

Applied Mathematics and Informatics Program

Multivariate Symbolic Aggregate Approximation for ECG Analysis

Moritz M. Konarski

A Thesis Submitted to the Applied Mathematics and Informatics Program of American
University of Central Asia in Partial Fulfillment of the Requirements for the Degree of
Bachelor of Arts

Author

Moritz M. Konarski

Certified by Thesis Supervisor

Professor Taalaibek M. Imanaliev

Accepted by

Sergey N. Sklyar

Head of Applied Mathematics and
Informatics Program, AUCA

May 18, 2021

TODO Don't forget to put simply May, 2021 or something Bishkek, Kyrgyz Republic

ABSTRACT

Electrocardiograms are the most common tool used to diagnose heart diseases, which claim more lives each year than any other disease. Since their invention, **CITE** cite when, and who electrocardiograms needed to be analyzed by a trained professional like a cardiologist. Since **CITE** cite when it became a thing, they can be analyzed using computers and computer-assisted methods. These methods can be more accurate, faster, and more versatile than humans. One discord discovery method is the Symbolic Aggregate Approximation, which transforms an electrocardiogram into a shorter, symbolic form. This form is faster and simpler to analyze.

Multivariate Symbolic Aggregate Approximation takes more than one electrocardiogram lead into account and should thus be more accurate when it comes to discord discovery using Heuristically Ordered Time series using Symbolic Aggregate Approximation.

This paper shows, with **TODO** insert significance level, that Multivariate Symbolic Aggregate Approximation increases the sensitivity of HOTSAX compared to Symbolic Aggregate Approximation.

Keywords: acute cardiac ischemia, ECG, mathematical modeling

TODO - relevance in 1 line if possible

- result summary
- major implications
- quantitative
- no abbr, citations

1. what did I do
2. why did I do it? which questions
3. how was it done? methods
4. what are the results?
5. why do the results matter? + at least 1 application

ACKNOWLEDGEMENTS

This research project is part of a Faculty Research Grant **TODO** insert specifications or ID here. I also want to thank my Supervisor, Professor Taalaibek M. Imanaliev, for his invaluable help and precise feedback.

TODO thank Imanaliev; mention the grant?

DRAFT

TABLE OF CONTENTS

1	Introduction	6
2	Background and Related Work	9
2.1	Time Series and Time Series Analysis	9
3	State of Computerized ECG Analysis	11
4	Methods	14
4.1	Mathematical Foundations	14
4.2	SAX	14
4.3	MSAX	15
4.4	HOTSAX	15
4.5	Statistical Analysis of Results	15
4.6	Implementation	15
4.6.1	ECG acquisition	15
4.6.2	preprocessing	16
4.6.3	SAX	16
4.6.4	MSAX	16
4.7	Statistical Evaluation	16
5	Results	17
5.1	First Run	17
5.2	Second Run	17
5.3	SAX	17
5.4	MSAX	17
5.5	MSAX vs SAX	18
6	Discussion	19
6.1	MSAX vs SAX	19
7	Conclusion	20
	References	21

TODO table of figures? **TODO** table of tables? **TODO** table of abbreviations?

TODO Sections (not explicit):

- motivation
- objectives and contribution
- falsifiable hypothesis
- introduce relevance of sensitivity here

TODO write last;

- hook to get reader interested
- brief overview of current research state
- why was this work required?
- statement of hypothesis
- enough background to clear up goal

- only stuff immediately relevant to thesis and goal
- scope of work: what will I and won't I look at?
- verbal table of contents / roadmap
- make clear what is new / novel

TODO Just start from scratch:

- introduction with who data for ihd, arrh
- prevention and detection is important
- ECG is THE method for that
- how an ECG works; pros and cons of an ECG
- how an ECG is also a multivariate time series
- where does automated ECG analysis come into play?
- time series methods to ECG
- discords that then lead to next point: connect to the 5 steps, signal which I will look at, where is SAX, HOTSAX, where is the cardiologist...
- the approach of sifting out the abnormal beats and having a cardiologist look at them instead of doing everything internally
- SAX, MSAX, HOTSAX and my hypothesis
- use definitions etc

1 INTRODUCTION

TODO UPDATE: In the year 2016, over 9.4 million people worldwide died of ischemic heart disease (IHD). IHD is responsible for 16.6% of all deaths, making it the most common cause of death globally. All forms of cardiovascular disease make up 31.4% of all deaths (17.9 million). Death caused by IHD disproportionately affects people over 50 years of age, with 91% of deaths for men and 95% of deaths for women occurring in that age range. In Kyrgyzstan, 13% of all deaths in 2016 were caused by IHD [who2018].

