

Universitatea de Științe Agricole și Medicină Veterinară Cluj-Napoca
Școala Doctorală
Facultatea de Horticultură

Lorentz JÄNTSCHI

Rezumat al tezei de doctorat

Algoritmi genetici și aplicații ale acestora

Conducător Științific:
Prof. Univ. Dr. Radu E. SESTRĂȘ

Cluj-Napoca
2010

Cuprins

<u>Introducere</u>	2
<u>Probleme de optimizare a relațiilor structură-activitate</u>	3
<u>Simularea evoluției cu algoritmi genetici</u>	4
<u>Cadrul cercetării, scop și obiective</u>	6
<u>Definirea problemei de optimizare a relației structură-activitate</u>	7
<u>Definirea problemei genetice și proiectarea algoritmului</u>	9
<u>Definirea experimentului de simulare a evoluției</u>	13
<u>Rezultate și discuții</u>	15
<u>Analiza variabilității și diversității</u>	19
<u>Interpretarea rezultatelor din observațiile pe observabile calitative</u>	24
<u>Analiza obiectivului evoluției folosind un eșantion întâmplător de generații</u>	27
<u>Analiza legii de distribuție a obiectivului evoluției folosind un studiu sistematic în cadru generalizat</u>	29
<u>Ce lege urmează momentele de apariție a evoluției?</u>	33
<u>Ce distribuție urmează numărul de evoluții?</u>	36
<u>Concluzii și recomandări</u>	38
<u>Lucrări reprezentative publicate</u>	40

Introducere

Teza “Algoritmi genetici și aplicații ale acestora” este un demers de cercetare fundamentală interdisciplinar, având ca scop simularea evoluției cu algoritmi genetici în probleme de optimizare a relațiilor structură-activitate.

Problematica principală avută în vedere o constituie problemele dificile (cele cu complexitate exponențială) de optimizare a relațiilor cantitative între structura compușilor chimici și activitatea lor biologică.

Modelele informatice elaborate prin intermediul algoritmilor genetici emulează modelele biologice evoluționiste, asigurând rezolvarea unor probleme concrete de optimizare sau căutare în experiențele de genetică și ameliorare a plantelor. Prin intermediul elementelor individuale, reprezentate sub forma șirurilor binare, și a operatorilor de natură biologică definiți asupra populației și a modelului molecular, algoritmi genetici manipulează cele mai promițătoare șiruri, evaluate conform unei funcții obiectiv, căutând soluții mai bune, tinzând în esență spre soluția “optimă”, dar acceptând în final una apropiată de optim.

Prezenta lucrare demonstrează că algoritmi genetici, ca tehnici adaptive de căutare euristică, bazate pe principiile geneticii și selecției naturale, pot fi eficient utilizați în simularea procesului biologic al evoluției și în cel de ameliorare a plantelor. În acest sens, a fost conceput cadrul necesar pentru construcția și aplicarea unui algoritm genetic care să rezolve problemele de optimizare, prin elaborarea unui algoritm genetic într-un cadru definit. Respectivul algoritm genetic a fost implementat într-un program evolutiv și aplicat pe un set de date experimentale, fiindu-i urmărită evoluția.

S-a realizat un design experimental cu scopul de a face trecerea de la problema de optimizare la o problemă de simulare, și anume simularea evoluției în diferite strategii de selecție și de supraviețuire. A fost creată o contingentă de 3x3 strategii distincte pentru selecție și supraviețuire (proporțional, în turnir și deterministic) și s-a urmărit evoluția pe parcursul a 20000 de generații în mod repetat de 46 de ori în fiecare strategie în parte. S-au analizat inferențele statistice în observabile calitative și cantitative ale procesului de evoluție controlată de diferitele strategii de evoluție, cu ajutorul diferitelor variabile pe care programul evolutiv a fost configurat să le înregistreze.

Informațiile și rezultatele obținute au în mare măsură un caracter fundamental. Analiza statistică a rezultatelor obținute din simularea proceselor de evoluție a permis obținerea unor răspunsuri la întrebări precum: *Care este legea de distribuție a obiectivului evoluției? Care este legea de distribuție a momentelor evoluției? Care este legea de distribuție a numărului de evoluții? Cum este influențată variabilitatea și diversitatea genotipică de alegerea unei strategii de evoluție? Cât de timpuriu se produc evoluțiile în raport cu strategia de evoluție aleasă? Cât de frecvent se produc evoluțiile în raport cu strategia de evoluție aleasă? Cât de dispers (și respectiv cât de predictibil) se produc evoluțiile în raport cu strategia de evoluție aleasă? Care sunt similaritățile și deosebirile între evoluțiile care au loc urmând diferite strategii?*

Au fost obținute și o serie de rezultate cu caracter aplicativ, cum sunt: implementarea algoritmului genetic într-un program evolutiv capabil să rezolve o problemă dificilă de optimizare a unei relații structură-activitate folosind familii de descriptori de structură; implementarea unor module de calcul automate pentru optimizarea geometriei moleculare; implementarea unor programe pentru calculul statisticii Anderson-Darling de agrement între model și observație; implementarea procedurii Grubbs de identificare și eliminare a observațiilor în eroare față de un model.

Teza oferă și soluții de transfer tehnologic, cuprinzând răspunsuri la o serie de probleme de optimizare în domeniul horticul în care evoluția către un obiectiv de ameliorare sau de planificare este influențată de o serie de parametri specifici materialului genetic și/sau ai arealului și în care obiectivul fixat este influențat de strategia aleasă; în acest cadru experimentul de simulare a evoluției, teza oferă soluții evidențiate statistic cu privire la influența strategiilor de selecție și supraviețuire.

Probleme de optimizare a relațiilor structură-activitate

Scurt istoric

Abordarea matematică a SAR (Relații Structură-Activitate) pentru BAC (Compuși Biologic Activi), începută în secolul nouăsprezece, s-a concretizat prin apariția conceptului de relații cantitative structură-activitate (QSAR = Quantitative Structure-Activity Relationships), metodă matematică care permite identificarea legăturii cantitative dintre structura chimică și activitatea biologică a compușilor investigați - ([Hammett, 1935](#)). Observații de SAR au fost publicate în literatura de specialitate încă din 1868, când Crum-Brown & Fraser au stipulat ideea că activitatea compușilor este o funcție a structurii și compoziției chimice ([Crum-Brown & Fraser, 1868](#)), însă au trecut aproape patruzeci de ani de când paradigma relații cantitative structură-activitate (QSAR) și-a dovedit utilitatea practică în agrochimie, chimie farmaceutică, toxicologie etc. ([Hansch & Leo, 1979](#)).

Ce sunt relațiile structură-activitate

Activitatea biologică sau bioactivitatea este termenul comun pentru efectul benefic sau advers al unui compus (sau amestec de compuși chimici) asupra materiei vii.

Manifestarea și cuantificarea calitativă și/sau cantitativă a activității biologice a unui anumit compus chimic este un proces extrem de complex prin natura foarte variată a efectelor pe care un compus chimic le poate avea asupra diferitelor organisme vii. Procedura de determinare a activității acestuia asupra organismului este standardizată (ex. Schema ADMET, Tabelul 1 - [Anexa 2-1](#) din Teză). Activitatea biologică este evaluată prin intermediul unor proceduri specifice, supuse standardizării (ex. Tabelul 2 - [Anexa 2-1](#) din Teză redă două astfel de activități biologice). O serie de procese biologice sunt referite distinct atunci când se exprimă activitatea biologică a unui compus chimic (ex. Tabelul 3 - [Anexa 2-1](#) din Teză redă definiții specifice mediului acvatic). Standardizarea care se referă la determinarea toxicității acute ([UNE-CE-4, 2009](#)) stabilește că aceasta trebuie determinată folosind una din metodele prezentate în Tabelul 4 - [Anexa 2-1](#) din Teză.

Compușii biologic activi (BAC) au o largă utilizare în domeniul agricol și horticol. Tabelul 5 - [Anexa 2-1](#) din Teză redă o clasificare a regulatorilor de creștere, în acord cu [Societatea Americană de Științe Horticole](#).

Seriile de compuși sunt alcătuite din compuși (congeneri) înrudiți, atât în ceea ce privește structura, cât și proprietățile fizico-chimice și/sau activitățile biologice. Atunci când se supune observației o serie de compuși, se pornește de la ipoteza că aceștia au în comun atât elemente de structură, cât și elemente de proprietate/activitate ce fac ca valorile acestora să fie relativ apropiate.

Pentru ca rezultatele observației să capete consistență în ceea ce privește interpretarea statistică, trebuie asumată și ipoteza de convergență la normalitate asupra valorilor observate în eșantionul seriei de compuși, spațiul complet al acestora fiind în acest caz un exemplu tipic de populație finită distribuită normal.

În aceste ipoteze, de înrudire a compușilor atât sub aspect structural, cât și sub aspect al proprietății/activității măsurate, și de distribuire normală a valorilor observate, se pot formula și verifica (cu ajutorul testelor statistice) ipoteze de inferență (dependență) între structură și activitatea/prorietatea măsurată. Relațiile structură-activitate (SAR) și respectiv relațiile structură-prorietate (SPR) stabilesc legături funcționale între structura compușilor chimici și proprietățile măsurate de natură biologică (SAR) și fizico-chimică (SPR) ale acestora.

Relațiile cantitative (q) care se stabilesc între structură și activitate (qSAR) sau respectiv proprietate (qSPR) se exprimă prin intermediul unor ecuații care au un domeniu de aplicabilitate definit cel mai frecvent de seria de compuși pe care au fost obținute și de proprietatea sau activitatea supusă observației.

Elaborarea și valorificarea de relații structură-activitate

Fluxul de informații de specialitate face numeroase referiri la metodologia de obținere a noilor compuși biologic activi. În acest sens, monografia Diudea ([Diudea & alții, 2001](#)) este cuprinzătoare.

Simularea evoluției cu algoritmi genetici

Scurt istoric

“Hard inheritance” (“Moștenirea dură”) ([Weismann, 1893](#)) și “Soft inheritance” (“Moștenirea ușoară”) ([Lamarck, 1809](#)), selecția și supraviețuirea ([Darwin, 1859](#)), genele și recombinarea genică ([Morgan & alții, 1915](#)), transmiterea caracterelor ([Mendel, 1866](#)) constituie problematici îndelung dezbătute și disputate de-a lungul secolului al XIX-lea ([Fisher, 1954](#)), toate contribuind la fundamentarea geneticii moderne de azi ([Ayala & alții, 1994](#)), și oferind sursele de inspirație ale algoritmilor genetici.

Primele simulări ale evoluției se regăsesc în studiile lui Nils Aall BARRICELLI ([Barricelli, 1954](#)). Puțin mai târziu, Alex FRASER (1923-2002) a publicat o serie de lucrări despre simularea selecției artificiale a organismelor cu locuși multipli ce controlează o trăsătură măsurabilă. Simulările lui FRASER ([Fraser, 1957-1970](#)) includ toate elementele esențiale ale algoritmilor genetici moderni.

În ce situații sunt aplicabili algoritmi genetici

Instrumentul de dezvoltare a algoritmilor genetici îl constituie informatica. Astfel, uzual, în viața de zi cu zi, și la fel în cercetarea științifică, se operează cu *probleme*. În informatică și ramurile derivate ale acesteia (cum e cazul bio-informaticii și chemo-informaticii) o problemă are o semnificație precisă, foarte apropiată de cea ilustrată de *algoritm*. Un algoritm este în esență o rețetă, specificând ce trebuie făcut în anumite condiții, pentru a obține un anumit obiectiv. Un algoritm necesită două resurse pentru a *rezolva* o problemă, și anume: timp (cu sensul de timp de execuție, mărime corelată cu numărul de instrucțiuni elementare) și spațiu (pentru stocarea datelor de intrare și a variabilelor). Nu toate problemele sunt de aceeași *complexitate*, și același lucru este valabil și pentru algoritmi de rezolvare. Astfel, unele probleme au complexitate exponențială, ceea ce înseamnă că cel mai bun algoritm rezolvă problema într-un timp de execuție ce crește exponențial în funcție de dimensiunea (volumul, mărimea) datelor de intrare. Acest tip de probleme sunt numite *dificile*, deoarece chiar și cel mai bun algoritm (care există, sau ar putea exista) va fi probabil nepractic cu date de intrare din practică ([Falkenauer, 1998](#)). Dacă o problemă este dificilă, atunci căutarea *optimului* frecvent iese în afara timpului disponibil pentru aplicațiile reale. Chiar dacă există această problemă, există totuși o serie de probleme întâlnite în practică când obținerea optimului nu este necesară (obligatorie). De cele mai multe ori, o *soluție bună* este suficientă.

Ce sunt algoritmi genetici

Deoarece întotdeauna cercetătorii s-au confruntat cu mai multe probleme dificile, de foarte mult timp s-a încercat rezolvarea acestora, unul sau mai mulți *euristici* fiind de-a lungul anilor concepuți în acest sens. Aceștia sunt seturi de reguli gândite pentru a rezolva o problemă anume, uzual bazați pe bunul simț (în ceea ce privește soluția așteptată) prin evitarea erorilor grosolane, dar care nu sunt gândiți pentru a produce întotdeauna soluția cu exactitate și, respectiv, să fie capabili să producă o soluție pentru orice valori de intrare. Chiar dacă cei mai mulți euristici sunt foarte mult ad-hoc și dependenți de problema dată, odată cu dezvoltarea informaticii, cercetătorii au reușit să formuleze trei euristici care sunt foarte generali, și anume aplicabili la o mare varietate de probleme dificile. Din cauza generalității pe care o presupun, aceștia au căpătat numele de *meta-euristici*. Toți trei sunt stocastici în natura lor (*A fi stocastic: Implicând sau conținând una sau mai multe variabile aleatoare, implicând șansa sau probabilitatea*), doi dintre aceștia (SA și GA) fiind bazați pe procese naturale care au loc în jurul nostru din totdeauna. Împreună cu *călirea simulată* (SA - Simulated Annealing) și *căutarea tabu* (TS - Tabu Search) sunt și *algoritmi genetici* (GA - Genetic Algorithm). Chiar dacă primele studii în care au apărut algoritmi genetici se situează în anul 1954 ([Barricelli, 1954](#)), studii de amploare ale acestora au apărut după 1970 ([Bosworth & alții, 1972](#); [Holland, 1975](#)), ei fiind re-invențați ceva mai târziu ([Davis, 1991](#); [Holland, 1992](#)) odată cu dezvoltarea tehnicii de calcul.

