

Primzahlregression mit Hilfe genetischer Algorithmen

Björn Boss Henrichsen und Simon Budinsky
Betreut von Herrn Christoph Fischer

6. Februar 2017

Inhaltsverzeichnis

1	Einleitung	3
1.1	Problemstellung	3
1.2	Ziel des Projektes	3
2	Relevanz von Primzahlen	4
2.1	RSA-Verschlüsselung	4
2.1.1	Funktionsweise RSA	4
2.1.2	Der RSA-Wettbewerb	4
2.1.3	Folgen einer geknackten Verschlüsselung	5
2.2	Primzahltests	5
2.3	Finden der n ten Primzahl	5
2.4	Die größte Primzahl	6
3	Das Sieb des Eratosthenes	7
3.1	Allgemein	7
3.2	Funktionsweise	7
3.3	Offizielle Variante	8
3.3.1	Programmablaufplan	8
3.3.2	Komplexität	8
3.4	Optimierungsmaßnahmen	9
3.4.1	Weglassen der eins	9
3.4.2	Quadrieren statt Radizieren	9
3.4.3	Zahlen im Vorraus streichen	9
3.5	Messungen	12
3.5.1	Methodik	12
3.5.2	Vergleich zwischen offizieller und optimierter Variante	15
3.5.3	Messung mit Ausgabe aller Primzahlen	18
4	Genetischer Algorithmus für Regressionen	19
4.1	Definition genetischer Algorithmus	19
4.2	Funktionsweise eines genetischen Algorithmus	19
4.3	Unsere Implementierung	22
4.3.1	Die Funktionssyntax	22
4.3.1.1	Beispiel	23
4.3.2	Definitionen	24
4.3.2.1	Headerdatei primereg.h	24
4.3.2.2	Erläuterung	25
4.3.3	Wahl des Funktionsparameter Datentyps	26
4.3.4	Funktionsweise	26

4.3.4.1	Programmablaufplan	26
4.3.4.2	Erläuterung	28
4.4	Vor-/Nachteile gegenüber herkömmlichen Regressionsmethoden	31
4.4.1	Vorteile	31
4.4.2	Nachteile	31
5	Ergebnisse	32
6	Fazit	33
A	Quellen	34

1 Einleitung

1.1 Problemstellung

In der Mathematik gibt es einige charakteristische Zahlen, die auf Grund der Komplexität ihrer Bestimmung eine besondere Anwendung finden. Eine solche Zahl stellen unter anderem Primzahlen dar. Unter einer Primzahl versteht man eine natürliche Zahl größer eins, welche nur durch sich selbst und eins teilbar ist. Somit ist die Mächtigkeit der Teilermenge einer Primzahl zwei.

Über 2000 Jahre lang konnte man keinen praktischen Nutzen aus dem Wissen über die Primzahlen ziehen. Dies änderte sich erst mit dem Aufkommen elektronischer Rechenmaschinen, bei denen die Primzahlen beispielsweise in der Kryptographie eine zentrale Rolle spielen.

Primzahlen sind hierfür deshalb so gut geeignet, weil es bisher keine Formel gibt, welche Primzahlen effizient berechenbar generiert.

1.2 Ziel des Projektes

Das Ziel unseres Projektes besteht darin, eine Folge zu finden, die Primzahlen möglichst gut beschreibt. Dazu setzen wir einen genetischen Algorithmus ein, welcher zu vorgegebenen Datenpunkten eine Regressionsfunktion bestimmen soll. Diese Datenpunkte sollen Primzahlen sein, welche wir mit Hilfe einer optimierten Variante des Siebes des Eratosthenes erzeugen.

2 Relevanz von Primzahlen

2.1 RSA-Verschlüsselung

In der Kryptographie gibt es Algorithmen, welche auf der Primfaktorzerlegung basieren. Ein modernes Beispiel hierfür ist der asymmetrische Verschlüsselungsalgorithmus RSA. Da bisher offiziell keine effiziente Methode gefunden wurde, Primfaktorzerlegung zu betreiben, ist eine Verschlüsselung mit RSA sicher.

2.1.1 Funktionsweise RSA

$\varphi(n)$: Eulersche Phi-Funktion

$$\varphi(n) := |\{a \in \mathbb{N} | 1 \leq a \leq n \wedge \text{ggT}(a, n) = 1\}|$$

1. Nimm zwei (sehr, sehr große) Primzahlen p und q
2. Bilde das RSA-Modul $N = p * q$
3. Bestimme $\varphi(N) = \varphi(p * q) = \varphi(p) * \varphi(q) = (p - 1) * (q - 1)$
4. Wähle eine Zahl e mit $1 < e < \varphi(N)$ mit $\text{ggT}(e, \varphi(N)) = 1$
5. Bestimme d mit $e * d \equiv 1 \text{ mod } \varphi(N)$

Jetzt bilden (e, N) den öffentlichen und (d, N) den privaten Schlüssel. Mit Hilfe des privaten Schlüssels lässt sich ein Klartext T zum Geheimtext G verschlüsseln:

$$G = T^e \text{ mod } N$$

Der Geheimtext lässt sich mit dem öffentlichen Schlüssel wieder entschlüsseln:

$$T = G^d \text{ mod } N$$

2.1.2 Der RSA-Wettbewerb

Aktuell gibt es einen RSA-Wettbewerb. Ziel ist es, die zwei Primfaktoren p und q einer RSA Zahl N zu finden. Denn findet man diese, so ist die Verschlüsselung geknackt.

Dies wird mit einem entsprechenden Preisgeld belohnt. Die bisher größte zerlegte RSA Zahl hat 232 Dezimalstellen (768 Bits, RSA-768) und wurde

am 12. Dezember 2009 von Thorsten Kleinjung, Paul Zimmermann u.a. für ein Preisgeld von \$ 50000 USD faktorisiert. Hierfür wurde ein Rechencluster verwendet, welches etwa tausend Mal schneller als ein 2,2GHz AMD Opteron single-core Prozessor ist. Trotzdem dauerte die Faktorisierung zwei Jahre. Dabei muss bedacht werden, dass heutzutage eher längere RSA-Zahlen verwendet werden. Häufig verbreitet sind Zahlen im Bereich von 309 bis 617 Dezimalstellen (1024-2048 Bits (das Preisgeld für eine RSA-2048-Zahl beträgt aktuell \$ 200000 USD (02.01.17))).

2.1.3 Folgen einer geknackten Verschlüsselung

Zusätzlich sollte die Verschlüsselung in ein paar Minuten, am besten wenigen Sekunden, erfolgen. So könnte das Passwort für den Bankaccount per Man-In-The-Middle-Angriff einfach ausgelesen werden. Dies ist natürlich nur eine von den vielen fatalen Folgen, wenn RSA geknackt werden würde. Der Mensch könnte seine Privatsphäre verlieren und zum „gläsernen Menschen“ werden.

2.2 Primzahltests

Eine primzahlbeschreibende Folge wird vor allem in der Kryptographie benötigt. Damit RSA funktioniert, werden zunächst zwei Primzahlen erzeugt. Diese müssen zunächst gefunden werden.

Hierfür werden Test wie der Rabin-Miller-Test verwendet, der allerdings nur über eine Zahl zwei Aussagen treffen kann: zusammengesetzt oder wahrscheinlich prim.

Er liefert also nicht mit hundertprozentiger Sicherheit Zahlen, die anschließend für RSA verwendet werden. Ist eine Zahl nicht prim, funktioniert der RSA-Algorithmus nicht und kann geknackt werden.

Dies kann verhindert werden, indem stets mit hundertprozentiger Sicherheit über eine Zahl ausgesagt werden kann, ob sie prim ist. Hierfür könnte eine Folge für Primzahlen verwendet werden und somit die Primzahltest ablösen.

2.3 Finden der n ten Primzahl

Ein weiteres Problem ist folgendes: Finde die n te Primzahl.

