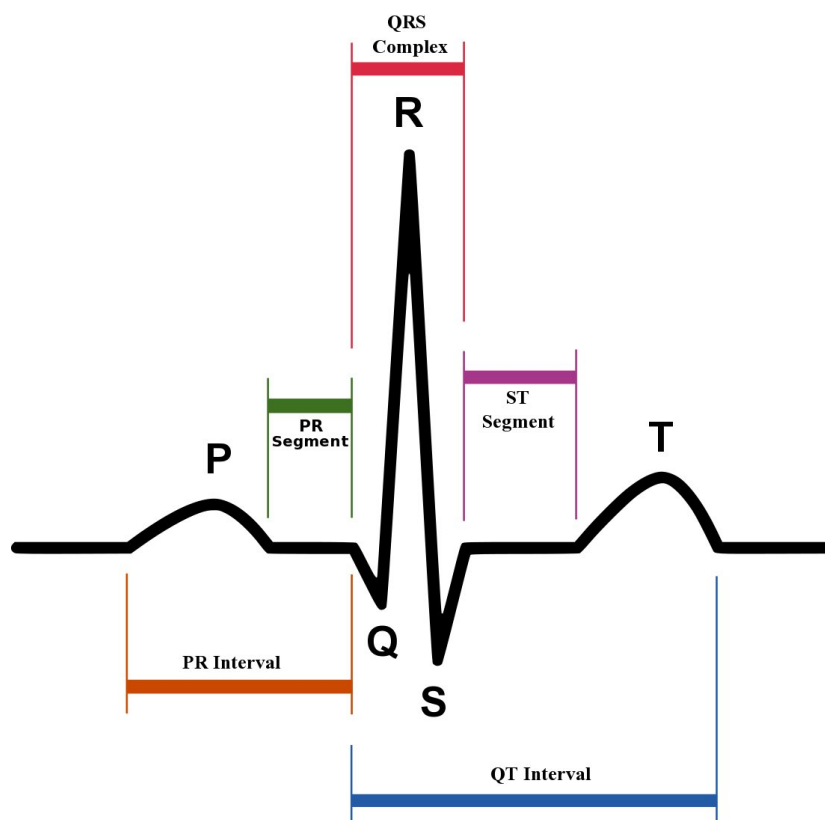
| ECG Components | Normal Parameters | Abnormal Parameters | Causes of Abnormal Parameters |
|---|---|---|---|
| P Wave | Upright in most leads including lead II. Duration: < 0.11 seconds Amplitude: 0.5-2.5 sec. | Inverted<br><br>Notched or tall | Junctional Rhythm<br><br>Atrial rhythm, atrial hypertrophy |
| PR Interval | Duration: 0.12 - 0.20 sec. | Duration: shorter or longer than normal | Junctional rhythm, Wolff-Parkinson-White syndrome |
| Q Wave | Duration: <0.04 seconds Amplitude: <25% the amplitude of the R wave | Duration: 0.04 sec. or longer Amplitude: at least 25% the amplitude of the R wave | Myocardial infarction |
| QRS Complex | Upright, inverted or biphasic waveform Duration: < 0.11 seconds Amplitude: 1 mm or more | Duration: 0.11 second or more | Bundle branch block, ventricular ectopic i.e. PVC |
| QT Interval | Duration: less than 1/2 the width of the R-R interval | Duration: at least 1/2 the R-R interval | Long QT syndrome, cardiac drugs, hypothermia, subarachnoid hemorrhage<br><br>Short QT associated with hypercalcemia |
| ST Segment | In line with PR or TP segment (baseline) Duration: shortens with increased heart rate | Deviation of 0.5 mm or more from baseline | Cardiac ischemia or infarction, early repolarization, ventricular hypertrophy, digoxin dip, pericarditis, subarachnoid hemorrhage |
| T Wave | Upright, asymmetrical and bluntly rounded in most leads Duration: 0.10-0.25 sec. Amplitude: less than 5 mm | Peaked, inverted, biphasic, notched, flat or wide waveforms | Cardiac ischemia or infarction, subarachnoid hemorrhage, left-sided tension pneumothorax, left bundle branch block, hyperkalemia, hypokalemia |
| U Wave | Upright Amplitude: < 2 mm | Peaked or Inverted Amplitude: > 2 mm | Hypokalemia, cardiomyopathy, ventricular hypertrophy, diabetes, digoxin, quinidine |

Table 3.1: *Outlines the parameters that define normal and abnormal ECG components*

# Chapter 4

# Clustering methods

## 4.1 Introduction

Cluster analysis or clustering is the assignment of a set of observations into subsets (called clusters) so that observations in the same cluster are similar in some sense. Clustering is a method of unsupervised learning, and a common technique for statistical data analysis used in many fields, including machine learning, data mining, pattern recognition, image analysis and bioinformatics.

Besides the term clustering, there are number of terms with similar meanings, including automatic classification, numerical taxonomy, botryology and typological analysis.

Many clustering algorithms require specification of the number of clusters to produce in the input data set, prior to execution of the algorithm. Barring knowledge of the proper value beforehand, the appropriate value must be determined, a problem for which a number of techniques have been developed.

In this thesis, two algorithms were implemented and tested for clustering: Ant Colony Optimization Clustering (ACOC) and K-means.

## 4.2   Ant Colony Optimization

ACO is a probabilistic technique for solving computational problems which can be reduced to finding good paths through graphs. It is part of the swarm intelligence algorithms category.

Initially proposed by Marco Dorigo in 1992 in his PhD thesis[19], the first algorithm was used for finding an optimal path in a graph, based on the behavior of ants seeking a path between their nest and a source of food. The original idea has since been diversified to solve a wider class of numerical problems[20].

### 4.2.1   Ant Colony Optimization Algorithm

In nature, ants initially wander randomly, and upon finding food return to their colony while laying down pheromone trails. Pheromone is a chemical substance for communicating information between ants.  The pheromone trails are reinforced by other artificial ants and validated over time (ACO has been applied successfully to a range of different combinatorial optimization problems recently)[20]. If other ants find a pheromone trail, they are likely not to keep travelling at random, but to instead follow it, returning to the nest and reinforcing it if they eventually find food.
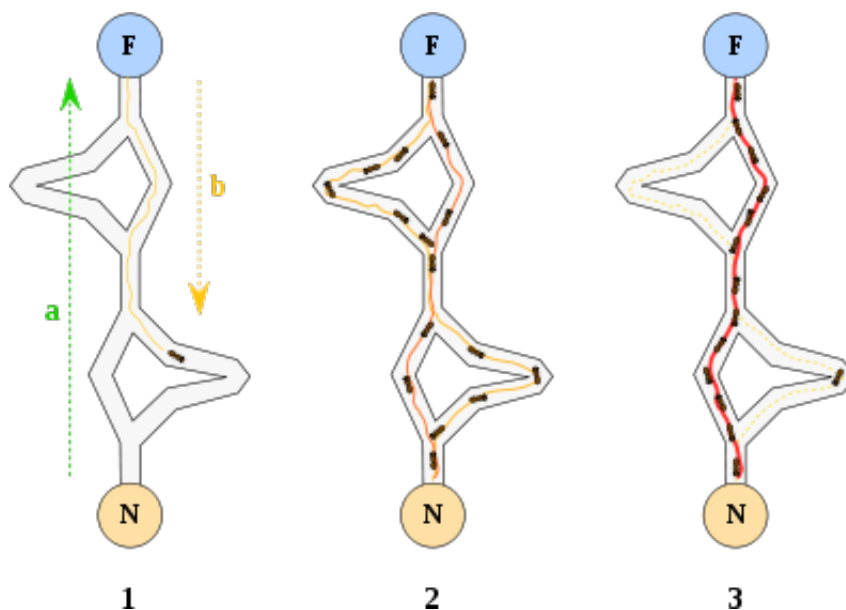
Over time, however, pheromone begins to evaporate, making the current path less attractive to ants.  The more time it takes for an ant to travel

down a path and back again, the more pheromone evaporates. A short path, by comparison, is crossed in a faster rate, and the associated pheromone level remains high because it is laid on the path as fast or faster than it can evaporate. Pheromone evaporation provides the advantage of avoiding the convergence to a locally optimal solution. If pheromone did not evaporate, the path selected by the first ant would tend to be increasingly more attractive to the rest of them. In that case, the exploration of the solution space would be constrained to a local extreme.

Thus, when one ant finds a short path from the colony to a food source, other ants are more likely to follow that path, and eventually all ants follow that path. The idea of the ACO is to mimic this behavior with "simulated ants" walking around a graph that is a representation of the problem to solve.

In a series of experiments on a colony of ants with a choice between two unequal length paths leading to a source of food, biologists have observed that ants tended to use the shortest route. A model explaining the behavior of ants while searching for food can be described as follows:

1. Each ant initially runs more or less randomly around the colony.

2. If it discovers a food source, it returns directly to the nest, depositing pheromone on the path it uses.

3. Pheromone attracts nearby ants making them follow, more or less directly, the track.

4. Returning to the colony, these ants will strengthen the route.

Figure 4.1: **Ant movement.**

5. If two different routes can be used to reach the same food source, the shorter one will be, in the same time, traveled by more ants than the longer one.

6. Pheromone on the shortest route will remain at a high level, attracting more ants.

7. Pheromone on the longer route will eventually diminish.

8. Finally, all the ants will be obliged to use the shortest route.

Ants use their environment as a medium of communication. They exchange information indirectly by depositing pheromones, detailing the status of their search. This provides positive feedback (the deposit of pheromone attracts other ants that will strengthen it themselves) and negative feedback (dissipation of pheromone by evaporation prevents ants from following

a path). Theoretically, if the quantity of pheromone remained the same on all
edges over time, no route would be finally selected. However, because of the
feedback, a slight variation on an edge will be amplified and thus propagate
across all possible routes. The algorithm will move from an unstable state in
which no edge is more attractive than another, to a stable state where edges
with high pheromone levels form the solution of the problem[20].

The Ant colony optimization algorithm have been applied to many combi-
natorial optimization problems, ranging from the quadratic assignment prob-
lem to folding proteins or vehicle routing. A lot of derived methods have been
adapted to solving dynamic problems with real variables, stochastic prob-
lems, multi-targets and parallel implementations. It has also been used to
produce near-optimal solutions to the travelling salesman problem[21]. This
method has an advantage over simulated annealing and genetic algorithm
approaches of similar problems where the graph may change dynamically;
the ant colony algorithm can be run seamlessly and adapt to changes in
real time. This is of interest in network routing and urban transportation
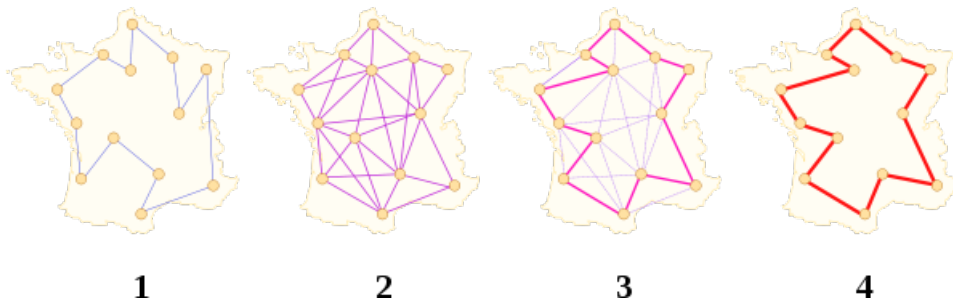systems.



Figure 4.2: *ACO applied on TSP problem (four steps).*

## 4.2.2 Ant Colony Optimization Clustering

In the ACOC algorithm, a weighted graph (V,E) is built where the vertexes represent the data to be clustered (V) and the weight of the edges between vertexes (E) is the acceptance rate between two data points. The ants traverse the graph and update the pheromone on the paths they cross. The graph is modified by gradually omitting some edges whose pheromone values are below a threshold. The strong connected components of the updated graph are computed every ten iterations to form the data clusters. The algorithm finally selects the clustering with the best performance.

Several quantities are defined to implement ACOC[23].

- $\tau_{ij}(t)$ : represents the amount of pheromone between stage $i$ and $j$ at time $t$,

- $n_{ij}$ : a problem-dependent heuristic function that evaluates the quality of local path $(i, j)$,

- $p_{ij}^k(t)$: the transition probability for ant $k$ to select the path $(i, j)$ at time $t$, which depends on the amount of pheromone on trail $(i, j)$ $\tau_{ij}(t)$ and the value of heuristic function$(n_{ij})$.

- $\alpha, \beta$: two parameters which decide respectively the effect proportion of $\tau_{ij}(t)$ and $n_{ij}$ on ants selecting path.

**Set of data items**

We use $S = (O, A)$ to denote a set of $n$ data items,where:

- $O = \{O_1, O_2, \ldots O_n\}$ represents the data set objects,

- $A = \{A_1, A_2, \ldots A_r\}$ represents the attributes of data objects, where

  $\forall i, i \in (1, 2, \ldots n), \exists a_i k, k \in (1, 2, \ldots r)$ denotes the attribute $A_k$ of $O_i$,

  therefore $A_k$ could be denoted as a r-dimensional vector $(a_{i1}, a_{i2}, \ldots a_i r)$,

  $i \in \{1, 2, \ldots n\}$.

**Difference between data items**

The difference between two data items $O_i$ and $O_j$ is defined as:

$$diff(i, j) = \sum_{k=1}^{r} dist(a_{ik}, a_{jk}), i, j = 1, 2, \ldots n, \qquad (4.1)$$

where dist is the euclidean distance between $a_{ik}$ and $a_{jk}$ (4.2).

$$dist(a_{ik}, a_{jk}) = \sqrt{\sum_{k=1}^{r} (A_{ik} - A_{jk})^2}. \qquad (4.2)$$

**Similarity between data items**

For two data items $O_i$ and $O_j$, we use $Sim(i, j)$ to denote their similarity:

$$sim(i, j) = 1 - \frac{diff(i, j)}{\max diff} \qquad (4.3)$$

The closer two items $a_{ik}$ and $a_{jk}$ are, the smaller the $diff(i, j)$ is and the larger the $sim(i, j)$ is.

Here

$$\max diff = \max_{1 \le i,j \le n} diff(i,j)$$

denotes the largest difference among the data items

We use $avesim(i)$ and $\max sim(i)$ to denote the average and maximum similarity of $O_i$ with all the other data items:

$$avesim(i) = \frac{1}{n-1} \sum_{j=1}^{n} sim(i,j) \qquad (4.4)$$

$$\max sim(i) = \max_{1 \le j \le n} sim(i,j) \qquad (4.5)$$

**Acceptance rate between data items**

We use $accept(i,j)$ to denote the acceptance rate of data item $O_i$ to $O_j$:

$$accept(i,j) = sim(i,j) - \frac{1}{2}[avesim(i) + maxsim(i)] \qquad (4.6)$$

Using $accept(i,j)$ as the weight of the $edge(i,j)$, we can form a weighted digraph as the initial pheromone graph. These values will be updated in every iteration via an ant crossing by any given edge. The value of $accept(i,j)$ may be negative or zero which means $O_i$ rejects $O_j$. In this case, this edge is treated as an invalid one and it will not be included in $E$ ($E$ is the set of valid edges of graph). $E$ is updated in each iteration and some of the edges with pheromone less than a certain threshold are removed from the graph[22].