Complexitatea algoritmică

O problemă importantă legată de *complexitatea algoritmică* este reprezentată de teorema *inexistenței mesei pe gratis* (NFLT - No Free Lunch Theorem; Wolpert & Macready, [1995](#) & [1997](#);

(English, 1996), teoremă care, utilizând trei criterii de evaluare a calității unui algoritm (viteză, precizie și scop) sugerează că toți algoritmi sunt strict echivalenți. În esență, aceasta înseamnă că pentru doi algoritmi A și B, pentru fiecare set de date pentru care A performează mai bine decât B, există un set de date pentru care B performează mai bine decât A.

Construcția algoritmilor genetici

÷ Se operează asupra unei populații de reprezentări abstracte numite (după elementele genetice pe baza cărora au fost imaginate) cromozomi sau genotipuri ale unui **genom**; la rândul său, fiecare reprezentare abstractă a unui **cromozom** este compusă din **gene**.

÷ Fiecare **generație** este compusă dintr-o populație de șiruri de caractere (sau alte forme de reprezentare abstractă) analog cu cromozomii ADN-ului. Fiecare element al populației reprezintă un punct în spațiul de căutare și în același timp o soluție posibilă.

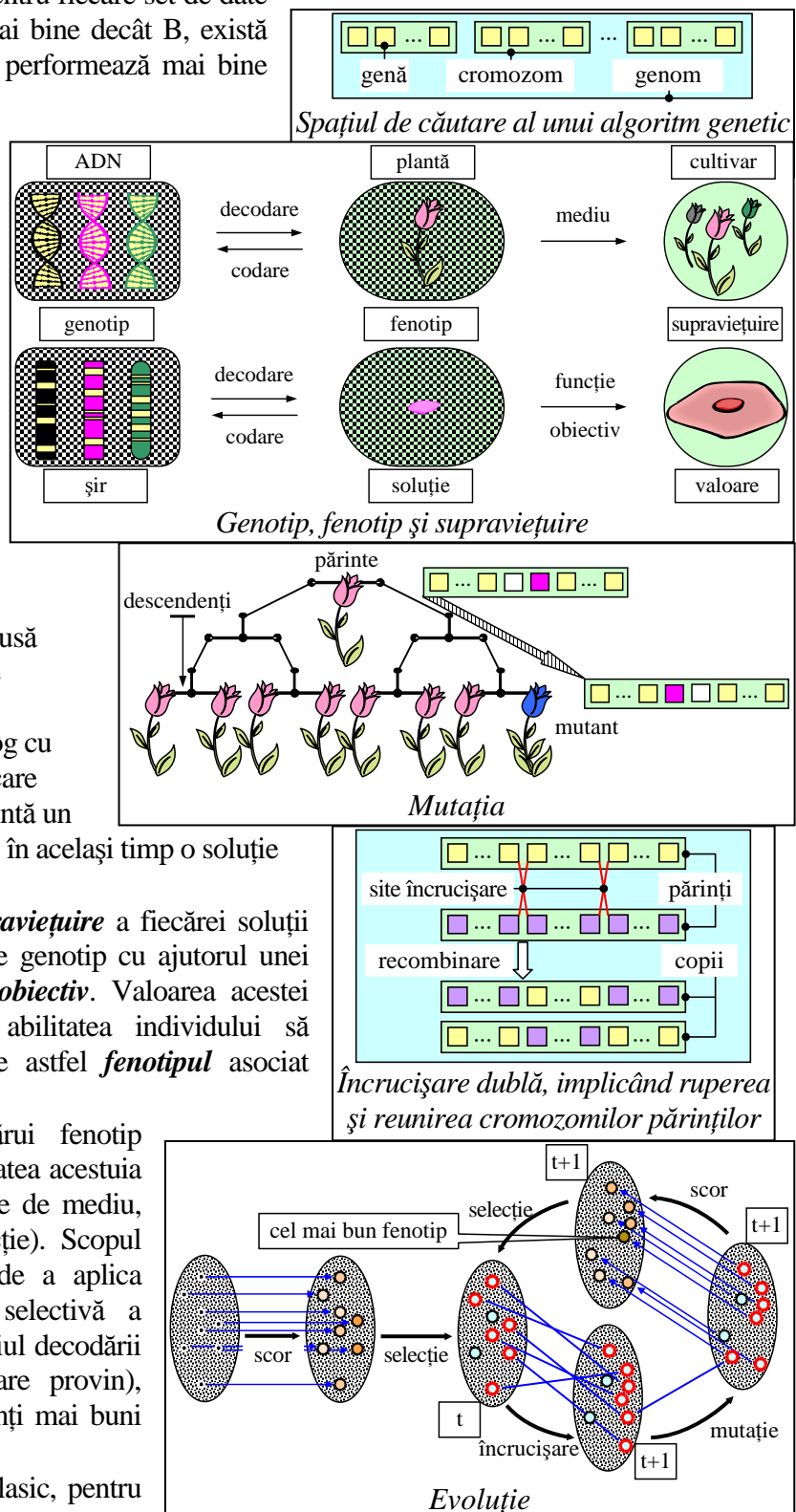
÷ Un scor sau **șansă de supraviețuire** a fiecărei soluții este calculată pentru fiecare genotip cu ajutorul unei funcții, numită și **funcție obiectiv**. Valoarea acestei funcții este asociată cu abilitatea individului să supraviețuiască și definește astfel **fenotipul** asociat genotipului.

÷ Scorul este asociat fiecărui fenotip (soluție) reprezentând abilitatea acestuia de a concura pentru resurse de mediu, pentru **supraviețuire** (selecție). Scopul algoritmului genetic este de a aplica **încrucișarea** și **mutația** selectivă a fenotipurilor (prin intermediul decodării lor în genotipurile din care provin), pentru a produce descendenți mai buni decât părinții lor.

÷ Într-un algoritm genetic clasic, pentru a rezolva o problemă, se generează întâmplător sau se inițiază cu valori predefinite o populație de un volum dat de genotipuri și evoluția se iterează prin repetiția selecției, mutației și încrucișării până când cel mai bun fenotip al populației satisface o condiție impusă (condiție care reprezintă condiția de sfârșit a algoritmului).

Elaborarea și valorificarea algoritmilor genetici

Algoritmii genetici servesc în clasificarea filogenetică (Jäntschi & alții, 2008-PTA), analiza



secvențelor de gene ([Jäntschi & alții, 2009-GSA](#)), probleme dificile de dinamica proceselor ([Jäntschi & alții, 2009-DPA](#)) și în orice altă categorie de probleme dificile de decizie, clasificare, optimizare sau simulare ([Falkenauer, 1998](#)).

Cadrul cercetării

Dezvoltarea continuă a depozitelor de cunoștințe de genul celor administrate de NIH, cum sunt PubMed, PubChem, Genome etc., accentuează necesitatea de a poseda instrumente eficiente de a relaționa aceste cunoștințe, iar relațiile structură-activitate reprezintă unul dintre aceste instrumente. Problema de simulare aleasă în studiu, și anume a evoluției (prin intermediul diferiților parametri ce caracterizează eșantionul supus evoluției) este o problemă insuficient explorată în literatura de specialitate al cărui subiect sunt algoritmi genetici.

Studii asupra altor operatori esențiali pentru evoluție sunt axate pe eficiența algoritmică (viteza cu care se atinge obiectivul și apropierea de maximul global). O colecție de lucrări de acest tip este reprezentativă în acest sens ([Martin & Spears, 2001](#)). Astfel subiectul îl constituie diferiții operatori de încrucișare ([Prügel-Bennett, 2001](#)), mutația și încrucișarea ([Spears, 2001](#)), sau alți parametri dinamici ([Droste & alții, 2001](#)).

Studiile sunt adesea concentrate spre rezolvarea problemelor dificile cu ajutorul algoritmilor genetici, uneori abordându-se direcționat eficiența acestora (ca timp de execuție, resurse de memorie necesare), dar foarte puțin influența diferitelor strategii de evoluție asupra obiectivului urmărit. În acest din urmă caz, se are în vedere în special eficiența algoritmului, și aproape niciodată parametrii ce caracterizează eșantionul supus evoluției.

Datorită potențialului de valorificare a rezultatelor pe care îl au, algoritmi genetici au depășit demult granițele domeniului informatică. Teze de doctorat având ca obiectiv proiectarea de algoritmi genetici, implementarea de programe evolutive și realizarea de studii cu ajutorul lor, se regăsesc practic în toate domeniile de cercetare. Astfel, în domeniul agricultură și-au găsit utilizarea la planificarea culturilor ([Matthews & Kraw, 2001](#)), evaluarea riscului de eroziune a solului ([Osman & McManus, 2007](#)), în bioinginerie la controlul eficient al poluării la nivelul unui bazin hidrografic ([Veith & Wolfe, 2002](#)), în chimie la designul proceselor controlate senzorial ([Dai & Lodder, 2007](#)), în economie la probleme de optimizare cu opțiuni multiple ([Aickelin & Dowsland, 1999](#)), în management la modelarea proceselor multi-scală ([Sastry & alții, 2007](#)), în mecanică la optimizarea structurilor compozite ([Gantovnik & Gürdal, 2005](#)) și în mediu la alegerea strategiei pentru controlul calității apei ([Tufail & Ormsbee, 2006](#)). În domeniul biologie, se desprind două direcții principale în ceea ce privește elaborarea și utilizarea algoritmilor genetici: în probleme de evoluție ([Suzuki & Iwasa, 1998](#)) și în studii filogenetice ([Zwickl & Hills, 2006](#)). În privința caracterului practic, de utilizare a algoritmilor genetici în domeniul agricol și horticol, algoritmi genetici au o largă aplicabilitate, de la studii de creștere ([Venard & Vaillancourt, 2006](#)), la clasificări taxonomice ([Sarmiento-Monroy & Sharkey, 2006](#)) sau analiza diversității genetice ([Zhang & Ghabrial, 2006](#)).

Scop și obiective

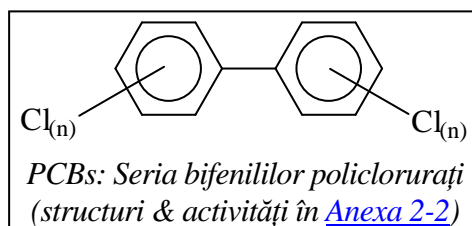
Scop: Simularea evoluției cu algoritmi genetici în probleme de optimizare a relațiilor structură-activitate. Proiectarea, implementarea și analiza statistică a influenței diferitelor metode de selecție și de supraviețuire asupra evoluției unui algoritm genetic utilizat pentru obținerea de relații structură-activitate în serii de compuși chimici biologic activi.

Obiective:

1. Elaborarea algoritmului genetic (Definirea problemei dificile, cu complexitate exponențială; Formularea problemei în termeni genetici; Proiectarea algoritmului genetic; Implementarea și documentarea programului evolutiv bazat pe algoritmul genetic);
2. Simularea evoluției (Evoluție: definirea observabilelor; Definirea contingenței selecție vs. supraviețuire; Proiectarea experimentului statistic; Realizarea experimentului statistic);
3. Analiza și interpretarea rezultatelor: Observabile calitative; Obiectivul evoluției - r^2 - observabilă cantitativă);

Definirea problemei de optimizare a relației structură-activitate

Setul de molecule ales pentru studiu este seria bifenililor policlorurați (PCBs), formată din 209 compuși, a cărui studiu este de o deosebită importanță pentru impactul acestora asupra ecosistemului.



Chiar dacă producția de PCB a fost stopată în 1970 datorită toxicității ridicate a celor mai mulți congeneri ai familiei PCB, efectele acestora sunt încă prezente în mediu, datorită faptului că PCBs sunt poluanți organici persistenti - clasificați ca atare, care se acumulează plante și animale. Coeficientul de partiție octanol/apă (K_{ow}) este raportul concentrațiilor unui compus chimic între octanol și apă aflate în contact la o anumită temperatură. Este un parametru adimensional (fiind un raport de concentrații) care frecvent se exprimă pe scară logaritmică ($\log K_{ow}$ sau mai simplu l_{kow}). Această proprietate fizico-chimică este utilizată în multe studii de mediu în determinarea efectului compușilor chimici în mediu, un exemplu fiind utilizarea acestuia pentru prezicerea magnitudinii de bioacumulare în pești ([U.S. Geological Survey, 2008](#)).

O serie de statistici au fost aplicate pentru verificarea ipotezei de normalitate pe seria de activități biologice observate (pe 206 din 209 compuși).

Ipoteza de normalitate a l_{kow} măsurate, 206 compuși (descrierea statisticilor în [Anexa 3-3](#) din Teză)

Statistică	Valoare	Probabilitate de observație	Concluzie
KS	0.03348	96.91%	Nu se respinge ipoteza de normalitate
AD	0.44432	27.2%; 25.2%; 19.2%	Nu se respinge ipoteza de normalitate
CS	11(df=7)	13.8%	Nu se respinge ipoteza de normalitate
WS	0.98709	5.8%	Nu se respinge ipoteza de normalitate
$Z_{Skewness}$	1.48	14%	Nu se respinge ipoteza de normalitate
$Z_{Kurtosis}$	2.51	1.2%	Se respinge ipoteza de normalitate
JB	7.577	2.3%	Se respinge ipoteza de normalitate

Nota: testul JB nu este afectat de valori pereche (tied); KS și AD sunt afectate

Ipoteza de normalitate a l_{kow} măsurate, 205 compuși (206\{PCB209(l_{kow} =9.603)\})

Statistică	Valoare	Probabilitate de observație	Concluzie
KS	0.03579	94.68%	Nu se respinge ipoteza de normalitate
AD	0.37878	40.3%; 39.5%; 21.0%	Nu se respinge ipoteza de normalitate
CS	8.64(df=7)	27.9%	Nu se respinge ipoteza de normalitate
WS	0.98709	47.8%	Nu se respinge ipoteza de normalitate
$Z_{Skewness}$	1.48	79.2%	Nu se respinge ipoteza de normalitate
$Z_{Kurtosis}$	2.51	41.5%	Nu se respinge ipoteza de normalitate
JB	0.56146	75.5%	Nu se respinge ipoteza de normalitate

Valoarea eliminată (9.603) a fost identificată folosind procedura Grubbs elaborată pe baza testului Grubbs ([Grubbs, 1969](#)). În setul format de cei 205 compuși, se remarcă nu numai agrementul între măsurat și model (Gauss) dar și agrementul între statistici.