Beim Sieb des Eratosthenes zum Beispiel müssen zunächst alle Primzahlen vor der n ten gefunden werden. Gäbe es eine Folge für Primzahlen, könnte dieses Problem wesentlich effizienter gelöst werden. Folglich wären Algorithmen zum Finden von Primzahlen wie das Sieb des Eratosthenes obsolet.

2.4 Die größte Primzahl

Vor etwa 2300 Jahren hat Euklid bewiesen: Es gibt unendlich viele Primzahlen. Dies wurde 1737 durch Leonhard Eulers Beweis bestärkt, der folgendes zeigte:

$$\sum_{p \in \mathbb{P}} \frac{1}{p} = \infty$$

Heutzutage versucht man stets größere Primzahlen zu finden. Daraus wurde eine Art Wettbewerb. Die Aktuell größte Primzahl (Stand Januar 2017) hat 22.338.618 Dezimalstellen und wurde 2016 vom *Great Internet Mersenne Prime Search*-Projekt gefunden.

Bei der Suche kann jeder teilnehmen. Hierfür muss lediglich eine kostenlose Software heruntergeladen werden, welche nach Primzahlen sucht. Findet der Computer eine neue „größte“ Primzahl, erhält man \$3000 USD.

Mit Hilfe einer Folge für Primzahlen könnte diese direkt gefunden werden.

3 Das Sieb des Eratosthenes

3.1 Allgemein

Das Sieb des Eratosthenes ist ein Algorithmus, mit welchen man Primzahlen findet kann.

Er wurde nach dem griechischen Mathematiker Eratosthenes von Kyrene benannt, welcher im 3. Jahrhundert v. Chr. lebte.

Laut wikipedia.org habe dieser selbst nicht das Verfahren entdeckt. Eratosthenes habe lediglich die Bezeichnung „Sieb“ eingeführt.

Für Primzahlen kleiner 10.000.000 zählt der Algorithmus zu den effizientesten Methoden. Für größere Zahlen empfehlen sich effizientere Algorithmen. Beispiele sind der Lehmann-Test oder der Rabin-Miller-Test, welcher auch in modernen Verschlüsselungsalgorithmen (RSA) verwendet wird.

3.2 Funktionsweise

Zu Beginn werden alle natürlichen Zahlen von 2 bis einer Grenze n aufgeschrieben. Anfangs sind all diese nicht gestrichen und somit potentielle Primzahlen. Nun durchläuft man das Intervall. Trifft man auf eine ungestrichene Zahl, so ist diese eine Primzahl. Jetzt werden alle Vielfachen dieser gestrichen, die größer oder gleich deren Quadrat sind. Ist man am Ende des Intervalls angekommen, liegen sowohl gestrichene als ungestrichene Zahlen vor. Dabei ist jede ungestrichene Zahl prim.

Es genügt, die Vielfachen von Primzahlen zu streichen, die kleiner oder gleich der Wurzel der Grenze n sind. Denn ist ein Primfaktor einer zusammengesetzten Zahl immer kleiner gleich der Wurzel der zusammengesetzten Zahl.

Außerdem muss beim Streichen der Vielfachen erst mit dem Quadrat der Primzahl begonnen werden, da alle kleineren Vielfachen bereits getrichen wurden.

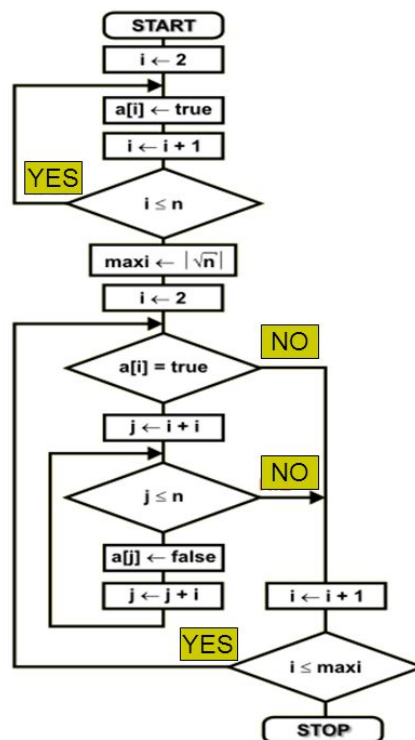
3.3 Offizielle Variante

3.3.1 Programmablaufplan

Sieve of Eratosthenes flowchart

$a[i] = \text{true}$ - prime number

$a[i] = \text{false}$ - composite number



3.3.2 Komplexität

Die asymptotische Laufzeitkomplexität O in Abhängigkeit von einer oberen Grenze n ergibt sich wie folgt:

Sei $p_j \in \mathbb{P}$

So werden für jede Primzahl

$$p_j \leq \sqrt{n}$$

höchstens $\frac{n}{p_j}$ Zahlen gestrichen. Somit erhält man für die Anzahl an „Streichungen“:

$$\frac{n}{2} + \frac{n}{3} + \frac{n}{5} + \dots = \sum_{p_j \leq \sqrt{n}} \frac{n}{p_j} = n * \sum_{p_j \leq \sqrt{n}} \frac{1}{p_j}$$

Nun muss die asymptotische Ordnung O vom zweiten Faktor gefunden werden. Diese ist $O(\log \log n)$ (Der Beweis hierfür ist nicht trivial und würde den Rahmen der Facharbeit sprengen). Daher ergibt sich eine Gesamtkomplexität von $O(n \log \log n)$.

3.4 Optimierungsmaßnahmen

3.4.1 Weglassen der eins

Da die eins keine Primzahl ist, kann man sie vorab ignorieren. Somit wird die Länge des primes-Puffers um ein Byte verringert (Speichereffizienz). Zusätzlich muss zu Beginn ein Byte weniger initialisiert werden. In der äußeren `for()`-Schleife wird ebenfalls ein Durchlauf weniger benötigt (Zeiteffizienz).

3.4.2 Quadrieren statt Radizieren

Des Weiteren ist das Quadrieren an einem Computer effizienter als das Radizieren. Dies erklärt sich darin, dass für das Radizieren eine Bisektion (Software) benötigt werden kann. Dahingegen ist das Quadrieren lediglich eine Multiplikation, welche direkt in der ALU (Hardware) ausgeführt wird. Somit kann man schreiben:

i : Aktuell betrachtete Zahl

$$i \leq \sqrt{n} \Rightarrow i * i \leq n$$

3.4.3 Zahlen im Vorraus streichen

Der Algorithmus kann weitergehend optimiert werden, indem man Zahlen vorab „streicht“. So haben wir alle Vielfachen von zwei bereits „gestrichen“. Daraus ergibt sich folgende Speichereinsparung:

a : Belegter RAM-Speicher mit allen Zahlen $[0;n]$

b : Belegter RAM-Speicher ohne eins und geraden Zahlen

n : Obere Grenze des Siebes

k : Proportionalitätsfaktor

$$a = k * b \tag{1}$$

$$\Leftrightarrow n = k * \left(\frac{n}{2} - 1\right) \tag{2}$$

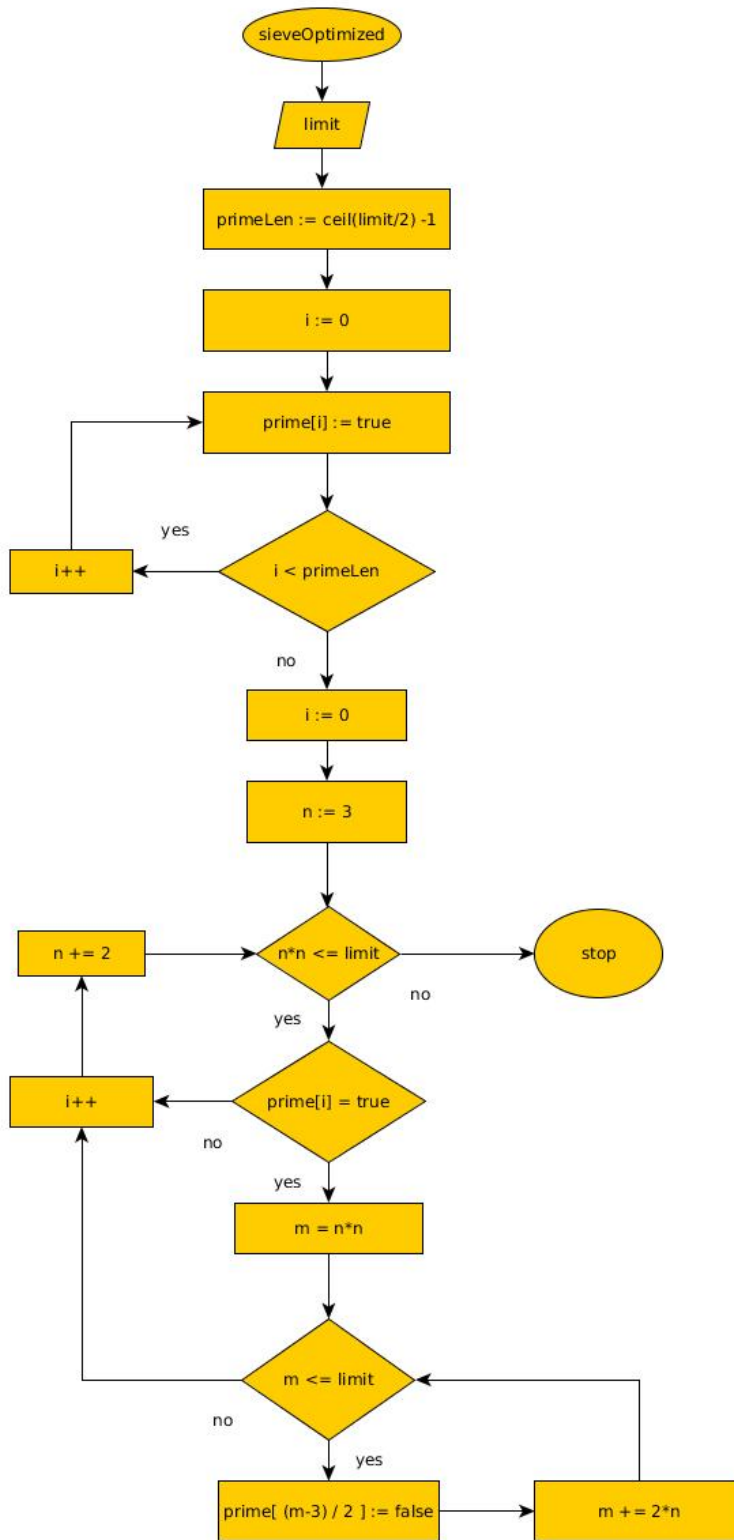
$$\Leftrightarrow k = \frac{n}{\left(\frac{n}{2} - 1\right)} \tag{3}$$

$$= \frac{2n}{n - 2} \tag{4}$$

Folglich verwendet die optimierte Variante $\frac{n-2}{2n}$ Mal so viel Speicher, im Vergleich zur offiziellen. Somit können mehr Primzahlen auf dem gleichen Rechner gefunden werden.

Auch zeigen unsere Messungen, dass, wenn man alle Vielfachen von zwei weglässt, der Algorithmus mehr als doppelt so schnell terminiert (Zeiteffizienz).

Nach den oben genannten Optimierungen könnte ein Programmablaufplan so aussehen:



Nun könnte man immer mehr Zahlen im Vorraus „streichen“. Hierfür gibt es allerdings andere Algorithmen. Das *Sieb des Atkin* zum Beispiel schließt vorher Zahlen bestimmter Restklassen aus und hat eine lineare Laufzeitkomplexität $O(n)$.

3.5 Messungen

3.5.1 Methodik

Um die Laufzeit des Siebes für verschiedene Grenzen *limit* zu messen, wurde das Programm für jedes *limit* jeweils neu kompiliert und ausgeführt. Nachdem die Primzahlen gefunden sind, gibt das Programm die benötigte Zeit zwischen Speicherallokation für den Primzahl-Puffer und dem Finden der Primzahlen in Sekunden aus. Diese Ausgabe wurde in eine .csv-Datei per Linux-Bash umgeleitet, um somit später mit LibreOfficeCalc Diagramme zu erzeugen. Dabei besteht diese Datei aus drei Spalten. Die erste Spalte repräsentiert die aktuelle Grenze *limit*, die zweite die gemessene Zeit in Sekunden und in der dritten Spalte steht die Anzahl gefundener Primzahlen.

Dafür verwendeten wir folgendes Bash-Script:

```
1  #!/bin/bash
2  #i: limit
3  for i in {1000000000..6000000000..1000000000}
4  do
5      echo -n "$i," >> runtimes.csv
6      gcc -std=c11 -DLIMIT=$i sieveRuntime.c -o sieveRuntime -lm \
7          && ./sieveRuntime >> runtimes.csv
8  done
```

Wir implementierten folgende Algorithmen in C zur Zeitmessung:

Sieb ohne die eins:

```
1  #include    <inttypes.h>
2  #include    <stdbool.h>
3  #include    <stdlib.h>
4  #include    <string.h>
5  #include    <sys/time.h>
6  #include    <stdio.h>
7
8  int main( void ) {
9      struct timeval start, end;
10     gettimeofday(&start, NULL);
11
12     bool* prime = (bool*)malloc(sizeof(bool) * (LIMIT-1));
13     if( prime == NULL ) return 1;
14     memset(prime, true, LIMIT-1);
15
16     for(uint64_t i = 2; i*i <= LIMIT; i++) {
17         if(prime[i-2]) {
18             for(uint64_t j = i*i -2; j <= LIMIT; j += i)
19                 prime[j] = false;
20         }
21     }
22
23
24     gettimeofday(&end, NULL);
25
26     // Count primes
27     uint32_t nPrimes = 0;
28     for( uint32_t i = 0; i < LIMIT-1; i++)
29         if( prime[i] ) nPrimes++;
30
31     printf("%.10f,%" PRIu32 "\n",
32           end.tv_sec + end.tv_usec/1e6 - start.tv_sec - start.tv_usec/1e6,
33           nPrimes);
34
35     return 0;
36 }
```

Sieb ohne die eins und geraden Zahlen:

```
1  #include    <inttypes.h>
2  #include    <stdbool.h>
3  #include    <stdlib.h>
4  #include    <string.h>
5  #include    <sys/time.h>
6  #include    <stdio.h>
7  #include    <math.h>
8
9  int main( void ) {
10     struct timeval  start, end;
11     gettimeofday(&start, NULL);
12
13     uint32_t primeLen = ceil(LIMIT / 2) - 1;
14     bool* prime      = (bool*)malloc(sizeof(bool) * primeLen);
15     if( prime == NULL ) return 1;
16     memset(prime, true, primeLen);
17
18     for(uint64_t i = 0, n = 3; n*n <= LIMIT; i++, n += 2) {
19         if(prime[i]) {
20             for(uint64_t m = n*n; m <= LIMIT; m += 2*n)
21                 prime[ (m-3) / 2 ] = false;
22         }
23     }
24
25     gettimeofday(&end, NULL);
26
27     // Count primes
28     uint32_t nPrimes = 1;
29     for( uint32_t i = 0; i < primeLen; i++)
30         if( prime[i] ) nPrimes++;
31
32     printf("%.10f,%" PRIu32 "\n",
33         end.tv_sec + end.tv_usec/1e6 - start.tv_sec - start.tv_usec/1e6,
34         nPrimes);
35     return 0;
36 }
```

[Anmerkung:]

Es wird bewusst nicht die Zeit zum Ausgeben (printf()) mit einbezogen, da diese allgemein einen großen Anteil an Laufzeit beansprucht. So ergibt sich bei unserem optimierten Algorithmus für $n = 1.000.000$ eine Zeit von etwa 1389.054 Sekunden mit Ausgabe aller Primzahlen und etwa 8.0484920487 Sekunden ohne.