**Heuristic function**

$\eta_{ij}$ indicates the "attractivity" of an edge $(i, j)$ towards the ants. The more similar two objects are, the higher preference an edge should have. Therefore, $\eta_{ij}$ is defined as:

$$\eta_{ij} = sim(i, j) \qquad (4.7)$$

In every iteration, the movement of ants causes changes in the graph. Because of these changes, each ant may not be allowed to travel to a specific vertex. Thus $\eta_{ij}$ is calculated on every iteration with respect to the subset of edges that an ant can cross. This is done in order to avoid rejecting edges of the graph too early.

**Probability function**

A probability function is used to facilitate the selection of the "right" path when an ant reaches a crossing between three or more edges. The next node $j$ is selected using the following formula:

$$j = \begin{cases} arg\{\max_{u \in allowed_k}[\tau_{iu}^{\alpha}(t)\eta_{iu}^{\beta}]\}, & \text{when } q \leq q_0 \\ \\ selected\ by\ probability\ p_{ij}^k(t), & \text{otherwise.} \end{cases} \qquad (4.8)$$

where $allowed_k$ is the set of vertexes that can be selected by the $k$-th ant, $q_0$ is a threshold for the vertex connected by the edge with the largest amount of pheromone to be selected. In each iteration, a random number is generated and compared with $q_0$. $q_0$ value is typically 0.9 meaning that the

next path will be selected according to the amount of pheromone on it for 90% of the time mostly, and 10% of the time according to $p_{ij}^k(t)$ function. When $q \geq q_0$ , data item $j$ is selected by the probability function defined as:

$$p_{ij}^k(t) = \begin{cases} \dfrac{\tau_{ij}^\alpha(t)\eta_{ij}^\beta(t)}{\sum_{r \in allowed_k} \tau_{ir}^\alpha(t)\eta_{ir}^\beta(t)} & j \in allowed_k \\ \\ 0 & \text{otherwise.} \end{cases} \quad (4.9)$$

otherwise the ant selects the vertex connected by the edge with the largest amount of pheromone[23].

**Pheromone updating**

After each iteration, the pheromone on the edges of the graph is updated according to the following formula:

$$\tau_{ij}(t+1) = (1 - \rho) \cdot \tau_{ij}(t) + \sum_{k=1}^{m} \Delta\tau_{ij}^k \quad (4.10)$$

$\rho \in (0, 1)$ is a constant called coefficient of evaporation. At each iteration the pheromone on the edges of the graph will be evaporated with a rate of $\rho$. The increment of $\tau_{ij}$ by ant $k$ is denoted as $\Delta\tau_{ij}^k$ and can be found using:

$$\Delta\tau_{ij}^k = \begin{cases} Q \cdot sim(i, j), & \text{if ant } k \text{ passes path}(i, j) \\ \\ 0 & \text{otherwise.} \end{cases} \quad (4.11)$$

Q is a constant with a typical value of 10. This value is used to increase the pheromone value on the path that was crossed by each ant.

It can easily be seen from 4.11 that the more ants pass through an edge, the more pheromone is deposited on it, and the edge will have higher probability to be included in the final graph.

**Update of the parameters $\alpha, \beta$ adaptively**

In 4.9 $\alpha, \beta$ determine the relative influence of the trail strength $\tau_{ij}$ and the heuristic information $\eta_{ij}$. At the initial stage of the algorithm, the pheromone value on each edge is relatively small. The ants should select the path mainly according to the heuristic information $\eta_{ij}$. Therefore, the value of $\beta$ should be relatively large. After some iterations, the pheromone values on the edges are increased, and as a result their influence becomes more and more important. Therefore, the value of $\alpha$ should become relatively large. Since the adjustment of the values of $\alpha$ and $\beta$ is based on the distribution of the pheromone on edges, in 4.12 $\tau_{ave}$ is defined as the average amount of pheromone on the pheromone graph and in 4.13 $\psi$ is defined as the variance of pheromone levels on the graph[22].
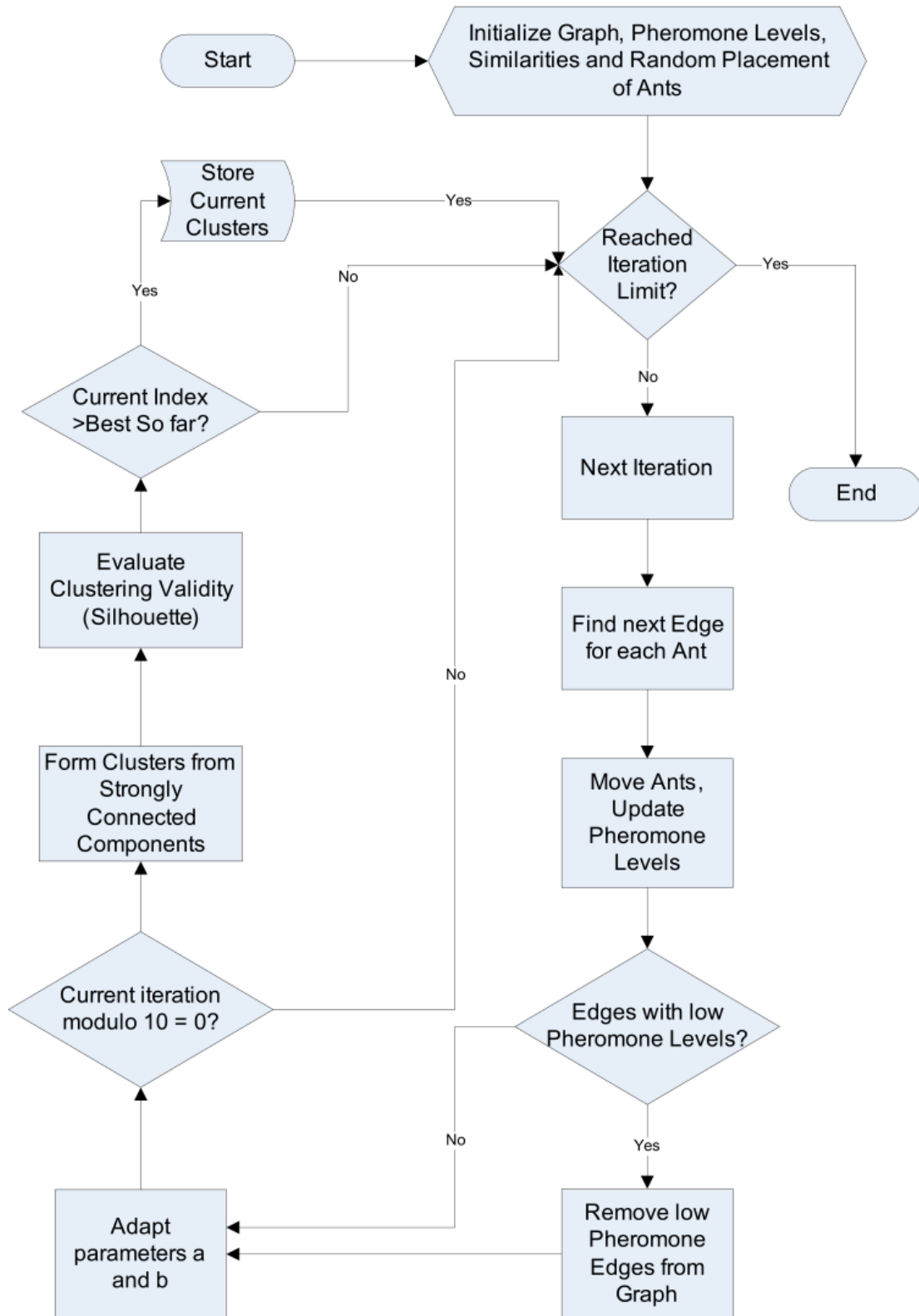
$$\tau_{ave} = \frac{\sum_{(i,j)\in E} \tau_{ij}}{|E|} \tag{4.12}$$

$$\psi = \frac{1}{|E|} \left[ \sum_{(i,j)\in E} (\tau_{ave} - \tau_{ij})^2 \right]^{\frac{1}{2}} \tag{4.13}$$

Using the pheromone distributing weight , the algorithm updates the value of $\alpha, \beta$ as follows:

$$\alpha = e^{-\psi}, \ \beta = \frac{1}{\alpha} \tag{4.14}$$

Acceleration of convergence and avoidance of local minima can be achieved by adapting $\alpha, \beta$. Furthermore, since the amount of pheromone is an important measure for data clustering, the pheromone distributing weight $\psi$ is also a critical factor to terminate the algorithm[22].

Figure 4.3: **Flow chart of ACOC algorithm.**

## 4.3   K-means Clustering

In statistics and machine learning, K-means clustering is a method of cluster analysis which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean. It is similar to the expectation-maximization algorithm for Gaussian mixtures in that they both attempt to find the centers of natural clusters in the data.

A non-hierarchical approach to forming good clusters is to specify a desired number of clusters, say, k, then assign each case (object) to one of k clusters so as to minimize a measure of dispersion within the clusters. A very common measure is the sum of distances or sum of squared Euclidean distances from the mean of each cluster. The problem can be set up as an integer programming problem but because solving integer programs with a large number of variables is time consuming, clusters are often computed using a fast, heuristic method that generally produces good (but not necessarily optimal) solutions. The Lloyd K-means algorithm is one such method[24].

The k-means clustering algorithm is commonly used in computer vision as a means for image segmentation. The results of the segmentation are used to aid border detection and object recognition. In this context, the standard Euclidean distance is usually insufficient in forming the clusters. Instead, a weighted distance measure utilizing pixel coordinates, RGB pixel color and/or intensity, and image texture is commonly used.

In addition, k-means is used in image recognition and therefore in pattern recognition as it is fast and mostly accurate after many iterations.
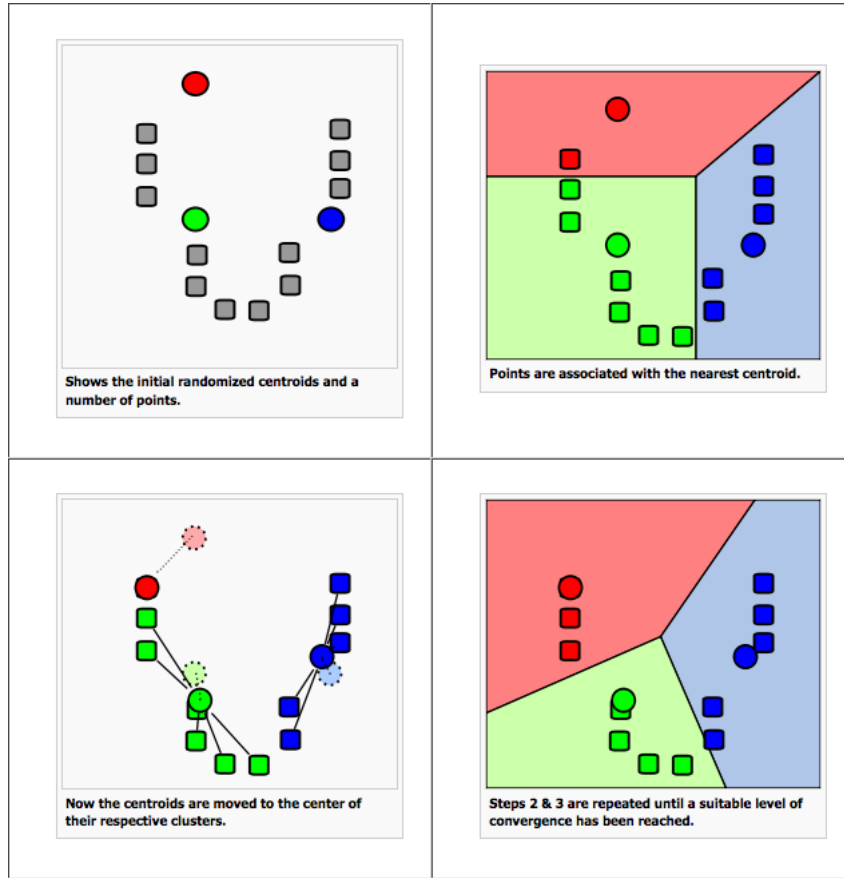
Figure 4.4: **Steps of K-means (Centroids are marked as circles and data as squares. Different colors are used for each cluster).**

Some modern applications of k-means are implemented in decision-making systems and data mining in complex web environments[24].

Given a set of observations $(x_1, x_2, \ldots, x_n)$, where each observation is a $d$-dimensional real vector, representing the various features of a datapoint, then $k$-means clustering aims to partition the $n$ observations into $k$ sets $(k < n)$ $S = \{S_1, S_2, \ldots, S_k\}$ so as to minimize the within-cluster sum of

squares(WCSS):

$$\min_{S} \sum_{i=1}^{k} \sum_{x_j \in S_i} ||x_j - \vec{\mu_i}||^2 \tag{4.15}$$

where $\vec{\mu_i}$ is the mean of $S_i$.

The most common algorithm uses an iterative refinement technique. Due to its ubiquity it is often called the k-means algorithm; it is also referred to as Lloyd's algorithm, particularly in the computer science community.