Modelarea structurii moleculare este condiția obligatorie pentru o analiză structură-activitate. Realizarea unui model tridimensional (3D) se impune în situațiile în care descriptorii calculați izează de geometria moleculară, aspect valabil pentru cazul prezent. Obținerea modelului 3D se poate realiza folosind un program de modelare moleculară cum ar fi: HyperChem, Spartan, Gaussian, Molecular Modelling Pro, Mopac. În modelarea structurii PCBs s-a folosit programul de modelare moleculară HyperChem (licență [v. 8.0/2007](#)). Obținerea informației de structură 3D presupune parcurgerea unei serii de pași, care cuprinde definirea unui model de mecanică moleculară (a fost definit AMBER), optimizarea geometriei moleculare până la convergență folosind un algoritm de optimizare (a fost folosit POLAK-RIBIERE), definirea unei metode semi-empirice de calcul energetic (a fost definită AM1) și obținerea unei serii de parametrii energetici ce

caracterizează optimul împreună cu sarcinile electrice parțiale.

Pentru ca problema de optimizare aleasă să poată fi asociată unei probleme de evoluție cu ajutorul unui algoritm genetic, este obligatoriu ca descriptorii moleculari să provină dintr-o populație care posedă un cod genetic (de aici denumirea de familie). Cu alte cuvinte, posibilitatea de a asocia un cod genetic populației de descriptori este cea care creează oportunitatea execuției unui algoritm genetic.

Algoritmul genetic a fost elaborat având ca model patru familii de descriptori de concepție proprie. Descrierea detaliată a familiilor de descriptori moleculari este redată în [Anexa 4](#) din Teză. Pentru analiza structură-activitate implicând PCBs și activitatea măsurată a acestora Ikw s-a ales familia MDF (în [Anexa 4-2](#) din Teză) să definească informația structurală în asociere cu activitatea biologică, această familie fiind cea mai bine documentată și susținută de rezultate (cu ajutorul acesteia au fost analizați peste 50 de seturi de compuși).

Relația structură-activitate căutată este compusă din patru descriptori membrii MDF și definește o problemă dificilă de optimizare datorită volumului populației MDF de descriptori (787968 membri).

Definirea problemei genetice și proiectarea algoritmului

Fiecare *genă* [1] codifică câte un operator folosit în construcția *cromozomului* [2] unui descriptor molecular ([Tabelul 5-2](#) în Teză). Fiecare descriptor al unei familii de descriptori este un *genotip* [3] și toți împreună constituie *materialul genetic* [4] al familiei respective.

Secvență din Tabelul 5-2. Spațiul de căutare al MDF

Familie	Gene	Genom
MDF (©2005)	D _M	t g
	A _P	C H M E G Q
	I _D	D d O o P p Q q J j K k L l V E W w F f S s T t
	I _M	r R m M d D
	F _C	m M D P
	S _M	m M n N S A a B b P G g F f s H h I i
L _O	I i A a L l	

Numărul de valori pe care le codifică fiecare genă variază de la 2 valori (în cazul genei ce codifică tipul de metrică de distanță - topologică și geometrică - D_M pentru FPIF și MDF și D_O pentru MDFV și SAPF) până la 58 de valori în cazul descriptorului de interacțiune I_D al familiei MDFV. Volumul materialului genetic variază, [Tabelul 5-3](#) sumarizând aceste variații:

Tabelul 5-3. Volumele familiilor de descriptori moleculari

Familie	Gene								Volum (N)
FPIF (©2000)	I _M :2	D _M :2	A _P :4	P _D :8	F _C :6	S _M :5	M _I :4	L _O :3	46080
MDF (©2005)	D _M :2	A _P :6	I _D :6	I _M :24	F _C :4	S _M :19	L _O :6		787968
MDFV (©2008)	D _O :2	A _P :7	I _D :58	S _F :7	S _M :7	I _T :10	E _U :2	L _O :3	2387280
SAPF (©2009)	C _F :3	D _O :2	A _P :5	D _P :6	P _P :6	O _M :2	M _P :6	L _O :6	77760

Metodologia de lucru a algoritmilor genetici presupune prelevarea inițială (întâmplător sau deterministic) a unui *eșantion* [5] de cromozomi din materialul genetic format dintr-un șir de descriptori X₁, ..., X_p care este supus procesului de *evoluție* [6] în *cultivar* [7]. *Algoritmul genetic* [8] operează astfel asupra eșantionului care suferă modificări în fiecare *generație* [9]. Fiecare mulțime de `n` descriptori distincți reprezintă un punct în *spațiul de căutare* [10] și în același timp o *soluție posibilă* [11]. Operatorii de bază ai unui algoritm genetic sunt încrucișarea și mutația. *Încrucișarea* [12] a două genotipuri presupune alegerea unei porțiuni de încrucișat de-a lungul

[1] genă = una din valorile de pe coloana Gene a Tabelului 36; ex. I_M pentru FPIF

[2] cromozom = secvența de gene a unei familii în Tabelul 36; ex. D_MA_PI_DI_MF_CS_ML_O pentru MDF

[3] genotip = o concretizare posibilă a valorilor fiecărei gene a unui cromozom; ex. TCJtAAfDI pentru MDFV

[4] material genetic = mulțimea tuturor combinațiilor posibile de valori de pe coloana Genom în tabelul 36; ex. {D, P, C} × {T, G} × {C, H, M, E, A} × {I, E, H, G, A, Q, S} × {S, M} × {I, E, H, G, A, Q, S} × {I, E, H, G, A, Q, S} × {I, A, S, T, Q, R, L} pentru SAPF

[5] eșantion = submulțime a materialului genetic al familiei de descriptori moleculari; ex. {DTCIISII, DTCIESII, DTCGISII} reprezintă un eșantion de volum 3 al materialului genetic al SAPF

[6] evoluție = proces genetic complex care implică procese de selecție, încrucișare și mutație

[7] cultivar = spațiu (de memorie, virtual) în care genotipurile sunt transformate în fenotipuri prin aplicarea operatorilor definiți de valorile genelor pentru întreg setul de `m` molecule supus studiului; fenotipul asociat genotipului este astfel un șir de `m` valori numerice (câte una pentru fiecare moleculă a setului)

[8] algoritm genetic = algoritm care descrie prin instrucțiuni procesul de evoluție asupra eșantionului

[9] generație = una din iterațiile algoritmului genetic

[10] spațiul de căutare = mulțimea posibilităților de selecție a `n` descriptori din `V_S` posibilități (relația 20)

[11] soluție posibilă = o ecuație de regresie cu `n` descriptori distincți definită de relația (1) sau (2)

[12] încrucișarea = procesul prin care o porțiune a materialului genetic al unui cromozom este înlocuită de porțiunea corespunzătoare a altui cromozom și viceversa; încrucișarea este făcută în speranța că dacă se recombina porțiuni de genotipuri de succes, atunci acest proces este probabil să producă descendenți chiar mai buni decât părinții din care provin

șirului de gene (întâmplător sau deterministic), valorile celor două porțiuni de șiruri fiind schimbate între ele la descendenți. **Mutația** [13] unui genotip presupune modificarea unei valori a unei gene a cromozomului cu o altă valoare din lista valorilor posibile pentru gena respectivă. Rezultatul încrucișării și al mutației este obținerea de **descendenți** [14] sau urmași, cu genotipuri diferite. **Selecția** [15] genotipurilor este operația prealabilă necesară încrucișării și mutației și se face pe baza unui **scor de selecție** [16]. Cel puțin o parte a descendenților reprezintă descriptori **viabili** [17] putând face atunci parte din soluții candidate ale generațiilor următoare. Descendenții viabili înlocuiesc o parte corespunzătoare a indivizilor din eșantion în urma unui proces de **supraviețuire** [18] aplicat indivizilor din eșantion pe baza unui **scor de supraviețuire** [19]. Un alt parametru al algoritmului genetic îl reprezintă **obiectivul evoluției** [20] care este urmărit pe baza unei **funcții obiectiv** [21]. Urmărirea obiectivului evoluției se face odată la fiecare generație selectând din eșantion acei indivizi care maximizează sau, după caz, minimizează valoarea funcției obiectiv (acei indivizi care fac parte din cea mai bună ecuație de regresie obținută pe baza indivizilor din cultivar). Se poate opta ca indivizii care ating obiectivul evoluției într-o generație să fie păstrați în eșantion, caz în care acestora nu li se mai aplică procesul de supraviețuire, ei fiind automat incluși în eșantionul generației următoare.

Așa cum rezultă din aplicarea procesului de supraviețuire, nu toți indivizii unei generații supraviețuiesc și sunt incluși în generația următoare. Motivul acestui fapt este păstrarea unui număr constant de genotipuri în eșantionul dezvoltat în cultivar, astfel încât numărul de indivizi înlocuiți din eșantion este egal cu numărul de descendenți viabili obținuți în urma selecției, încrucișării și mutației.

Întrucât selecția și supraviețuirea au la bază scorurile de selecție și de supraviețuire, se realizează folosind o **modalitate de selecție și supraviețuire** [22].

[13] Mutație = operatorul care introduce modificări noi (inexistente în eșantionul unei generații); ceea ce este caracteristic în general mutației și implicit și operatorului acesteia corespondent în algoritmi genetici este că ea se petrece cu o probabilitate scăzută, fiind deci aplicată cu o probabilitate scăzută

[14] Descendenți = genotipurile obținute din încrucișarea și eventual mutația indivizilor din eșantion

[15] Selecție = operatorul cu ajutorul căruia se extrage din eșantion mai mulți indivizi care participă la înmulțire

[16] Scor de selecție = valoare numerică asociată individului din eșantion calculată pe baza (sau exprimată din) tăriei fenotipului în cultivar

[17] Viabilitatea (unui descriptor molecular) = referă potențialul acestuia de a fi folosit în regresii; un descriptor este viabil dacă (ceea ce urmează referă fenotipul acestuia, viabilitatea referind astfel manifestarea fenotipică) cel puțin are valori reale și finite pentru toate moleculele din set și nu are toate valorile identice; suplimentar i se pot impune și alte condiții, cum ar fi o variabilitate rezonabilă (prin intermediul unui coeficient de variație), o abatere de la normalitate rezonabilă (prin intermediul unui test de normalitate cum este Jarque-Bera) și o capacitate de explicare a proprietății măsurate rezonabilă (prin intermediul coeficientului de determinare din regresia liniară simplă cu proprietatea măsurată).

[18] Supraviețuire = operatorul cu ajutorul căruia se extrage din eșantion mai mulți indivizi care vor fi înlocuiți în eșantion de către descendenți

[19] Scor de supraviețuire = valoare numerică asociată individului din eșantion care poate fi o valoare obținută atât pe baza caracteristicilor genotipice ale individului (asociind o măsură a similarității acestuia cu alte genotipuri din cultivar în scopul menținerii diversității materialului genetic) cât și pe baza trăsăturilor fenotipice ale individului (asociind o măsură a similarității acestuia cu alte fenotipuri din cultivar în scopul menținerii diversității fenotipice)

[20] Obiectivul evoluției = parametrul sau caracteristica (unei ecuații de regresie) care constituie subiectul optimizării (minimizare - ex. suma pătratelor diferențelor între (erorilor) valoarea măsurată și cea explicată de model; maximizare - ex. coeficientul de determinare)

[21] Funcția obiectiv = algoritmul (procedura) de calcul al parametrului sau caracteristicii ce constituie obiectivul evoluției eșantionului.

[22] modalitate de selecție și modalitate de supraviețuire = metodă de extragere a unui individ din eșantion ce folosește drept parametru valorile scorurilor (de selecție și respectiv de supraviețuire) ale indivizilor ce compun eșantionul; ceea ce Tabelul 38 prezintă în mod formal exprimă faptul că se folosesc trei alternative de selecție (proporțional, deterministic și turnir) care se pot aplica valorilor scorurilor sau rangurilor scorurilor (când selecția e bazată pe valorile rangurilor în locul valorilor scorurilor); de asemenea, valoarea scorului poate fi supusă unui proces de normalizare care corectează (ajustează relativ) scorurile indivizilor din cultivar în raport cu două valori (una de minim și alta de maxim) care se actualizează global în fiecare generație pe parcursul întregii evoluții a eșantionului în cultivar

÷ Până când (condiție care reprezintă condiția de sfârșit a algoritmului)

- Se satisface o condiție impusă cu privire la valoarea funcției obiectiv (atinge o valoare impusă) sau se realizează un număr dat de iterații (evoluții).

Programul evolutiv [24] rezultat a fost gândit astfel încât să poată lucra cu oricare familie de descriptori moleculari (FPIF, MDF, MDFV, SAPF date în [Anexa 4](#) din Teză) și să poată fi parametrizat; soluția de implementare aleasă este crearea și utilizarea de fișiere de configurare. Prezentarea și documentarea programului evolutiv ce implementează algoritmul genetic realizat este redată în [Anexa 5](#) din Teză.

Procesul de optimizare a relației SAR sub formă de MLR cu 4 variabile MDF pentru predicția l_{kw} a PCBs are ca scop obținerea unei cât mai bune relații, care să posede bune capacități de estimare (reziduuri teoretic vs. experimental cât mai mici) corectat/reformulat și de predicție (pentru molecule care aparțin de aceeași clasă de compuși, dar care nu au fost incluși în analiză).

Identificarea celui mai bun instrument cu ajutorul căruia se poate obține acest lucru??? este o discuție deschisă în literatura de specialitate.

Se poate alege să se minimizeze reziduurile, să se maximizeze determinarea și lista de posibilități poate continua (de exemplu, suma pătratelor reziduurilor poate fi înlocuită cu suma modulelor reziduurilor, determinarea în estimare poate fi înlocuită cu determinarea în validarea încrucișată, corelația Pearson cu corelația Spearman, ș.a.m.d.) - ([Bolboacă & Jäntschi, 2006-PSK](#)).

Soluția aleasă a fost de a maximiza coeficientul de determinare în expresia sa clasică (pătratul coeficientului de corelație multiplă Pearson).

[24] program evolutiv = în accepțiunea generală este un program ce implementează un algoritm genetic

Definirea experimentului de simulare a evoluției

Fișierul de configurare a execuției algoritmului genetic, excluzând opțiunile de afișare, numără 30 de parametri (redați în Tabelul 5 - Anexa 5 din Teză), dintre care 19 parametri ordinali și 12 parametri cu valori dintr-o listă de valori impuse (finită și definită). Atenție: $19+12=31!$ Performanța algoritmului genetic a fost evaluată prin prisma rezultatelor colectate în fișierele de ieșire. Configurația parametrilor de ieșire este redată în Tabelul 6-2 din Teză.

