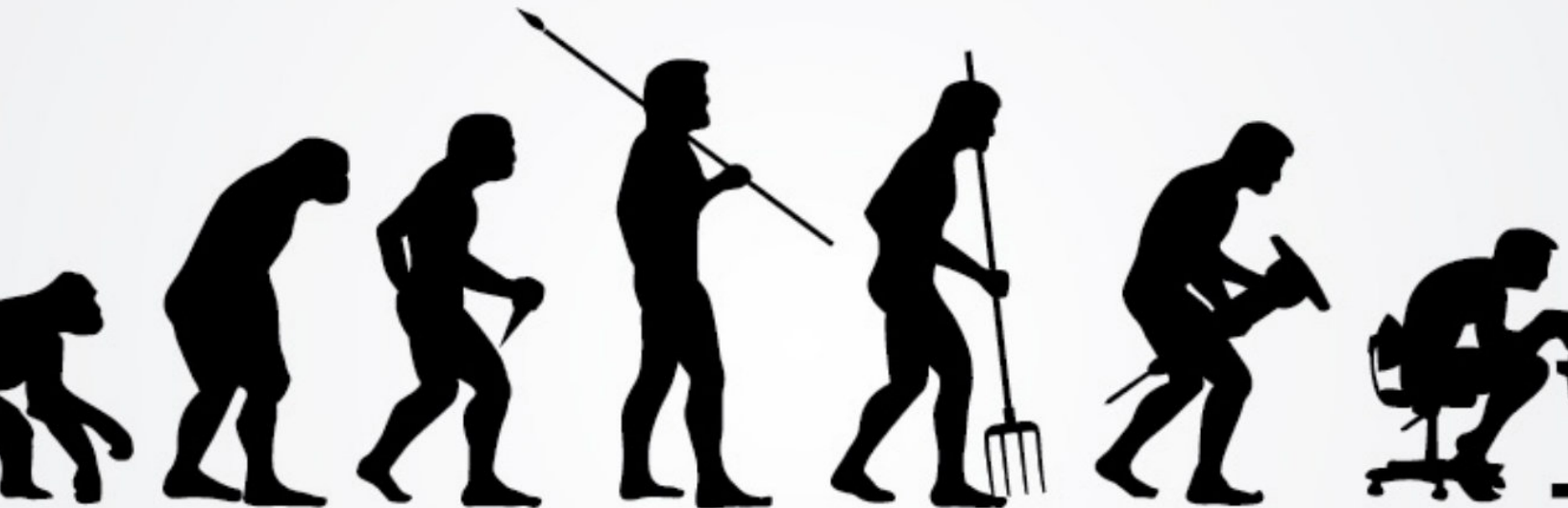


02122 Software Technology Project

Detecting Structural Breaks in Time Series via Genetic Algorithms

Group 3: Markus B. Jensen, s183816



Contents

1	Statusrapport 1.5 (Fra april 2021, med som reference)	1
1.1	Revideret projektplan	1
1.1.1	Skift i fokus	1
1.1.2	Revideret tidsplan	1
1.2	Statusrapport	1
2	Statusrapport 2 (Deadline: 7. maj 2021)	3
2.1	Tidsplan	3
2.2	Algoritmen	3
2.3	Fremadrettet	3
3	Introduction	4
4	Goals	5
4.1	Non-project goals	5
5	Terminology	5
6	Problem analysis and design	6
6.1	The algorithm	6
7	Implementation	6
8	Test	6
9	Project Management	6
10	Conclusion	6
	References	7

1 Statusrapport 1.5 (Fra april 2021, med som reference)

En ekstraordinær statusrapport som følge af gruppemedlem er sprunget fra.

1.1 Revideret projektplan

I den første projektplan blev der lagt vægt på, at Johan havde bedre overblik over matematikken og derfor stod for størstedelen af algoritmen. Han fik skrevet en del på algoritmen, men han nåede aldrig at skrive den helt færdig. I stedet for at færdiggøre hans algoritme besluttede jeg mig for at researche genetiske algoritmer selv og implementere algoritmen helt på ny. Dette har givet mig en rigtig god forståelse for genetiske algoritmer og bedre overblik over koden.