Die Messungen fanden in folgender Umgebung statt:

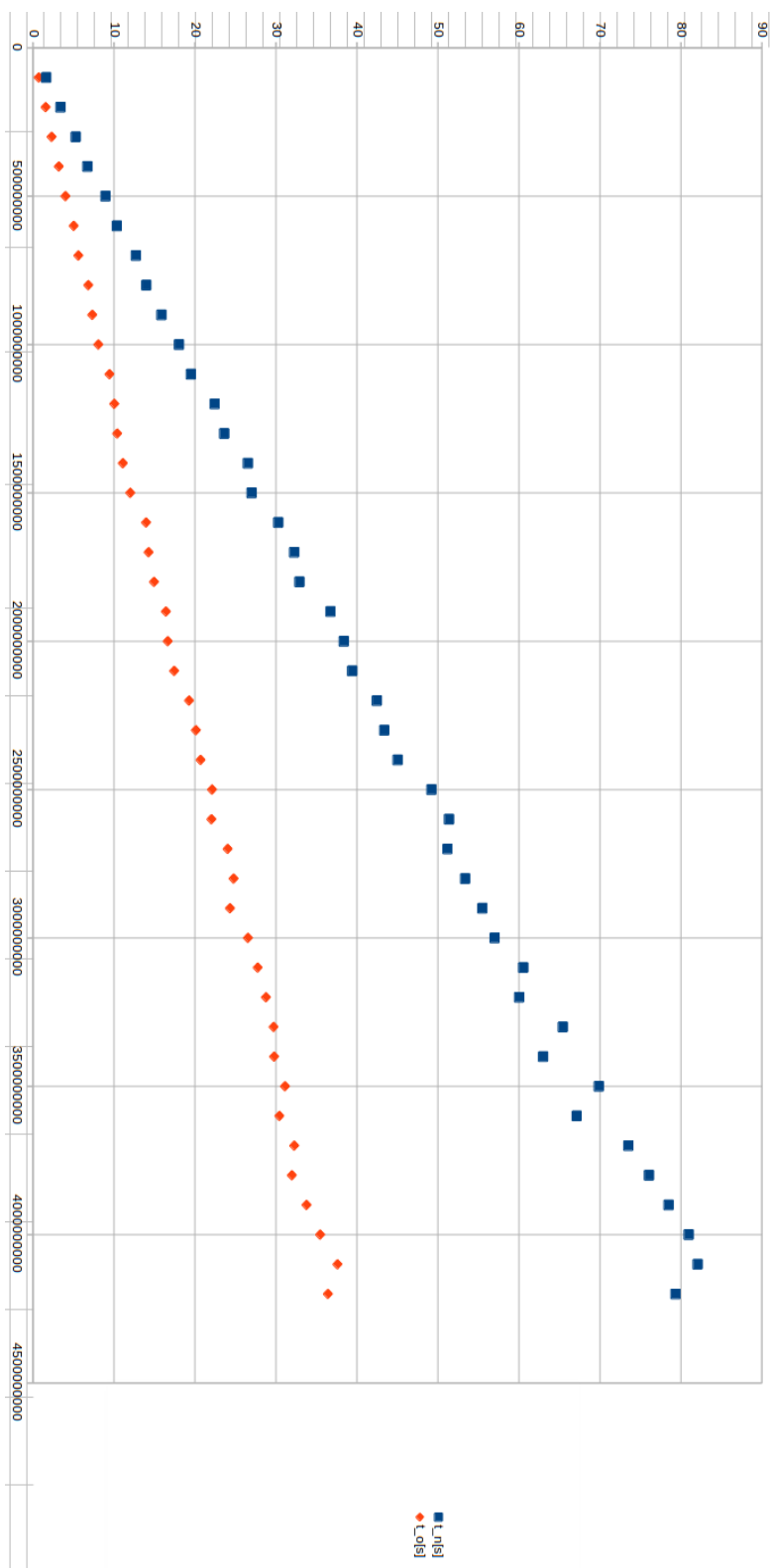
Operating System	<i>Linux Mint 18.1 Cinnamon 64-Bit</i>
Cinnamon Version	<i>3.2.7</i>
Linux Kernel	<i>4.4.0-53-generic</i>
Processor	<i>Intel Core i5-4440 CPU @ 3.10GHz x 4</i>
Memory RAM	<i>7.7GiB</i>

3.5.2 Vergleich zwischen offizieller und optimierter Variante

t_g : Benötigte Zeit in Sekunden, offizielle Variante (blau)

t_u : Benötigte Zeit in Sekunden, optimierte Variante (rot)

<i>limit</i>	t_g [s]	t_u [s]	$\frac{t_g}{t_u}$
1000000000	1,6142570771	0,6927529358	2,3302060427
2000000000	3,3796610065	1,5407490117	2,193518205
3000000000	5,2564000593	2,2783710453	2,3070869296
4000000000	6,7034089816	3,1714290402	2,1136872043
5000000000	8,9381100896	3,9901179841	2,2400616035
6000000000	10,3455499889	5,011878017	2,0642062624
7000000000	12,6841129465	5,5832778916	2,271803982
8000000000	13,9744109066	6,8093409411	2,0522413296
9000000000	15,8358931078	7,3026679844	2,1685078853
10000000000	18,0069118877	8,0484920487	2,2373025629
11000000000	19,4879198822	9,4263940965	2,0673780114
12000000000	22,4203640433	10,0174689825	2,2381266248
13000000000	23,5954410898	10,3775840185	2,2736930915
14000000000	26,5114490883	11,0751651153	2,393774613
15000000000	26,9755919041	11,9925879208	2,2493553587
16000000000	30,2582220339	13,9449501191	2,1698336513
17000000000	32,2235829942	14,2521269333	2,2609666013
18000000000	32,8872118869	14,924550968	2,2035645801
19000000000	36,7206629502	16,3945279494	2,2398121534
20000000000	38,3611690887	16,617560026	2,3084718231
21000000000	39,3934800971	17,4123429244	2,2623882534
22000000000	42,4355878925	19,252283056	2,2041847073
23000000000	43,353328041	20,0979319791	2,1571039292
24000000000	45,0081739355	20,6734950301	2,1770955453
25000000000	49,1815101048	22,0877529756	2,2266416217
26000000000	51,3432890466	22,012063034	2,3325069062
27000000000	51,1586249138	24,0191830281	2,129906952
28000000000	53,3409920051	24,7522308979	2,1549973505
29000000000	55,47018401	24,3075119497	2,282018173
30000000000	56,9847639827	26,5199970692	2,148746994
31000000000	60,5142700262	27,7222269376	2,1828791086
32000000000	60,012045889	28,7487220151	2,0874683006
33000000000	65,3936210892	29,68840601	2,2026652784
34000000000	62,9704729115	29,7541698854	2,116357914
35000000000	69,8409569373	31,0977290192	2,2458539302
36000000000	67,1172831177	30,4034039054	2,2075581842
37000000000	73,4939509822	32,2375419811	2,2797628624
38000000000	76,0375679339	31,9458291063	2,3802033023
39000000000	78,455664996	33,7550660629	2,3242634113
40000000000	80,9324689933	35,4401829989	2,2836357531
41000000000	82,0384289596	37,5744409737	2,1833572725
42000000000	79,314230075	36,3976849106	2,1791308935



3.5.3 Messung mit Ausgabe aller Primzahlen

<i>limit</i>	t_g [s]	t_u [s]	Δt	<i>anzahl</i>
100.000	0,003	0,002	0,001	9592
1.000.000	1,216	1,205	0,011	78498
10.000.000	14,463	14,387	0,076	664579
100.000.000	143,18	142,439	0,741	5761455
1.000.000.000	1397,5	1389,054	8,446	50847534

4 Genetischer Algorithmus für Regressionen

4.1 Definition genetischer Algorithmus

Ein genetischer Algorithmus, oder auch Evolutionärer Algorithmus, ist eine Klasse von Algorithmen, welche an die Funktionsweise der Natur angelehnt sind.