Given an initial set of k means $m_1^{(1)}, \ldots, m_k^{(1)}$, which may be specified randomly, by some heuristic, or be $k$ items from the data set in question, the algorithm proceeds by repeating these two steps:
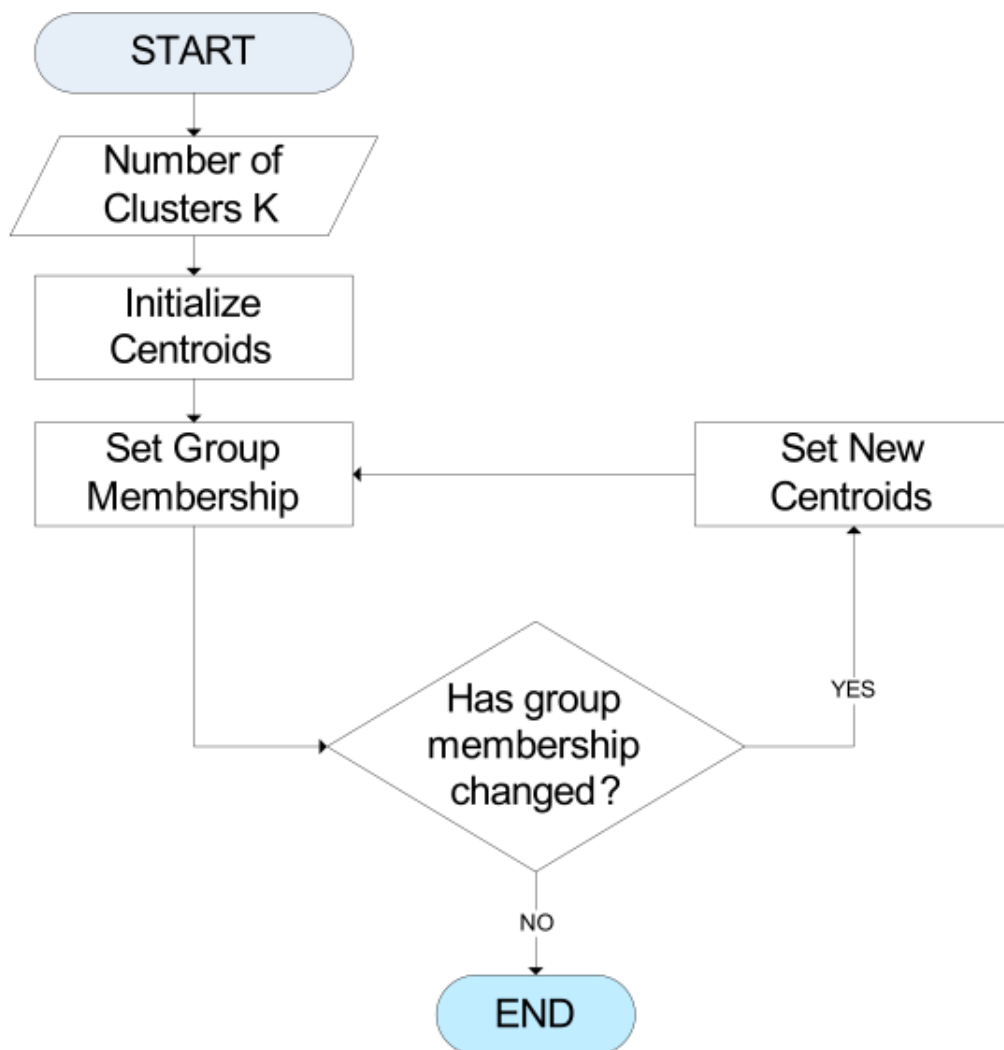
- Assignment step: each observation is assigned to the cluster with the closest mean.

$$S_i^{(t)} = \{x_j : ||x_j - m_i^{(t)}|| \leq ||x_j - m_{i*}^{(t)}|| \; \forall i^* = 1, \ldots, k\} \tag{4.16}$$

- Update step: calculate the new means of each cluster to be the centroid of the observations.

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j \tag{4.17}$$

The algorithm converges when the assignments on clusters no longer change.

Figure 4.5: **Flow chart of K-means algorithm.**

## 4.4 Comparison

In order to compare the two algorithms and finally select the most appropriate one, they were both tested in five synthetic data sets. All data sets consisted of two-dimensional data that belong to a normal distribution:

$$\mathcal{N}(\mu, \sigma^2). \tag{4.18}$$

where $\mu$ is the mean and $\sigma$ is the standard deviation.

### Silhouette

In order to evaluate the performance of the two algorithms the silhouette function was used. Silhouette refers to a method of interpretation and validation of clusters of data. The technique provides a graphical representation of how well each object lies within its cluster[25].

This technique calculates the silhouette width for each sample, average silhouette width for each cluster and overall average silhouette width for a total data set. The average silhouette width is applied for evaluation of clustering validity and also could be used to decide how good is the number of selected clusters.

For each datum, $i$ let $a(i)$ be the average "dissimilarity" of $i$ with all other data within the same cluster. Any measure of dissimilarity can be used but distance measures are the most common. We can interpret $a(i)$ as how well matched $i$ is to the cluster it is assigned (the smaller the value, the better the matching). Then the average dissimilarity of $i$ with the data of another

single cluster is found. The process is repeated for every cluster that $i$ is not a member of. The cluster with the lowest average dissimilarity to $i$ is denoted by $b(i)$. This cluster is said to be the neighbouring cluster of $i$ as it is, aside from the cluster $i$ is assigned, the cluster $i$ fits best in.

$$s(i) = \frac{b(i) - a(i)}{\max{(a(i), b(i))}} \tag{4.19}$$

which can be written as:

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases} \tag{4.20}$$

From the above definition it is clear that $-1 \le s(i) \le 1$.

For $s(i)$ to be close to one we require $a(i) << b(i)$. A small value of $a(i)$ means $i$ is well matched to its cluster. Furthermore, a large $b(i)$ implies that $i$ is badly matched to its neighbouring cluster. Thus an $s(i)$ close to one means that the datum is well clustered. If $s(i)$ is close to negative one, then by the same reasoning it can be seen that $i$ would be more appropriate if it was clustered in its neighbouring cluster. An $s(i)$ near zero indicates that the datum is on the border of two natural clusters.

The average $s(i)$ of a cluster is a measure of how tightly all data are grouped in the cluster they are. Thus the average $s(i)$ of the entire dataset is a measure of how appropriately the data has been clustered. If there are too many or too few clusters some of the clusters will display much narrower

silhouettes than the rest.

## Data Set Preparation

All data sets were mapped to $[0, 1]$ using a logarithmic transformation. The Euclidean distance metric used in both ACOC and K-means gives erroneous results when data sets are scaled differently across the features of the data set. The following steps were used:

An offset was determined:

$$offset = 1 - \min(V), \qquad (4.21)$$

where $V$ is the given data set.

Addition of the $offset$ to the data set

$$V_{new} = V + offset, \qquad (4.22)$$

where $V_{new}$ is the new data set.

Finally the new data set was created,

$$V_{final} = \frac{\log V_{new}}{\log\left(ceil(max(V))\right)}, \qquad (4.23)$$

where $ceil$ is the next largest integer.

## Synthetic data sets

### Data set 1

The first data set consists of five types of data which are:$[\mathcal{N}(2, 1^2), \mathcal{N}(2, 1^2)]$, $[\mathcal{N}(6, 1^2), \mathcal{N}(6, 1^2)], [\mathcal{N}(11, 1^2), \mathcal{N}(11, 1^2)], [\mathcal{N}(11, 1^2), \mathcal{N}(2, 1^2)], [\mathcal{N}(2, 1^2),$ $\mathcal{N}(11, 1^2)]$, which respectively belong to normal distribution (4.18).

In this data set it is obvious (Fig.4.6) that the five data types are clearly distinguished and as a result the outcome of the clustering of both algorithms (Fig.4.6) is characterized by high silhouette values (Table 4.2).

Both algorithms construct the same number of clusters but ACOC algorithm converges in the solution much earlier than k-means (Table 4.1).Graph before and after application of ACOC is presented (Fig.4.7).
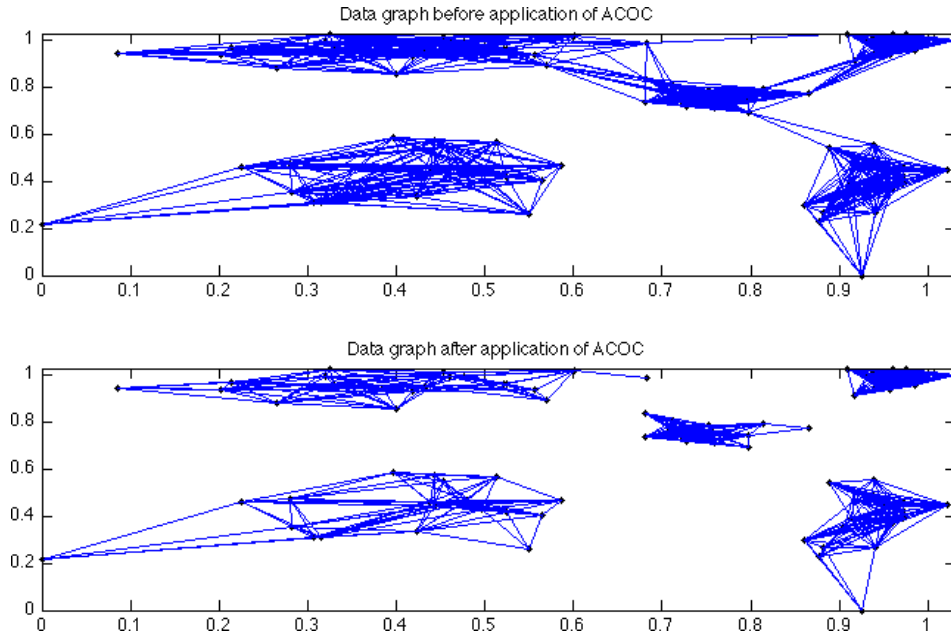


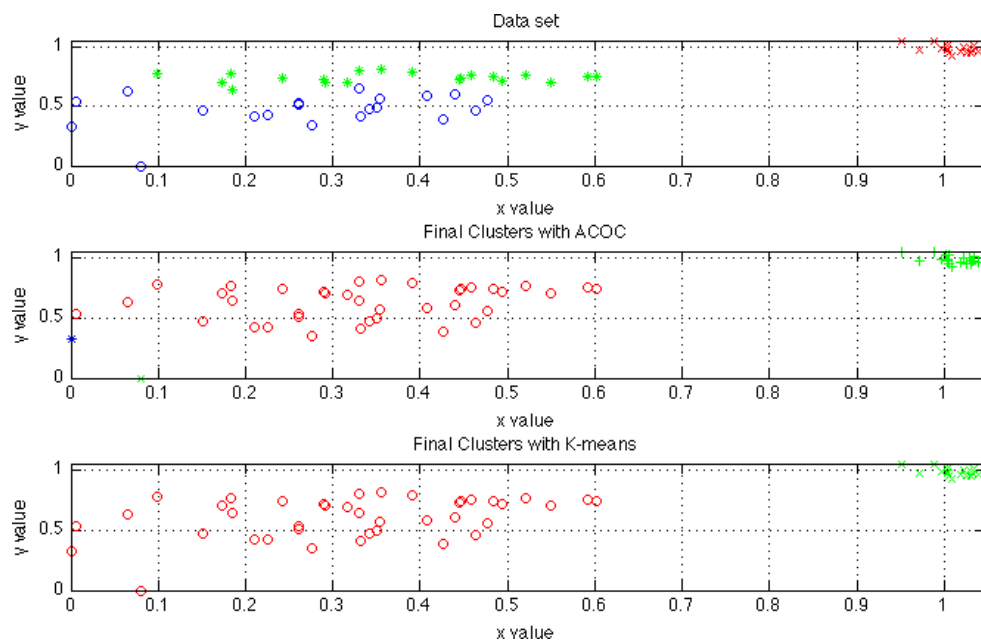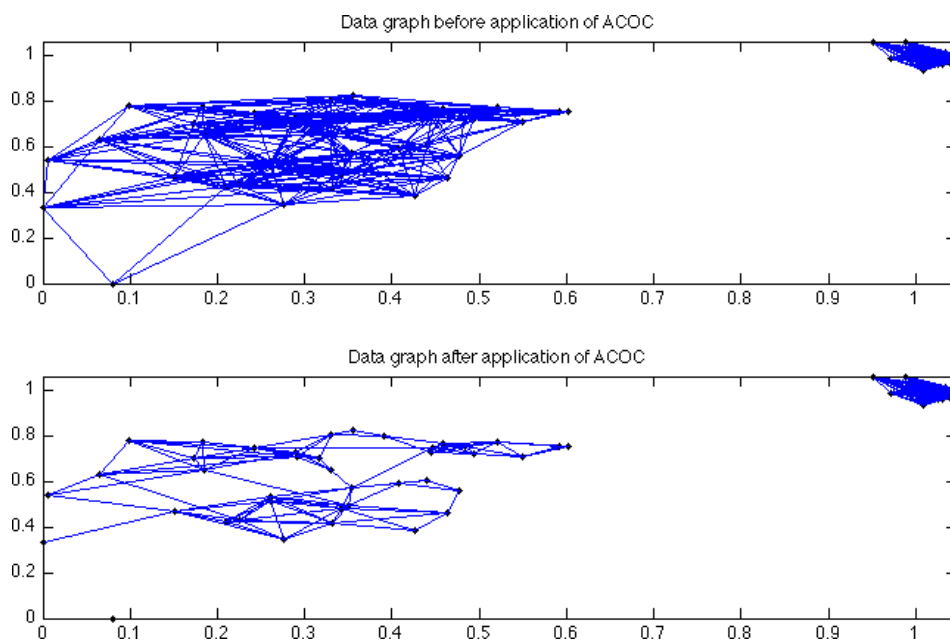Figure 4.6: **Data set 1, final clusters with ACOC and K-means.**

Figure 4.7: **Graph before and after ACOC.**

**Data set 2**

The second data set consists of three types of data which are: $[\mathcal{N}(1, 1^2),$
$\mathcal{N}(1, 1^2)], [\mathcal{N}(12, 1^2), \mathcal{N}(12, 1^2)], [\mathcal{N}(1, 1^2), \mathcal{N}(5, 1^2)]$, which respectively be-
long to normal distribution (4.18).

In this data set it is obvious (Fig.4.8) that the two of three data types
are close to each other and we expect that it will be difficult to cluster them
efficiently. ACOC and k-means converge to a solution with two clusters.

The performance (silhouette value) is relatively the same although ACO
converges much earlier (Table 4.2). Graph before and after application of
ACOC is presented (Fig.4.9).

Figure 4.8: **Data set 2, final clusters with ACOC and K-means.**



Figure 4.9: **Graph before and after ACOC.**

**Data set 3**

The third data set consists of three types of data which are: $[\mathcal{N}(1, 1^2),$
$\mathcal{N}(1, 1^2)], [\mathcal{N}(7, 1^2), \mathcal{N}(3, 1^2)], [\mathcal{N}(1, 1^2), \mathcal{N}(1, 3^2)]$, which respectively belong
to normal distribution (4.18).

In this data set it is obvious (Fig.4.10) that the two of three data types
are so close to each other so that it is extremely difficult to decide whether
some of them belong to the one data type or the other and we expect that
it will be difficult to cluster them efficiently.

