

AARHUS UNIVERSITET

INSTITUT FOR FYSIK OG ASTRONOMI

Eksamensprojekt i
EKSPERIMENTELFYSIK OG STATISTISK DATAANALYSE

Robuste Fitemetoder: Genetiske Algoritmer

Af:
MADS B. CARLSEN
EMIL HANSEN
NIKOLAJ KLINKBY

13. maj 2020

Resumé

To typer af genetiske algoritmer (GA'er) bliver præsenteret: binære og kontinuerte GA'er, og den teoretiske baggrund og implementeringen af disse bliver gennemgået. Dernæst ses på benchmarking af algoritmerne ift. minimering af Rastrigin-funktionen, og der ses på hvordan algoritmen kan bruges til at fitte til eksperimentel data via minimering af χ^2 , på både et simpelt og avanceret problem. Til avanceret fitting foreslås at den bedste løsning er en hybridløsning: at bruge GA'en til at finde et godt startgæt og derefter bruge SciPy's `curve_fit` til at finde nøjagtigt minimum samt kovariansmatricen.

1 Indledning

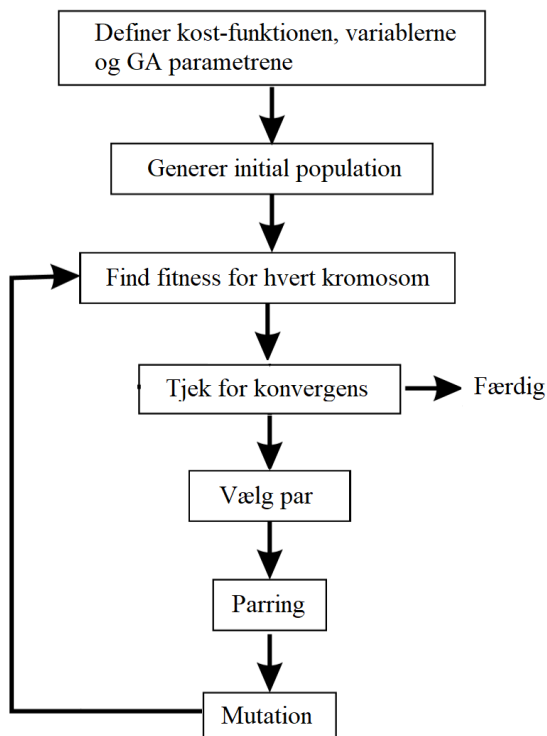
I alle problemstillinger, vi har mødt gennem kurset, har det været tilstrækkeligt at bruge forskellige udgaver af en lokal minimeringsalgoritme (som f.eks. gradient descent eller variationer af Newton's method)^{5;4} til at minimere funktioner og fitte til data. En enkelt gang eller to har det dog krævet lidt ekstra tid, fordi metoderne ofte behøver et godt start gæt for at finde et tilfredsstillende minimum. Dette problem kan overkommes ved at anvende en global optimerings algoritme, der i stedet for at skulle bruge et godt gæt (og derfra finde det nærmeste minimum) skal kende det område, den skal søge i og finder det globale minimum i området. Umiddelbart virker det mærkeligt, hvorfor man så ikke altid bruger disse algoritmer, men oftest kræver de meget længere tid om at finde et minimum og skal i mange tilfælde også skrives specifikt til det problem, som de skal behandle.³

En type af globale optimerings algoritmer er den såkaldte genetiske algoritme. Idéen bag genetiske algoritmer - og evolutionære algoritmer generelt - kommer fra Charles Darwins princip om "*survival of the fittest*". Grundlæggende går en genetisk algoritme ud på, at først danne en tilfældig samling af potentielle løsninger til et problem. De bedste af disse (typisk dårlige) potentielle løsninger 'avles', for at danne en endnu bedre. Dette gentages i flere generationer, til en optimal løsning er nået.⁴

Lidt mere specifikt indeholder en GA en *befolkning* af N_{pop} mulige løsninger til et problem, også kaldet *kromosomer*. Hver kromosom indeholder information af de N_{var} parametre for problemet, såkaldte *gener*. Hvordan selve generne er udtrykt varierer de forskellige typer GA imellem, som vil blive tydeligt i de følgende afsnit. Typisk vil hele populationen dog repræsenteres vha. en matrix \mathbf{P} , hvor hver række er et kromosom. Hvor god en løsning de enkelte kromosomer er, evalueres gennem en såkaldt *kostfunktion*. Herved tildeles hvert kromosom en værdi for, hvor gode de er, kaldet *fitness-værdien*. Oftest ønskes at minimere en matematisk funktion $f : \mathcal{S} \rightarrow \mathbb{R}$, hvor $\mathcal{S} \subset \mathbb{R}^n$. Målet er således at finde et kromosom repræsenterende et $\mathbf{x}^* \in \mathcal{S}$, hvor

$$f(\mathbf{x}^*) \leq f(\mathbf{x})$$

for alle $\mathbf{x} \in \mathcal{S}$.³ Dvs. de bedste løsninger i dette tilfælde er de løsninger med lavest fitness-værdi.⁴



Figur 1. Generelt flowchart for genetiske algoritmer (inspireret af *Practical Genetic Algorithms*⁴).

Herefter udvælges den X_{keep} bedste procentdel af den inertielle befolkning, baseret på deres fitness-værdi. De $N_{parents} = N_{pop} X_{keep}$ fitteste individer, rundet op til næste lige tal, får lov til at danne nye individer, der erstatter de ikke-udvalgte (de dårligste løsninger). Inden denne *parring* skal løsninger sammensættes i par. Der findes forskellige metoder til at udvælge parrene på. Det kan både gøres hvor parrene vælges tilfældigt og uafhængig af fitness-værdien, men også hvor hvert individs fitness-værdi bliver vægtet, så de bedste løsninger har størst sandsynlighed for at føre deres gener videre. De udvalgte par skal nu danne $N_{offspring} = N_{pop} - N_{parents}$ afkom, så den nye population har samme antal som den oprindelige. Igen er der mange måder at gøre dette på, men den simpleste, som er direkte inspireret af genetikken, er at vælge ét eller flere *crossover*-punkter i kromosomerne, hvor

kromosomerne krydser over og bytter gener⁴. Denne metode bliver også illustreret under binære algoritmer.