În evaluarea algoritmului genetic s-a considerat de importanță teoretică și practică compararea performanțelor obținute pentru două elemente esențiale în procesul de evoluție, și anume: strategia de selecție a indivizilor generatori de descendenți prin încrucișare și mutație și strategia de selecție a indivizilor pentru înlocuire de către descendenți (supraviețuire).

Principiul parsimoniei [25] este esența care stă la baza asocierii *Optimizare (SAR)* → *Evoluție (Observabile)*. Principiul a fost aplicat în simularea evoluției controlate către obiectivul evoluției, folosind diferite strategii de selecție și supraviețuire. Principiul a avut ca scop evidențierea diferențelor în valorile observabilelor pe parcursul evoluției.

Conform acestui principiu, pe toată durata experimentului s-au definit și păstrat aceeași parametri intrinseci execuției programului evolutiv și evoluției algoritmului genetic.

Parametri de configurare în experimentul de execuție a programului evolutiv

Clasă	Parametru	Valoare
Topologia familiei de descriptori moleculari	Genes	mp/fc/oi/id/ap/dm
	Addre	fc/ap/id/oi/dm/mp
	mp	mMnNSPsAaBbGgFfHhIi
	fc	mMDP
	oi	RrMmDd
	id	DdOoPpQqJjKkLlVvEwWwFfSsTt
	ap	CHMEGQ
	dm	gt
Topologia infrastructurii informaționale	Mydb	MDFSARs
	TabE	PCB_ikow_data
	TabM	PCB_ikow_tmpx
Material genetic și cultivar	sn0_SAMPLE_Size	12
	a_v_ADAPT_Variance	0.1
	ajb_ADAPT_JarqueBera	0.1
	a_c_ADAPT_Correlation	0.1
	g_r_GENERATIONS_first_rich	Yes
	b_k_RUNS_kepp_best_in_sample	Yes
	b_f_RUNS_get_best_from_file	No
Înmulțire	cn0_CROSSOVER_Pairs	2
	m_m_MUTATION_Genes	2
	mpp_MUTATION_Parent_probability	5%
	mcp_MUTATION_Child_probability	5%
Obiectivul evoluției	m0_REGRESSION_Multiple	4
	b_p_SELECTION_parameter	r2
	b_o_SELECTION_objective	max
Evoluție	eIn_GENERATIONS_max	20000
	eOn_RUNS_number	46
Selecție	sfn_FITNESS_normalized	No
	sfr_FITNESS_ranks	No
	sfa_FITNESS_accuracy	10000
	sff_FITTEST_function	r2_min
	sfo_FITTEST_objective	max
	fr2_FITTEST_r2_p	1.0
	fse_FITTEST_se_p	1.0
	fMt_FITTEST_Mt_p	1.0
	fHr_FITTEST_Hr_p	1.0
Supraviețuire	v_p_SURVIVAL_phenotyping_p	1.0
	v_g_SURVIVAL_genotyping_p	1.0
	vfr_SURVIVAL_ranks	No

Doi parametri: strategia de selecție (*sfs_FITNESS_strategy* în fișierul de configurare - vezi

[25] parsimonie - adoptarea celor mai simple presupuneri în formularea teoriei sau interpretarea datelor, în special în acord cu regula lamei de ras a lui Ockham (principiu atribuit logicianului William of OCKHAM, care subliniază că trebuie eliminate toate acele presupuneri care nu fac nici o diferență în predicțiile observate ale ipotezelor explicatoare sau teoriei); în latină: *lex parsimoniae - entia non sunt multiplicanda praeter necessitatem*.

Tabelul 6-1 în Teză) și strategia de supraviețuire (*vfs_SURVIVAL_strategy* în fișierul de configurare - vezi Tabelul 6-1 în Teză) au luat pe rând valorile: *proportional* (pentru strategie proporțională), *deterministic* (pentru strategie deterministă) și *tournament* (pentru strategie în turnir). S-au proiectat astfel nouă execuții ale programului evolutiv, în fiecare din cele nouă execuții dând valori diferite celor doi parametri ce definesc cele două strategii, ceilalți parametri rămânând cu aceleași valori pe toată durata experimentului.

Modalități de selecție și supraviețuire: design experimental de execuție program

Supraviețuire Selecție	Proportional (P) <u>SV_Pro</u>	Deterministic (D) <u>SV_Det</u>	Turnir (T) <u>SV_Tur</u>
Proportional (P) - <u>SL_Pro</u>	P:P (1)	P:D (2)	P:T (3)
Deterministic (D) - <u>SL_Det</u>	D:P (4)	D:D (5)	D:T (6)
Turnir (T) - <u>SL_Tur</u>	T:P (7)	T:D (8)	T:T (9)

Rezultate și discuții

Strategiile implementate pentru selecție și supraviețuire

Pentru a realiza o selecție bazată pe un scor de selecție se impune parcurgerea unei serii de pași. Algoritmul *FS* (Algoritm 1 în Teză - pentru scorurile de selecție) parcurge o astfel de serie de pași și realizează astfel premisele pentru aplicarea unei strategii de selecție. Algoritmul *PS* (Algoritm 2 - pentru strategia proporțională) realizează o selecție proporțională folosind un șir de scoruri de selecție și dă o șansă de selecție proporțională cu scorul de selecție. Algoritmul *DS* (Algoritm 3 - pentru strategia deterministă) realizează o selecție deterministă folosind un șir de scoruri de selecție extrăgând cele mai mari *N_Sel* scoruri de selecție. Algoritmul *TS* (Algoritm 4 - pentru strategia în turnir) realizează o selecție în turnir folosind un șir de scoruri de selecție extrăgând acele valori care se califică în urma unui turnir între două valori candidate repetat de un număr de ori dat de *N_Sel*. Scorul de supraviețuire (*VS*) este o valoare compozită, menită să asigure deopotrivă diversitatea genotipică și cea fenotipică în cultivar. Pentru aceasta, două măsuri de similaritate între doi descriptori intră în expresia scorului de supraviețuire: o măsură de similaritate fenotipică (*VSP* - dată de diferența între valorile scorurilor de selecție) și o măsură de similaritate genotipică (*VSG* - dată de diferența între codurile genetice ale celor doi descriptori). Tabelul 7-1 din Teză redă expresiile de calcul folosite. Odată obținute valorile *VS* ale scorurilor de supraviețuire în modalitatea expusă în tabelul de mai sus (în care scara de similaritate individuală are acum același obiectiv - cele mai mari valori dau cei mai potenți candidați pentru înlocuire în materialul genetic) pentru șirul de indivizi reprezentați în cultivar, algoritmul proiectat pentru pregătirea scorurilor de selecție *FS* este legitimat a fi folosit și pentru scorul de supraviețuire *VS*, ceea ce a și fost făcut (motiv pentru care nu există descris un algoritm separat pentru *VS*). În mod similar, algoritmi ce implementează strategiile proporțională, deterministă și în turnir sunt perfect legitimați pentru a fi folosiți și pentru șirul scorurilor de supraviețuire (*PV = PS*, *DV = DS*, *TV = TS*); nici pentru aceștia nu există descriere separată.

Algoritm 2. Algoritmul *SP* ce realizează o strategie de selecție proporțională

Scorul de selecție se calculează astfel:
└ Date de intrare:
÷ *FS_Array* - șirul scorurilor de selecție pentru fiecare genotip (*n_of_sample* valori)
÷ *FSD_Array* - șirul valorilor distincte de scor de selecție (cel mult *n_of_sample* valori)
÷ *FSC_Array* - șirul numărului de apariții pentru fiecare scor distinct de selecție
÷ *N_Sel* - număr de selecții pentru încrucișare și mutație;
└ Inițializează *Selected_Genotypes_Array* la \emptyset (mulțimea vidă);
└ Pentru fiecare selecție (de la 1 la *N_Sel*)
└ Calculează suma scorurilor pentru genotipurile încă neselectate în *FS_Sum*;
└ Generează întâmplător (distribuție uniformă) un număr *FS_Freq* între 0 și *FS_Sum* (inclusiv);
└ Găsește primul (indice) *Group* din *FSD_Array* pentru care $FS_Freq \leq \sum_{i \leq Group} FSD_Array[i] * FSC_Array[i]$;
└ Generează întâmplător (distribuție uniformă) un număr *FSD_Next* între 1 și *FSC_Array[Group]* (inclusiv)
└ Aadaugă a *FSD_Next*-a valoare egală cu *FSD_Array[Group]* neselectată încă din *FS_Array* în *Selected_Genotypes_Array* și scade o unitate din *FSC_Array[Group]*;
└ Sfârșit 'Pentru'
└ Date de ieșire: *Selected_Genotypes_Array* - șirul genotipurilor selectate (în număr de *N_Sel*)

Algoritm 3. Algoritmul *SD* ce realizează o strategie de selecție deterministă

Scorul de selecție se calculează astfel:
└ Date de intrare:
÷ *FS_Array* - șirul scorurilor de selecție pentru fiecare genotip (*n_of_sample* valori)
÷ *FSD_Array* - șirul valorilor distincte de scor de selecție (cel mult *n_of_sample* valori)
÷ *FSC_Array* - șirul numărului de apariții pentru fiecare scor distinct de selecție
÷ *N_Sel* - număr de selecții pentru încrucișare și mutație;
└ Inițializează *Selected_Genotypes_Array* la \emptyset (mulțimea vidă);
└ Inițializează *Already_Selected* la 0;
└ Inițializează *Group* la *n_of_sample*;
└ Cât timp $Already_Selected + FSC_Array[Group] \leq N_Sel$
└ Pune primul indice din *FS_Array* egal cu *FSD_Array[Group]* în *Selected_Genotypes_Array*;

```

└─ Dacă  $FSC\_Array[Group] > 0$  atunci  $dec(FSC\_Array[Group])$  altfel  $dec(Group)$ ;
└─ Sfârșit 'Cât timp'
└─ Cât timp  $Already\_Selected \leq N\_Sel$  (au mai rămas câteva scoruri identice într-un ultim grup din care trebuie făcută o selecție):
└─ Generează întâmplător (distribuție uniformă) un număr  $FSD\_Next$  între 1 și  $FSC\_Array[Group]$  (inclusiv)
└─ Aduagă a  $FSD\_Next$ -a valoare egală cu  $FSD\_Array[Group]$  neselectată încă din  $FS\_Array$  în  $Selected\_Genotypes\_Array$  și scade o unitate din  $FSC\_Array[Group]$ ;
└─ Sfârșit 'Cât timp'
└─ Date de ieșire:  $Selected\_Genotypes\_Array$  - șirul genotipurilor selectate (în număr de  $N\_Sel$ )

```

Algoritm 4. Algoritmul *ST* ce realizează o strategie de selecție în turnir

```

Scorul de selecție se calculează astfel:
└─ Date de intrare:
├─  $FS\_Array$  - șirul scorurilor de selecție pentru fiecare genotip ( $n\_of\_sample$  valori)
├─  $FSD\_Array$  - șirul valorilor distincte de scor de selecție (cel mult  $n\_of\_sample$  valori)
├─  $FSC\_Array$  - șirul numărului de apariții pentru fiecare scor distinct de selecție
├─  $N\_Sel$  - număr de selecții pentru încrucișare și mutație;
└─ Inițializează  $Selected\_Genotypes\_Array$  la o permutare întâmplătoare (distribuție uniformă) a mulțimii  $\{1..n\_of\_sample\}$ 
└─ Pentru fiecare  $i\_Sel$  de la 2 la  $N\_Sel$  (primele  $N\_Sel$  genotipuri din permutare concurează în turnir):
└─ Dacă  $FS\_Array[i\_Sel] \leq FS\_Array[i\_Sel-1]$  atunci
└─ Dacă  $FS\_Array[i\_Sel] = FS\_Array[i\_Sel-1]$  atunci dacă  $Random(\{0,1\}) = 0$  atunci continuă de la începutul iterației 'Pentru';
└─ Permută în  $FS\_Array$  valorile de pe pozițiile  $i\_Sel$  &  $i\_Sel-1$ ;
└─ Sfârșit 'Dacă'
└─ Sfârșit 'Pentru'
└─ Dacă  $N\_Sel < n\_of\_sample$  atunci (ultimul nu a participat în turnir și încă mai sunt genotipuri cu care să concureze în eșantion)
└─ Generează întâmplător (distribuție uniformă) un număr  $i\_Sel$  între  $N\_Sel + 1$  și  $n\_of\_sample$ ;
└─ Dacă  $FS\_Array[N\_Sel] \leq FS\_Array[i\_Sel]$  atunci
└─ Dacă  $FS\_Array[N\_Sel] = FS\_Array[i\_Sel]$  atunci
└─ Dacă  $Random(\{0,1\}) = 0$  atunci Stop (turnir complet);
└─ Permută în  $FS\_Array$  valorile de pe pozițiile  $N\_Sel$  &  $i\_Sel$ ;
└─ Sfârșit 'Dacă'
└─ Sfârșit 'Dacă'
└─ Date de ieșire:  $Selected\_Genotypes\_Array$  - șirul genotipurilor selectate (în număr de  $N\_Sel$ )

```

Fișiere rezultat

Execuția programului evolutiv s-a făcut pe calculatoare din generația P6 (Dual P5) în perioada Ianuarie - Februarie 2009 și rezultatele au fost salvate de program în fișierele date în tabelul de mai jos.

Fișiere rezultat (configurare și evoluție) după designul experimental

Selecție	Supraviețuire	Configurare	Evoluție
Proporțional	Proporțional	PCB_4044_cfg.txt	PCB_4044_evo.txt
Proporțional	Deterministic	PCB_2441_cfg.txt	PCB_2441_evo.txt
Proporțional	Turnir	PCB_9878_cfg.txt	PCB_9878_evo.txt
Deterministic	Proporțional	PCB_5108_cfg.txt	PCB_5108_evo.txt
Deterministic	Deterministic	PCB_6369_cfg.txt	PCB_6369_evo.txt
Deterministic	Turnir	PCB_6690_cfg.txt	PCB_6690_evo.txt
Turnir	Proporțional	PCB_5828_cfg.txt	PCB_5828_evo.txt
Turnir	Deterministic	PCB_4872_cfg.txt	PCB_4872_evo.txt
Turnir	Turnir	PCB_1758_cfg.txt	PCB_1758_evo.txt

Disponibile pentru descărcare de la adresa:

http://l.academicdirect.org/Horticulture/GAs/MLR_MDF_selection_vs_survival/

Verificarea datelor: Testul Benford

S-au cumulat frecvențe din 46 de execuții independente ale numărului (num_obs) și aparițiilor (sum_obs) genotipurilor viabile reprezentate în cultivar în generațiile ce au produs evoluție pentru fiecare asociere de selecție (Sel) și supraviețuire (Srv): Srv, Sel ∈ {P, T, D} ([Tabelul 7-3](#) în Teză).