Both algorithms converge in solutions with more clusters than the real
ones (Fig.4.10). However ACOC is faster than k-means (Table 4.1). Graph
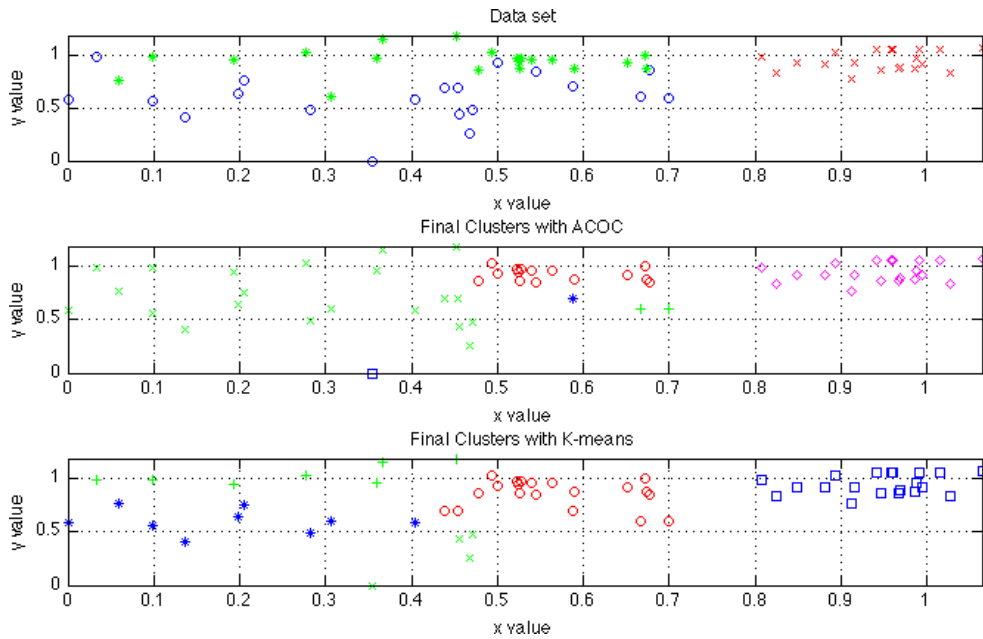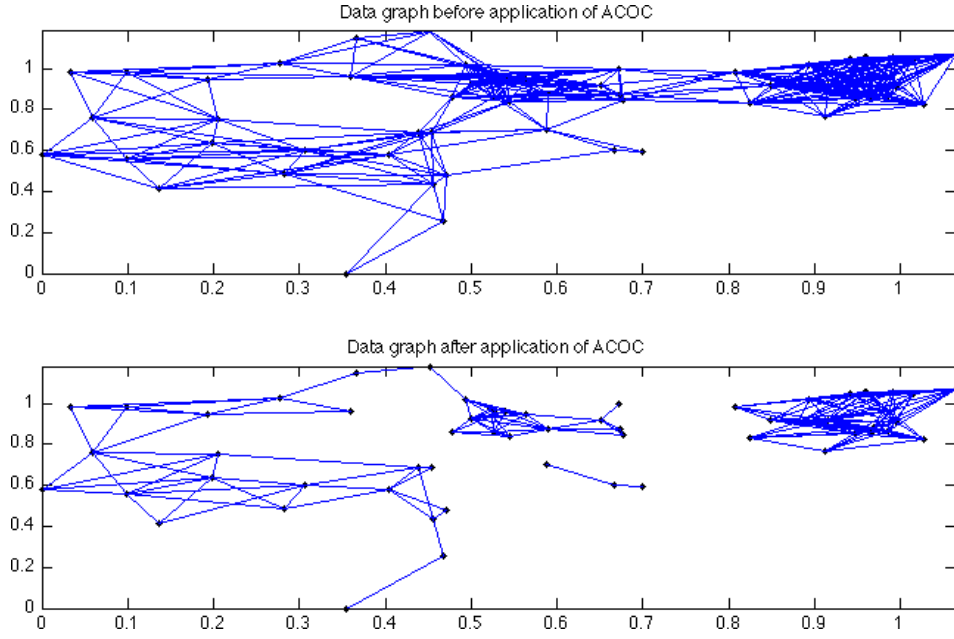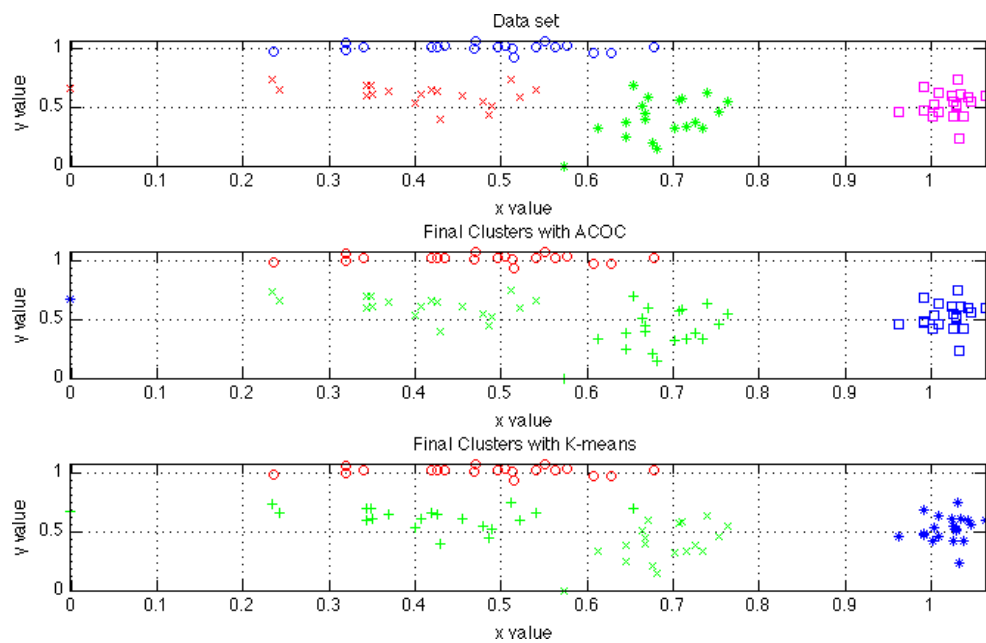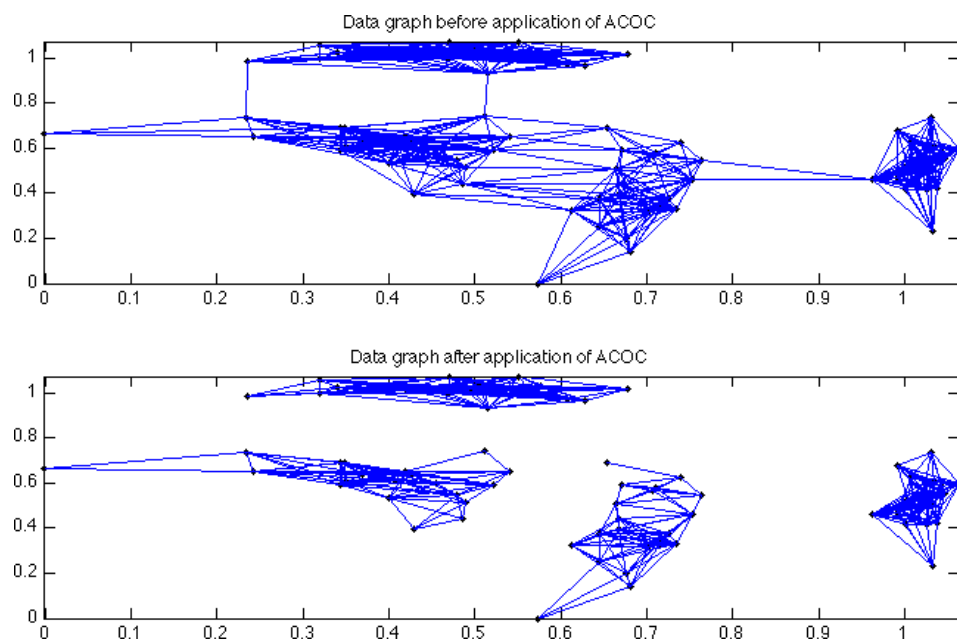before and after application of ACOC is presented (Fig.4.11).



Figure 4.10: **Data set 3, final clusters with ACOC and K-means.**

Figure 4.11: *Graph before and after ACOC.*

**Data set 4**

The fourth data set consists of four types of data which are: $[\mathcal{N}(1, 1^2),$
$\mathcal{N}(11, 1^2)], [\mathcal{N}(.5, 1^2), \mathcal{N}(2, 1^2)], [\mathcal{N}(4, 1^2), \mathcal{N}(.5, 1^2)], [\mathcal{N}(11, 1^2), \mathcal{N}(1, 1^2)],$
which respectively belong to normal distribution (4.18).

In this data set it is obvious (Fig.4.12) that the three of four data types are
relatively close to each other. ACOC algorithm clusters the mixed data into
five groups and k-means converges to a solution with four clusters (Fig.4.12).
However ACOC is faster than k-means and converges earlier (Table 4.1).
Graph before and after application of ACOC is presented (Fig.4.9).

Figure 4.12: **Data set 4, final clusters with ACOC and K-means.**



Figure 4.13: **Graph before and after ACOC.**

| Set | No.of data | Clusters | | | Cpu time (sec.) | |
|---|---|---|---|---|---|---|
| | | Real | ACOC | K-means | ACOC | K-means |
| 1 | 100 | 5 | 5 | 5 | 5.53 | 22.73 |
| 2 | 60 | 4 | 2 | 2 | 4.04 | 18.96 |
| 3 | 60 | 3 | 6 | 5 | 3.74 | 22.93 |
| 4 | 80 | 4 | 4 | 4 | 5.69 | 20.49 |

Table 4.1: ***Clusters, cpu time ACOC vs K-means (cpu time and clusters are average values).***

| Set | No.of data | Silhouette | | Iterations | |
|---|---|---|---|---|---|
| | | ACOC | K-means | ACOC | K-means |
| 1 | 100 | 0.80 | 0.80 | 200 | 2000 |
| 2 | 60 | 0.70 | 0.65 | 200 | 2000 |
| 3 | 60 | 0.50 | 0.55 | 200 | 2000 |
| 4 | 80 | 0.72 | 0.70 | 200 | 2000 |

Table 4.2: ***Silhouette, iterations ACOC vs K-means (silhouette and iterations are average values).***

In all synthetic data sets figures (4.6, 4.8, 4.10, 4.12), data are transformed in $[0, 1]$. This happens because a logarithmic mapping technique was applied to the synthetic data sets (see section 4.4).

## Selection

In order to select the more efficient algorithm for this thesis the following factors were taken into consideration :

- Time complexity,

- Number of parameters needed to set[1],

---

[1]Nine parameters for ACOC four for K-means

- Perfomance of the algorithms for certain data sets using silhouette (4.4),

- Precision in clustering.

Despite the fact that K-means is fast there is no guarantee that it will converge to the global optimum, and the result may depend on the initial clusters. As the algorithm is usually very fast, it is common to run it multiple times with different starting conditions.

ACOC is a heuristic algorithm too, but it converges to the final clustering (with sufficient precision) fast enough with its first initialization, without the need of running multiple times.

In the above synthetic data sets ACOC was always faster than K-means. Also, ACOC has a great advantage in precision of clustering. This is proved by the ability of ACOC to discriminate groups of data that are close to each other as different clusters. In addition to the little time that ACOC needs to converge in the solution, the stability of the algorithm is a factor that has to be mentioned. The lack of "*a priori*" knowledge of the groups of data and as a result the parametrs' adjustment is also a great advantage of ACOC against K-means in this type of data. If we assume "*a priori*" knowledge of the data groups K-means is always faster that ACOC but there is a need for precise numbering of the clusters that will be formed. This is a senior drawback of the K-means algorithm which limits its usage for large data sets without knowing the exact number of clusters. ACOC's only disadvantage is the sensitivity of its parameters. Anyone who lacks knowledge of the algorithm cannot adjust its parameters easily enough.

# Chapter 5

# QRS Complex detection

## 5.1 Introduction

The duration, amplitude and morphology of the QRS complex is significant for the diagnosis of cardiac arrhythmias, conduction abnormalities, ventricular hypertrophy, myocardial infarction and other disease states. In this chapter the collection and the transformation of the real data sets (ECG) are described. The problem of QRS complex detection is a relatively difficult one because of the variety of different morphologies that can be encountered. There was a need for transforming the ECG to another form that is characterized by less noise, holds the morphology of the initial data and finally is suitable for clustering. Many techniques were used for this purpose that are analyzed extensively.

## 5.2    Method proposed

The period of QRS complexes in an ECG is not stable because of many reasons. In pathological cases, like Atrial Fibrillation, the interval RR varies[26]. An important factor that intensifies this change is the state of the patient, because they might "be stressed" or move during the recording of the ECG.

Thus, the algorithm that was used for the detection of QRS complexes, in the ECG that were studied, is based on the statistical characteristics of the ECG and not on likely periodicities.

In order to become noise free, Fourier transformation is applied on the ECG. In the inverse Fourier transformation, the components that correspond to the 80% of total energy of signal are used while all remainder become zero. The result is a more smooth varying time series.
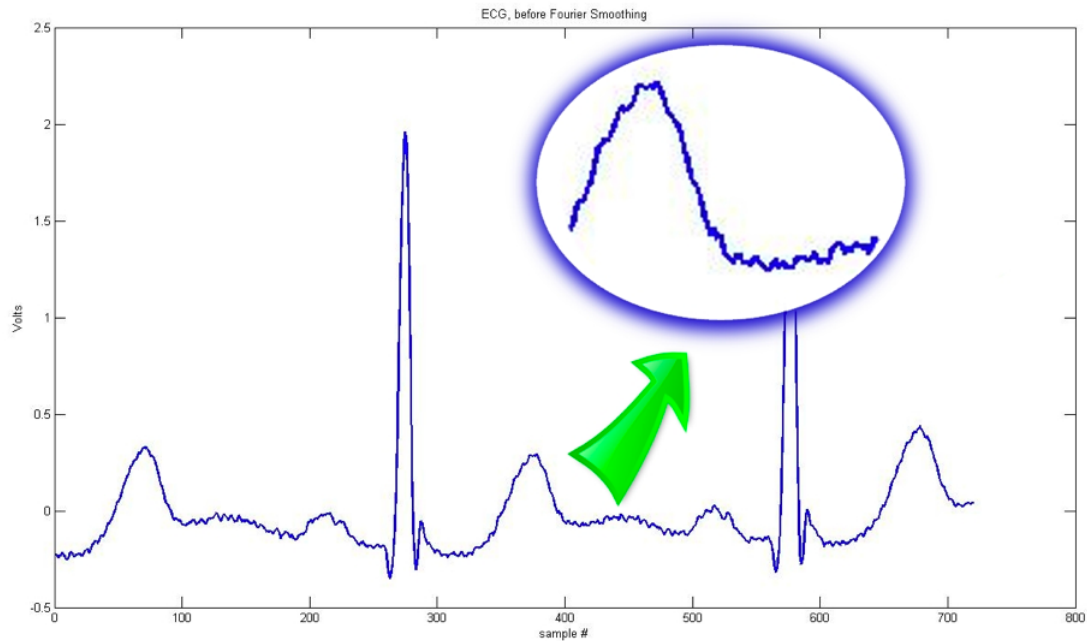


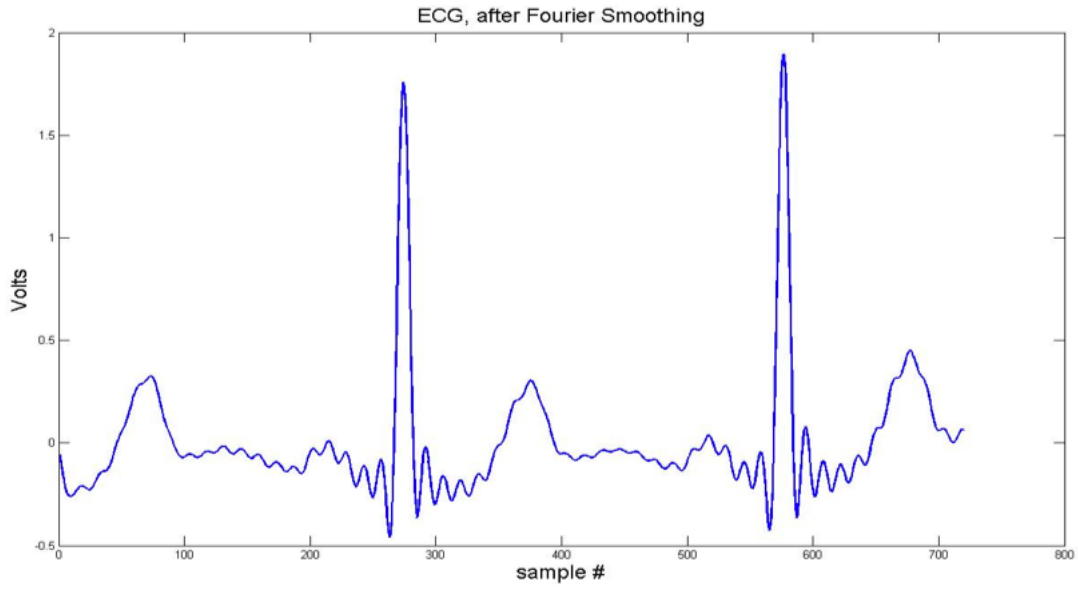Figure 5.1: *ECG before smoothing with Fourier transform.*