TODO FOCUS MORE ON MY ACTUAL TOPIC:
Ischemic heart disease is characterized by restricted blood flow to an area of the heart, causing it to not receive enough blood and oxygen. Blood flow restriction is caused by a blockage (or narrowing) in a blood vessel supplying the heart muscle. An artery can be blocked by a blood clot, but the most common cause is plaque buildup, which is called atherosclerosis. If the circulation to the heart is completely blocked, the cells in the heart muscle begin to die. This is called myocardial infarction, more commonly known as a heart attack. The deprivation of oxygen the heart experiences leads to the characteristic chest pain commonly associated with heart attacks [iom2010].

TODO mention arrhythmia too, it is what the mit database looks at
IHD can be diagnosed before it leads to a heart attack. The diagnosis can be performed based on a patient's medical history, pharmacologically induced stress, or stress induced by physical exercise. During an exercise stress test, an electrocardiograph (sometimes combined with other methods) records the patient's heart activity, resulting in an electrocardiogram (ECG) [iom2010].

TODO make this its own paragraph section

- TODO**
- heart's electrical activity
 - up to 12 leads
 - common medical diagnostic tool
 - electricity is what causes the contraction
 - this can be measured on the skin
 - a bit on ECG theory
 - specific electrodes and positions
 - mention ion flow

The ECG is a diagnostic tool used to evaluate patients with (suspected) heart problems. It is a non-invasive, real-time, and cost-effective method that may be used to diagnose IHD. It is the most common tool used for cardiac analysis and diagnosis [alghatrif2012, kligfield2007, xie2020]. The most common form of the ECG is the 12-lead variant. The 12-lead ECG consists of 6 leads connected to the limbs and 6 leads connected to the torso of the patient. The leads record the differences in electrical potential between the places on the body that they are attached to. This reflects the differences in

voltage that the heart experiences with each heart beat because those voltage differences are conducted by the body. The measurements are taken in millivolts (mV). The ECG represents the state of the heart; a recorded ECG has the shape of a wave (the ECG wave) [kligfield2007, xie2020].

TODO datasets are available online, the most significant leads tend to be included

If the state of the heart beat changes as the result of a disease like IHD (changing the measurable potentials or their occurrence over time), the ECG is able to record these changes.

The characteristic shape of an ECG for two heart beats is shown in Figure 1.1; the figure is taken from [wasilewski2012]. The figure has been annotated to show the significant features of an ECG. The peaks (or waves) P, Q, R, S, T, and U, as well as the segments between them, are the focus of ECG analysis. Multiple points together form what is called a complex; **TODO** FIX: the QRS complex is a good example of this. Using these waves, the heart activity can be described and analyzed. In an ECG, the P-wave is the result of the atria depolarizing, which is the process of blood entering the heart as the first step in a heart beat. The QRS complex represents ventricular depolarization, the contraction of the heart causing it to pump blood. The T-wave is the return of the ventricle to its polarized state. The U-wave is only present in roughly 25% of the population and may be caused by mechanical-electric feedback. The RR interval can be used to calculate the heart rate because it represents one complete heart beat [wasilewski2012]. The shape of the P, Q, R, S, T, and U waves as well as the duration of various intervals between them are used as indicators of cardiac diseases.