Det sidste skridt er *mutation*. Her skal vælges endnu en vigtig faktor, *mutationsraten* R_{mut} , som også tager værdier mellem 0 og 1, men normalt ligger på en værdi på omkring 0.2.^{3;4} Den mest simple tilgang til mutation er at tilfældigt udvælge indgange i \mathbf{P} , og så erstatte indgangen med en tilfældig ny passende værdi (den nye værdi skal stadig være en mulig værdi for parameteren). Antallet af mutationer der skal foretages beregnes ved $N_{\text{mutations}} = \lceil N_{\text{pop}} R_{\text{mut}} \rceil$. En meget vigtig pointe ved mutation er at sørge for ikke at mutere populationens bedste kromosom, for at forhindre at man smider den bedste løsning væk.⁴

Ovenstående serie af skridt gentages nu som illustreret på Fig. 1, indtil en passende løsning er opnået, eller et maximalt antal iterationer nåes.

2 Genetiske algoritmer i praksis

2.1 Binær genetisk algoritme

Den binære GA er den oprindelige udgave af de genetiske algoritmer. Inspireret af hvordan DNA oplagerer information gennem rækkefølgen af baser i DNA'et, lader man hver kromosom i algoritmen bestå af en bit-streng - dvs. en liste af 1'ere og 0'ere.⁴ Hvis man fx havde en genetisk algoritme der arbejder med tre parametre, kunne et kromosom have formen:

$$[101\ 010\ 110]$$

Hvert gen i kromosomet består altså af N_{gen} bits. Generne skal omregnes fra binær til reelle tal, for at kunne evalueres i værdi-funktionen. Dette kan gøres vha. følgende ligninger:⁴

$$q = 2^{-(N_{\text{gen}}+1)} + \sum_{m=1}^{N_{\text{gen}}} \text{gen}[m] \cdot 2^{-m}$$

$$p = q(p_{hi} - p_{lo}) + p_{lo}$$

Hvor p er den reelle værdi af genet, $\text{gen}[m]$ er den m 'te bit i genet og p_{hi} og p_{lo} er den højeste hhv. laveste værdi for det interval som de binære gener skal repræsentere. Grundet dets binære opbygning, kan hvert gen kun repræsentere $2^{N_{\text{gen}}}$ værdier på dette interval. Hvis fx $N_{\text{gen}} = 8$ giver dette kun en opløsning på 256 forskellige værdier. Den diskrete natur af den binære algoritme gør dog samtidig, at denne er ideel til at løse diskrete problemstillinger.⁴ Bemærk nemlig, at man ikke er begrænset til at lade generne repræsentere reelle tal, hvis man ikke forsøger at optimere en matematisk funktion.

Som udgangspunkt følger den binære genetiske algoritme de trin og metoder, der er beskrevet i forrige afsnit. Der findes mange metoder at parre sine individer, men den mest

simple er *single-point crossover*. Her vælges et bit i forældre-kromosomerne og to nye individer dannes ved at ombytte bit-strengene før og efter dette punkt de to forældre imellem. Dette forstås bedst ved et eksempel. Lad følgende være to forældre:

$$\begin{array}{c} [010110101] \\ [100111001] \end{array}$$

Det 4. bit vælges som overkrydsningspunktet og to nye individer dannes:

$$\begin{array}{c} [100110101] \\ [010111001] \end{array}$$

På denne måde får man dannet nye individer til populationen, der er en blanding af genetisk materiale fra forældrene.

Mutation i den binære GA er en ret simpel procedure. Her udvælges blot et antal bits, givet ved $N_{mut} = \lceil N_{pop} X_{mut} N_{gen} \rceil$, i hele populationen, og disse bits inverteres. Hvis fx et udvalgt bit har værdien 0 inverteres denne således til 1.⁴

2.2 Kontinuert genetisk algoritme

De kontinuerlige GA'er tager samme form som de binære GA'er, men kromosomerne er opbygget forskelligt. I det kontinuerlige tilfælde tildeler vi ikke hver variabel en bit streng, men beholder dem direkte, som en værdi i kromosomet. Dvs. at en N_{var} -dimensional funktion med $p_1, p_2, \dots, p_{N_{var}}$ variabler har et kromosom givet ved,⁴

$$\mathbf{x} = (p_1 \ p_2 \ \cdots \ p_{N_{var}}).$$

Fordelen ved dette er at vi ikke længere skal tænke over mængden af bits pr. variable eller i det hele taget konvertere bits.

Herfra løber processen som før og der skal laves en initial population, hvilket gøres ved at lave en $N_{pop} \times N_{var}$ -matrix \mathbf{P} med tilfældige værdier mellem 0 og 1 i hver indgang, og er på formen

$$\mathbf{P} = \begin{pmatrix} - & \mathbf{x}_1 & - \\ - & \mathbf{x}_2 & - \\ & \vdots & \\ - & \mathbf{x}_{N_{pop}} & - \end{pmatrix}$$

Dette skaber N_{pop} kromosomer i populationen. Hver gang vi ønsker at finde værdien for hvert kromosom de-normere vi det, ved at tage de øvre og nedre grænser for variablerne, som vi

ønsker at undersøge (vi udvælger et område af funktionen at undersøge). For alle variablerne bruger vi da at⁴

$$p_{\text{denorm}} = (p_{\text{max}} - p_{\text{min}})p_{\text{norm}} + p_{\text{min}}.$$

De de-normerede kromosomer kan vi indsætte direkte i kostfunktionen for at evaluere kromosomerne, men vi beholder de normede kromosomer, som er dem vi arbejder videre med. Hvis vi vælger at parre dem på samme måde, som den binære GA, med et crossover-punkt/punkter, så introducerer vi aldrig nogen ny information, men bytter bare rundt på værdierne for vores variabler. I stedet kan vi bruge en anden metode, nemlig at introducere en tilfældig afvigelse fra enten moder kromosomet p_m eller fader kromosomet p_f , afhængigt af både moderen og faderen. Dette opnås ved at udvælge en tilfældig, positiv, værdi β så,⁴

$$\begin{aligned} p_{\text{ny1}} &= p_m - \beta(p_m - p_f) \\ p_{\text{ny2}} &= p_f + \beta(p_f - p_m) \end{aligned}$$