Hierbei werden aus einer Datenmenge, die dafür bestimmt ist, ein bestimmtes Problem zu lösen, die effizientesten ausgewählt und dafür verwendet die Anzahl an Datenmengen erneut zu erweitern mit den effizientesten als Grundgerüst.

Somit „lernt“ der Algorithmus aus seinen Fehlern, da nur die effizientesten Datenmengen „überleben“.

4.2 Funktionsweise eines genetischen Algorithmus

Zu Beginn wird erst einmal eine Darstellungsweise für das zu lösende Problem aufgestellt.

Anschließend wird eine Population generiert mit einer zufällig aufgebauten, oder bereits leicht geordneten, Struktur dieser Darstellungsweise. Anschließend muss jedes Element dieser Population von einer bewertenden Funktion, meist der Fitness-Funktion, die Fitness des Elements berechnen/bewerten.

Diese Funktion ist notwendig da es sonst, wie in der Natur auch, keinen Grund gäbe sich zu verändern, wenn es nicht für das Überleben beziehungsweise das Wohlergehen notwendig ist. Die Fitness-Funktion wird basierend auf das Problem festgelegt und ist in ihrer Struktur meistens einfacher. So kann diese beispielsweise bei einem Algorithmus, welcher ein simuliertes Auto zu fahren lernen soll, einfach nur eine Division von der zurückgelegten Strecke durch die benötigte Zeit sein.

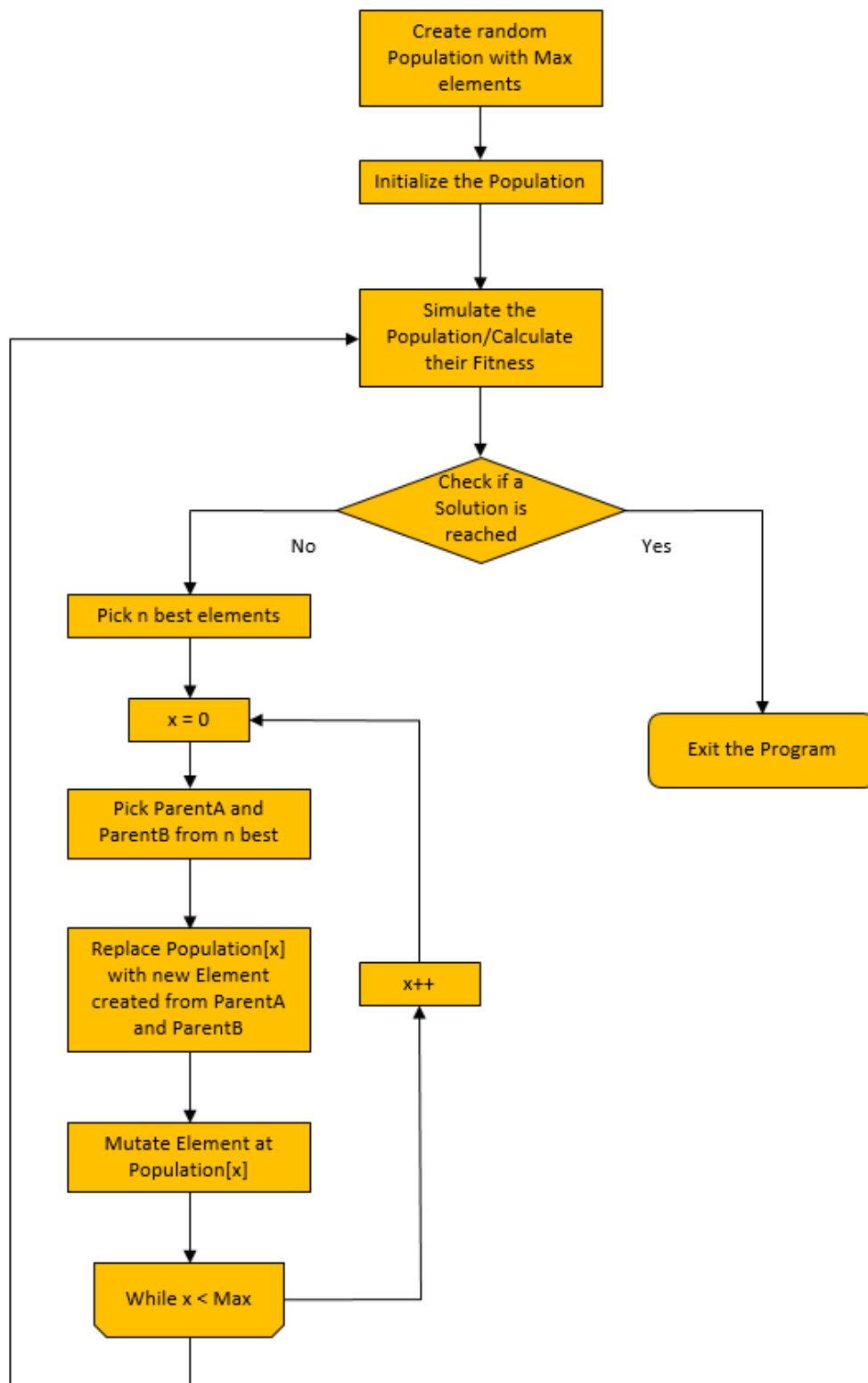
Nachdem jedes Element der Population bewertet wurde und jedem ein Fitness Wert zugeordnet worden ist müssen die Elemente die am besten abgeschnitten haben aussortiert werden. In unserem Beispiel von oben wären dies die Elemente, welche die höchste Fitness erreicht haben. Dies ist aber nicht zwingend. So kann es auch sein, dass eine kleinere Fitness besser, als eine größere ist. Das hängt von der Fitness-Funktion und deren Funktionsweise ab.

Diese besseren ausgewählten Elemente repräsentieren nun die Eltern. Diese werden nun aus der Population entfernt und außerhalb gespeichert. Anschließend werden alle Elemente der Population gelöscht.

Nun wird für jeden Platz in der Population durch Zufall n , vorher festgelegte, Eltern aus den beiseitegelegten Elementen ausgewählt. Von diesen n

Eltern werden nun, meist ebenfalls durch Zufall, Teile aus der Datenmenge entnommen und in das neue leere Element der Population getan.

Anschließend kann es noch zu einer geringen Wahrscheinlichkeit, die meist als 1 bis 2 Prozent angelegt ist, zu einer Mutation kommen. Diese Mutation ist notwendig, damit sich die Population nicht auf den am Anfang durch Zufall erstellten Datenmengen festfährt, sondern immer was Neues hinzukommt und was Altes entfernt wird. Anschließend können die beiseitegelegten Elemente ebenfalls entfernt werden, da man für diese keinen Nutzen mehr hat. Nun kann der gesamte Prozess von vorne beginnen mit der neu erstellten Population, die auf der alten Population aufbaut und die statistisch betrachtet nun bessere Fitness Werte erzielen wird.



4.3 Unsere Implementierung

4.3.1 Die Funktionssyntax

Damit der Algorithmus mit mathematischen Funktionen umgehen kann, bedarf es einer speziellen Funktionssyntax, welche von diesem interpretiert werden muss.