Figure 5.2: **ECG after inverse Fourier transform.**

Then, a Possible QRS complexes is selected which corresponds to a QRS complex and this complex is used as a template in order to the remainder clusters be recognized. This process in detail:

**Step 1**

Local maxima are collected from the ECG using a simple first derivative check:

$$ecg(i) > ecg(i-1), ecg(i) > ecg(i-2)$$

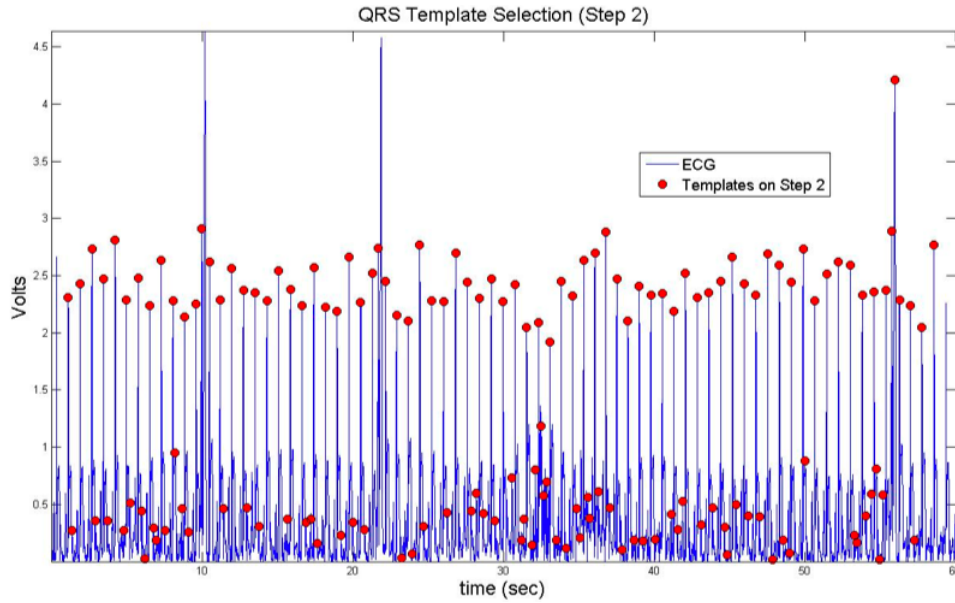$$\text{and}$$

$$ecg(i) > ecg(i+1), ecg(i) > ecg(i+2)$$

Figure 5.3: ***Possible QRS complexes remaining after step one.***

## Step 2

The voltage range recorded in the ECG around the points found in step 1 is calculated for a duration of 0.125 seconds.

## Step 3

Since the QRS complex is the most intense variation on an ECG, the part of the ECG that contains it is expected to vary greatly in voltage. All possible QRS complexes from the second step are classified per range. The top 25% of them are selected.
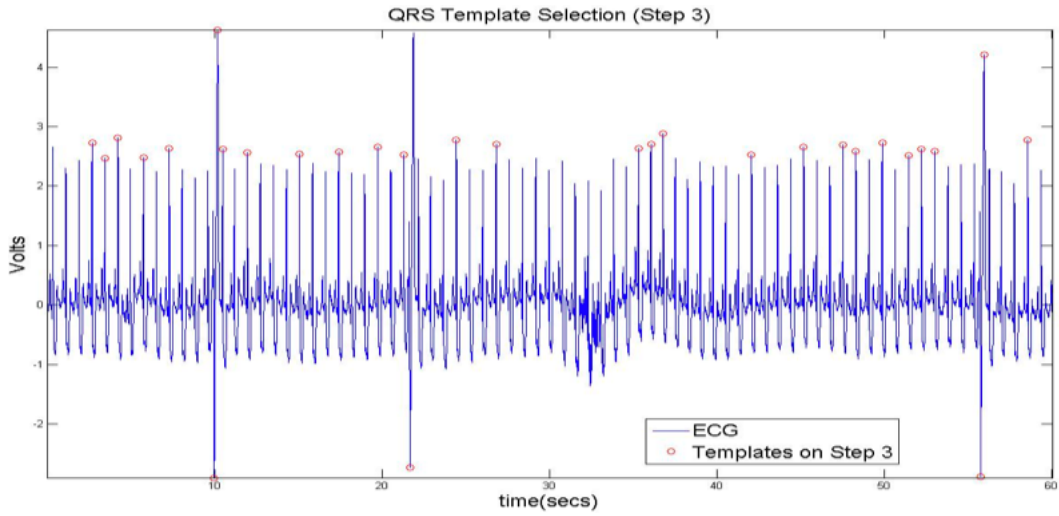
Figure 5.4: **Possible QRS complexes remaining after step three. The ECG comes from a patient with supra ventricular arrhythmia**

## Step 4

Waveforms of 0.125 sec length were taken from the ECG around each possible QRS complex. By logarithmic transformation, the values of these waveforms change in the interval $[0, 1]$. This technique has been described in chapter 4.

This process maintains the morphology of each waveform, changing only its width so higher QRS complexes are not rejected by the process of selection. The degree of correlation of dissimilar QRS complexes decreases .
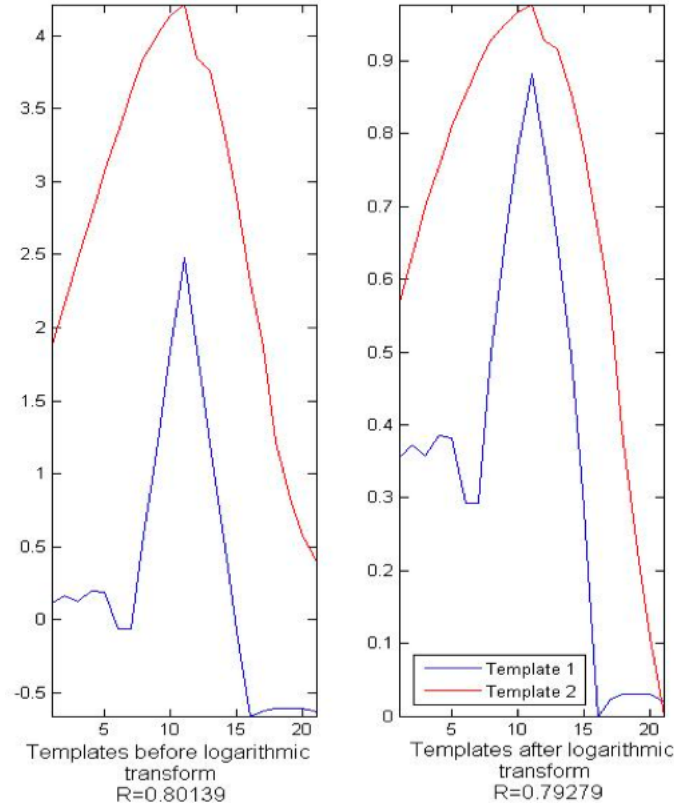
Figure 5.5: **Two possible QRS complexes before and after logarithmic transformation. The red line is an ectopic heart beat which is not evaluated during diagnosis. The blue line corresponds to a normal complex. The degree of correlation decreases.**

**Step 5**

For each waveform of the fourth step, its correlation with the rest is calculated, using the following formula:

$$R(i,j) = \frac{C(i,j)}{\sqrt{C(i,i) \cdot C(j,j)}} \qquad (5.1)$$

where $C$ is the covariance matrix of the two possible QRS complexes and $R$ is a $2x2$ matrix.
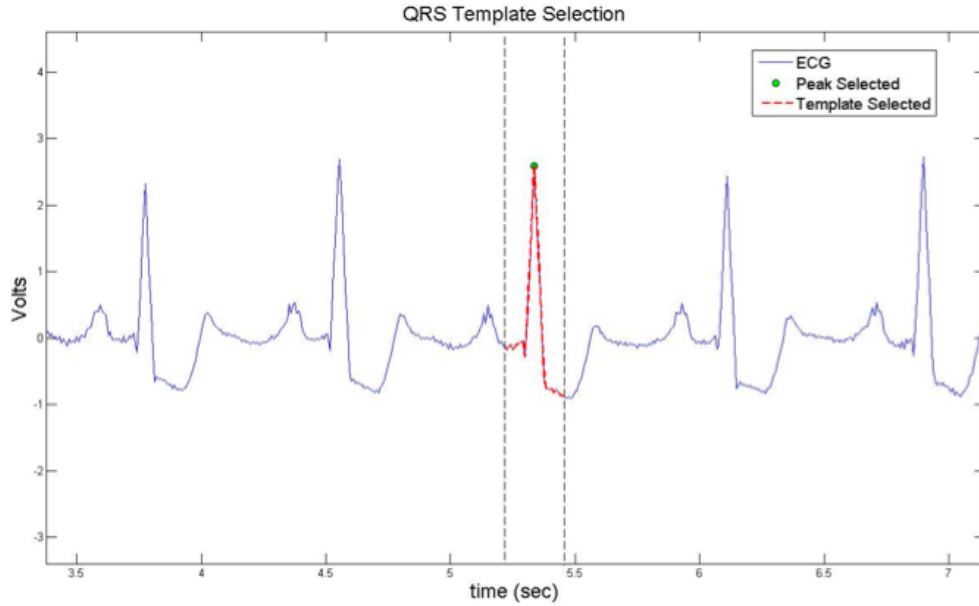


Figure 5.6: *The QRS complex that finally was selected for the scanning of the ECG*

The waveform that mostly "describes" the form of the QRS complexes presented in the ECG, is chosen as the one with the highest sum of correlation within all of the waveforms from step 4. This is called a QRS template. After the selection of the template, the whole ECG is scanned and the correlation between the ECG and the template is calculated (Fig.5.6). A typical result of this process is shown in the following figure.
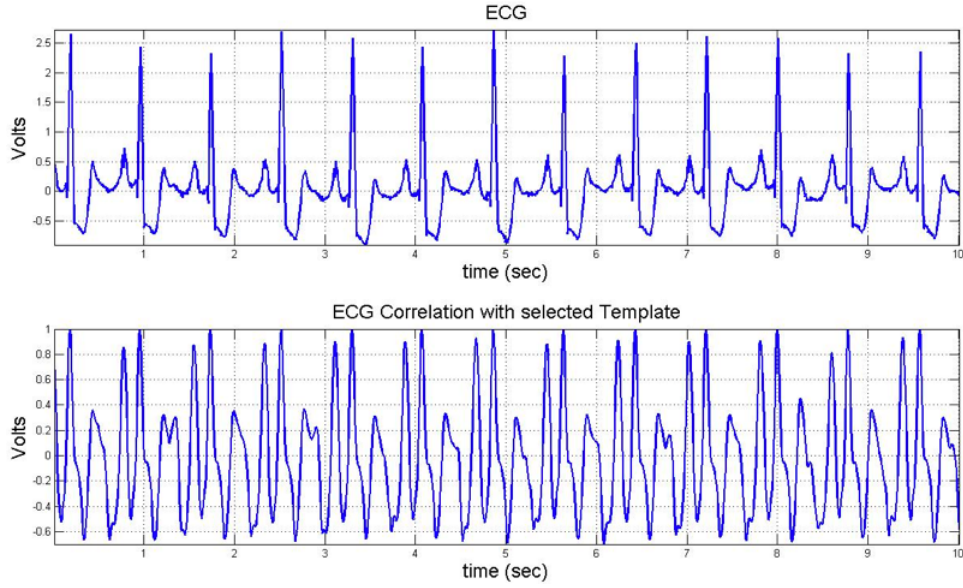
Figure 5.7: ***The correlation of the ECG with the template QRS complex. Maximum values are shown in the rest of the QRS complexes***

Finally, the problem of detecting the presence of QRS complexes has been converted to a problem of finding greatest values of correlation. The newly acquired signal is much smoother than the initial ECG. Sweeping the time series of the correlation, its derivative is examined along with a simple threshold value rejection rule. The waveforms that fulfill the following requirements are defined as candidate complexes:

$$|Correlation(i)| \geq 0.3$$
$$Correlation(i) > Correlation(i-1)$$
$$Correlation(i) > Correlation(i-2)$$
$$Correlation(i) > Correlation(i+1)$$
$$Correlation(i) > Correlation(i+2)$$

where $i$ is the *ith* window of 0.125 seconds in the ECG

The possible QRS complexes that have remained are converted in pairs of values, where the x-component is the correlation with the template and the y-component is the normalized to $[0, 1]$ voltage range around the possible QRS complex.
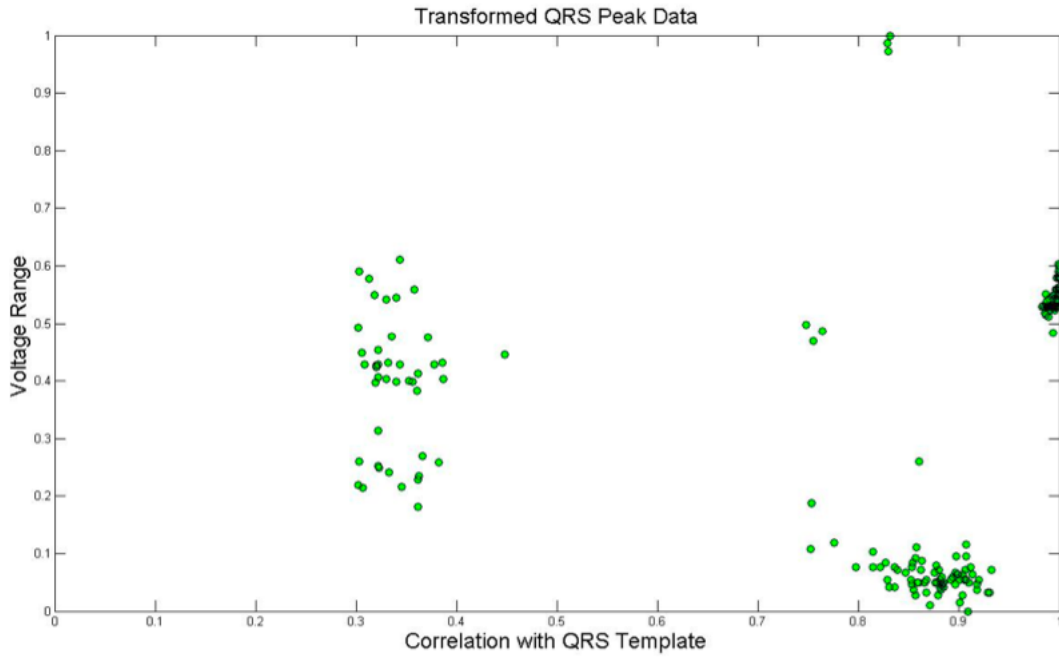


Figure 5.8: *Figure of correlation with the template QRS complex in comparison with the range around each possible QRS complex.*

So by applying a clustering algorithm to these transformed data, the true QRS complexes are found by forming a cluster that contains the template QRS complex.