På den måde kan vi introducere variabler udenfor det område, hvor fadere og moderen ligger. Hvis vi vil udvide det meget, så kan vi lade $0 \leq \beta \leq a$, for en værdi $a > 1$. Oftest vil vi dog bare vælge $a = 1$, hvorved vi ender med at undersøge variabler ud fra moderen og faderen i intervaller $[p_m, p_f]$ og $[p_f, 2p_f - p_m]$. Hvis dette gøres for alle variabler og man lader hvert par lave to afkom, så kan afkommet skrives som⁴

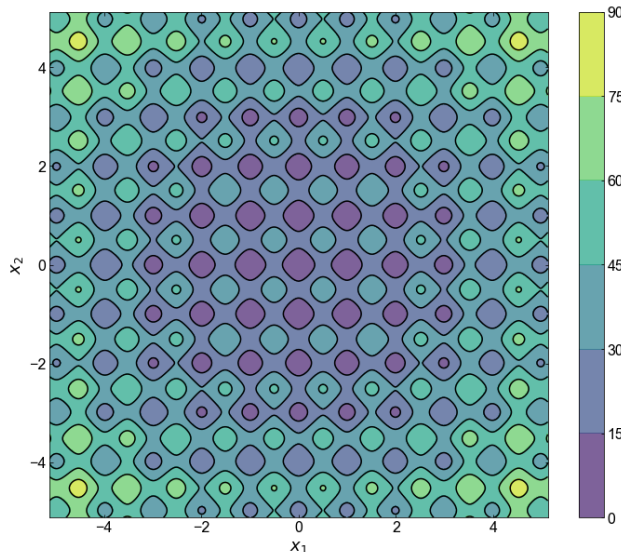
$$\begin{aligned} \text{afkom}_1 &= \text{moder} - [\beta_1(p_{m,1} - p_{f,1}), \beta_2(p_{m,2} - p_{f,2}), \dots, \beta_{N_{\text{var}}}(p_{m,N_{\text{var}}} - p_{f,N_{\text{var}}})] \\ \text{afkom}_2 &= \text{fader} + [\beta_1(p_{f,1} - p_{m,1}), \beta_2(p_{f,2} - p_{m,2}), \dots, \beta_{N_{\text{var}}}(p_{f,N_{\text{var}}} - p_{m,N_{\text{var}}})] \end{aligned}$$

Den nye population af forældre og afkom, får udvalgt hvilke indgange der skal muteres og da det hele er normeret, så indsættes bare en ny værdi mellem 0 og 1 på den muterede plads. Herefter er det vigtigt at tjekke alle indgange stadig ligger indenfor 0 og 1. Hvis de ikke gør kan man enten sætte dem til den nærmeste værdi 0 eller 1, eller lade dem køre i ring i intervallet.

Herfra har man nu en ny population, der er klar til at blive de-normeret og evalueret.

3 Resultater

3.1 Testkørsler af genetiske algoritmer



Figur 2. Contour plot af funktion (1) med $n = 2$. Farvesøjlen til højre giver sammenhæng mellem farve og funktionsværdier.

Til at teste genetiske algoritmer findes flere standardfunktioner. Her vælges følgende funktion, kendt som Rastrigin funktionen, der også benyttes til at teste andre typer af minimeringsalgoritmer:⁶

$$f(\mathbf{x}) = nA + \sum_{i=1}^n [x_i^2 - A \cos(2\pi x_i)] \quad -5.12 \leq x_i \leq 5.12 \quad (1)$$

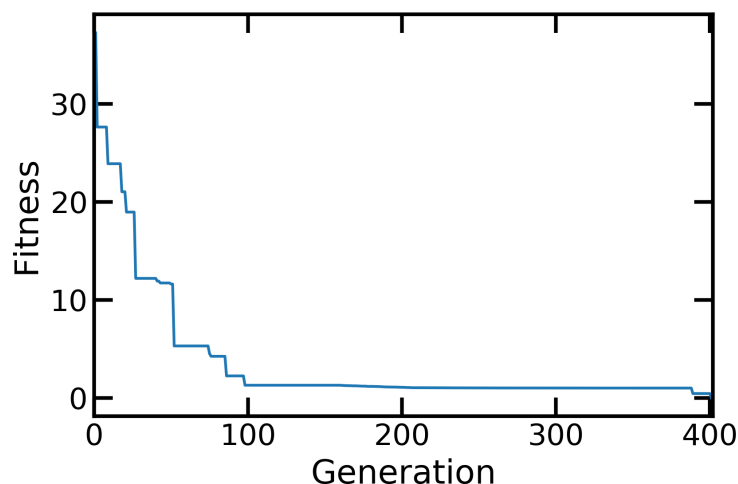
Her sættes $A = 10$ hvorved funktionen vil have globalt minimum i $\mathbf{x} = \mathbf{0}$ med funktionsværdien $f(\mathbf{0}) = 0$. Funktionen er valgt, da det er let at opskalere dimensionen, og dermed øge kompleksiteten af problemet, der skal løses. Derudover har funktionen mange lokale minima (se evt. Fig. 2), som de lokale minimeringsmetoder vil konvergere imod, hvis deres udgangspunkt er i nærheden af disse. Det vil derfor være utroligt svært at finde det globale minimum af (1) med metoder som f.eks. gradient descent, hvorimod det skulle kunne være muligt med en genetisk algoritme.

For at teste dette er skrevet 2 forskellige genetiske algoritmer, fundet i bilag A1 og A2. Algoritmen i A1 er kontinuert, mens den anden er binær. I det følgende vil algoritmerne i bilag A1 og A2 testes og sammenlignes. For den binære algoritme benyttes en population

på $N_{pop} = 8$, en mutationsrate på $R_{mut} = 0.03$ og $N_{gen} = 12$ bits til at repræsentere hver parameter. Algoritmen i A2 udvælger par til parring, ved at vægte de bedste løsninger højest, (såkladt rank-weighting⁴). Algoritmen søger indtil den bedste løsning har en fitness-værdi på under 0.01 ((1) antager kun denne værdi i brønden ved det globale minimum), da den binære GA ikke vil kunne ramme en funktionsværdi på præcis 0.0 grundet dens endelige opløsning.