S-au aplicat testele Chi-Square și Kolmogorov-Smirnov pentru a verifica dacă numerele urmează legea Benford. Tabelele de mai jos redau această analiză.

Testul χ^2 aplicat frecvențelor observate pentru respingerea ipotezei că prima, a doua și a treia cifră semnificativă a numerelor urmează [distribuția Benford](#)

Cifra			Frecvență așteptată			Frecvență observată			$(O_i - E_i)^2$			$(O_i - E_i)^2 / E_i$		
d ₀	d ₁	d ₂	d ₀	d ₁	d ₂	d ₀	d ₁	d ₂	d ₀	d ₁	d ₂	d ₀	d ₁	d ₂
0	0	0	-	40	19	-	28	25	-	144	36	-	3.60	1.89
1	1	1	108	38	19	117	41	18	81	9	1	0.75	0.24	0.05
2	2	2	63	37	18	72	37	18	81	0	0	1.29	0.00	0.00
3	3	3	45	35	18	48	33	20	9	4	4	0.20	0.11	0.22
4	4	4	35	34	18	33	34	11	4	0	49	0.11	0.00	2.72
5	5	5	29	33	18	17	42	15	144	81	9	4.97	2.45	0.50
6	6	6	24	32	18	16	34	18	64	4	0	2.67	0.13	0.00
7	7	7	21	31	18	18	30	18	9	1	0	0.43	0.03	0.00
8	8	8	18	30	18	19	31	16	1	1	4	0.06	0.03	0.22
9	9	9	17	29	18	20	29	23	9	0	25	0.53	0.00	1.39
Σ	Σ	Σ	360	339	182	360	339	182	402	244	128	11.0	6.60	7.00

Numărul gradelor de libertate ale legii de distribuție Benford este 1 (baza de numerație, 10)
d₀: $X^2=11 < 14.7 = \chi^2(9-2, 5\%)$; d₁: $X^2=6.6 < 15.5 = \chi^2(10-2, 5\%)$; d₂: $X^2=7 < 15.5 = \chi^2(10-2, 5\%)$;

Testul K-S aplicat frecvențelor observate pentru respingerea ipotezei că prima, a doua și a treia cifră semnificativă a numerelor urmează [distribuția Benford](#)

Cifra			Frecvență cumulată așteptată și observată						Diferență			Diferență		
d ₀	d ₁	d ₂	d _{0a}	d _{1a}	d _{2a}	d _{0o}	d _{1o}	d _{2o}	d ₀	d ₁	d ₂	d ₀	d ₁	d ₂
0	0	0	0	40	19	0	28	25	0	12	-6	0	12	6
1	1	1	108	78	38	117	69	43	-9	9	-5	9	9	5
2	2	2	171	115	56	189	106	61	-18	9	-5	18	9	5
3	3	3	216	150	74	237	139	81	-21	11	-7	21	11	7
4	4	4	251	184	92	270	173	92	-19	11	0	19	11	0
5	5	5	280	217	110	287	215	107	-7	2	3	7	2	3
6	6	6	304	249	128	303	249	125	1	0	3	1	0	3
7	7	7	325	280	146	321	279	143	4	1	3	4	1	3
8	8	8	343	310	164	340	310	159	3	0	5	3	0	5
9	9	9	360	339	182	360	339	182	0	0	0	0	0	0
Σ	Σ	Σ	-	-	-	-	-	-	-66	55	-9	82	55	37

d₀: $D\sqrt{n} = \frac{21}{360}\sqrt{9} = \frac{14}{80} < \frac{31}{80} = K(9, 5\%); \frac{14}{80} = K(9, 90.82\%)$
d₁: $D\sqrt{n} = \frac{12}{339}\sqrt{10} \approx \frac{17}{152} < \frac{56}{152} \approx K(10, 5\%); \frac{17}{152} = K(10, 95.16\%)$
d₂: $D\sqrt{n} = \frac{7}{182}\sqrt{10} \approx \frac{62}{510} < \frac{188}{510} \approx K(10, 5\%); \frac{62}{510} = K(10, 94.60\%)$

Valorile χ^2 obținute arată că pentru fiecare dintre primele trei cifre semnificative nu se poate respinge semnificativ statistic ipoteza de distribuție după [legea Benford](#) la un nivel de semnificație de 5%. Pentru ca să se poată asigura semnificativ statistic distribuția cifrelor după [legea Benford](#), ar fi fost necesar ca valoarea X^2 să fi fost de cel mult $2.17 = \chi^2(7, 95\%)$ pentru prima cifră, și de cel mult $2.73 = \chi^2(8, 95\%)$ pentru următoarele două cifre semnificative ale numerelor.

Se impune o măsură prealabilă în aplicarea [testului \$\chi^2\$](#) între observat și așteptat la datele primare ([Tabelul 7-4](#) din Teză), și anume asigurarea apropierii de normalitate a distribuției pătratelor diferențelor ([Fisher, 1920-Accuracy](#)). Astfel, dacă valoarea $b_2 = m_4/m_2^2$ este apropiată de 3, atunci

distribuția pătratelor erorilor se poate aproxima de legea de distribuție normală și pătratele diferențelor constituie o statistică suficientă ([Fisher, 1922-Estimation](#)); dacă b_2 este apropiată de 6, atunci distribuția pătratelor erorilor se poate aproxima de legea de distribuție dublu exponențială și suma diferențelor modulelor constituie o statistică suficientă. Pentru diferențele din Tabelul 7-4 din Teză valorile β_2 sunt: $\beta_2(d_0) = 2.13$; $\beta_2(d_1) = 4.60$; $\beta_2(d_2) = 2.71$ cu o valoare medie de 3.15, deci aplicarea [testului \$\chi^2\$](#) pentru pătratele diferențelor este consistentă cu legea de distribuție a acestora.

Concluzia analizei este că nu se poate respinge ipoteza [distribuției Benford](#) cu o probabilitate de 95% (la nivelul de semnificație de 5%). Mai mult, valorile probabilităților din [distribuția Kolmogorov](#) arată că un număr relativ scăzut de date experimentale concordă mai bine cu legea de [distribuție Benford](#) (9.18% pentru d_0 ; 4.84% pentru d_1 ; 5.4% pentru d_2).

Analiza variabilității

Frecvența de apariție a genotipurilor în eșantion de-a lungul evoluțiilor permite aprecieri cu privire la capacitatea de adaptare a acestora, și în același timp reprezintă o măsură a variabilității materialului genetic al eșantionului pe care o induce metoda de selecție și metoda de supraviețuire.

Numărul de genotipuri viabile (reprezentate fenotipic în cultivar) s-a urmărit folosind schema de contingență {Top23, Total}X{NGD, NTG, Part}, unde Top23 - referind cele mai frecvente 23 iar Total - toate genotipurile reprezentate în cultivar în 46 de execuții independente, NGD - numărul de genotipuri distincte, NTG - suma numărului de genotipuri și NGR - numărul de genotipuri ale fenotipurilor participante în regresii (valori date în Tabelul 8-1 în Teză).

Pentru testarea independenței între metodele de selecție și supraviețuire în ceea ce privește numărul de genotipuri s-a folosit testul χ^2 aplicat la o tabelă de contingență de 3x3 pentru fiecare serie de valori numerice din Tabelul 8-1 din Teză și rezultatele sunt în Tabelele 8-2, ...8-7 din Teză.

(Tabelul 8-2 în Teză): Sunt independente metoda de selecție față de metoda de supraviețuire în ceea ce privește numărul de genotipuri distincte din cultivar în generațiile ce produc evoluție? - NU

(Tabelul 8-3 în Teză): Sunt independente metoda de selecție față de metoda de supraviețuire în ceea ce privește numărul total de genotipuri din cultivar în generațiile ce produc evoluție? - NU

(Tabelul 8-4 în Teză): Sunt independente metoda de selecție față de metoda de supraviețuire în ceea ce privește genotipurile participante la regresii ce produc evoluție? - NU

(Tabelul 8-5 în Teză): Sunt independente metoda de selecție față de metoda de supraviețuire în ceea ce privește numărul de genotipuri distincte din Top 23 în generațiile ce produc evoluție? - NU

(Tabelul 8-6 în Teză): Sunt independente metoda de selecție față de metoda de supraviețuire în ceea ce privește numărul total de genotipuri din Top 23 în generațiile ce produc evoluție? - NU

(Tabelul 8-7 în Teză): Sunt independente metoda de selecție față de metoda de supraviețuire în ceea ce privește genotipurile din Top 23 participante la regresii ce produc evoluție? - NU

Confidența în dependența de strategia de selecție și supraviețuire crește în ordinea: număr de genotipuri distincte; număr total de genotipuri; număr de genotipuri participante la regresii în același timp cu faptul că numărul de observații nu crește în aceeași ordine. În baza dependenței remarcate în strategia de selecție și de supraviețuire la toți parametrii ce caracterizează numărul de genotipuri pe parcursul evoluției, s-a impus o caracterizare a acestei dependențe.

La întrebarea “Există legătură între cele trei serii de numere de genotipuri?” se răspunde calculând coeficientul de corelație (calcul dat în tabelul de mai jos).

Există legătură între numărul de genotipuri distincte (NGD), numărul total de genotipuri (NTG) și numărul de genotipuri participante la regresii (NGR)? - DA

Serii	Coeficient de determinare	Valoare F; probabilitate de a greși
NGD vs. NTG	0.982 (y=ax)	924; 10^{-15}
NGD vs. NGR	0.982 (y=ax)	951; 10^{-15}
NTG vs. NGR	0.999 (y=ax)	16110; 10^{-25}

S-au calculat media și deviația standard (în ipoteza că distribuția de eșantionare induce o distribuție normală a acestor statistici ale eșantionului) în jurul parametrului statistic (al populației) asociat. Teorema Limită Centrală asigură faptul că se pot folosi în această analiză statistică valoarea medie (m) și abaterea standard (s), prezentate în [Tabelul 3 & Tabelul 4 - Anexa 3-1](#) din Teză, și pe baza cărora se poate exprima intervalul de încredere al acestora din distribuția Student t. S-a folosit aplicația [Statistica](#) pentru a calcula aceste valori.

Analiza valorilor medii a permis formularea următoarelor concluzii:

÷ Selecția deterministă (D) face ca:

- Indiferent de metoda de supraviețuire, numărul total de genotipuri distincte să scadă semnificativ statistic;
- Folosind supraviețuirea turnir (T) sau proporțională (P), se remarcă scăderea semnificativă statistic la toți parametrii observați (Top 23 și Total; Distincți, Apariții și Participări în regresii), în timp ce, folosind supraviețuirea deterministă (D), se remarcă creșterea

- semnificativă statistic numai în ceea ce privește cele mai frecvente genotipuri pentru toți parametrii (Distincți, Apariții și Participări în regresii);
- ÷ Supraviețuirea deterministă (D) face ca:
 - Folosind supraviețuirea turnir (T) sau proporțională (P), să mărească semnificativ numărul total de genotipuri pentru toți parametrii (Distincți, Apariții și Participări în regresii).

Analiza abaterilor standard a permis formularea următoarelor concluzii:

- ÷ Supraviețuirea deterministă (D) îmbogățește semnificativ statistic grupul celor mai frecvente genotipuri (Top23) din generațiile ce produc evoluție, în timp ce selecția deterministă (D) sărăcește semnificativ statistic numărul total al genotipurilor din generațiile care produc evoluție;
- ÷ Practic, fiecare metodă de selecție definește câte o populație genotipică în generațiile care produc evoluție; argumentul este că oricare ar fi parametrul urmărit pentru numărul total de genotipuri (Distincți, Apariții și Participări în regresii), și luând pentru exemplificare numărul de genotipuri distincte (Num), se obține:
 - Varianța totală: 1583^2 cu intervalul de încredere de 95%: $[1069^2, 3033^2]$;
 - Varianța populației produse de selecția proporțională (P): $656^2 < 1069^2$;
 - Varianța populației produse de selecția turnir (T): $731^2 < 1069^2$;
 - Varianța populației produse de selecția proporțională (P): $523^2 < 1069^2$.
- ÷ Nu aceeași concluzie se poate trage cu privire la metoda de supraviețuire, pentru care se produce segregare populațională doar pentru supraviețuirea deterministă (D), care creează o populație cu un număr mediu de genotipuri semnificativ statistic mai mare decât supraviețuirea proporțională (P) și respectiv turnir (T).

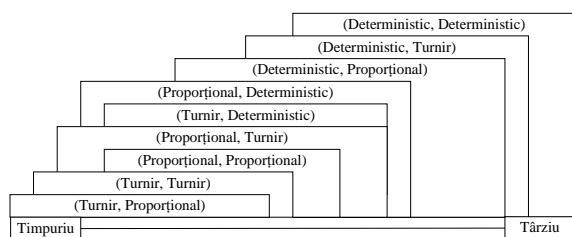
Pe baza rezultatelor experimentale obținute, poate fi interpretat și un alt parametru important al evoluției: numărul de generații care produc evoluție în cursul execuției cu număr impus de generații (20000), ca măsură a capacității de adaptare definită de combinația celor două metode (de selecție și supraviețuire), precum și valoarea medie a numerelor generațiilor care produc evoluție ca măsură a vitezei de adaptare.

Rezultatele au fost prelucrate (în Tabelul 8-11 din Teză) astfel: pentru fiecare execuție independentă a algoritmului genetic s-au consemnat numerele consecutive ale generațiilor care au produs îmbunătățirea valorii funcției obiectiv, și anume s-a obținut o ecuație de regresie validă (toți coeficienții sunt semnificativi statistic) cu un coeficient de determinare mai mare decât cel obținut în generațiile anterioare. S-a calculat apoi valoarea medie a acestui număr (care întotdeauna este mai mic decât numărul maxim de generații al unei execuții) și numărul de evoluții distincte (numărul de valori), informații care sunt prezentate în Tabelul 8-11 din Teză. Întrucât ambele valori (media și numărul de valori) au fost obținute printr-o repetare (de 46 de ori) a experimentului pentru fiecare pereche de metode (selecție, supraviețuire) valorile obținute aproximează distribuția de eșantionare, astfel încât s-a putut presupune aproximația la normalitate a acestora (atât cele ca valori ale eșantionului de 46 de observații) cât și populația din care provin, care așa cum s-a dovedit mai sus este caracteristică (distinctă) cel puțin după metoda de selecție. Pentru a se realiza compararea perechilor de metode (selecție, supraviețuire) informațiile din Tabelul 8-11 din Teză au fost supuse unei analize statistice descriptive, care a inclus calcularea valorilor medii și a deviațiilor standard, împreună cu intervalele de încredere la un nivel de semnificație de 95%, rezultate care sunt redată în tabelul de mai jos. Datele din tabelul de mai jos au servit pentru caracterizarea genotipurilor reprezentate în cultivar (numărul și media acestora). Caracterizarea s-a făcut pentru fiecare parametru statistic calculat, atribuind semnificațiile fizice măsurilor folosite și interpretând valorile obținute. Rezultatele acestei caracterizări sunt redată în continuare.