TODO UPDATE TO TIKZ FIGURE:

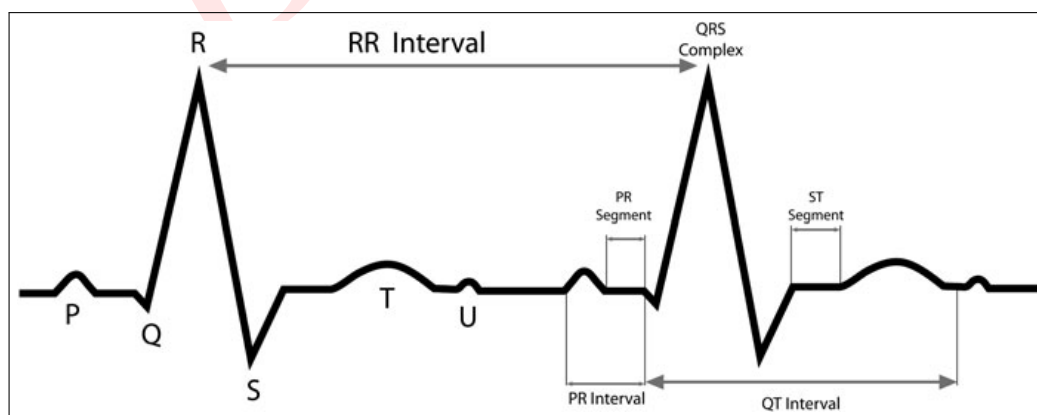


Figure 1.1: A schematic of an ECG waveform, annotated; from [wasilewski2012]

TODO also: lots of data, cannot be analyzed quickly

Using an ECG to diagnose a cardiac condition is difficult in practice. Small changes in the components of the ECG can be indicators of diseases and those changes can be overlooked, even by trained and specialized physicians. The chance to make a mistake is even higher for non-specialized physicians and trainees [alghatrif2012, xie2020].

TODO make it about assisting in their diagnosis by pointing out important segments

-> discords

For the diagnosis of IHD, changes in the ST-segment and T-wave are of particular interest. An elevation of the ST-segment compared to a normal heart beat is one of the main indications of IHD and myocardial infarction. A downward depression of the ST-segment, especially in combination with chest pain, is another indication of IHD. The changes in the ST-segment are thought to be caused by current flow between healthy heart muscle and ischemic heart muscle [rautaharju2009, wasilewski2012].

TODO after introducing idea of discord discovery, introduce ECGs as time series:

- ECGs are multivariate: 2-12 leads
- discrete: measured at discrete points
- ordered sequences: come one after another after equal time segments

TODO then, time series analysis methods are ready to be applied to ECGs

TODO what is required of the method:

- fast
- accurate
- adaptable
- versatile

The diagnosis of IHD on the basis of an ECG is time sensitive. If a patient has IHD or suffers from a heart attack, treatment has to be started as soon as possible. Some forms of treatment are most effective in the first 3 hours after symptom onset and lose most of their effectiveness after 9 to 12 hours. The diagnosis required for treatment to begin should thus be as quick as possible. The ECG delivering information in real-time is an advantage here, even though there are more time consuming methods that can deliver more accurate results than an ECG [herring2006].

TODO WHY automated ECG analysis:

The widespread use of ECGs and the time-sensitive nature of their application as diagnostic tools makes errors, delays, or inconsistencies in their interpretation unacceptable. A recent approach to minimizing this problem is the application of computer technology in ECG recording, storage, and analysis. The main steps of computerized ECG analysis are [kligfield2007]

TODO fix this one, adapt to my writing (1) signal acquisition and filtering, (2) data transformation or preparation for processing, (3) waveform recognition, (4) feature extraction, and (5) classification or diagnosis.

This research will investigate steps (3) and (4) through the use of different feature extraction algorithms. The ECG data will be retrieved from the European ST-T Database. This database provides ECG recordings that can be used as trial data to test feature extraction algorithms. The European ST-T Database contains annotations made by cardiologists indicating the ST-segment, T-wave, and their changes. They also include information about the suspected disease [physionet, taddei1992]. This information can be used to determine the effectiveness of the feature extraction algorithms.

2 BACKGROUND AND RELATED WORK

This section provides background information on time series and ECGs, as well as their analysis.