I den kontinuerte GA fundet i A1 bruges at størrelsen af populationen er givet ved $N_{pop} = 8n$, hvor n er antallet af parametre. Dette gøres grundet en større population gør algoritmen mere effektiv når der skal findes minimum i højere dimensioner.⁴ Derudover er $X_{keep} = 0.5$ og $R_{mut} = 0.2$. At sætte mutationsraten eller X_{keep} meget højere eller lavere resulterer som oftest i markant længere beregningstid. Udvalgelsen er tilfældig blandt de N_{keep} bedste for at bevare nogle mindre gode løsninger, så algoritmen ikke bliver for elitær.

En tilfældig kørsel med GA'en fra bilag A1, på (1) er illustreret på Fig. 3. Her ses, at det lykkedes algoritmen at opnå en løsningen efter ca. 400 generationer. Det bemærkes dog også, at der til start sker et hurtigt fald, mens at udviklingen i den bedste løsning flader ud hen mod de høje generationer.

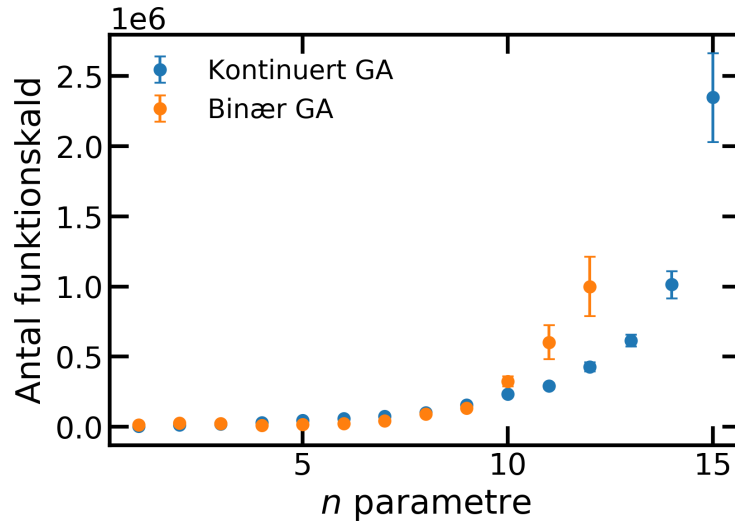


Figur 3. *Eksempel på en kørsel med den kontinuære GA i bilag A1 på (1) med 4 parametre. Fitness-værdien for den bedste løsning er plottet som funktion af antallet af generationer.*

Fordi GA'er er baseret på tilfældighed (både start population, avl og mutation), vil det også variere en del hvor mange generationer, det tager at opnå en løsning. For at undersøge kompleksiteten af beregningerne udføres derfor 15 kørsler for hver dimension. Dette lave antal kørsler skyldes den store mængde tid involveret i beregningerne. Istedet for at måle antallet af generationer, måles istedet hvor mange gange kostfunktionen kaldes (dette gøres ofte,⁶ da dette mål er uafhængig af algoritmens opbygning). Fordelingen af antal kald til

kostfunktionen for hver dimension viser sig ikke at være Gauss-fordelt, heller ikke ved højere antal kørsler. Grundet kombinationen af få datapunkter og ikke Gauss-fordelt data vælges usikkerheden på dataet at findes vha. bootstrap metoden.⁵ Dvs. for hvert samling af 15 målinger, dannes et stort antal samples, bestående af tilfældige udtrukne værdier fra den oprindelige samling. Hver sample indeholder samme antal elementer som den oprindelige, dvs. 15. For hver sample beregnes nu en middelværdi, og denne samling af middelværdier fra de forskellige sample fordelinger afspejler nu fordelingen af den oprindelige samling af elementer. Ved at sortere middelværdierne kan 2/3 konfidensintervallet herefter findes for fordelingen, og dermed for dataet fundet vha. GA'en.

Ovenstående udføres for både den binære GA fra bilag A2 og den kontinuerte GA fra bilag A1. Resultatet er plottet i Fig. 4.



Figur 4. Antallet af funktionskald for den kontinuerte og binære GA fra bilag A1 til minimering af (1). Algoritmen stopper og returnerer resultatet, hvis $f(\mathbf{x}) \leq 0.01$.

Det er altså lykkedes de fremstillede algoritmer, at løse (1) med op til hhv. 12 og 15 parametre. Det ses umiddelbart at den kontinuerte algoritme klarer sig bedre end den binære algoritme ved de højere dimensioner, mens de følges ad op til omkring 9 dimensioner. Der er kun 12 datapunkter for den binære GA, da beregningstiden for højere dimensioner var meget stor. I dataet ses desuden en tendens til at antallet af funktionskald stiger eksponentielt. Læg desuden mærke til, at skalaen på akse er 10^6 . Dvs. at når der skal optimeres et stort antal parametre på én gang, benytter GA'erne op mod flere millioner funktionskald.

3.2 Fitting med genetiske algoritmer

Her vil det ses på, hvordan det er muligt at bruge en kontinuert genetisk algoritme til at udføre et fit af en funktion $y(x, \mathbf{a})$, hvor \mathbf{a} er en vektor med funktionens uafhængige parametre, til et datasæt bestående af n datapunkter (x_i, y_i) for $i = 1, 2, \dots, n$. Der sammenlignes i afsnittet med fitte-rutinen `curve_fit` fra Python's Scipy. Denne benytter som standard Marquardt-Levenberg metoden til minimering, som er en blanding af gradient descent og Gauss-Newton metoden (en videreudvikling af Newton metoden).^{5;1} For at udføre et fit med en genetisk algoritme minimeres χ^2 -funktionen givet ved,⁵

$$\chi^2(\mathbf{a}) = \sum_i^n \left(\frac{y_i - y(x_i, \mathbf{a})}{\alpha_i} \right)^2, \quad (2)$$

hvor α_i er usikkerheden på det i 'te datapunkt. Målet er således at finde den optimale sammensætning af \mathbf{a} , nemlig \mathbf{a}^* hvorom det gælder at $\chi_{\min}^2 = \chi^2(\mathbf{a}^*)$.