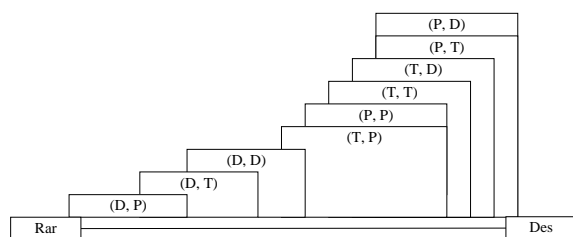
Media generațiilor care produc evoluții este o măsură a ***cât de timpuriu / târziu se produc evoluțiile***. Figura (Figura 8-2 în Teză) reprezintă valorile mediei și intervalului său de încredere; analizând, se remarcă: selecția deterministă (și cu atât mai mult însoțită de supraviețuirea deterministă) produce cele mai târzii evoluții; selecția turnir (și cu atât mai mult însoțită de supraviețuirea proporțională sau turnir) produce cele mai timpurii evoluții.

Statistici ale generațiilor ce produc evoluție în funcție de strategia de selecție și supraviețuire

Parametru	Medie	CI(95%,Medie)	Deviație	CI(95%,Deviație)
m(D,D)	4120	3518 4722	2027	1681 2553
m(D,T)	3907	3326 4488	1957	1623 2465
m(D,P)	3714	3032 4396	2296	1904 2892
m(P,D)	3335	2631 4039	2369	1965 2984
m(T,D)	3307	2748 3866	1882	1561 2371
m(P,T)	3214	2520 3908	2338	1939 2945
m(P,P)	3196	2671 3722	1770	1468 2229
m(T,T)	2929	2400 3458	1781	1478 2244
m(T,P)	2916	2322 3510	2001	1660 2520
n(P,D)	32.0	29.0 35.1	10.1	8.4 12.8
n(P,T)	31.8	29.0 34.6	9.3	7.7 11.7
n(T,D)	31.1	28.3 33.8	9.2	7.7 11.6
n(T,T)	30.0	27.2 32.7	9.3	7.7 11.7
n(P,P)	29.4	26.4 32.4	10.2	8.4 12.8
n(T,P)	28.7	25.2 32.1	11.5	9.5 14.5
n(D,D)	23.6	20.8 26.3	9.2	7.6 11.6
n(D,T)	21.7	19.1 24.2	8.5	7.1 10.8
n(D,P)	18.5	16.2 20.8	7.9	6.5 9.9

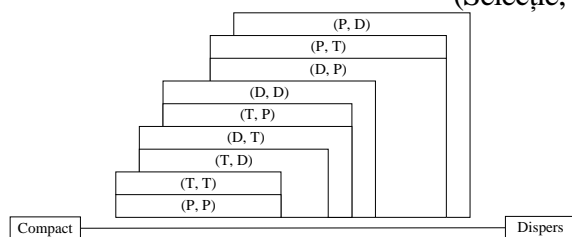


[CI(95%,Medie)_{n=46}(generație medie)]
Cât de timpuriu se produc evoluțiile?

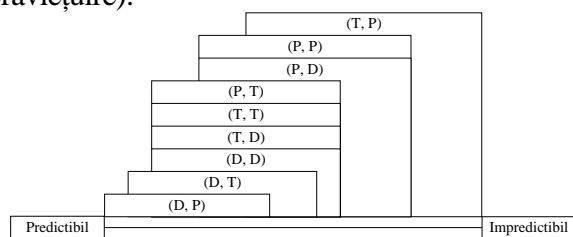


[CI(95%,Medie)_{n=46}(număr evoluții)]
Cât de frecvent se produc evoluțiile?

(Selecție, Supraviețuire):



[CI(95%,Deviație)_{n=46}(generație medie)]
Cât de dispers se produc evoluțiile?



[CI(95%,Deviație)_{n=46}(număr evoluții)]
Cât de predictibil se produc evoluțiile?

Numărul de evoluții dintr-un număr dat de generații este o măsură a **cât de frecvent se produc evoluțiile**. Figura (Figura 8-3 în Teză) reprezintă valorile mediei și intervalului său de încredere; analizând, se remarcă: selecția proporțională (și cu atât mai mult însoțită de supraviețuirea deterministă sau turnir) produce cele mai multe evoluții, în timp ce selecția deterministă (și cu atât mai mult însoțită de supraviețuirea proporțională sau turnir) produce cele mai rare evoluții; în ceea ce privește frecvența evoluțiilor, o selecție proporțională aproape că nu face diferența între supraviețuirea deterministă și supraviețuirea în turnir (medii 32 și 31.8, cu intervale de încredere rotunjite la întreg egale) în timp ce viteza evoluțiilor aceeași selecție proporțională de verificat/corectat! nu face diferența între valoarea medie pentru supraviețuirea proporțională și turnir (diferența de 150 între valorile medii reprezentând cel mult 30% din lărgimea intervalului de încredere la oricare dintre ele).

Variabilitatea momentului evoluției pe parcursul generațiilor este o măsură a **cât de compact / dispers se produc evoluțiile**. Figura (Figura 8-4 în Teză) reprezintă valorile deviației

standard și intervalului său de încredere; analizând, se remarcă: cea mai bună compactitate de evoluție o au selecția proporțională asociată cu supraviețuirea proporțională și selecția în turnir asociată cu supraviețuirea în turnir; cea mai mare împrăștiere în timp a evoluției se observă la selecția proporțională asociată cu supraviețuirea deterministă, urmată îndeaproape de selecția deterministă asociată cu supraviețuirea proporțională.

Variabilitatea numărului evoluțiilor este o măsură a *cât de predictibil / impredictibil se produc evoluțiile*. Figura (Figura 8-5 în Teză) reprezintă valorile deviației standard și intervalului său de încredere; analizând, se remarcă: patru asocieri de selecție și supraviețuire au rezultate similare sub aspectul predictibilității evoluției: (selecție deterministic, supraviețuire deterministic), (selecție în turnir, supraviețuire deterministic), (selecție în turnir, supraviețuire în turnir) și (selecție proporțională, supraviețuire în turnir); predictibilități extreme ale evoluției obțin selecția deterministă asociată cu supraviețuirea proporțională având cea mai mare predictibilitate (cea mai mică variabilitate) și selecția în turnir asociată cu supraviețuirea proporțională având cea mai mică predictibilitate (cea mai mare variabilitate) de evoluție.

Analiza diversității

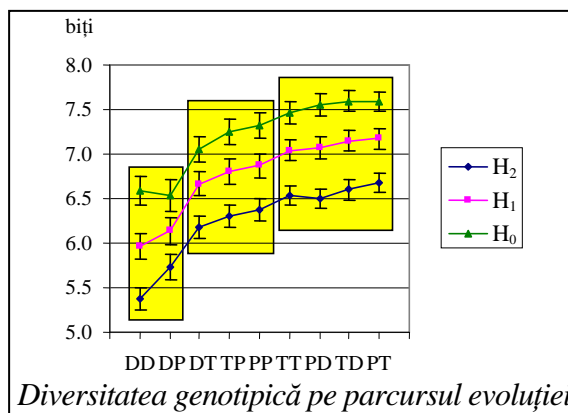
Diversitatea genotipurilor reprezentate în cultivar în momentele evoluției se poate cuantifica folosind entropia informațională, având la dispoziție o familie de măsuri entropice, date de expresia entropiei generalizate (sau Rényi) - (Rényi, 1961) - $H(p_1, p_2, \dots, p_n; \alpha) = H_\alpha(p_1, p_2, \dots, p_n)$ și unde are loc relația:

$H(\{p\};\infty)$ = $H_\infty(\{p\})$	\leq	$H(\{p\};2)$ = $H_2(\{p\})$	\leq	$H(\{p\};1)$ = $H_1(\{p\})$	\leq	$H(\{p\};0)$ = $H_0(\{p\})$	\leq	$2H(\{p\};\infty)$ = $2H_\infty(\{p\})$
H_∞ : entropia min; H_2 : logaritmul cu semn schimbat al diversității Simpson; H_1 : entropia Shannon; H_0 : entropia Hartley (entropia max);								

Măsurile H_0 , H_1 și H_2 sunt frecvent folosite în literatura de specialitate pentru a descrie gradul de dezordine sau diversitatea. Astfel, entropia Hartley (H_0) - (Hartley, 1928) - este aplicată în informatică la reconcilierea informației, entropia Shannon (H_1) - (Shannon, 1948) - este aplicată în fizică și chimie la caracterizarea stării materiei, iar entropia derivată din indicele de diversitate Simpson (H_2) - (Simpson, 1949) - în biologie și științele conexe la caracterizarea diversității populațiilor de organisme. Dacă logaritmul din expresia entropiei se calculează în baza 2 atunci unitatea de măsură a entropiei este biți, dacă se calculează în baza e (logaritm natural) atunci măsura este nați, iar dacă se calculează în baza 10 atunci măsura este diți.

Fișierele de rezultate au fost prelucrate pentru calculul măsurilor entropice H_0 , H_1 și H_2 :

- ÷ Au fost extrase pentru fiecare execuție, fiecare asociere (selecție, supraviețuire) și fiecare evoluție lista genotipurilor reprezentate în cultivar în momentul evoluției (coloanele Gen0..Gen11 în fișierele de rezultate);
- ÷ S-au concatenat listele de genotipuri pentru toate evoluțiile unei execuții și unei asocieri;
- ÷ S-au sortat și apoi numărat aparițiile acestor genotipuri pentru fiecare execuție și asociere;
- ÷ Rezultatul obținut reprezintă frecvențe genotipuri la o execuție și o asociere;
- ÷ S-au calculat probabilitățile $\{p\}$ din frecvențe; s-au aplicat formulele de mai sus pentru calculul lui H_0 , H_1 și H_2 ; valorile calculate sunt redată în Anexa 6;
- ÷ Pentru valorile calculate ale entropiilor H_0 , H_1 și H_2 și pentru fiecare asociere de strategie de selecție și supraviețuire s-au calculat valorile medii și lărgimea intervalelor de încredere pentru medii; rezultatele sunt redată în Tabelul 8-13 din Teză și s-au reprezentat grafic în ordinea crescătoare a valorii medii împreună cu intervalul de încredere în figura alăturată.



Graficul arată că cele nouă strategii de evoluție se grupează în 3 grupe după diversitatea genotipică care o produc în cultivar:

- ÷ Grupul (DD, DP) - selecție deterministă, supraviețuire deterministă sau proporțională este semnificativ statistic diferit în diversitate de toate celelalte strategii în toate trei măsurile entropice; produc cea mai scăzută diversitate genotipică;
- ÷ Pentru grupul (DT,TP,PP) - (selecție deterministă, supraviețuire în turnir) împreună cu (supraviețuire proporțională, selecție în turnir sau proporțională) se poate evidenția statistic că este în medie semnificativ statistic diferit de grupul de ce produce mai mare diversitate (TT, PD, TD, PT) chiar dacă nu fiecare metodă a grupului este semnificativ statistică diferită de metodele celui alt grup - de exemplu pentru H_2 calculul mediei arată că $Medie(DT,TP,PP) = 6.29 \pm 0.07 (df=138)$ în timp ce $Medie(TT,PD,TD,PT) = 6.58 \pm 0.05 (df=184)$;
- ÷ Grupul (TT, PD, TD, PT) este compus din cele mai favorabile strategii pentru păstrarea diversității: (selecție și supraviețuire în turnir), (selecție proporțională și supraviețuire deterministă), (selecție în turnir și supraviețuire deterministă), și respectiv selecție proporțională și supraviețuire în turnir.

O remarcă generală cu privire la strategia de selecție se poate face concatenând observațiile de la cele trei strategii de supraviețuire și calculând din nou valoarea medie împreună cu intervalul său de încredere și în mod similar pentru strategia de supraviețuire. Tabelul următor cumulează aceste rezultate pentru H_1 (tabelul de mai jos, Tabelul 8-14 în Teză).

Entropia Shannon (H_1) pentru selecție și pentru supraviețuire

Strategie	Medii și intervale de încredere la riscul de a fi în eroare de 5%	
Selecție	D(6.26±0.10)	
	T(7.00±0.07)	
	P(7.04±0.07)	
Supraviețuire	P(6.61±0.10)	
	D(6.73±0.12)	
	T(6.96±0.08)	
D: strategie deterministă; T: strategie în turnir; P: strategie proporțională		

Analiza statistică din tabel arată extrem de clar că influență decisivă asupra reducerii dramatice a diversității în cultivar o are strategia de selecție deterministă; diferența între diversitatea produsă de strategia de selecție proporțională și cea în turnir nu este semnificativă statistic la riscul de a fi în eroare de 5%. Strategia de supraviețuire are, de asemenea, o influență decisivă asupra diversității genotipurilor. Strategia de supraviețuire în turnir produce o diversitate mai mare (fapt evidențiat statistic cu un risc de a fi în eroare de 5%) decât strategia deterministă și respectiv proporțională; acestea din urmă nu se disting statistic la riscul de a fi în eroare de 5%, cu toate că observația arată că strategia deterministă produce o mai mare diversitate.

Interpretarea rezultatelor din observațiile pe observabile calitative

Analiza numărului de genotipuri reprezentate în cultivar în generațiile ce au produs evoluție

O ipoteză importantă poate fi formulată pe baza datelor prezentate în Tabelul 8-1, și anume: dacă numărul mediu de genotipuri viabile (ce rezultă din împărțirea la numărul de generații ce au produs evoluție observate num_obs a numărului de genotipuri viabile observate în aceste generații sum_obs) este sau nu independent de strategia de selecție și strategia de supraviețuire.

Pentru această analiză (care să răspundă la întrebarea: “În ce măsură numărul mediu de genotipuri viabile este sau nu independent de strategia de selecție și strategia de supraviețuire?”) s-a implementat și folosit un program pentru analiza de distribuție a perechilor de eșantioane folosind [testul Anderson-Darling](#) descris în [Anexa 3-3](#).