2.1 Time Series and Time Series Analysis

TODO insert some formulas and symbols, make sure that they are consistent

A time series is a set of values recorded at specific times. A common form of time series are discrete-time time series (often simply called discrete time series). Discrete time series are time series whose values are recorded at discrete points in time, the most common example of this are time series with values recorded at fixed intervals. Continuous-time time series are time series that are recorded continuously over an interval [1]. Time series are used in many disciplines to record information on time-dependent processes, e.g. stock prices in economics, the sun's activity in physics, and the heart's activity in medicine. Time series can be recorded digitally, physically, or physically recorded and then digitized. The recorded data can then be used to gain insight into the underlying processes. To gain insight from time series, the relevant information needs to be extracted from the time series, a process that is often called data mining. Data mining of time series is a vast discipline that includes the visualization, forecasting, indexing, clustering, anomaly detection, classification, and summarization of time series. Time series indexing is the **TODO** insert some symbols here process of finding the most similar time series to a given time series. The division of time series into groups of similar time series is called clustering. Identifying parts of a time series that are not "normal" (don't fit certain parameters) constitutes anomaly detection. To classify a time series means to assign it a label based on its features, e.g. "sick" or "healthy". Finally, summarization is the reduction of the complexity (often length) of a time series while preserving its main features [2, 3].

Challenges for time series analysis include the often very large data sets that are difficult for humans to analyze and take up considerable digital storage space. Analyzing very large data sets requires a large amount of computational power because most time series analysis algorithms do not scale well for large data sets [2]. To mitigate this issue, time series dimension reduction (also known as dimensionality reduction and time series representation) is used. Dimension reduction transforms a "raw" (unmodified) time series into a representation that is simpler but nonetheless resembles the raw time series. This can be achieved by either using a method that reduces the number of values in a time series, or by extracting only the relevant features from the time series [3]. According to [3], there are four types of dimension reduction methods: 1. non-data adaptive, 2. data adaptive, 3. model-based, and 4. data dictated. **TODO** fix this next segment, add more citations, maybe short descriptions **TODO** turn these into subsections? **TODO** also cite shieh2008 here, has table and explanations too **TODO** maybe just mention the categories, explain what they mean, then do a "in this research the focus is on SAX, a ..." Non-data adaptive

methods operate on time series segments with a fixed size to reduce the dimension, they are useful for comparing multiple time series with each other. These methods include wavelets, the Discrete Wavelet Transform, and the Piecewise Aggregate Approximation (PAA). Data adaptive methods use non-fixed size segments and aim to fit the raw data more closely. Examples of data adaptive methods are the Piecewise Polynomial Approximation, Piecewise Linear Approximation, and Piecewise Constant Approximation. Model based methods use stochastic methods such as Markov Models and Hidden Markov Models, including the Auto-Regressive Moving Average (ARMA). All three methods mentioned above have varying compression ratios, based on user-given input. Data dictated methods derive their compression ratios from the data automatically, e.g. Clipped [3].

TODO to put in:

- previous research in this area
- current research in the area
- show why more work is still important
- enough background to really clear up the problem
- stuff relevant to question, and some background

TODO Sections:

- overview of methods for ECG analysis
- main elements
- current foci
- use the confusion matrix at all?
- support assertions made in introduction
- arrive at natural conclusion that SAX/MSAX/HOTSAX should be investigated
- describe all 4 statistical measures in some detail

TODO Structure:

- follow the same overall structure as the introduction, order of topics
- background on history of ECG
- how an ECG works in more technical terms
- how do cardiologists detect heart diseases?
- move the annotated graph down here and leave the general graph in intro?
- which time series methods are being applied to ECGs? strengths and weaknesses
- what are current hot topics in this field?
- more research on sax, msax, hotsax and what people have done with it
- ecg applications; why is it relevant
- RESEARCH MORE ABOUT "ecg discord discovery algorithm"
- use some of their arguments to support my choice
- try to make it flow so that my choices seem to come from the literature analysis
- maybe add a section that talks about how to evaluate these types of algorithms

3 STATE OF COMPUTERIZED ECG ANALYSIS

Recent advances in computer technology have enabled the use of computers in every aspect of ECG acquisition, processing, analysis, and storage. In light of these developments, the American Heart Association published recommendations for the interpretation and standardization of the ECG. They recommend that the low-frequency cutoff for low-frequency filtering of an ECG should be 0.05 Hz or 0.67 Hz for filters that do not exhibit phase distortion. For high-frequency filtering they recommend a cutoff of at least 150 Hz. For the storage of digital ECG samples (at 500 samples per second), it is recommended use use compression with an error of less than 10 microvolt [Kligfield2007].