Først ses på et simpelt eksempel; at fitte til en trigonometrisk funktion, i dette tilfælde en sinusfunktion på formen

$$y(x) = A \sin(\omega x + \phi). \quad (3)$$

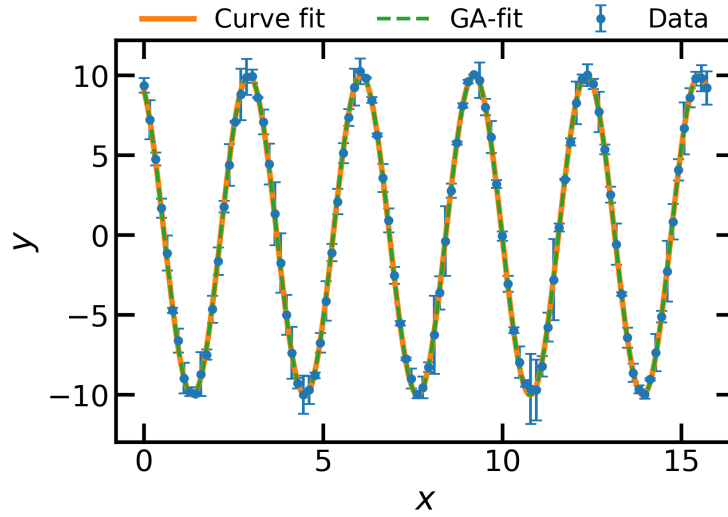
Det, der oftest volder problemer i denne forbindelse at få et ordentligt gæt på vinkelfrekvensen ω , da fitte-rutinen oftest er meget sensitiv overfor denne.

På Fig. 5 ses resultatet af et fit med både `curve_fit` og en kontinuær GA (bilag A1) på noget tilfældigt genereret sinusdata med støj. Visuelt er det tydeligt, at de to metoder finder samme resultat. På Tabel 1 finder vi de returnede parametre fra hhv. `curve_fit` og GA'en, hvor vi ser afvigelsen mellem de to er ganske lille, men grundet i at GA'en ikke nødvendigvis finder minimum lige så nøjagtigt som `curve_fit`. Men mens `curve_fit` skulle have et relativt nøjagtigt gæt for ω (ca. intervallet fra 1.8-2.3) for at konvergere til ovenstående løsning, var GA'en indstillet til blot at undersøge ω i intervallet 0-50. Det lykkedes altså GA'en at søge over et større område, og frasortere de lokale minima, som `curve_fit` ellers godt kan falde i. Til gengæld benyttede GA'en 11.200 funktionskald mens `curve_fit` benyttede omkring 30 (alt efter hvor nøjagtig startgættet var).

Metode	A	ω	ϕ	χ^2	$P(\chi^2; \nu)$
<code>curve_fit</code>	10.003	1.998	2.018	87.751	0.738
GA-fit	10.002	1.999	2.017	87.803	0.737

Tabel 1. Funde parametre til fittene illustreret på Fig. 5.

Eftersom det lykkedes at udføre et simpelt fit, forsøges i det følgende med et lidt mere ambitiøst. Der forsøges at fitte til data fra Fizeau's eksperiment, hvor der ses på intensiteten



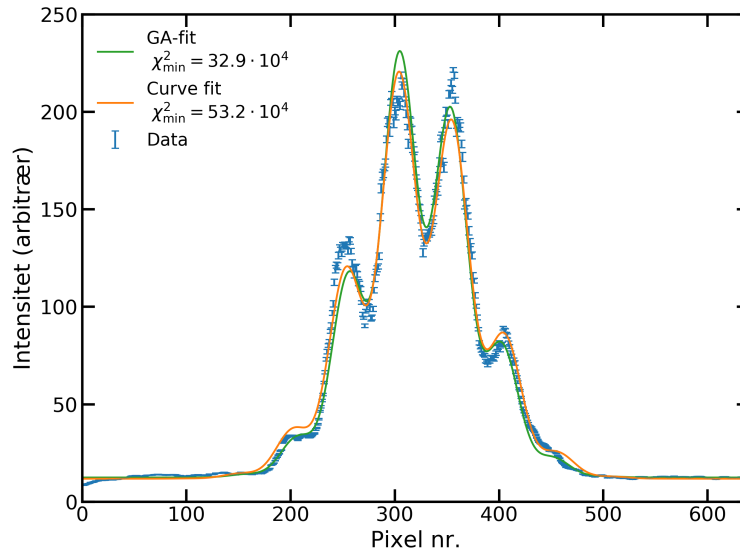
Figur 5. Generet sinusdata med normalfordelt støj (blå) hvortil der er lavet to fits: Fit med den genetiske algoritme (grøn) og med SciPy's `curve_fit` (orange). Dataet er genereret ud fra (3) med $A=10$, $\omega = 2$ og $\phi = 2$. Bemærk desuden, at der ikke regnes med usikkerheder på datapunkterne, da interessen er hvor godt selve fittet er.

på tværs af et interferensmønster. Der skal gøres opmærksom på, at dataet ikke er optimalt, og at fittet derfor ikke vil blive det uanset. Interessen ligger i, om GA'en kan finde et evt. bedre bud på en løsning end `curve_fit`, og uden at bruge lang tid på at finde et start-gæt. Funktionen, der skal fittes til dataet er givet ved ligningen²

$$f(x) = \sin^2(kx + \delta) \left[b \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} - c \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} \right] + c \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} + a, \quad (4)$$

med parametrene $[a, b, c, k, \delta, \mu, \sigma]$. Her skal (2) altså minimeres for 7 variable på samme tid. At have så mange variable og en så kompliceret funktion gør, at man skal være ret præcis med sit startgæt, hvis man skal opnå et brugbart fit med `curve_fit`.

Hvis vi betragter Fig. 6 kan vi se, at der er blevet fundet en løsning vha. den genetiske algoritme selvom søgerummet er stort. Som forventet har begge fits har en høj χ^2 -værdi og er ikke gode fits (faktisk har begge fits en $P(\chi^2_{\min}; \nu)$ -værdi på 0), men vi ser dog at χ^2 -værdien er lavest for fittet fra den genetiske algoritme. Pointen er, at der med GA'en er fundet et sæt af parametre (omend ikke godt) uden at skulle bruge lang tid på at gætte på startværdier til `curve_fit`.