Wir haben uns für eine Assembler-Artige Darstellung von mathematischen Funktionen entschieden, wobei ein „Term“ wie folgt aufgebaut ist:

Funktion *Parameter_a* optional *Parameter_b* *Solution*

Dabei repräsentieren konstante Bytes mathematische Funktionen, um den Speicherverbrauch möglichst gering zu halten (im Vergleich zur Verwendung von Zeichenketten):

1	VAR	= 0x00,	<i>//Usage of Variable x</i>	
2	SLA	= 0x01,	<i>//Temp. Solution A</i>	
3	SLB	= 0x02,	<i>//Temp. Solution B</i>	
4	NUM	= 0x03,	<i>//following is a double</i>	
5	ADD	= 0x10,	<i>//double add(double a, double b)</i>	<i>-> a + b</i>
6	SUB	= 0x11,	<i>//double subtract(double a, double b)</i>	<i>-> a - b</i>
7	MUL	= 0x12,	<i>//double multiply(double a, double b)</i>	<i>-> a * b</i>
8	DIV	= 0x13,	<i>//double divide(double a, double b)</i>	<i>-> a / b</i>
9	POW	= 0x14,	<i>//double power(double a, double b)</i>	<i>-> a ^ b</i>
10	LOG	= 0x15,	<i>//double logarithm(double a, double b)</i>	<i>-> log(a) / log(b)</i>
11	RND	= 0x16,	<i>//double round(double a)</i>	<i>-> (int64_t)(a + 0.5)</i>
12	ABS	= 0x17,	<i>//double absolute(double a)</i>	<i>-> +a</i>
13	SIN	= 0x18,	<i>//double sin(double a)</i>	<i>-> sin(a)</i>
14	COS	= 0x19,	<i>//double cos(double a)</i>	<i>-> cos(a)</i>
15	TAN	= 0x1a,	<i>//double tan(double a)</i>	<i>-> tan(a)</i>
16	ASIN	= 0x1b,	<i>//double arcsin(double a)</i>	<i>-> asin(a)</i>
17	ACOS	= 0x1c,	<i>//double arccos(double a)</i>	<i>-> acos(a)</i>
18	ATAN	= 0x1d,	<i>//double arctan(double a)</i>	<i>-> atan(a)</i>
19	EXP	= 0x1e,	<i>//double exponent(double a)</i>	<i>-> e ^ a</i>
20	RET	= 0xff	<i>//end of function</i>	

4.3.1.1 Beispiel

$$(|\sin(x) + \cos(x)|)^{\tan(6)}$$

Dies sieht schließlich wie folgt aus:

```
SIN VAR SLA  
COS VAR SLB  
ADD SLA SLB SLA  
ABS SLA SLA  
TAN NUM 6.0 SLB  
POW SLA SLB SLA  
RET SLA
```

Das würde folgendes bewirken:

Zuerst wird der Sinus von der Variable berechnet und im temporären Register A gespeichert.

Anschließend wird der Cosinus der unabhängigen Variable berechnet und im temporären Register B gespeichert.

Die Werte dieser Register werden nun miteinander addiert und wieder im Register A gespeichert.

Von dem Wert in Register A wird nun der Betrag genommen und wieder in Register A zurückgespeichert.

Nun wird der Tangens der Dezimalzahl 6.0 berechnet und in Register B gespeichert.

Danach wird der Wert von Register A hoch den Wert von Register B gerechnet und dies im Register A gespeichert.

Zum Schluss wird Register A als das mit dem Ergebnis anerkannt und die Funktion wird beendet.

4.3.2 Definitionen

4.3.2.1 Headerdatei primereg.h

```
1  #pragma once
2
3  #define N_FUNCTIONS          256 //Number of Functions per
   ↪ generation
4  #define N_PARENTS           32  //Number of parents to pick
5  #define N_NUMBERS           32  //Number of Numbers needed
6  #define N_PRINTDATA         256 //Number of loops till print
   ↪ of data
7  #define PERC_MUTATION       5   //Chance in percent of a
   ↪ mutation of a function
8  #define PERC_INHERIT        60  //Chance of child taking
   ↪ over data from parent B
9  #define FUNCTION_START      0x10 //Number of first function
10 #define FUNCTION_ONEARG     0x16 //Number of first function
   ↪ to take one argument
11 #define FUNCTION_END        0x1e //Number of last function
12 #define MAX_RAND_DOUBLE     512 //Maximum double to be
   ↪ picked as random number
13 #define MAX_RAND_DOUBLE_DIG 5    //Digits of
   ↪ MAX_RANDOM_DOUBLE
14 #define PERC_ADD_NEW         60  //Percentage of adding a new
   ↪ function on creation
15
16 #define DEBUG                //If defined, every
   ↪ generation results will be printed to stdout
17 #ifdef DEBUG
18 #undef N_PRINTDATA
19 #define N_PRINTDATA 1
20 #endif // DEBUG
```

4.3.2.2 Erläuterung

Diese stellen alle wichtigen Variablen des Quellcodes dar und können somit leicht verändert werden.

Hierbei steht das N_ für eine Nummer/Anzahl, das PERC_ für ein Prozentwert und die anderen für vordefinierte Werte.

Falls DEBUG definiert ist, dann werden alle Funktionen nach einem Loop ausgegeben. Sonst werden sie nur alle N_PRINTDATA Durchläufe ausgegeben.

N_FUNCTIONS gibt an, wie viele Elemente in unserer Population sein sollen.

N_PARENTS gibt an, wie viele der besten Elemente in der Population erhalten bleiben sollen und als die Grundgerüste der neuen Elemente gewertet werden sollen.

N_NUMBERS hält die Anzahl an Nummern, mit denen die der Elemente abgeglichen werden sollen.

N_PRINTDATA gibt an, nach wie vielen Durchläufen des primären Loops die Daten ausgegeben werden sollen.

PERC_MUTATION ist ein Prozentwert, der die Wahrscheinlichkeit einer Mutation des Elementes angibt. Dieser Wert sollte gering gehalten werden, da es sonst sehr schnell zu dem Fall kommen kann, dass die neuen Elemente keine Ähnlichkeit zu den N_PARENTS Elementen mehr haben.

PERC_INHERIT ist ein Prozentwert, der angibt, wie groß die Wahrscheinlichkeit ist, dass ein neues Element die Eigenschaft eines der besseren Elemente erbt. Dieser Wert sollte nicht zu groß werden, da die neuen Elemente sonst sehr schnell größer und komplexer werden, da sie dann nahezu alles ihrer Eltern übernehmen werden.

FUNCTION_START gibt den Wert der ersten Funktion in unserem Syntax an.

FUNCTION_ONEARG gibt den Wert der ersten Funktion an, die nur ein Argument und nicht zwei nimmt.

FUNCTION_END gibt den Wert der letzten Funktion in unserem Syntax an.

MAX_RAND_DOUBLE stellt die maximale Größe eines durch Zufall generierten Zahlenwertes dar.

MAX_RAND_DOUBLE_DIG stellt die maximale Anzahl an Nachkommastellen der durch Zufall generierten Zahlenwerte dar.