Xie2020 provide an overview of the current approaches to computerized ECG analysis. The standard approach to using computerized methods in ECG analysis is comprised of four steps (1) denoising of the raw ECG signal(s), (2) feature engineering, (3) dimensionality reduction, and (4) classification. To denoise an ECG, digital filters are often used. Their drawbacks are that they only filter out very specific frequencies. Because noisy ECGs contain different types of contaminations, digital filters can be inaccurate. Using wavelet transforms for denoising has the advantage that noise can be more precisely targeted and the clean signal reconstructed afterwards. Choosing appropriate wavelet parameters can be challenging and methods to optimize this process have been proposed. Empirical mode decomposition is the third option generally employed to denoise an ECG. It does not require the user to set parameters but it can lead to a mixing of oscillations of different time scales.

After the signal has been appropriately denoised, feature engineering is performed. Feature engineering is the process of extracting features that are relevant for diagnosis from the many points the ECG signal contains. The main features targeted for extraction are the PQRST features mentioned in the introduction. The fast Fourier Transform provides a way of analysing the frequency domain of the ECG signal, enabling the detection of the QRS complex and other features. The missing time information in the fast Fourier Transform can lead to difficulties in detecting time-dependent features. The short-time Fourier Transform adds time information to the fast Fourier Transforms data. This can increase the accuracy of the feature extraction. This transform has the drawback that there is a tradeoff between the time and frequency resolutions. Wavelet transforms can also be used for feature extraction. They have the advantage that they are suitable for all frequency ranges. Choosing the right wavelet base for the desired application can be a challenge. The discrete wavelet transform is the most widely used wavelet transform, thanks to its computational efficiency. Statistical methods are also used to extract features from ECGs; those methods are generally less affected by noise in the signal.

After the features of the ECG have been extracted, it is often necessary to reduce the number of features. The reason for this is that a large number of features, despite their high accuracy, require a high amount of computation to classify. This lengthy computa-

tion can negate the advantages gained by high accuracy. This process sacrifices a certain amount of information and sometimes precision, but significantly speeds up the classification. Feature selection is a process that attempts to select a subset of the original data that adequately describes the whole data. Feature selection can be performed by a filter that filters out unnecessary attributes based on some metric. This method is relatively simple, but the filtering process removes data and thus negatively impacts the precision of further steps. Feature extraction on the other hand uses dimensionality reduction methods to keep as much of the original information as possible. Principal component analysis preserves as much of the variance in the original data as it can. Other algorithms focus on separating classes of data, pattern recognition, or retaining the structure of the original data.

The final stage of the ECG processing is the classification stage. In this stage judgements are made based on the prepared input data and the result should be a disease diagnosis. In the early stages of computerized ECG analysis classification was performed by algorithms based on human actions when reading an ECG. Those algorithms were basic and not particularly accurate. Currently, the classification at the end of the preparation process is performed by a machine learning algorithm. Such models include the k-nearest-neighbors model which classifies points into groups but which is very expensive to calculate for high-dimensional data. Support vector machines are used for pattern recognition and are able to work with small samples. Artificial neural networks are robust and can work with complex problems, they are generally more accurate than support vector machines. The newest approach is to forego the stages discussed here and use a single neural network to perform all the required tasks "end-to-end". These networks are fed raw data and the denoising, feature extraction, selection, and classification is performed internally by the model [xie2020].