Figur 6. Dataene fra forsøget (blå) samt fits lavet hhv. SciPy's `curve_fit` (orange) og vores genetiske algoritme (grøn). Vi ser her, at χ^2_{\min} -værdien er lavest for fittet genereret med den genetiske algoritme. Se koden i bilag B1.

4 Diskussion

Dataene fra Fig. 4 ligger op til en diskussion om hvilken slags GA er bedst. Med den begrænsede data vi har produceret ligner det at den kontinuerlige GA klarer sig bedre, men det kan være svært at sammenligne dem, da de bygger på to forskellige måder at lave en GA på. Som det også bliver nævnt i [Haupt & Haupt], så kræver en problemstilling ofte en skredersyet GA, for optimal løsning. Vores data fra Fig. 4 peger på at den kontinuerlige algoritme i bilag A1 er bedre til at løse en Rastrigin funktion end den binære algoritme i bilag A2, men det betyder ikke at den binære algoritme ikke ville være bedre til at løse et andet problem. En kontinuerlig GA er, i tilfælde med kontinuerlige funktioner, bare mest praktisk. Som nævnt i afsnit 3.1, så arbejder den binære algoritme fx meget langsommere end den kontinuerlige. En undersøgelse af koden viste, at størstedelen af tiden (og beregningerne) gik til konvertering fra bit til reelle tal. De kontinuerlige algoritmer har den klare fordel at de slipper for disse konverteringer.

Parametrene benyttet for GA'erne i afsnit 3.1, var fundet som de umiddelbart bedste. En pointe er, at de overhovedet ikke er ens for de to algoritmer. Generelt kan man ikke fastsætte nogle parametre, der fungerer bedst, det viser sig i stor grad at afhænge af hvordan ens GA er opbygget. Vælges fx en alternativ måde at parre sine kromosomer, vil dette have stor indflydelse på effektiviteten af algoritmen og hvilke parametre, der er mest optimale. Da

en GA ikke er en matematisk formuleret algoritme, er det desuden meget svært at forudsige hvilke parametre og metoder, der vil være gode. En advarsel er således, at man hurtigt vil kunne bruge lige så lang tid på at finde de optimale parametre og opsætning for ens GA, som man kan på at finde et godt start gæt for en lokal minimeringsrutine. Endvidere vil nogle problemstillinger kun kunne løses med den helt rette sammensætning af metoder i ens GA.

Hvorvidt det er relevant at bruge en GA til en problemstilling ser ud til at afhænge stærkt af problemet. I tilfældet med den trigonometriske funktion ville man hurtigt kunne komme med et passende gæt til `curve_fit` og fordelten ved den genetiske algoritme kommer nok først i det, man ikke længere ønsker at være begrænset af at skulle gætte. For eksempel i et tilfælde hvor man skal fitte mange funktioner til mange datasæt, hvor alle datasættene ikke kan bruge samme gæt for `curve_fit`. Til gengæld skal man være opmærksom på, at en GA kræver væsentlig mere beregningskraft, som set i afsnit 3.2.

En stor fordel ved `curve_fit` er også, at den giver én covarians-matricen. Man kunne bruge `curve_fit` på sin GA's bedste løsning for at få dette. En kombination af en GA og en lokal minimeringsalgoritme kaldes også for en hybridalgoritme. Ideen er, at man med den genetiske algoritme kan finde den rigtige brønd, hvori det globale minimum ligger, hvorefter en lokal minimeringsalgoritme kan tage over og hurtigt konvergere til en præcis løsning. Man kan også aktivere den lokale minimeringsalgoritme, hvis forbedringen i værdien for kost-funktionen for den GA'en begynder at aftage med en hvis værdi.⁴ Et godt eksempel er Fig. 3, hvor det ville have været mere optimalt at bruge en hybrid algoritme, der blev aktiveret efter 100 generationer.

Det skal nævnes, at GA'erne skrevet her, ikke er optimale eller specielt hurtige (jævnfør Fig. 4). Skriver man dem rigtigt (og tilpas avanceret), kan de dog være meget effektive og løse (1) med 400 parametre på ganske kort tid.⁶ Det skal også nævnes, at der findes andre globale minimeringsalgoritmer, der i mange tilfælde arbejder hurtigere og mere effektivt, fx differentiell-evolution algoritmer.³ Hvis man programmerer i Python, så har biblioteket `scipy.optimize` allerede en differentiell evolution algoritme til rådighed.¹

5 Konklusion

Vi har i dette projekt betragtet og udviklet to forskellige genetiske algoritmer. Først en binær genetisk algoritme, som egner sig bedst til diskrete problemer (men som kan bruges til kontinuerte problemer) og dernæst en kontinuert genetisk algoritme, som modsat egner sig bedst til kontinuerte problemer. Algoritmerne fandt succesfuldt det globale minimum på rastrigin-funktionen, op til 15 dimensioner. Illustreret er også hvordan man kan bruge den kontinuerte GA til at minimere χ^2 på både et simpelt og et mere avanceret problem. Konklusionen er, at den bedste måde at udnytte GA'en på, er til at finde et godt startgæt til `curve_fit`, og dermed bruge kombinationen af de to som en hybridalgoritme til fitting.

Litteratur

- [1] Optimization and root finding (scipy.optimize), Dec 19, 2019. URL <https://docs.scipy.org/doc/scipy/reference/optimize.html>.
- [2] Mads Slot Bertelsen. A reenactment of the fizeau experiment. Master's thesis, Aarhus University, 2015.
- [3] Andries P. Engelbrecht. *Computational Intelligence An Introduction*. John Wiley & Sons, Inc, second edition edition, 2007.
- [4] Randy L. Haupt and Sue Ellen Haupt. *Practical Genetic Algorithms*. John Wiley & Sons, Inc, second edition edition, 2004.
- [5] Ifan G. Huges and Thomas P.A. Hase. *Measurements and their Uncertainties*. Oxford University Press, 2010.
- [6] H. Mühlenbein, M. Schomisch, and J. Born. The parallel genetic algorithm as function optimizer. *Parallel Computing*, 17(6):619 – 632, 1991. ISSN 0167-8191. doi: [https://doi.org/10.1016/S0167-8191\(05\)80052-3](https://doi.org/10.1016/S0167-8191(05)80052-3). URL <http://www.sciencedirect.com/science/article/pii/S0167819105800523>.