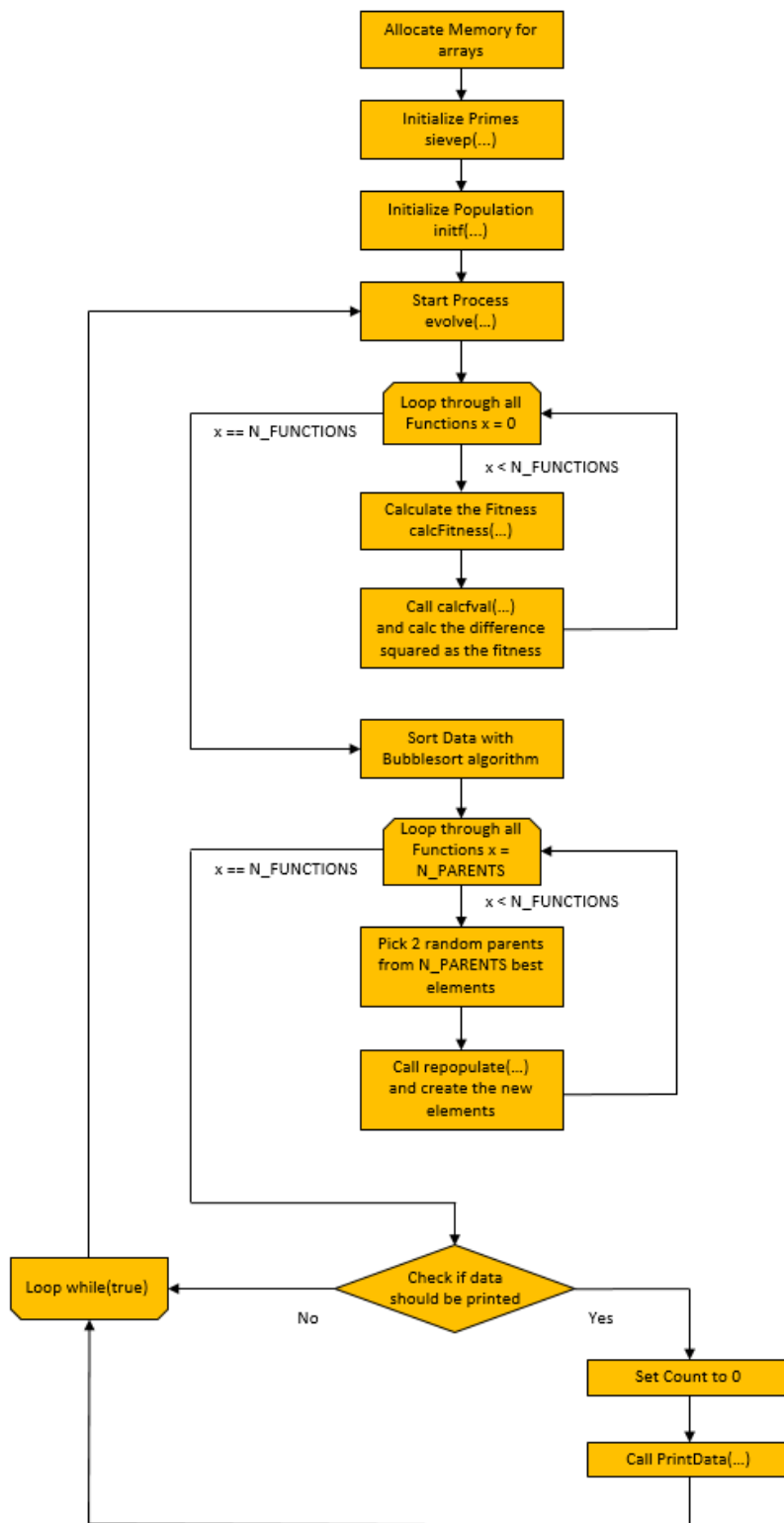
PERC_ADD_NEW gibt in Prozent an, wie groß die Wahrscheinlichkeit ist, dass eine Funktion bei der Initialisierung eine weitere mathematische Funktion erhält. (Alle erhalten mindestens eine)

4.3.3 Wahl des Funktionsparameter Datentyps

Wir haben uns dazu entschieden, den Datentypen `double` als den primären Datentypen für unsere Funktionen zu wählen, da dieser eine hohe Präzision gewährleistet, jeder x86 und x86-64 Prozessor mit diesem umgehen kann und er nur 8 Byte pro Variable aufnimmt. Zudem stellt die Mathematik-Bibliothek, welche wir für die mathematischen Funktionen verwendet haben, Funktionen bereit, welche mit `double` arbeiten.

4.3.4 Funktionsweise

4.3.4.1 Programmablaufplan



4.3.4.2 Erläuterung

Unser Algorithmus startet damit, den Speicher für die Population zu allokalieren. Anschließend allokiert er den Speicher für die Primzahlen, deren Erzeugung noch beschrieben wird. Danach wird der Zufallsgenerator mit der vergangenen Zeit seit dem 1. Januar 1970 in Millisekunden geseedet, damit wir möglichst zufällige Ergebnisse erhalten.

Nun werden die einzelnen Elemente der Population initialisiert. Hierfür wird die erste mathematische Funktion durch Zufall gewählt und deren Parameter ebenfalls als entweder VAR oder NUM. Anschließend wird geschaut, ob $\text{rand()} \bmod 101 \neq \text{PERC_ADD_NEW}$ gilt (Mit $\text{rand()} \bmod 101$ erhält man einen Zufallswert zwischen 0 und 100). Wenn dies gilt, dann wird eine weitere zufällige mathematische Funktion hinzugefügt. Falls nicht, dann wird geschaut, ob beide Register einen Wert beinhalten. Falls ja, dann wird noch eine zufällige Funktion, die zwei Argumente nimmt, gewählt und die nimmt dann als Parameter die beiden Register. Somit ist gewährleistet, dass kein Wert verloren geht und man nur noch ein volles Register am Ende hat, welches dann als das mit dem Ergebnis gewählt wird.

Falls als Parameter NUM gewählt wird, so heißt das, dass eine zufällige Fließkommazahl generiert werden muss. Diese wird wie folgt erzeugt: Zuerst wird ein HighDouble Wert erzeugt. Dieser liegt innerhalb der Grenzen 0 und MAX_RAND_DOUBLE. Anschließend wird ein LowDouble Wert erzeugt, welcher in den Grenzen 0 und 10*MAX_RAND_DOUBLE_DIG liegt. Dieser wird nun solange in einer Schleife durch 10 geteilt, bis er kleiner als 1 ist. Der Generierte Wert setzt sich dann zusammen aus dem HighDouble und dem LowDouble, die miteinander addiert werden. Zum Schluss wird noch geschaut, ob der Wert negativ sein soll, was durch ein Aufruf der $\text{rand}()$ Funktion geschieht.

Als Beispiel:

MAX_RAND_DOUBLE 512

MAX_RAND_DOUBLE_DIG 5

So werden als Zufallszahlen Zahlen zwischen -511.9999 und 511.9999 erzeugt.

Nachdem die gesamte Population mit validen Funktionen gefüllt worden ist, geht es in die Haupt-Schleife (*Main-Loop*). Hier wird zum einen die evolve Funktion aufgerufen, welche im Grunde den genetischen Algorithmus beinhaltet, und dann wird noch eine Zählervariable gehalten, welche angibt, wann die Daten alle ausgegeben werden sollen, und eine, die die derzeitige Gene-

ration mit sich trägt.

In der *evolve* Funktion wird zuerst die Güte (*Fitness*) der einzelnen Elemente der Funktion berechnet. Dabei wird diese anhand der Summe der Fehlerquadrate ermittelt: Diese Fitnessfunktion berechnet erst den Funktionswert des Elements für einen gegebenen x-Wert und anschließend seine Abweichung von der tatsächlichen Primzahl und quadriert diesen Wert. Diese Ergebnisse werden dann von $x = 0$ bis $x = N_NUMBERS - 1$ berechnet und aufaddiert. Der dann erhaltene Wert repräsentiert die Fitness. Umso geringer dieser Wert ist, desto besser ist die Funktion.

Wichtig ist, dass wir mathematische Fehler als vollständigen Eliminierungsfaktor ansehen. So werden Funktionen, die beispielsweise versuchen durch 0 zu teilen, mit dem Fitnesswert -1.0 versehen, was sie als nicht gültig kennzeichnet.

Sobald die Fitness für jede Funktion berechnet wurde, werden nun die *N_PARENTS* besten Elemente ausgewählt. Dies erfolgt mit Hilfe des Bubble-sort Algorithmus, welcher dafür sorgt, dass die besten *N_PARENTS* Elemente die ersten der Population sein werden. Wir haben uns dafür entschieden, die *N_PARENTS* besten Elemente nicht für die neue Population zu verwerfen, da wir der Ansicht waren, dass es besser wäre, diese zu behalten. Somit werden in jedem Durchlauf lediglich $N_FUNCTIONS - N_PARENTS$ Elemente der Population als neue Elemente aus den *N_PARENTS* Elementen generiert. Dies hat ebenfalls zur Folge, dass eine Funktion, die die Fitness 0 hat und somit perfekt ist, nicht durch Zufall verändert wird, zum Schlechteren und dann rausgeworfen wird. Dadurch ist gewährleistet, dass die *N_PARENTS* ersten Elemente der Population immer gute Werte erzielen werden.