## 5.3   Data (ECG) selection

The following four ECG databases were used for testing the algorithms:

- **Normal Sinus (NS)**. This database is located in http://www.physionet.org/physiobank/database/nsrdb/.

- **Atrial Fibrillation (AF)**. This database is located in http://www.physionet.org/physiobank/database/afdb/.

- **ST Change (STC)**. This database is located in http://www.physionet.org/physiobank/database/stdb/.

- **Supra ventricular arrhythmia (SVA)**. This database is located in http://www.physionet.org/physiobank/database/svdb/.

# Chapter 6

# Results

## 6.1  Introduction

The problem of detecting QRS complexes has been converted to a problem
of finding the greatest values of correlation and then the data clustering
algorithms that are described in chapter 4 are applied in the transformed
data. The group with the higher correlation with the template QRS complex
has to be a separate cluster and as a result this cluster finally contains the
QRS complexes of the ECG.

## 6.2  Results

Two records from each disease that were selected are presented in the results.
The QRS complexes belong in the same group with the template, while the
rest of them are not real QRS complexes.

Therefore, the clustering performance can be judged by how well the QRS

complexes cluster is isolated from the rest of the data set.

## Normal Sinus (NS)

### Record: 100 Lead: 2

In the first record the possible QRS complexes form three groups which are away from each other. This leads to an efficient clustering with both ACOC and K-means (Fig.6.1). The result of the application of the two algorithms in this record is almost perfect (Table 6.1). However, ACOC is faster than K-means (Table 6.2).
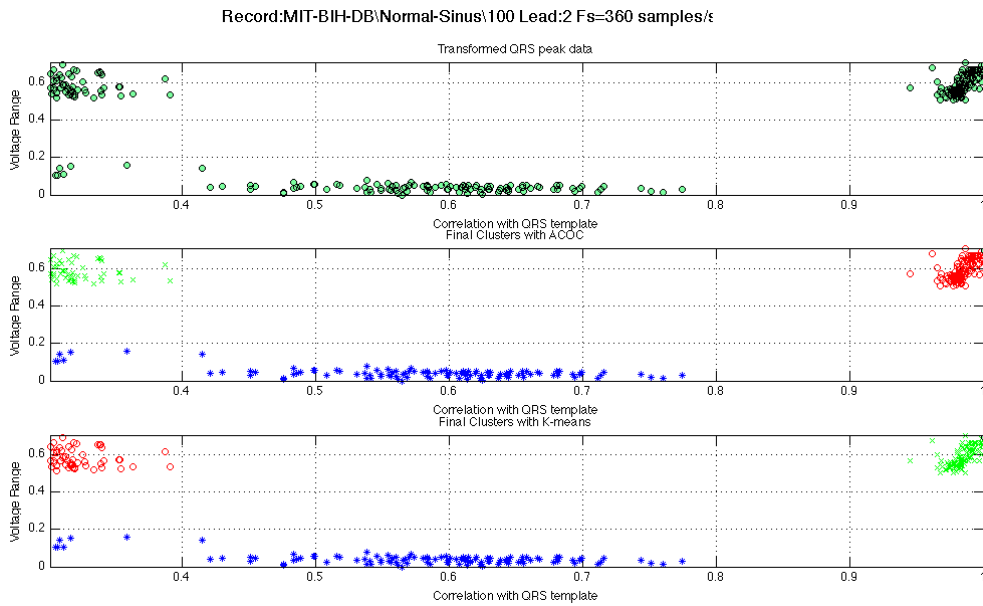


Figure 6.1: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*
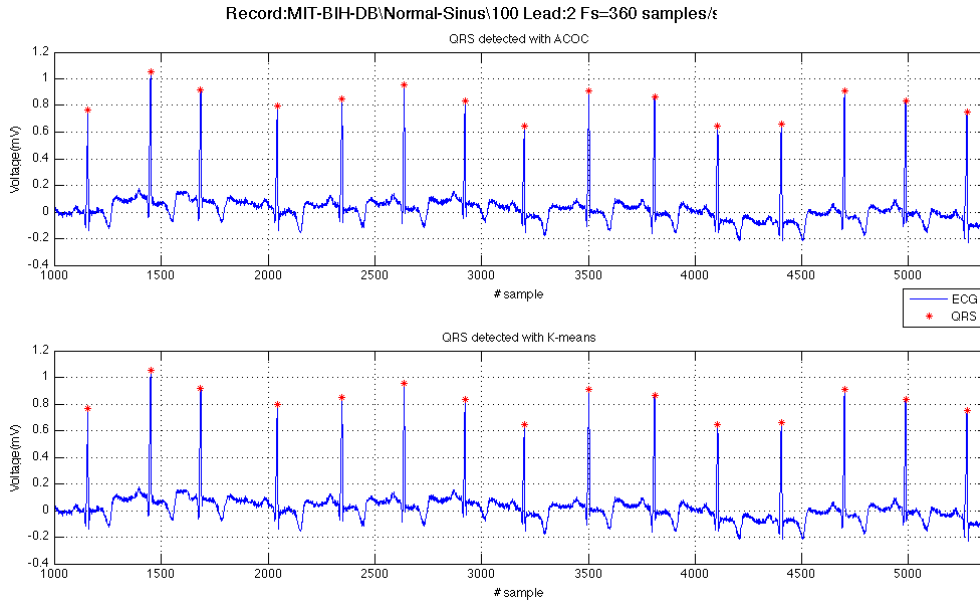
Figure 6.2: **Possible QRS complexes detection in the ECG (Only a small part is presented).**

**Record: 102 Lead: 1**

In this record the possible QRS complexes form three or four groups, which are relatively close to each other. This leads to a clustering with more errors by ACOC algorithm than by K-means (Fig.6.3). The precision of the QRS complex detection is lower with ACOC than K-means. However the percentage of the real QRS detected with ACOC is about 88% which is satisfactory (Table 6.1). Time cost of K-means is much more than of ACOC (Table 6.2). In Fig.6.4 a part of the ECG is presented.
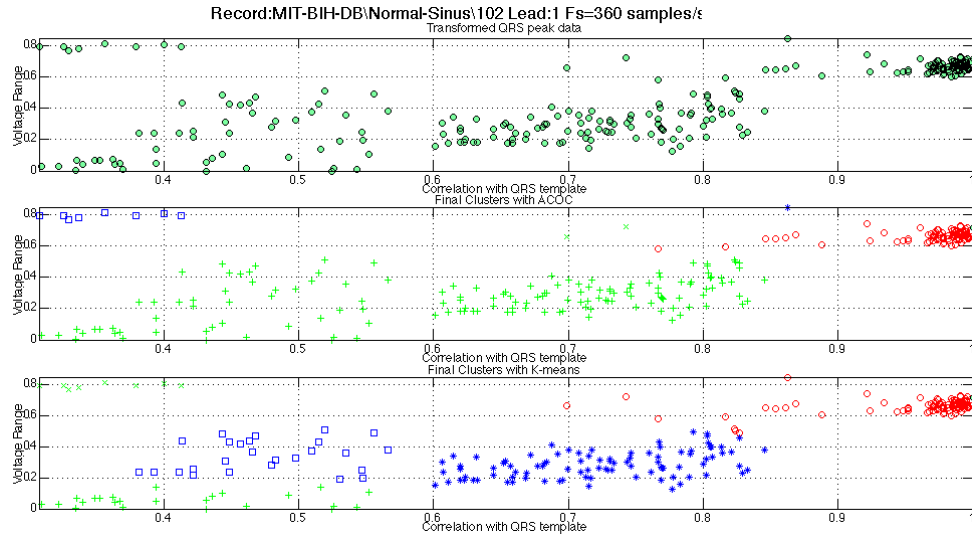
Figure 6.3: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*
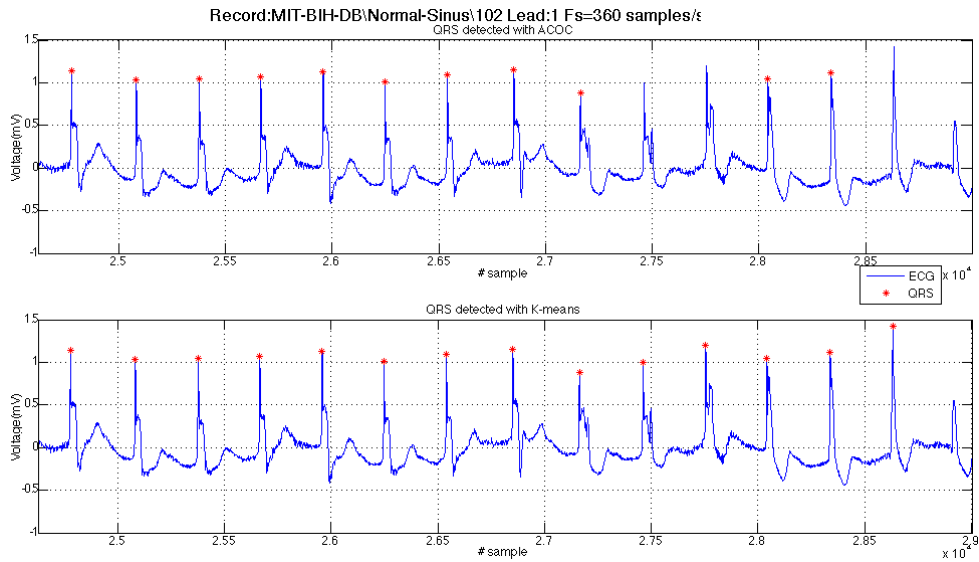


Figure 6.4: *Possible QRS complexes detection in the ECG (Four peaks are not detected with ACOC vs one with K-means).*

## Atrial Fibrillation (AF)

### Record: 04015 Lead: 1

In this record the possible QRS complexes form three or four groups, which are distinct (Fig.6.5). This leads to a satisfactory clustering with both algorithms. Their precision in QRS complex detection is almost identical with a light advantage of K-means (Table 6.1). However, in time cost ACOC has a great advantage (Table 6.2).
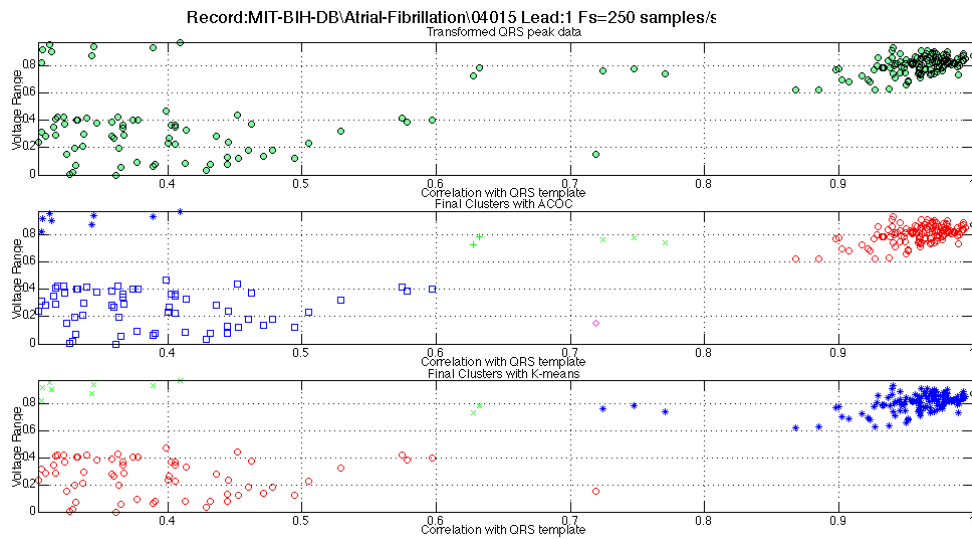


Figure 6.5: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*

Figure 6.6: **Possible QRS complexes detection in the ECG (Two peaks are not detected with ACOC).**

### Record: 04015 Lead: 2

In this record the possible QRS complexes form three or four groups, which are not distinct. However, the group with the template is clearly distinct. Although ACOC behaves like a binary classifier, a satisfactory clustering is achieved by both algorithms (Fig.6.7). Their precision in QRS complex detection is almost identical with a light advantage of ACOC (Table 6.1). ACOC is better in time cost (Table 6.2).