A Genetiske algoritmer

A.1 Kontinuert algoritme 1

Her ses koden til den ene af de to kontinuerte genetiske algoritmer, der er vedhæftet. Under selve algoritmen findes et eksempel på at finde minimum på Rastrigin-funktionen i tre dimensioner, da udpakningen af værdierne fra funktionen ikke er trivial.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 ### ALGORITMEN ###
5 def ga(cost_function, constraints, mutrate=0.2, X_keep=0.5, N_pop=100,
6     generations=1000, terminate=False, criteria=0.1):
7     # Konstanter
8     N_var = len(constraints)
9     N_keep = int(2*np.ceil(N_pop*X_keep/2))
10    N_offspring = int(N_pop - N_keep)
```

```

11     N_mutations = int(N_pop*mutrate)
12
13     # Counter til antal funktionskald
14     func_call_count = 0
15
16     # Befolkning med N_pop kromosomer
17     pop = np.random.rand(N_pop, N_var)
18
19     # Constraints
20     min_c, max_c = np.asarray(constraints).T
21     diff = np.fabs(min_c - max_c)
22
23     # Beregn forste fitness
24     pop_denorm = min_c + np.multiply(pop, diff)
25     fitness = np.asarray([cost_function(chromosome) for chromosome in
pop_denorm])
26     func_call_count += len(fitness)
27
28     # Udvaelgelse
29     best_indx = np.array(fitness).argsort()[:N_keep]
30     parents = pop[best_indx]
31     pop = parents
32
33     for i in range(generations):
34
35         # Parring
36         offspring = np.zeros((N_offspring, N_var))
37         for _ in range(N_offspring)[0::2]:
38             parent1 = parents[np.random.randint(N_keep)]
39             parent2 = parents[np.random.randint(N_keep)]
40             alpha = np.random.rand(N_var)
41             beta = np.random.rand(N_var)
42             offspring1 = np.clip(parent1 + alpha*(parent1 - parent2), 0, 1
)
43             offspring2 = np.clip(parent2 - beta*(parent1 - parent2), 0, 1)
44             offspring = np.vstack([offspring, offspring1, offspring2])
45
46         pop = np.vstack([parents, offspring])
47
48         # Mutationer
49         rows = np.random.randint(1, N_pop, N_mutations)
50         cols = np.random.randint(0, N_var, N_mutations)
51         for i, j in zip(rows, cols):
52             pop[i, j] = np.random.rand()
53
54         # Beregn fitness og find bedste kromosom

```

```

55     pop_denorm = min_c + np.multiply(pop, diff)
56     fitness = np.asarray([cost_function(chromosome) for chromosome in
pop_denorm])
57     func_call_count += len(fitness)
58     best_idx = np.argmin(fitness)
59     best_individual = pop_denorm[best_idx]
60
61     # Udvaelgelse
62     best_idx = np.array(fitness).argsort()[:N_keep]
63     parents = pop[best_idx]
64     pop = parents
65
66     # Stop algoritme, hvis kriteriet er naaet
67     if terminate:
68         if fitness[best_idx] <= criteria:
69             yield best_individual, fitness[best_idx], func_call_count
70             break
71
72     yield best_individual, fitness[best_idx], func_call_count
73
74 ### EKSEMPEL - minimum paa rastrigin ###
75
76 def rastrigin(x):
77     A = 10
78     dim = len(x)
79     return (A*dim + sum(x**2 - A*np.cos(2*np.pi*x)))
80
81 # Antal dimensioner
82 N_dim = 3
83
84 # Kor algoritmen og udpak vaerdier
85 # Terminer algoritmen, hvis der faaes en fitness paa 0.0
86 params, fitness, func_calls = zip(*list(ga(rastrigin,\
87     [[-5.12, 5.12]]*N_dim,\
88     N_pop=25*N_dim,\
89     generations=10000,\
90     terminate=True, criteria=0.0)))
91
92 # Print parametre og fitness
93 print(f'Parametre: {params[-1]}')
94 print(f'Fitness: {fitness[-1]}')
95 print(f'Funktionskald: {func_calls[-1]}')

```


A.2 Binær algoritme

Er man interesseret i at downloade og selv køre koden, findes den på følgende link: https://drive.google.com/open?id=1lljTx7TwO7A9aUP_V5wHpb4mLWVeA14

```
1 import numpy as np
2
3 def cost_func(p):
4     """Funktionen der forsoges at minimeres med GA'en. Her i from af
5     Rastrigin"""
6     n = len(p)
7     A = 10
8     return A*n + sum(x**2-A*np.cos(2*np.pi*x) for x in p)
9
10 class GA:
11     def __init__(self, cost_func, N_params, p_lo, p_hi, N_gene=8, N_pop=12
12     , mutation_rate = 0.2, selection_rate=0.5):
13         self.N_gene = N_gene
14         self.N_pop = N_pop if N_pop%4 == 0 else N_pop+N_pop%4
15         self.p_lo = p_lo #Nedre graense for bit
16         repraesentation
17         self.p_hi = p_hi #Ovre graense
18         self.N_params = N_params #Antal parametre
19         self.N_keep = int(2*np.ceil(selection_rate*self.N_pop/2)) #Antal
20         kromosomer beholdt, skal vaere lige
21         self.mutation_rate = mutation_rate
22         self.len_chromosomes = int(self.N_gene*self.N_params)
23         self.cost_func = cost_func
24
25         self.population = self.make_random_pop()
26         self.costs = self.eval_cost(self.population)
27         self.rank_prop = self.find_rank_prop()
28         self.N_news = self.N_pop - self.N_keep
29
30     def make_random_pop(self):
31         '''Danner den tilfaeldige inertielle population'''
32         population = np.random.randint(2,size=(self.N_pop,self.N_params,
33         self.N_gene))
34         return np.array([chromosome.flatten() for chromosome in population
35         ])
36
37     def eval_cost(self, eval_population):
38         '''Evaluerer cost-funktionen. Dette indebaerer at konvertere
39         bitrepraesentationen til reelle tal.'''
40         eval_population = np.array([np.split(chromosome,self.N_params) for
41         chromosome in eval_population])
```