Un calcul al numărului total de inferențe posibile arată că pentru media genotipurilor viabile (dată în Tabelul 9-1 din Teză) trebuie investigate un număr de $2^9 - C_9^0 - C_9^1 = 502$ inferențe statistice. Rezultatul analizei este redat în [Tabelul din Anexa 7-1](#).

Pe baza rezultatelor din [Tabelul din Anexa 7-1](#)) se desprinde că:

- ÷ Analiza grupurilor de câte două perechi de metode evidențiază statistic că:
 - Nu poate fi respinsă ipoteza apartenenței la populații identice pentru: DT și DP (cu un raport între valoarea critică și statistică $c/k = 1.2$); PP și PT ($c/k = 3.0$); PP și TT (3.2); TT și PT (2.1); TT și TP (1.3);
- ÷ Concluzia analizei grupurilor de câte două perechi de metode este:
 - Cu un risc de a fi în eroare de 5% următoarele provin din populații diferite: DD, PD, TD;
 - Nu poate fi respinsă ipoteza că DT și DP provin din populații identice;
 - Analiza populațională pentru PP, PT, TT și TP necesită investigația grupurilor de ordin superior;
- ÷ Analiza grupurilor de câte trei perechi de metode evidențiază statistic că:
 - Nu poate fi respinsă ipoteza că PP, PT și TT ($c/k = 2.2$) provin din populații identice;
 - Cu un risc de a fi în eroare de 5% toate grupurile de trei perechi de metode ce conțin metoda TP provin din populații diferite;
- ÷ Concluzia analizei grupurilor de câte cel mult trei perechi de metode este:
 - Cu un risc de a fi în eroare de 5% următoarele provin din populații diferite: DD, PD, TD;
 - Nu poate fi respinsă ipoteza că DT și DP provin din populații identice;
 - Nu poate fi respinsă ipoteza că TT și TP provin din populații identice;
 - Nu poate fi respinsă ipoteza că PP, PT și TT provin din populații identice;
- ÷ Analiza grupurilor de câte patru perechi de metode evidențiază statistic că:
 - Cu un risc de a fi în eroare de 5% în grupul de perechi de metode PP, PT, TP, TT cel puțin una provine din populație diferită ($c/k = 0.9$);
- ÷ Concluzia analizei grupurilor de perechi de metode este:
 - Cu un risc de a fi în eroare de 5% DD, PD, și TD provin din populații diferite;
 - Nu poate fi respinsă ipoteza că DT și DP provin din populații identice;
 - Nu poate fi respinsă ipoteza că TT și TP provin din populații identice;
 - Nu poate fi respinsă ipoteza că PP, PT și TT provin din populații identice;
 - Cu un risc de a fi în eroare de 5% se respinge ipoteza că PP, PT, TP și TT ar proveni din populații identice;

Analiza numărului de fenotipuri viabile din cultivar în generațiile ce au produs evoluție

Tabelul 9-3 din Teză cumulează rezultatele obținute pentru numărul de fenotipuri viabile observate (sum_obs) și media numărului acestora (avg_obs) pentru același număr de evoluții observate (num_obs) ale căror valori au fost redată în Tabelul 8-1 din Teză, fiind aici grupate după mii de generații (de la 0..1000 până la 19001-20000).

Ipoteza cu privire la numărul mediu de fenotipuri viabile (valori redată în Tabelul 9-3 din Teză) și anume “Strategia de selecție și strategia de supraviețuire produc populații distincte pentru numărul mediu de fenotipuri viabile?” s-a verificat folosind [testul Anderson-Darling](#) și rezultatul

analizei este redat în [Tabelul din Anexa 7-2](#).

Se observă din analiza statistică ([Tabelul din Anexa 7-2](#)) că interpretarea populațională a numărului mediu de fenotipuri din cultivar în generațiile ce produc evoluție este mult mai complexă decât interpretarea populațională ([Tabelul din Anexa 7-1](#)) a numărului mediu de genotipuri în aceleași evoluții.

Interpretarea rezultatelor din [Tabelul din Anexa 7-2](#) necesită identificarea celor mai mari grupuri de perechi de metode cu populații posibil identice (numită în continuare lista suspecților - suspecți de a proveni din populații identic distribuite) și eliminarea sub-grupurilor unice ale acestora de ordin inferior; astfel cele mai mari grupuri de ordin maxim sunt grupurile de ordin 5 (PP, PT, TT, TD, DD) - run 400 - și (PP, PT, PD, TD, DD) - run 448 - care intră în mod automat în lista suspecților; grupurile de ordin inferior sunt după cum urmează din aplicarea algoritmului de incluziune:

- ÷ Grupuri de ordin 5; lista suspecților: {(PP, PT, TT, TD, DD), (PP, PT, PD, TD, DD)}
 - (PP, PT, TD, DD) - run 384 - e simultan în (PP, PT, TT, TD, DD) și (PP, PT, PD, TD, DD) ceea ce face imposibilă decelarea apartenenței sale; se adaugă la lista suspecților;
 - (PT, TT, TD, DD) - run 145 - e doar în (PP, PT, TT, TD, DD); se elimină;
 - (PP, TT, TD, DD) - run 272 - e doar în (PP, PT, TT, TD, DD); se elimină;
 - (PP, PT, TT, DD) - run 392 - e doar în (PP, PT, TT, TD, DD); se elimină;
 - (PP, PT, TT, TD) - run 399 - e doar în (PP, PT, TT, TD, DD); se elimină;
 - (PP, PT, PD, DD) - run 440 - e doar în (PP, PT, PD, TD, DD); se elimină;
 - (PT, PD, TD, DD) - run 193 - e doar în (PP, PT, PD, TD, DD); se elimină;
 - (PP, PT, PD, TD) - run 447 - e doar în (PP, PT, PD, TD, DD); se elimină;
- ÷ Grupuri de ordin cel puțin 4; lista suspecților: {(PP, PT, TT, TD, DD), (PP, PT, PD, TD, DD), (PP, PT, TD, DD)}
 - (PT, TD, DD) - run 129, (PP, PT, DD) - run 376, (PP, TD, DD) - run 256, (PP, PT, TD) - sunt simultan în (PP, PT, TT, TD, DD), (PP, PT, PD, TD, DD) și (PP, PT, TD, DD); se adaugă;
 - (TT, TD, DD) - run 20, (PT, TT, TD) - run 144, (PT, TT, DD) - run 137, (PP, TT, TD) - run 271, (PP, TT, DD) - run 264 și (PP, PT, TT) - run 391 - sunt doar în (PP, PT, TT, TD, DD); se elimină;
 - (PT, PD, DD) - run 185, (PP, PT, PD) - run 439 și (PP, PD, DD) - run 312 - sunt doar în (PP, PT, PD, TD, DD); se elimină;
- ÷ Grupuri de ordin cel puțin 3; lista suspecților: {(PP, PT, TT, TD, DD), (PP, PT, PD, TD, DD), (PP, PT, TD, DD), (PT, TD, DD), (PP, PT, DD), (PP, TD, DD), (PP, PT, TD)}
 - (PT, DD) - run 121 - în run 376, 129, 384, 448, 400; se adaugă;
 - (TT, TD) - run 19, (PT, TT) - run 136, (TT, DD) - run 12, (PP, TT) - run 263 - în run 400; se elimină;
 - (PP, DD) - run 248 - în run 256, 376, 384, 448, 400; se adaugă;
 - (TD, DD) - 5 în 256, 129, 384, 448, 400; se adaugă;
 - (PT, TD) - 128 - în 383, 129, 384, 448, 400; se adaugă;
 - (PP, TD) - 255 - în 383, 256, 384, 448 și 400; se adaugă;
 - (DP, DT) - run 3 - nu se regăsește în grupurile de ordin superior; se adaugă;
 - (PP, PT) - 375 - în 383, 376, 384, 448, 400; se adaugă;
 - (PT, PD) - run 184 și (PD, DD) - în run 448; se elimină;
 - (TP, TT) - run 42 - nu se regăsește în grupurile de ordin superior; are însă cea mai mare susceptibilitate situându-se foarte aproape de confidența de 95 pentru respingere ($c/k = 1.09$); se elimină;
- ÷ Grupuri de ordin cel puțin 2; lista suspecților: {(PP, PT, TT, TD, DD), (PP, PT, PD, TD, DD), (PP, PT, TD, DD), (PT, TD, DD), (PP, PT, DD), (PP, TD, DD), (PP, PT, TD), (PT, DD), (PP, DD), (TD, DD), (PT, TD), (PP, TD), (DP, DT), (PP, PT)}; se impune continuarea includerii în listă cu populațiile distincte, urmând aceeași procedură:
 - PP se regăsește în mai multe grupuri de ordin superior; se adaugă;
 - PT: în mai multe grupuri de ordin superior; se adaugă;

- PD: într-un sigur grup de ordin superior: (PP, PT, PD, TD, DD); se elimină;
 - TP: nu există în grupuri de ordin superior; se adaugă;
 - TT: într-un sigur grup de ordin superior: (PP, PT, TT, TD, DD); se elimină;
 - TD: în mai multe grupuri de ordin superior; se adaugă;
 - DP: într-un sigur grup de ordin superior: (DP, DT); se elimină;
 - DT: într-un sigur grup de ordin superior: (DP, DT); se elimină;
 - DD: în mai multe grupuri de ordin superior; se adaugă;
- ÷ Toate grupurile distincte; lista finală: {(PP, PT, TT, TD, DD), (PP, PT, PD, TD, DD), (PP, PT, TD, DD), (PT, TD, DD), (PP, PT, DD), (PP, TD, DD), (PP, PT, TD), (PT, DD), (PP, DD), (TD, DD), (PT, TD), (PP, TD), (DP, DT), (PP, PT), PP, PT, TP, TD, DD};

La finele procesului, se poate întocmi lista posibilelor populații fenotipice împreună cu nivelul de încredere asociat. Pentru a obține nivelul de încredere asociat acestora, este necesară din nou inspecția valorilor conținute în [Tabelul din Anexa 7-2](#). Rezultatul acestei investigații este redat în [Tabelul 9-4](#) din Teză, clasificat după nivelul de confidență al rezultatului obținut.

Analiza numărului de asocieri (regresii) viabile din cultivar în generațiile ce au produs evoluție

Tabelul 9-6 din Teză cumulează rezultatele obținute pentru numărul de asocieri viabile (regresii cu parametrii semnificativ diferiți de zero statistic cu riscul de cel mult 5% de a fi în eroare) în ceea ce privește media numărului acestora (*avg_obs*) pentru același număr de evoluții observate (*num_obs*) fiind de asemenea grupate după mii de generații (de la 0..1000 până la 19001-20000).

Ipoteza cu privire la numărul mediu de asocieri viabile (valori redade în Tabelul 9-6 din Teză) și anume dacă metoda de selecție și strategia de supraviețuire produc populații distincte s-a verificat folosind [statistica Anderson-Darling](#) și rezultatul analizei este redat în [Tabelul din Anexa 7-3](#). [Tabelul din Anexa 7-3](#) arată că pentru un număr de 10 grupuri de perechi de metode nu s-a putut pune în evidență o diferență statistică semnificativă între legile de distribuție ale acestora. Interpretarea rezultatelor din [Tabelul din Anexa 7-3](#) se face în același mod în care s-au interpretat datele prezentate din [Tabelul din Anexa 7-2](#).

Cele mai mari grupuri de ordin maxim nediscriminate sunt grupurile de ordin 3 (PD, TD, DD) - run 66 - și (PP, PT, TD) - run 383 - care intră în mod automat în lista suspectilor; grupurile de ordin inferior sunt după cum urmează din aplicarea algoritmului de incluziune:

- ÷ Grupuri de ordin 2; lista suspectilor: {(PD, TD, DD), (PP, PT, TD)}
 - (PD, DD) - run 58, (TD, DD) - run 5, și (PD, TD) - run 65 - sunt doar în (PD, TD, DD); se elimină;
 - (TP, TT) - run 42 și (DP, DT) - run 3 - nu se regăsesc în nici unul din grupurile de ordin superior; se adaugă;
 - (PP, PT) - run 375, (PP, TD) - run 255, și (PT, TD) - run 128 - sunt doar în (PP, PT, TD); se elimină;
- ÷ Grupuri de ordin cel puțin 2; lista suspectilor: {(PD, TD, DD), (PP, PT, TD), (TP, TT), (DP, DT)}; se impune continuarea includerii în listă cu populațiile distincte, urmând aceeași procedură:
 - PP, PT se regăsesc într-un singur grup (PP, PT, TD); se elimină;
 - PD, DD: într-un sigur grup (PD, TD, DD); se elimină;
 - TP, TT: într-un sigur grup (TP, TT); se elimină;
 - DP, DT: într-un singur grup (DP, DT); se elimină;
 - TD: în două de ordin superior; se adaugă;
- ÷ Toate grupurile distincte; lista finală: {(PD, TD, DD), (PP, PT, TD), (TP, TT), (DP, DT), TD};

La finele procesului, se poate întocmi lista posibilelor asocieri fenotipice împreună cu nivelul de încredere asociat. Pentru a obține nivelul de încredere este necesară din nou inspecția valorilor conținute în [Tabelul din Anexa 7-3](#). Rezultatul acestei investigații este redat în Tabelul 9-7 din Teză, clasificat după nivelul de confidență al rezultatului obținut, care conține acele grupuri pentru care apartenența la populații identic distribuite nu a fost respinsă statistic cu o confidență de 95%.

Analiza obiectivului evoluției folosind un eșantion întâmplător de generații

Pentru această statistică, fiind variabilă continuă și producând variații doar la evoluție, s-a înlocuit ordonata *evoluție* (generație în care se produce o îmbunătățire a scorului funcției obiectiv) cu ordonata *generație*, cu ajutorul căreia salturile de evoluție se produc în medie mult mai lent, atenuând astfel variațiile observate, ceea ce permite observarea mult mai bună a funcției de distribuție a variabilei de interes (coeficientul de determinare). Procedura de transformare a datelor de la evoluții la generații este simplă, în generațiile intermediare evoluției statistica de interes având aceeași valoare cu cea pe care a avut-o în ultima sa evoluție.