1.1.1 Skift i fokus

Grundet jeg nu er alene (og man *undtagelsesvis* måtte være to i gruppen) må jeg begrænse arbejdet. Derfor dropper at implementere andre fitness-funktioner for i stedet at fokusere på at optimere algoritmens hastighed med *range trees*. Mit fokus vil således ligge på følgende:

1. Implementere algoritmen der virker.
2. Implementere GUI, hvor parametre kan ændres og graf kan ses.
3. Optimere algoritmens hastighed.

Sideløbende med ovenstående vil jeg skrive på rapporten, hvor listen ovenfor er så kort. Rapporten vil have høj prioritet.

1.1.2 Revideret tidsplan

Hermed en revideret tidsplan. Det er indforstået, at der arbejdes på rapporten sideløbende med alle de punkter, hvor det ikke udspecificeres.

- **Uge 9 (nu):** Lav "Statusrapport 1.5" og arbejd på at få algoritme til at virke.
- **Uge 10:** Færdiggør implementation af algoritme som i pseudo-kode.
- **Uge 11:** Arbejd på GUI.
- **Uge 12:** Færdiggør GUI, begynd at optimere algoritme.
- **Uge 13:** Fortsæt arbejde med at optimere algoritme. Dokumentér resultater. **Deadlines:** Send udkast af rapport til Paul (som aftalt) og aflever "Statusrapport 2".
- **3-ugers:** Fokus på rapport. Få arbejdet på ting fra 13-ugers perioden, som ikke virker tilfredsstillende endnu.

1.2 Statusrapport

Algoritmen er næste færdigimplementeret. Som følge af ikke at skulle implementere flere fitness-metoder er koden omstruktureret en smule til at gøre det en smule mere overskuelig. Regner med at have den på plads i næste uge.

GUI skal nok gå relativt hurtigt: Jeg har undersøgt, hvordan jeg kan sætte firkanter på grafer, og så laver jeg selve GUI med "Scene Builder", så det forhåbentlig ikke tager så lang tid.

Optimering har jeg ikke kigget på endnu. Det vil kræve en del research i første omgang og så skal der også bruges tid på implementeringen. Men det tror jeg sagtens der er tid til med skiftet i fokus.

Rapportskrivningen går fremad. Jeg har taget mange noter i et andet program, som bare skal skrives rent her, hvorfor dette dokument synes tomt. Fremadrettet vil jeg begynde at skrive ting direkte i dette dokument, så jeg ikke skal skrive ting to gange.

Mit største problem ligger i perfektionisme. Det har været fedt at få afklaret forventninger med Paul, og det har givet ro på. At gøre fokus meget snævert har også været en kæmpe hjælp. Jeg er sikker på, at jeg skal komme i mål uden at bruge for mange timer på det.

2 Statusrapport 2 (Deadline: 7. maj 2021)

2.1 Tidsplan

Tidsplanen fra min Statusrapport 1.5 (sektion 1) følges nogenlunde. På dette tidspunkt skulle GUI være nogenlunde færdigimplementeret, og jeg skulle være begyndt på at optimere algoritmen ved brug af *range trees*. Det tog en del længere tid at vise firkanter på en graf i JavaFX end jeg havde regnet med, så jeg er endnu ikke begyndt at optimere på algoritme. Så jeg er en smule bag tidsplanen.

Tidsplanen er dog ikke et kæmpe problem, da jeg kun har én eksamen i maj. Jeg har kun 15 ECTS-point dette 13-ugers, hvor 10 af dem er projektkurser. Derfor vil jeg have tid i eksamensperioden til at kigge lidt på fagprojekt.

2.2 Algoritmen

I Statusrapport 1.5 konkluderede jeg, at algoritmen spillede. Hvad jeg ikke vidste var, at *mutation*-metoden ikke blev kaldt nogensinde. Den skaber nu nogle problemer med at generere for mange break points. For nogle *time series* er "den bedste løsning" den *solution string*, hvor næsten alle punkter er markeret som et break point. Det er ikke særlig godt, og det kan måske bare rettes med at ændre fitness-funktionen: Implementere en anden udgave af $p(k)$ fra [Doerr et al., 2017]. Så algoritmen spiller ikke fuldstændig lige nu.