```

34     p_quant = np.array([[sum([gene[m-1]*2**(-m) for m in range(1,self.
N_gene+1))]+2**(-(self.N_gene+1)) \
35                         for gene in chromosome] for chromosome in
eval_population])
36     cost_func_params = p_quant*(self.p_hi-self.p_lo)+self.p_lo
37     return np.array([self.cost_func(params) for params in
cost_func_params])
38
39     def find_rank_prop(self):
40         '''Finder sandsynligheder til rank-selection. Kores kun en gang.
Ligning 2.11 i "Haupt & Haupt".'''
41         n_sum = sum(list(range(1,self.N_keep+1)))
42         rand_prop = [(self.N_keep-n+1)/n_sum for n in range(1,self.N_keep+
1)]
43         return np.cumsum(rand_prop)
44
45     def natural_selection(self):
46         '''Udforer en generation, dvs. samler alle funktionerne herunder.
,,,
47         population_to_keep, costs_of_old = self.find_survivors() #costs
retuneres, da den evt. kan bruges i mate
48         mates = self.rank_select(population_to_keep)
49         new_population_members = self.mate_one_cross(population_to_keep,
mates)
50         self.population = np.vstack((population_to_keep,
new_population_members))
51         self.mutate()
52         new_costs = self.eval_cost(self.population[1:])
53         self.costs = np.hstack((costs_of_old[0],new_costs))
54
55     def find_survivors(self):
56         '''Finder de bedste loesninger og sorterer saa bedste er overst'''
57         selection_index = np.argsort(self.costs)
58         sorted_costs = self.costs[selection_index]
59         sorted_population = self.population[selection_index]
60         return sorted_population[:self.N_keep], sorted_costs[:self.N_keep]
61
62     def rank_select(self, pop_to_mate):
63         '''Udvaelger hvilke par, der skal parrre ved rank-selection.
Bedste losninger har storst sandsynlighed. Side 39 i "Haput &
Haupt'''
64         rands = np.random.random((self.N_news//2,2))
65         mates = np.searchsorted(self.rank_prop,rands)
66         dupes_where = mates[:,0] == mates[:,1]
67         not_to_dupe = mates[dupes_where,0]
68

```

```

69     new_mate = [np.random.choice([x for x in range(self.N_genes//2) if
70 x != dupe]) for dupe in not_to_dupe]
71     mates[dupes_where,0] = new_mate
72     return mates
73
74     def mate_one_cross(self, population, mates):
75         '''Avler nye individer. Sker vha. single-point-crossover'''
76         crossover_points = np.random.randint(1, self.N_gene*self.N_params+1
77 ,size=len(mates))
78         new_population_members = np.array([[np.hstack((population[parent[0]
79 ][:point], population[parent[1]][point:])),
80 np.hstack((population[parent[1]][:point
81 ], population[parent[0]][point:]))]
82         for parent, point in zip(mates,
83 crossover_points)])
84         new_population_members = np.vstack((new_population_members[:,0],
85 new_population_members[:,1]))
86         return new_population_members
87
88     def mutate(self):
89         '''Mutation. Sker ved at udvalgte tilfældige bits og invertere
90 dem.'''
91         N_mutations = int(np.ceil(self.mutation_rate*(self.N_pop-1)*self.
92 len_chromosomes))
93         i_chromosome = np.random.randint(1, self.N_pop, size=N_mutations)
94         j_base = np.random.randint(self.len_chromosomes, size=N_mutations)
95         self.population[i_chromosome, j_base] = 1-self.population[
96 i_chromosome, j_base]
97
98 #EKSEMPEL:
99 max_generation = 1000
100 test = GA(cost_func=cost_func, N_params=2, p_lo=-5.12, p_hi=5.12, N_gene=1
101 2, N_pop=8, mutation_rate=0.03)
102 for gen in range(max_generation):
103     test.natural_selection()
104 print(min(test.costs))

```

B Fitting til χ^2

B.1 Fit til Fizeau's eksperiment data

Herunder er koden der er brugt til at udføre minimering af χ^2 til at fitte ligning 4 til dataene. Funktionen `ga` er funktionen fra appendix A1.

```

1  ### DATA ###
2  I = np.loadtxt('Dat_sren.txt')
3  I_err = np.loadtxt('Dat_fejl.txt')
4  pixels = np.arange(len(I))
5
6  ### FIT-FUNKTION ###
7  def fit(x, params):
8      a, b, c, k, delta, sigma, mu = params
9      return np.sin(k*x+delta)**2*(b*np.exp(-(x-mu)**2/(2*sigma**2))\
10         -c*np.exp(-(x-mu)**2/(2*sigma**2)))\
11         +c*np.exp(-(x-mu)**2/(2*sigma**2))+a
12
13  ## CHI^2 ##
14  def chi_square(params):
15      y_model = fit(pixels, params)
16      return sum(((I - y_model)/I_err)**2)
17
18  ## UDFoR FITTING ##
19  p_opt, fitness, function_calls = zip(*list(ga(chi_square, [[0, 15],\
20      [150, 230],\
21      [130, 250],\
22      [0, 1],\
23      [0, 30],\
24      [30, 70],\
25      [250, 330]],\
26      generations=100
27
28      00)))
29
30  ### PLOT ###
31  x = np.linspace(0, len(pixels), 600)
32
33  # Plot data
34  plt.errorbar(pixels, I, I_err, color='tab:blue', capsize=3, label='Data')
35
36  # Plot curve_fit-resultatet
37  cf_params = [1.18144971e+01, 1.22075435e+02, 2.18451351e+02,\
38      5.72908613e-02, -1.61860351e+00, 5.93860975e+01, 3.21019680e+
39      02]
40  plt.plot(x, fit(x, cf_params), color='tab:orange', label='Curve fit')
41
42  # Plot GA-resultat
43  plt.plot(x, fit(x, p_opt[-1]), color='tab:green', label='GA-fit')
44
45  plt.xlim(0, len(pixels))
46  plt.xlabel('Pixel nr.')
47  plt.ylabel('Intensitet')

```

```
45 plt.show()
```