Figure 6.7: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*
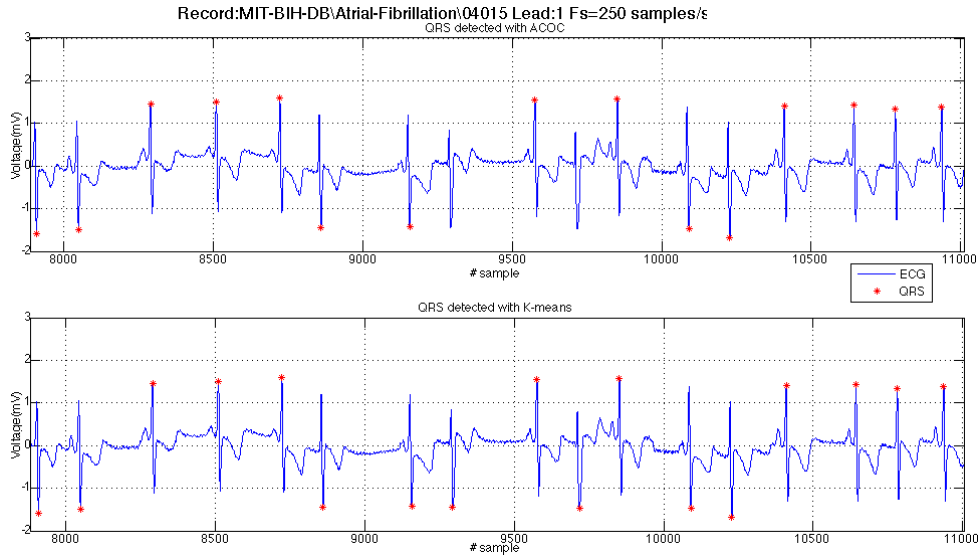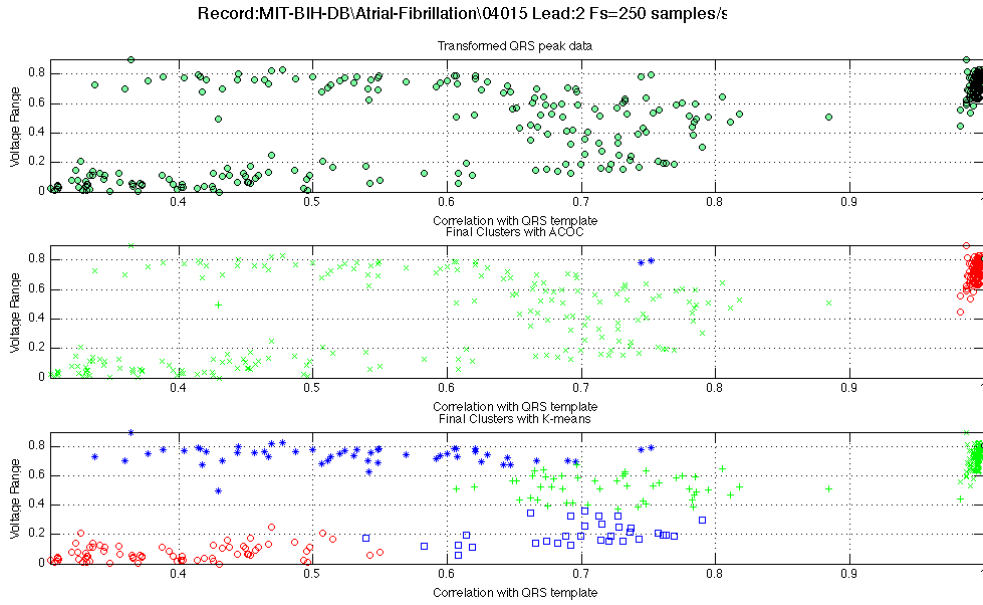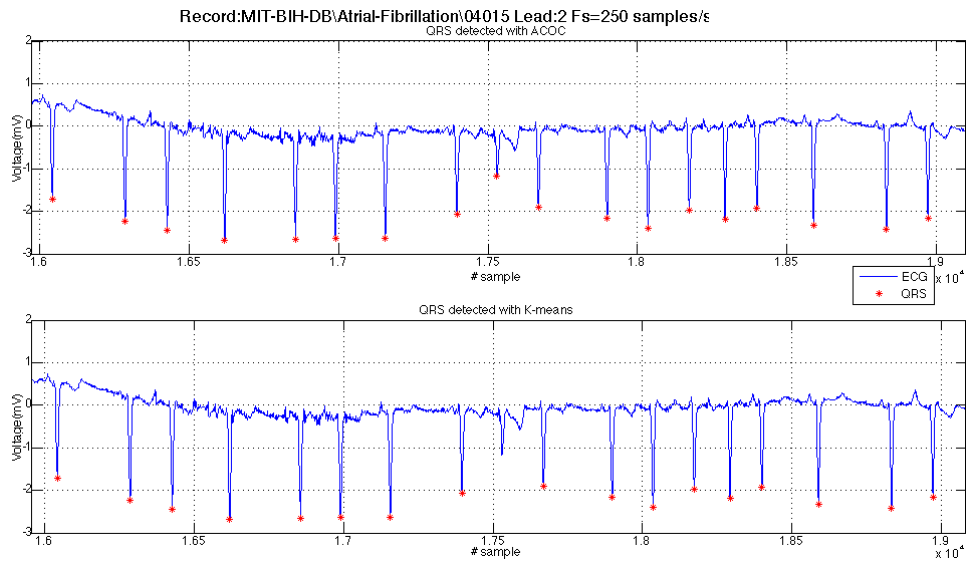


Figure 6.8: *Possible QRS complexes detection in the ECG (One peak is not detected with K-means).*

## ST Change (STC)

### Record: 302 Lead: 1

In this record the possible QRS complexes form two or three groups, which are distinct. A satisfactory clustering is achieved by both algorithms (Fig.6.9). Their precision in QRS complex detection is almost identical with a light advantage of ACOC (Table 6.1). ACOC is faster than K-means (Table 6.2). An erroneous QRS complex detection from both algorithms is presented. A false peak is detected by K-means (Fig.6.10).



Figure 6.9: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*

Figure 6.10: **Possible QRS complexes detection in the ECG (False detection of peak with K-means).**

**Record: 303 Lead: 2**

In this record the Possible QRS complexes form two or three groups, which are distinct. Although ACOC behaves like a binary classifier, a satisfactory clustering is achieved by both algorithms (Fig.6.11). Their precision in QRS complex detection is almost identical with a light advantage of K-means (Table 6.1). ACOC is faster than K-means (Table 6.2). However, in this part of the ECG the noise is sometimes greater than the QRS complex (Fig.6.12).
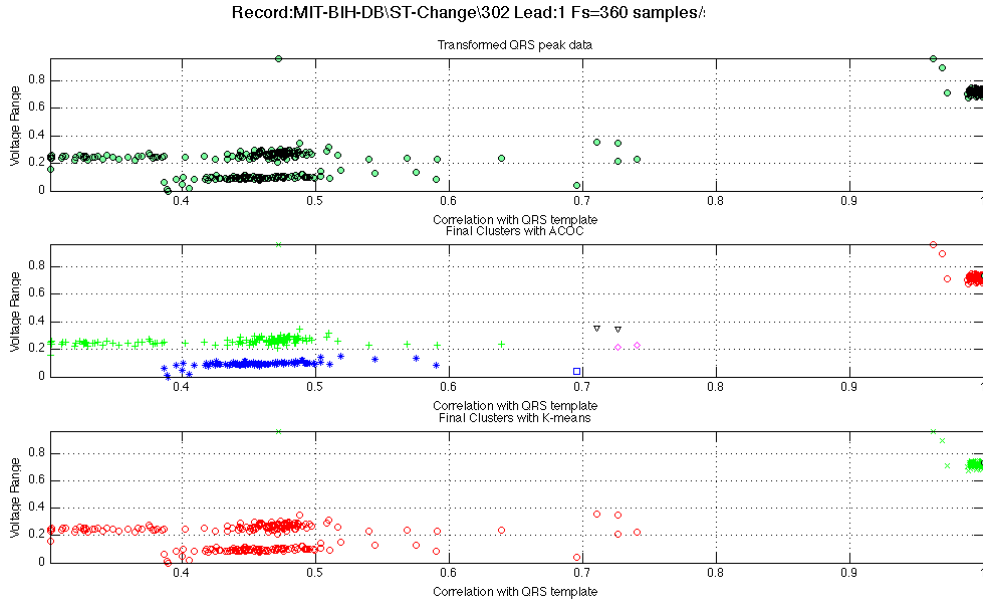
Figure 6.11: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*



Figure 6.12: *Possible QRS complexes detection in the ECG (Only a small part is presented).*

## Supra Ventricular Arrhythmia (SVA)

### Record: 800 Lead: 1

In this record the possible QRS complexes form two or three groups, which are distinct. Both algorithms behave like a binary classifier,although a satisfactory clustering is achieved (Fig.6.13). Their precision in QRS complex detection is identical (Table 6.1). ACOC is faster than K-means (Table 6.2).



Figure 6.13: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*
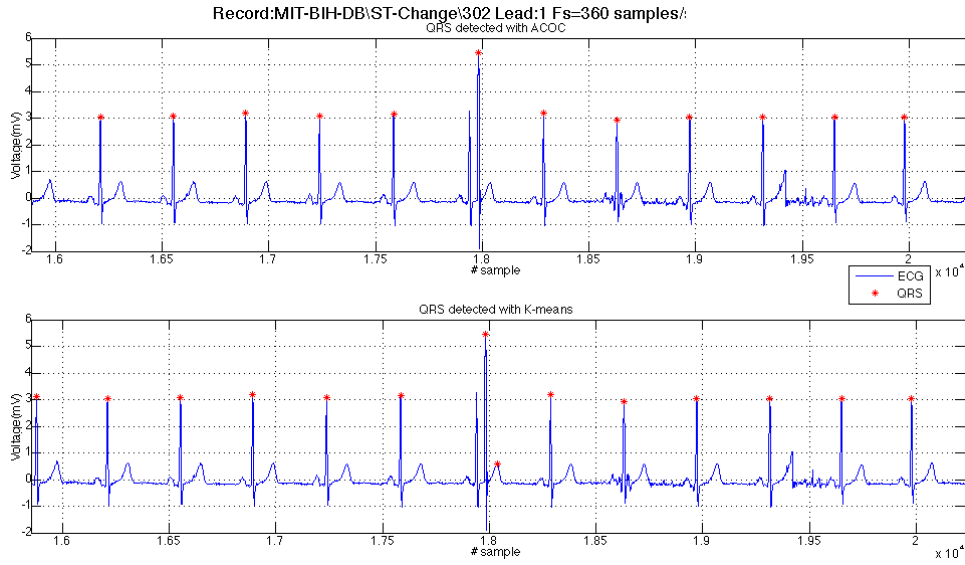
Figure 6.14: **Possible QRS complexes detection in the ECG (Only a small part is presented).**

**Record: 802 Lead: 2**

In this record the possible QRS complexes form two, which are distinct. A satisfactory clustering is achieved by both algorithms (Fig.6.15). Their precision in QRS complex detection is almost identical with a light advantage of K-means (Table 6.1). However, ACOC is faster than K-means (Table 6.2). An erroneous QRS complex detection from ACOC is presented (Fig.6.16).
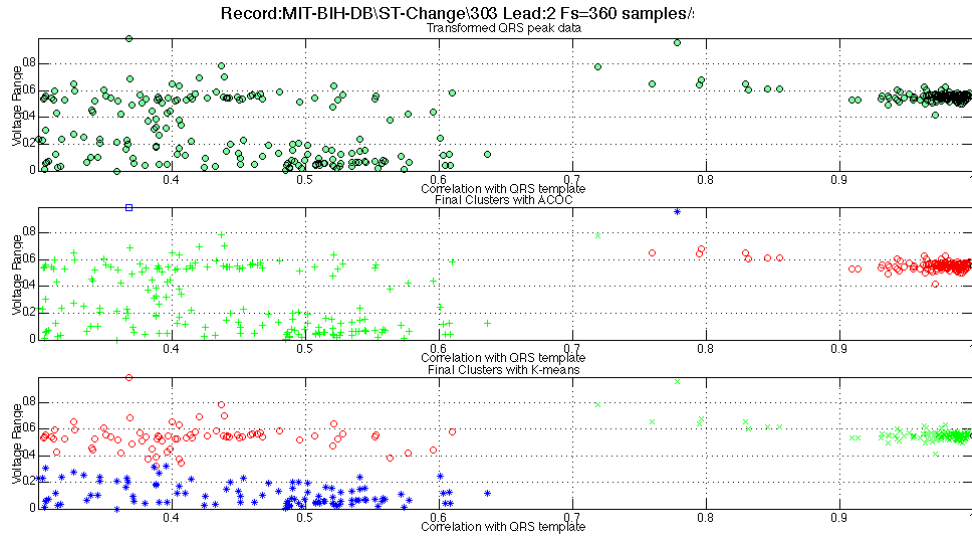
Figure 6.15: *Possible QRS complexes and clustering with ACOC and K-means (The green dot in the second and third plot is the template).*
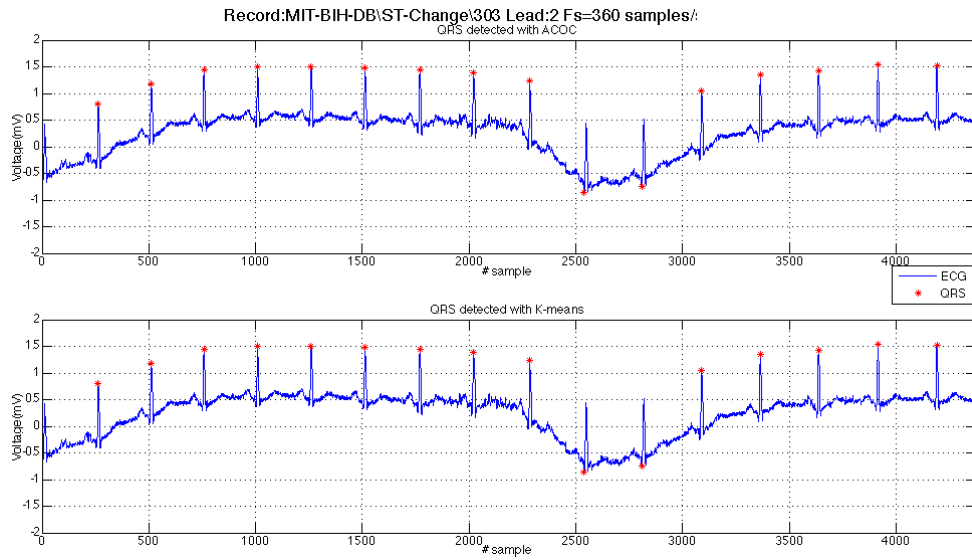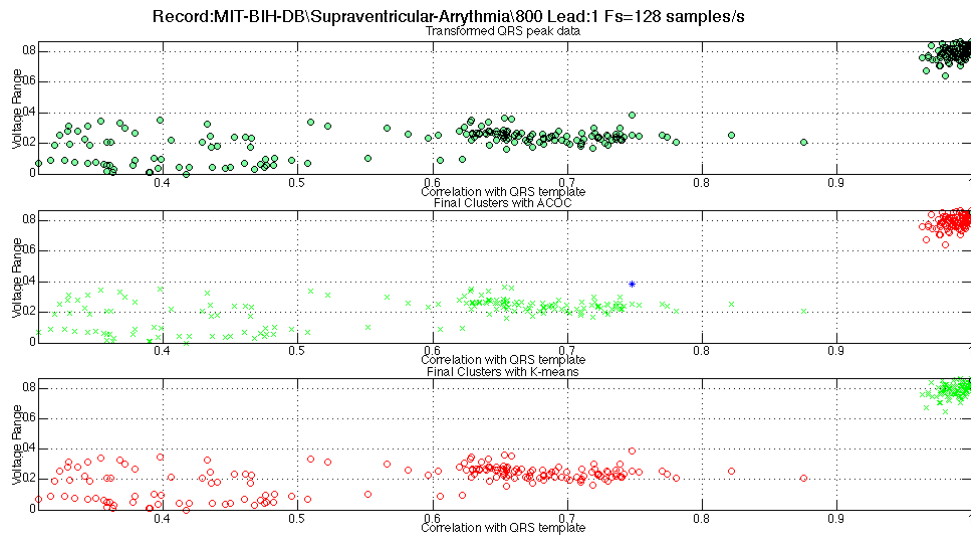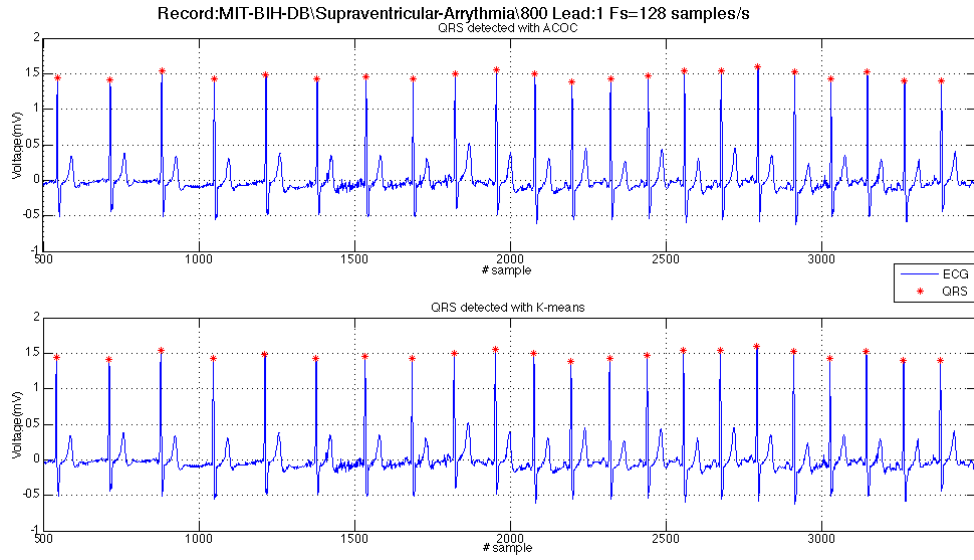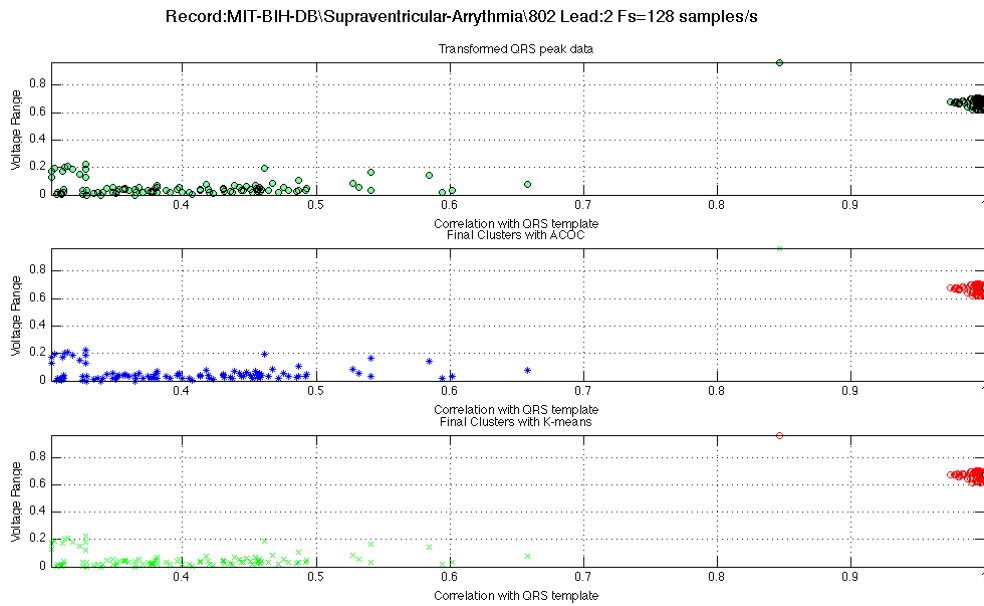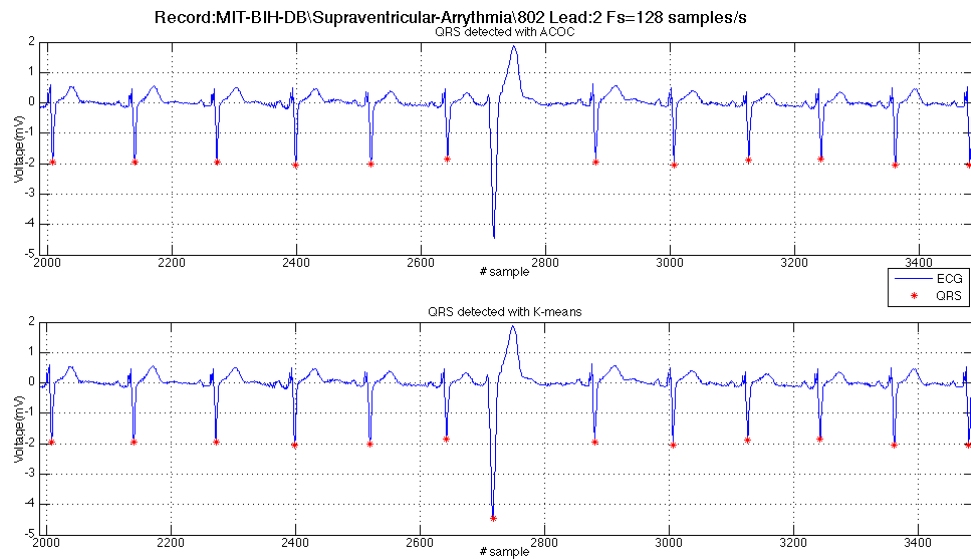