The end-to-end approach to ECG analysis is a relatively new development and is being actively researched. The more traditional method using denoising, features engineering, and classification as separate steps is also still relevant. The combination of denoising and feature extraction with a machine learning classifier can lead to very good results. **prasad2018** use the fast Fourier Transform to extract features from an ECG and then employ a multi-objective genetic algorithm to detect abnormal ECG signals with high accuracy. **vaneghi2012** compare 6 common feature extraction techniques with respect to their detection of ventricular late potentials. The compared methods are the autoregressive method, wavelet transform, eigenvector, fast Fourier Transform, linear prediction, and independent component analysis. **valupadasu2012** use the fast Fourier Transform to analyze the energy level in different frequencies in the ECG of patients with IHD. They find that the energy is distributed differently, allowing the distinction of ECGs with IHD from those without IHD. ^{kaur2016} **Kaur, Rajni, and Marwaha [4]** analyzed ECG signals with both the wavelet transform and principal component analysis. They found that the wavelet transform outperformed principal component analysis for the detection of heart beats in an ECG. Their model achieved an error rate of 0.221% of incorrectly classified heart beats

.

4 METHODS

TODO goals:

- reader can assess believability of results
- all information necessary to replicate the research
- describe all materials, procedure, ect
- all the formulae
- state all the limitations of the methods and the ones I impose myself
- analytical methods and languages

TODO answer questions:

- can someone else accurately replicate the study
- can the data be obtained again
- are all parts / instruments described with enough accuracy
- is the data freely available
- can the statistical analysis be repeated
- can the algorithms be replicated?

TODO Sections:

- general overview -> flowchart
- then explain each element of the flowchart one by one
- use formulae etc
- nice amount of tikz graphs
- section on implementation with details and the more important elements
- use another flow chart?
- use graphs to illustrate all important elements
- make a data description section that describes my process of data handling; which database
- explain the parameters that the methods have and what they mean
- describe how a got all the data

This section explains the methods used in this research. **TODO** create flow charts for all this shit to make it simpler. First methods section for the analytical methods in a mathematical way.

4.1 Mathematical Foundations

- section for the workings
- explain theoretical foundations of the approach
- what is it grounded in, who, what, when

4.2 SAX

- idea

- normalization
- dimensionality reduction
- discretization
- distance measure
- **TODO** all with graphs and formulas

4.3 MSAX

- idea
- normalization
- dimensionality reduction
- discretization
- distance measure
- **TODO** all with graphs and formulas
- **TODO** points out differences to SAX

4.4 HOTSAX

- what is hotsax
- its theoretical foundations
- advantages, disadvantages
- how does it work

4.5 Statistical Analysis of Results

- explain true positive, true negative, and so on
- explain recall, accuracy, precision, f1
- explain why recall was chosen and if that is fair
- introduce the correlations that we would expect to find if my hypothesis is true and also the ones that would disprove it
- which types of correlation, significance testing, and modeling will be used and why; what are the justifications

4.6 Implementation

How I implemented the above stuff. Languages, approaches, hurdles, all the details needed to reproduce this research. Also mention the simplifications I chose to make and why: no sliding window, only even divisors, only divisors within sampling frequency and cutting ECG to even multiple of sampling frequency.

4.6.1 ECG acquisition

flow chart for process

- where to download
- what exactly are the ECGs
- where do they come from
- technical parameters of them
- the physionet suite
- annotations, what they mean, how I can get them, etc

4.6.2 preprocessing

flow chart for process

TODO the codes and constants given for each thing

- how were they preprocessed
- physionet suite
- my script and what it does and why
- problems and limitations of this
- libraries used

4.6.3 SAX

TODO how was the whole data thing handled, how is the data created

flow chart for process

- how was sax implemented
- how does HOTSAX work here
- libraries used

4.6.4 MSAX

flow chart for process

- how was sax implemented
- how does HOTSAX work here
- **TODO** point out differences to SAX
- libraries used

4.7 Statistical Evaluation

- reading the data into R
- summarizing the data
- the summarized data files
- libraries used

5 RESULTS

TODO must include:

- present the actual results and findings
- range of validity
- DO NOT INTERPRET
- mention positive and negative results
- give enough information for others to make their own judgements
- use subheadings
- key results in clear statements at paragraph beginnings
- could be short

TODO Sections

- use confusion matrices for what is vs what was predicted [p. 44 anacleto2019]
- compare all the parameters and their influence

5.1 First Run

- parameters for this run
- why could I not let it continue
- what did this run indicate -> what did I change and modify for the next run
- keep in mind that the lower recall can be caused by the way I do the ECG checking, and that I did not want to assign data to segments that did not have it before out of fear that I would invent results.