Anschließend wird der Zufallsgenerator erneut geseedet, da es in dem von uns gewählten eine gewisse Wiederholung nach ungefähr 10^5 Zahlen gibt, die wir verhindern möchten (es werden keine *true random numbers* erzeugt). Dann kommt es zur Neubildung der Population. Hierfür werden immer zwei der *N_PARENTS* gewählt. Es kann auch vorkommen, dass diese beiden gleich sind. Das wirkt sich aber nicht negativ auf den Algorithmus aus und kann somit enthalten bleiben.

Die *repopulate*-Funktion ist die Zeitaufwändigste Funktion und zugleich die komplexeste. Zuerst geht sie durch die Funktion der Eltern und zählt alle mathematischen Funktionen, die in diesen aufgerufen werden. Anschließend allokiert sie genug Speicher, um einen Array von allen mathematischen Funktionen und deren Parametern zu erstellen, die bei den Eltern vorkommen. Nachdem sie diese mathematischen Funktionen der Eltern alle eingelesen hat, mutiert sie diese. Dies geschieht, indem die Funktion durch den Array geht

und, sobald $\text{rand}() \bmod 101 \neq \text{PERC_MUTATION}$ gilt, ersetzt sie entweder die mathematische Funktion oder einen der Parameter. Nachdem dieser Schritt getan ist, findet die Funktion ein verwendbares Erst-Element, dass keine Register als Parameter nimmt. Falls es keine mathematische Funktion findet, auf die dies zutrifft, verändert sie das erste Element des Arrays so, dass dieses keine abhängigen Parameter mehr nimmt. Zum Schluss geht es durch den Array von Funktionen und schaut, ob $\text{rand}() \bmod 101 \neq \text{PER_INHERIT}$ gilt. Wenn dies der Fall ist, erbt das neue Element diese Mathematische Funktion von denen der Eltern, falls dies von den Parametern her zutrifft. Sie kann aber nicht vererbt werden, wenn die neue mathematische Funktion als Parameter Register A nimmt und Register B, aber Register A zu diesem Zeitpunkt leer ist. Dann wird diese Funktion übersprungen. Nach Durchgehen des Arrays wird, wie bei der Initialisierung, geschaut, ob in beiden Registern ein Wert ist und falls ja, dann wird erneut eine mathematische Funktion zufällig ausgewählt, die beide Parameter nimmt und nur ein Register zurückgibt. Zum Schluss wird noch der Schluss, also RET, an die Funktion angehängen und der reservierte Speicher wird freigegeben beziehungsweise der von dem neuen Element verkleinert.

Diese Funktion wird nun für jedes Element der Population ausgeführt, bis die gesamte Population erneut gefüllt wurde. Dann beginnt die Schleife wieder von vorne und die Fitness der neuen Funktionen wird erneut berechnet.

4.4 Vor-/Nachteile gegenüber herkömmlichen Regressionsmethoden

4.4.1 Vorteile

Da unser Algorithmus lediglich die neu erstellten Funktionen mit denen der eingegebenen Funktionen abgleicht, ist er dazu in der Lage, eine Regression für jegliche Funktionstypen durchzuführen, sofern man ihm die richtigen Funktionswerte zum Vergleichen gibt. So kann der Algorithmus beispielsweise ohne Vorwissen über den möglichen Funktionstyp eine lineare oder quadratische oder exponentielle Funktion finden. In unserem Beispiel füllen wir den Array mit Primzahlen, um eine Funktion zu finden, welche diese beschreibt. Wir haben ihn aber auch mit den Zahlen einer quadratischen Funktion gestartet, welche er innerhalb von 3 Generationen, also 3 Aufrufen der *evolve*-Funktion, erfolgreich fand.

4.4.2 Nachteile

Der Speicherverbrauch und die Laufzeit lassen sich nicht einschätzen, da diese beiden Werte abhängig sind von der Rückgabe der `rand()` Funktion. Aber es lässt sich sagen, dass der Algorithmus für höhere `N_FUNCTIONS` oder höhere `N_PARENTS` mehr Zeit in Anspruch nehmen wird. Des Weiteren wird der Speicherverbrauch für höhere `PERC_ADD_NEW` und `PERC_INHERIT` ebenfalls schneller ansteigen. Da für höhere `PERC_INHERIT` die Tendenz des Speicherverbrauchs schnell steigend ist, muss man mit diesem Wert vorsichtig umgehen, da es sonst sehr schnell zu einem Speicherüberlauf kommen kann.

Auch liefert der Algorithmus keine Folge eines vorgegebenen Grades (wie etwa bei einer linearen, quadratischen, ... Regression), was für manche Sachverhalte unerwünscht sein kann.

5 Ergebnisse

Um Folgen einer ausreichenden Güte zu erhalten, muss der Algorithmus über längere Zeit laufen. Im Rahmen der Arbeitsprojektwoche stand uns leider nur sehr wenig Zeit zur Verfügung. Trotzdem möchten wir einige Durchläufe dokumentieren.

6 Fazit

Unser Algorithmus ist flexibel. Er versucht, alle ihm gegebene Datenpunkte möglichst gut zu approximieren, unabhängig davon, welchem Wachstum diese entstammen.

Allerdings muss man stets Glück haben, eine passende Folge zu finden. Dies ist darin zu begründen, dass diese per Zufall generiert werden. Somit ist das Verhalten des Algorithmus nicht vorhersehbar, weswegen konventionelle Regressionsmethoden schneller zu einem Ergebnis kommen.

Aber vielleicht hilft genau dieses Prinzip genetischer Algorithmen, Folgen für bisher unbekannte Wachstumsklassen, wie etwa Primzahlen, zu finden. Unsere bisherigen Messungen zeigen, dass unser Programm während der Laufzeit immer bessere Funktionen generiert. Somit könnte der Algorithmus irgendwann eine Folge für Primzahlen finden. Hierfür benötigen wir aber mehr Zeit und Rechenleistung.

Zusammenfassend konnte das volle Potential des Algorithmus im Rahmen der Arbeitsprojektwoche nicht völlig ausgeschöpft werden. Daher können wir zum aktuellen Zeitpunkt keine endgültige Aussage über die Eignung eines genetischen Algorithmus zur Regression treffen. Deshalb werden wir den Algorithmus weitergehend ausführen und sind auf dessen Ergebnisse gespannt.

A Quellen

- https://en.wikipedia.org/wiki/Formula_for_primes
- https://de.wikipedia.org/wiki/Sieb_des_Eratosthenes
- <https://de.wikipedia.org/wiki/Primzahl>
- <https://www.youtube.com/watch?v=ITLjfRiAc10>
- <http://www.iti.fh-flensburg.de/lang/algorithmen/asymp.htm>
- <http://www.programming-algorithms.net/article/41297/Sieve-of-Eratosthenes>
- <https://de.wikipedia.org/wiki/Miller-Rabin-Test>
- <http://www-cs-students.stanford.edu/~jl/Essays/ga.html>
- https://en.wikipedia.org/wiki/RSA_numbers#RSA-768
- https://en.wikipedia.org/wiki/RSA_Factoring_Challenge
- https://de.wikipedia.org/wiki/Eulersche_Phi-Funktion
- <https://www.youtube.com/watch?v=oXlY-yx1oIw>
- https://en.wikipedia.org/wiki/Sieve_of_Atkin
- https://en.wikipedia.org/wiki/Largest_known_prime_number
- https://en.wikipedia.org/wiki/Divergence_of_the_sum_of_the_reciprocals_of_the_primes
- https://de.wikipedia.org/wiki/Evolution%C3%A4rer_Algorithmus
- THE NATURE OF CODE by Daniel Shiffman