2.3 Fremadrettet

Jeg har kun brugt halvdelen af den normerede tid på kurset. Jeg har nogle gange en følelse af, at jeg er bagud, men jeg har jo indtil videre kun brugt halvdelen af den allokerede tid til kurset. Jeg har stadig masser af tid i 3-ugers til at færdiggøre projektet. Det skal jeg huske på.

Jeg glemmer at skrive på rapporten. Det er jo nogle gange sjovere at sidde og kode, end det er at skrive på rapporten. Men rapporten *er* virkelig vigtig. Jeg skal huske at have rapporten som en prioritet.

Jeg er lidt i tvivl om hvordan programmet skal behandle flere dimensioner. Jeg har ikke leget rundt med det endnu. Hvis jeg har tid, kunne det være super fedt at bruge det til noget, men det kan være, at jeg kun læser én dimension ind fra datafiler (og gør brugeren opmærksom på det).

3 Introduction

Detecting structural break points is an essential tool in analyzing time series. Structural break points are points on a time series where the pattern of the time series measurements changes in the amplitude. A simple example is shown in figure 1a where the break point is easily detected at $t = 500$. Structural breaks are where the pattern of a time series changes.

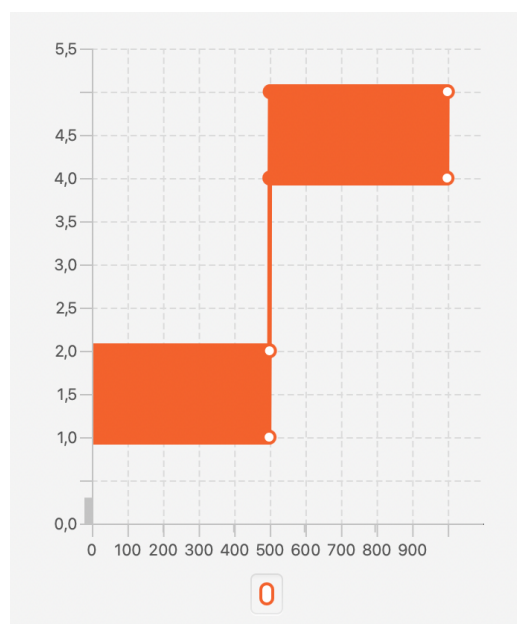
In the example above, the break point is easy to detect. As usual, the real world is far more messy. In figure 1b, which shows the number of patients hospitalized in Denmark due to Covid19, the break points are less obvious. This is where detecting break points is important: Being able to recognize when the pattern changes, so that a pandemic does not run amok or looking for fluctuations in the stock market.

In this project, the break points are found by using a so-called genetic algorithm. This algorithm mimics natural evolution: A number of individuals (solutions that indicate the positions of break points) mutate and mate during generations where survival of the fittest is the rule. In the end, the best solution will (hopefully) have found all the break points.

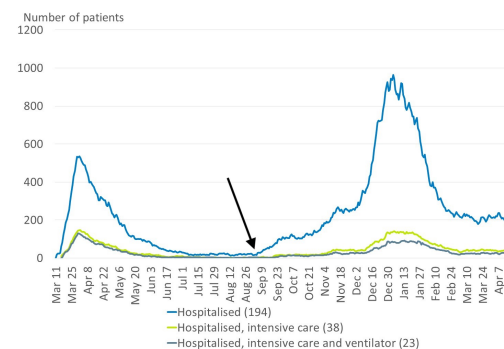
This project will focus on implementing the algorithm in the programming language Java and make a simple application. The application will give the user the ability to tweak some of the algorithm's parameters and see the break points directly on the time series graph. Further improvements to the speed of the algorithm will also be implemented "under the hood".

Please find another example than the stock market

Make the left figure more pretty!