5.2 Second Run

5.3 SAX

influence and significance of all the major parameters:

- k
- paa count
- subsequence count
- alphabet size
- which ones seem to be the best

5.4 MSAX

influence and significance of all the major parameters:

- k
- paa count

- subsequence count
- alphabet size
- which ones seem to be the best

5.5 MSAX vs SAX

Comparing SAX to MSAX is done using the recall value defined in **TODO** reference. Investigating the correlation between the methods (represented by a 1 for SAX and a 0 for MSAX), yields the correlation coefficient of -0.25. This coefficient indicates that for all investigated parameter combinations, the use of the MSAX method is weakly correlated with an increase in recall. When a specific set of parameters is selected and the correlation analysis is repeated, the correlation coefficient is -0.73, indicating a strong correlation. Here $k = -1$ and $paa_count = 12$.

- just the results that are gained directly from the data
- put results in graphs and tables to make them referencable

6 DISCUSSION

TODO what do your results mean?

- make sure to distinguish results and interpretation "I infer that"...
- start with summary of important results
- major patterns
- relations, trends, generalizations
- exceptions
- likely causes for the things above?
- agreement / disagreement with previous work
- interpret with respect to hypothesis
- hypothesis testing here?
- other questions this relates to?
- consider all possibilities
- what new insight have we gained?
- include the supporting evidence for each line of reasoning
- what is the significance of the current results
- cite related results

6.1 MSAX vs SAX

TODO use the results from the previous section to come to a conclusion

- do proper hypothesis testing of my hypothesis statement
- argue which sets of parameters are the most effective
- judge if I proved what I set out to prove
- what are the uses of this method
- which applications could this fit?
- what should be done in future research

7 CONCLUSION

TODO content

- strongest possible statement based on preceding sections
- what should people remember about the paper
- refer to research question, describe conclusions I reached
- summarize new observations, insights through this work
- broader implications of results
- paragraph on recommendations / in future research / how to deal with limitations of this paper

DRAFT

REFERENCES

- brockwell12016 [1] P. J. Brockwell and R. A. Davis, *Introduction to Time Series and Forecasting*, en, ser. Springer Texts in Statistics. Cham: Springer International Publishing, 2016, ISBN: 978-3-319-29852-8 978-3-319-29854-2. DOI: [10.1007/978-3-319-29854-2](https://doi.org/10.1007/978-3-319-29854-2).
- lin2003 [2] J. Lin, E. Keogh, S. Lonardi, and B. Chiu, “A symbolic representation of time series, with implications for streaming algorithms,” en, in *Proceedings of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery - DMKD '03*, San Diego, California: ACM Press, 2003, pp. 2–11. DOI: [10.1145/882082.882086](https://doi.org/10.1145/882082.882086).
- aghazorgi2015 [3] S. Aghabozorgi, A. Seyed Shirkhorshidi, and T. Ying Wah, “Time-series clustering – A decade review,” en, *Information Systems*, vol. 53, pp. 16–38, Oct. 2015, ISSN: 03064379. DOI: [10.1016/j.is.2015.04.007](https://doi.org/10.1016/j.is.2015.04.007).
- kaur2016 [4] I. Kaur, R. Rajni, and A. Marwaha, “ECG Signal Analysis and Arrhythmia Detection using Wavelet Transform,” en, *Journal of The Institution of Engineers (India): Series B*, vol. 97, no. 4, pp. 499–507, Dec. 2016, ISSN: 2250-2106, 2250-2114. DOI: [10.1007/s40031-016-0247-3](https://doi.org/10.1007/s40031-016-0247-3).