Figure 6.16: *Possible QRS complexes detection in the ECG (Peak not detected with ACOC).*

| Disease | Rec. | Lead | QRS Complexes | | | | |
| | | | Detected | | Real | Percentage (%) | |
| | | | ACOC | K-means | | ACOC | K-means |
|---|---|---|---|---|---|---|---|
| NS | 100 | 2 | 109 | 109 | 110 | 99.09 | 99.09 |
| | 102 | 1 | 95 | 100 | 108 | 87.96 | 92.59 |
| | 103 | 1 | 104 | 104 | 104 | 100 | 100 |
| | 103 | 2 | 104 | 104 | 104 | 100 | 100 |
| | 106 | 1 | 100 | 100 | 100 | 100 | 100 |
| Average | - | - | 102.4 | 103.4 | 105.2 | 97.41 | 98.34 |
| AF | 04015 | 1 | 128 | 130 | 132 | 96.97 | 98.48 |
| | 04015 | 2 | 128 | 127 | 132 | 96.97 | 96.21 |
| | 04043 | 1 | 156 | 156 | 157 | 99.36 | 99.36 |
| | 04043 | 2 | 154 | 154 | 155 | 99.35 | 99.35 |
| | 04048 | 2 | 107 | 107 | 107 | 100 | 100 |
| Average | - | - | 134.6 | 134.8 | 136.6 | 98.53 | 98.68 |
| ST | 302 | 1 | 96 | 95 | 97 | 98.97 | 97.94 |
| | 303 | 2 | 128 | 129 | 130 | 98.46 | 99.23 |
| | 304 | 1 | 79 | 79 | 80 | 98.75 | 98.75 |
| | 305 | 2 | 82 | 82 | 83 | 98.79 | 98.79 |
| | 306 | 1 | 95 | 95 | 95 | 100 | 100 |
| Average | - | - | 96 | 96 | 97 | 98.99 | 98.94 |
| SVA | 800 | 1 | 94 | 94 | 95 | 98.95 | 98.95 |
| | 802 | 2 | 89 | 90 | 90 | 98.95 | 100 |
| | 803 | 2 | 97 | 99 | 107 | 90.65 | 92.52 |
| | 804 | 1 | 108 | 108 | 109 | 99.08 | 99.08 |
| | 804 | 2 | 130 | 131 | 132 | 98.48 | 99.24 |
| Average | - | - | 103.6 | 104.4 | 106.6 | 97.22 | 97.96 |
| **Average** | - | - | **109.15** | **109.65** | **111.35** | **98.04** | **98.48** |

Table 6.1: ***QRS complexes detected, real and percentage of correct detection ACOC vs K-means (average values).***

In the transformed real data sets the performance of ACOC algorithm was similar to K-means. The clusters that finally formed by the two algorithms were not always the same but the cluster with the template QRS complex was always discriminated efficiently enough. This lead to an efficient QRS

complex detection with high percentage of real peaks found.  Both algorithms
has advantages and drawbacks.

| Disease | Rec. | Lead | Cpu time (s) | |
| --- | --- | --- | --- | --- |
| | | | ACOC | K-means |
| NS | 100 | 2 | 16.10 | 23.26 |
| | 102 | 2 | 20.58 | 24.56 |
| | 103 | 1 | 8.39 | 23.05 |
| | 103 | 2 | 6.36 | 23.65 |
| | 106 | 1 | 13.87 | 38.16 |
| Average | - | - | 13.06 | 26.54 |
| AF | 04015 | 1 | 8.71 | 42.68 |
| | 04015 | 2 | 20.72 | 38.29 |
| | 04043 | 1 | 5.49 | 23.72 |
| | 04043 | 2 | 4,44 | 25.23 |
| | 04048 | 2 | 5.83 | 29.13 |
| Average | - | - | 9.04 | 31.81 |
| ST | 302 | 1 | 18.25 | 22.84 |
| | 303 | 2 | 18.58 | 29.49 |
| | 304 | 1 | 5.07 | 27.13 |
| | 305 | 2 | 4.15 | 36.12 |
| | 306 | 1 | 6.97 | 23.61 |
| Average | - | - | 10.6 | 27.84 |
| SVA | 800 | 1 | 14.8 | 24.12 |
| | 802 | 2 | 8.41 | 20.11 |
| | 803 | 2 | 5.53 | 24.08 |
| | 804 | 1 | 3.75 | 30.61 |
| | 804 | 2 | 4.83 | 31.77 |
| Average | - | - | 7.46 | 26.14 |
| **Average** | - | - | **10.04** | **28.08** |

Table 6.2: ***Cpu time for various ECG (ACOC vs K-means).***

ACOC was faster in all data sets that was tested providing a suitable al-
gorithm for fast QRS complex detection and furthermore, for a tool in a real
time diagnostic system.  As an unsupervised method of clustering, there was

no need for providing ACOC with parameters that required previous knowledge of the data set. However, in large data sets with difficult to discriminate data, there was a need for adjusting the parameters of the algorithm so as to converge in an acceptable solution. Trial and error method was used for selecting the best parameters for such data sets which is a significant drawback because it requires complete knowledge of the formulas and the form of the algorithm.

On the other side, K-means is an efficient and fast algorithm that in most of the cases clusters the transformed ECG data quite well. The algorithm is usually very fast if the number of the clusters is provided. If not, K-means needs to run multiple times for different number of clusters till achieving the best convergence. This is the major drawback of the algorithm as time cost raises. Also K-means has to be run with a large value of iterations as its initialization is not reliable and can lead to a possible erroneous result.

# Chapter 7

# Conclusions and Future Work

ECG signals were processed, after being transformed and clustered, in order to detect the QRS complexes.

Most QRS complexes were detected using both ACOC and K-means while computationsl time was significantly lower using the ACOC algorithm.

According to the results and findings coming from the process, the complete method is considered as a reliable method of detecting QRS complexes in a typical ECG with high rates of real QRS detection achieved (over 90%).

Finally, in terms of future work there are many options. Implementation of the algorithm in high performance FPGA and its use in combination with an electrocardiograph can result in a complete near real time diagnostic tool for assisting medical diagnosis.

The adaptive adjustment of all the parameters of the ACOC algorithm could make clustering more efficient with the smaller time cost.

The overall method could find application in detecting more features on the ECG, for example, the T wave. This could be done by subtracting the

already detected QRS complexes, and reapplying the techniques used for QRS complex detection to the remaining signal.

# References

[1] B.U. Kohler, C. Henning, R. Orglmeister, "The principles of software QRS detection", IEEE Eng. Med. Biol. 21 (1), 42-57, 2002.

[2] R.M. Rangayyan, "Biomedical Signal Analysis: A Case-study Approach", Wiley-Interscience, New York, 2001, pp. 1828.

[3] R. Silipo, C. Marchesi, "Articial neural networks for automatic ECG analysis", IEEE Trans. Signal Process. 46 (5), 1998, 14171425.

[4] J. Pan, W.J. Tompkins, "A real-time QRS detection algorithm", IEEE Trans. Biomed. Eng. BME-32 (3), 1985, 230236.

[5] C.W. Li, C.X. Zheng, C.F. Tai, "Detection of ECG characteristic points using wavelet transforms", IEEE Trans. Biomed. Eng. 42(1) 1995, 2128.

[6] K.V. Suarez, J.C. Silva, Y. Berthoumieu, P. Gomis, M. Najim, "ECG beat detection using a geometrical matching approach", IEEE Trans. Biomed. Eng. 54 (4), 2007, 641650.

[7] X. Xu, Y. Liu, "ECG QRS complex detection using slope vector waveform (SVW) algorithm", Proceedings of the 26th International Conference of the IEEE EMBS, 2004, pp. 35973600.

[8] Kumar, Abbas, Fausto, Robbins and Cotran, "Pathologic Basis of Disease", 7th Ed. p. 556.

[9] Richard E. Klabunde, Cardiovascular Physiology Concepts, "Cardiac Anatomy", `http://www.cvphysiology.com/Heart%20Disease/HD001.htm`, accessed August 14, 2009.

[10] John A. McNulty, Ph.D., Loyola University Medical Center, "Heart Development ", `http://www.meddean.luc.edu/lumen/MedEd/GrossAnatomy/thorax0/heartdev/main_fra.html`, accessed August 13, 2009

[11] Guyton, A.C. Hall, "Textbook of Medical Physiology", J.E. 11th Ed. 2006 Philadelphia: Elsevier Saunder.

[12] Richard E. Klabunde, Cardiovascular Physiology Concepts, "Cardiac cycle", `http://www.cvphysiology.com/Heart%20Disease/HD002.htm`, accessed September 4, 2009.

[13] Health Sciences Library of the University of Utah, `http://library.med.utah.edu/kw/pharm/hyper_heart1.html` accessed August 28, 2009.

[14] TCHP Education Consortium, "*ECG Rhythm Interpretation Primer*", 2004-2007.

[15] Dubin D. "Cardiac Electrophysiology and Ion Channels", University of Maryland School of Medicine, Lecture 169, 1998

[16] Nursecom Educational Technologies, *"Six Second ECG Guidebook"*, 2003.

[17] K.P. Misra, *"A Primer Of Ecg"*, chapters 4-5, Orient Longman Limited, 2006.

[18] A. Ajith, G. Crina, R. Vitorino, *"Stigmergic Optimization"*, Studies in Computational Intelligence, volume 31, 2006.

[19] Marco Dorigo, *" Optimization, Learning and Natural Algorithms"*, PhD thesis, Politecnico di Milano, Italie, 1992.

[20] A. Colorni, M. Dorigo et V. Maniezzo, *"Distributed Optimization by Ant Colonies"*, Paris, France, Elsevier Publishing, 1991.

[21] Marco Dorigo, Thomas Stutzle, *"Ant Colony Optimization"*, "A Bradford book", The MIT Press, 2004.

[22] Ling Chen, Li Tu, Hong-Jian Chen *" Data clustering by ant colony on a digraph"*, Fourth International Conference on Machine Learning and Cybernetics, Guangzhou, August 2005.

[23] Marco Dorigo, Iridia *"Ant colony optimization"*, http://www.aco-metaheuristic.org/about.html, Universite Libre de Bruxelles, accessed April, 2009.

[24] MacQueen, J. B. (1967) "Some Methods for classification and Analysis of Multivariate Observations", Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability, pp. 281297.

[25] Peter J. Rousseeuw, " *Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis*", Computational and Applied Mathematics, 1987.

[26] D. T. Kremastinos, "Cardiology," Medical Publish P.X. Pasxalidis, p. 183