(a) Graph with break point at $t = 500$



(b) Number of people hospitalized due to Covid19 in Denmark. The arrow indicates a break point. Source: the Danish regions

Figure 1: Two figures illustrating break points in time series

4 Goals

A prioritized list of the functional goals for the project is listed below. While adding a lot of features seems very appealing, I will in this project focus on usability, stability, and speed. This leads to a somewhat conservative list of potential features, but the quality of the final product will reflect this decision.

1. Implement the algorithm using the rectangle-method from [Doerr et al., 2017] in Java.
2. Visualize two-dimensional time series graphs together with the rectangles produced by the algorithm in a simple graphical user interface (GUI). This GUI will allow user to load a time series data file and see the output of the algorithm on the time series. The user will also be able to tweak certain parameters of the algorithm.
3. Make the algorithm more flexible by allowing the user to alter the values of algorithm-parameters in the GUI.

A previous goal was also to implement further fitness functions (see explanation for *fitness functions* in section 5). Since this project ended up being a one-man-project, this goal was removed. In stead, a priority is to design the project so that it is easy to implement other fitness methods. This will be discussed further in section 6.

4.1 Non-project goals

As for non-project goals, there are a few:

- Learn the Maven project structure for Java. While previous courses have dealt with Maven a little bit, it has never been fleshed out. It seems to be a structure that is widely used and thus a good system to learn.
- Working on big project. This is the first big project I am working on. A big focus here is on the project management, report writing workflow and keeping track of sources in a bibliography. Especially the report writing workflow can become crucial, as I have a tendency to postpone it to the very last minute. Here, I will make it a part of my weekly work.

My ambition level is quite high; I am a perfectionist to the core. While I will attempt to keep the perfectionism to a minimum, I like working on bigger projects and making it work well. This will probably result in me working a lot on this project, purely because it will be fun, and I like improving my less-than-optimal solutions. I am aware that this course is only 5 ECTS points and will thus keep track of my hours spent on the project as to not overdo it.

5 Terminology

The terminology used in this report is specified below. The terminology differs a bit from [Doerr et al., 2017]. It takes inspiration from other sources ([Thede, 2004, Point,]) and follows a more biological narrative.

Individual An individual consists of a solution string. It is one possible solution to the problem.

Genome An individual consists of a genome. A genome is an array of genes. The genome *is* the solution string.

Allele A gene's value is called an allele. In this project, the allele is thus responsible for holding the information of whether or not a certain gene is a break point. An allele is the value at a certain gene in the genome/solution string.

Population A group of individuals.

Fitness function A function that measures how well the solution from an individual fits with the data. This can be any function relevant for a particular problem.

Parent and offspring individual Since genetic algorithms mimic evolution, procedures must include parent and offspring individuals. A parent individual is a parent to the offspring individual, meaning that any procedure must take one (or two) parent individual(s) and the procedure creates an offspring.

Crossover (Procedure) Crossover-procedures takes two parent individuals and creates an offspring from the genomes of the two parents. This project will incorporate *one-point crossover* and *uniform crossover*.

One-point crossover extracts the genes in the interval 0 to a random gene $i - 1$ from the first parent. Then, it appends the genome with the genes i to the last gene $n - 1$ from the second parent. Thus, the offspring's genome consists of the first i genes from parent one and $n - i$ last genes from parent two.

Uniform crossover is more simple. For each gene in the offspring's genome, there is a fifty-fifty chance whether it will be the gene from parent one or the gene from parent two.

Mutation (Procedure) Just like the corona virus, the genome of an individual can mutate. It means, that for any gene, there is a random chance that the allele will change. In this case, the change will be either creating break points or removing them.

6 Problem analysis and design

The main challenges are already hinted at in section 4. All these goals pose different problems. These problems (and especially their solutions) is best analyzed when also discussing the chosen design.

6.1 The algorithm

7 Implementation

8 Test

9 Project Management

10 Conclusion

References

- [Doerr et al., 2017] Doerr, B., Fischer, P., Hilbert, A., and Witt, C. (2017). Detecting structural breaks in time series via genetic algorithms. *Soft Computing*, 21(16):4707–4720.
- [Point,] Point, T. Genetic algorithms.
- [Thede, 2004] Thede, S. (2004). An introduction to genetic algorithms. *Journal of Computing Sciences in Colleges*, 20.