Obținerea eșantionului din distribuția uniformă discretă

O primă ipoteză verificată în acest context a fost dacă generatorul de numere aleatoare ce provin din distribuția uniformă discretă urmează într-adevăr această distribuție, fiind cunoscute problemele care pot apărea din generarea pseudoaleatoare [26]. În acest sens, simulări Monte-Carlo recente ([Steele & alții, 2005](#)) au arătat că cele mai potrivite teste statistice sunt [Pearson-Fisher \$\chi^2\$](#) și [Anderson-Darling AD](#), fiecare dintre acestea obținând performanțe superioare celeilalte pentru diferite cazuri de legi de distribuție, dar amândouă fiind superioare comparativ cu alte teste statistice. Din acest motiv, verificarea ipotezei statistice folosind ambele teste este necesară. Tabelul de mai jos ([Tabelul 10-1](#) în Teză) sumarizează această procedură.

Ipoteza de distribuție uniformă a numerelor aleatoare generate cu [Rand\(·,·\)](#) în PHP (numerele vor fi generații)

9 valori generate aleatoriu între 0 și 20000	9221; 4182; 14283; 15329; 8875; 4599; 994; 8620; 7404
Valorile ordonate crescător	994; 4182; 4599; 7404; 8620; 8875; 9221; 14283; 15329
Probabilități cumulate observate - pco	1/9; 2/9; 3/9; 4/9; 5/9; 6/9; 7/9; 8/9; 9/9
Probabilități cumulate așteptate (teoretice) - pca	.11; .22; .33; .44; .56; .67; .78; .89; 1.0
Testul Anderson-Darling	1AD = 0.9687; cAD _{teoretic} = 2.5024; c/k = 2.58 > 1 (5%)
Testul Pearson χ^2	La 10 clase de frecvență $X^2 = 8.5$; df = 7; $\chi^2(8.5, 7) = 29\% > 5\%$
Testul Kolmogorov - Smirnov	D = 0.317 (0.778-0.461); $D\sqrt{9} = 0.95 = K^{-1}(9, 13.31\%)$

Concluzia simultană a testelor Anderson-Darling, Pearson χ^2 și KS: există o asociere probabilă care depășește 5% între cele două distribuții (cea teoretică și cea observată) și nu poate fi respinsă ipoteza că cele 18 valori provin din [distribuția uniformă discretă](#) 0..20000.

S-a investigat dacă numerele ar putea proveni din altă distribuție. Analiza este redată în [Tabelul 10-2](#) în Teză, când s-a constatat că acceptarea ipotezei că datele ar putea proveni dintr-o lege de distribuție (se respinge ipoteza că cele două legi de distribuție sunt diferite cu un risc de a fi în eroare de 10%, 5%, 2.5% sau 1%) nu exclude posibilitatea ca aceleași date să provină din altă lege de distribuție. În același timp, raportul între statistica ajustată ([Anderson-Darling](#)) și valoarea sa la pragul de 10% arată că cea mai probabilă lege de distribuție este distribuția uniformă (cu 2.3) urmată de distribuția valorilor extreme (cu 2.2) și Pareto generalizată (cu 2.1).

Căutarea legii de distribuție a coeficientului de determinare al unei generații

A doua ipoteză care trebuie verificată este dacă valorile coeficientului de determinare obținut în cele 46 de execuții independente ale algoritmului genetic pentru fiecare pereche de strategie de selecție și supraviețuire provin dintr-o populație distribuită după o lege de distribuție (având la dispoziție o serie de alternative). Pentru aceasta, fiecare din cele nouă generații extrase din distribuția uniformă a generațiilor 0..20000 a constituit subiectul investigației, rezultatele fiind redare

[26] Mads HAAHR. 1998-2009. ©. <http://www.random.org/randomness/>

în Tabelele 10-3-DD..10-3-TT din Teză (folosind alternativele de distribuție: DE - dublu exponențială; LG - logistică; LN - log-normală; NO - normală; UN - uniformă (0,1); EX - exponențială; GU - valori extreme de tip I) pentru fiecare asociere de strategie de selecție și supraviețuire incluse în studiu. Tabelele 10-3-DD..10-3-TT din Teză arată genul de rezultat care nu ar fi fost de dorit de obținut; practic aproape fiecare tip de selecție și încrucișare determină propria lege de distribuție a determinării. Tabelul 10-3-SS din Teză cumulează aceste observații. Tabelul de mai jos ([Tabelul 10-4](#) din Teză) re-sintetizează rezultatele din [Tabelul 10-3-SS](#) din Teză. Se observă că pentru selecția deterministă însoțită de supraviețuirea deterministă nici o lege nu a întrunit acceptarea.

Legi de distribuție ale coeficientului de determinare într-o generație pentru diferite strategii de selecție și supraviețuire

Lege de distribuție	PP	PT	PD	TP	TT	TD	DP	DT	DD
Valori extreme de tip I (GU)	X		X	X		X	X	X	
Log-normală (LN)	X				X				
Dublu exponențială (DE)		X							X*
Logistică					X				

Se observă că legea valorilor extreme de tip I (GU) este legea de distribuție cea mai frecventă a coeficientului de determinare într-o generație a algoritmului genetic rulat cu diferite strategii de selecție și supraviețuire. Rezultatele sintetizate sugerează drept lege de distribuție pentru coeficientul de determinare legea de distribuție Gumbel, și respinge fără drept de apel ipoteza de distribuție normală (Gauss) a acestuia (cu nici măcar o singură intrare).

Media nu este o statistică suficientă

O problemă statistică importantă și-a găsit răspuns în secțiunea anterioară, și anume:

÷ *Valoarea medie a coeficientului de determinare pentru fiecare generație (ca valoare medie din cele 46 de execuții independente) este o **statistică suficientă** pentru acest număr? Varianța sau abaterea standard este o statistică suficientă?*

Răspunsul la această întrebare este NU, motivul fiind faptul că aproape pentru nici o pereche de strategii selecție - supraviețuire nu s-a putut accepta ipoteza distribuției normale a valorilor, care ar fi acceptat valoarea medie din eșantionul de 46 de valori drept estimator nedepășat al mediei populației tuturor execuțiilor independente iar varianța drept statistică suficientă.

Distribuția Gumbel nu este suficient de generală

Au fost făcute estimările parametrilor pentru fiecare generație (0...20000) folosind aplicația [DataPlot](#) și au fost salvate în fișiere distincte pentru fiecare asociere de selecție și supraviețuire. Rezultatele au fost apoi interpolate folosind aplicația [SlideWrite](#).

$\hat{Y}(X) = \frac{X^\gamma - 1}{\beta}$	eq.1	eq.2	$\hat{Y}(X) = 0.907 \left(1 - \frac{X^2 + a_3 X + a_2}{X^2 + a_1 X + a_0} \right)$
---	------	------	---

A fost definită o funcție putere (eq. 1) și parametrii de interpolare (α , β , și γ), intervalele de încredere și semnificația statistică t sociată au fost obținute pentru fiecare caz în parte în ceea ce privește tendința centrală. A fost definită o funcție rațională (eq. 2) și parametrii de interpolare (α , β , și γ), intervalele de încredere și semnificația statistică t sociată. De verificat: e corect t sociată? au fost obținute pentru fiecare caz în parte în ceea ce privește valorile extreme (minimumul și maximumul). Această ultimă analiză, a valorilor extreme, a arătat că distribuția Gumbel nu poate fi acceptată drept lege de distribuție, probabilitățile de apariție ale valorilor extreme calculate fiind inacceptabil de depărtate de valorile aceluiași probabilități observate.

Analiza cuprinzând valoarea medie observată ca estimator al tendinței centrale, valorile minime observate și valorile maxime observate este redată în [Anexa 8](#). În [Anexa 8](#) nu a fost cuprinsă analiza parametrilor estimați folosind distribuția Gumbel, și anume cei menționați cu asterix în tabelul prezentat dintr-un motiv foarte simplu: [concluzia studiului cuprins în Anexa 8](#) a fost respingerea ipotezei de distribuție după [legea valorilor extreme de tip I](#).

Analiza legii de distribuție a obiectivului evoluției folosind un studiu sistematic în cadru generalizat

Legea de distribuție Fisher-Tippett a obiectivului evoluției

Întrucât analiza din tabelul de mai sus (Tabelul 10-4 în Teză) relevă faptul că pentru majoritatea cazurilor observate legea de distribuție Gumbel este totuși acceptată, consecința rezultată este că legea de distribuție urmată de observații este o lege de distribuție mai generală decât legea de distribuție Gumbel, dar care o are pe aceasta (Gumbel) drept caz limită. Într-adevăr, există această lege de distribuție, numită Legea Generală a Valorilor Extreme, sau legea Fisher-Tippett ([Fisher & Tippett, 1928](#)), lege care are drept caz limită legea Gumbel (a valorilor extreme de tip I) ca frontieră de separație între alte două familii de legi de distribuție, legea Weibull și legea Fréchet. Expresiile distribuției de probabilitate (PDF) și funcției cumulative de probabilitate (CDF) ale legii Fisher-Tippett (FT) sunt:

$$\text{FT}_{\text{PDF}}(X) = \begin{cases} \frac{1}{\beta} \exp\left(-\left(1+k \frac{x-\lambda}{\beta}\right)^{-1/k}\right) \left(1+k \frac{x-\lambda}{\beta}\right)^{-1-1/k}, & k < 0 \quad \text{Weibull} \\ \frac{1}{\beta} \exp\left(-\frac{x-\lambda}{\beta} - \exp\left(-\frac{x-\lambda}{\beta}\right)\right), & k = 0 \quad \text{Gumbel} \\ \frac{1}{\beta} \exp\left(-\left(1+k \frac{x-\lambda}{\beta}\right)^{-1/k}\right) \left(1+k \frac{x-\lambda}{\beta}\right)^{-1-1/k}, & k > 0 \quad \text{Fréchet} \end{cases}$$

$$\text{FT}_{\text{CDF}}(X) = \begin{cases} \exp\left(-\left(1+k \frac{x-\lambda}{\beta}\right)^{-1/k}\right), & k < 0 \quad \text{Weibull} \\ \exp\left(-\exp\left(-\frac{x-\lambda}{\beta}\right)\right), & k = 0 \quad \text{Gumbel} \\ \exp\left(-\left(1+k \frac{x-\lambda}{\beta}\right)^{-1/k}\right), & k > 0 \quad \text{Fréchet} \end{cases} \text{Fisher - Tippett}$$

O primă problemă care se s-a cerut rezolvată a fost verificarea ipotezei de distribuție după legea Fisher-Tippett a observațiilor experimentale. Pentru a realiza această problemă, s-a folosit aplicația EasyFit [27]. Distribuția Fisher-Tippett a fost o alternativă dintr-un număr de peste 55 de legi de distribuție în total, număr care conține atât legi de distribuție continue cât și discrete, distribuții mărginite și respectiv nemărginite. O serie numeroasă dintre aceste distribuții nu s-au calificat din start sau după primele testări ca distribuții posibile. Setul a fost restrâns la distribuții mărginite și distribuții generalizate, set din care în studiu au fost cuprinse Beta, Johnson, Kumaraswamy, Pert, Putere, Reciprocă, Triangulară, Uniformă (distribuții mărginite) și Fisher-Tippett, Pareto, Log-Pearson tip III (distribuții generalizate). Trei statistici: Pearson-Fisher Chi-Square, Anderson-Darling și Kolmogorov-Smirnov au evaluat agrementul între observație și model. Analiza este redată în [Anexa 9](#).

O a doua problemă care s-a cerut rezolvată a fost identificarea parametrilor distribuției FT și anume forma (k), locația (λ) și scala (β) din cele 46 de observații independente ale valorii coeficientului de determinare pentru fiecare din cele 20000 de generații de evoluție supuse observației. Pentru a realiza această problemă s-a folosit din nou aplicația EasyFitXL [27].

Rezultatele comparației distribuției observate cu distribuțiile teoretice listate în [Tabelul 1 - Anexa 9](#) suportă analize multicriteriale, dintre care după generație și apoi după legea de distribuție, clasificare menită să pună în evidență influența generației asupra abaterii observate de la legea de distribuție teoretică, și după metoda de selecție și de supraviețuire și apoi după legea de distribuție să se pună în evidență dacă asocierea de selecție și supraviețuire generează o lege proprie de distribuție.

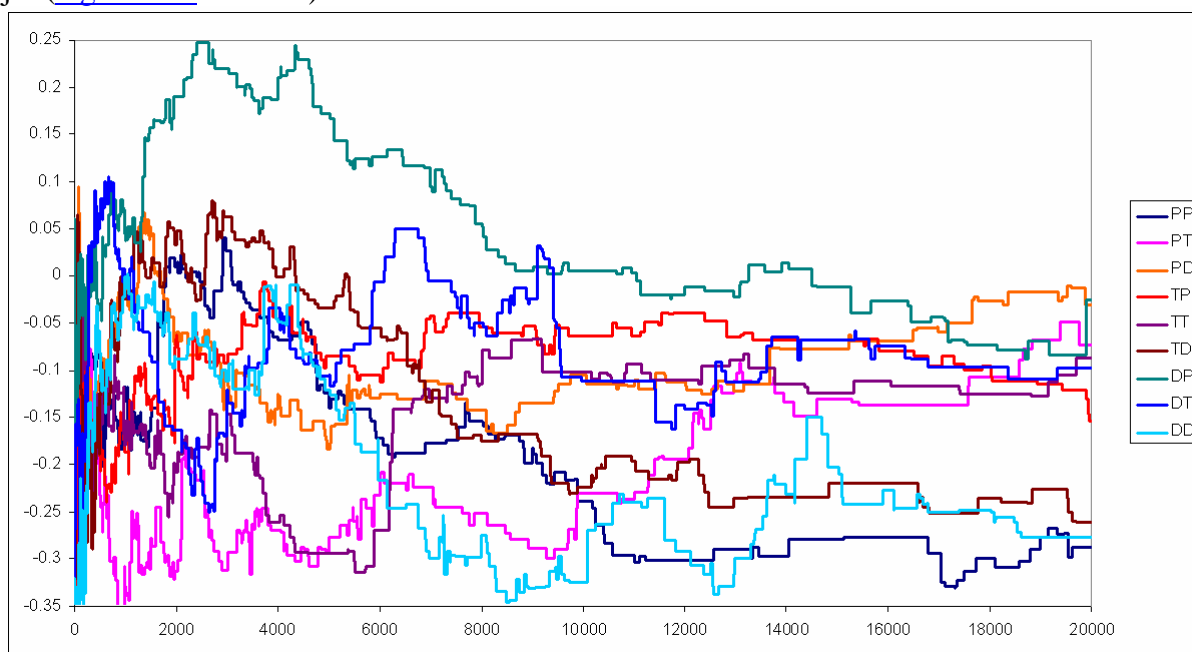
Însă cel mai important criteriu este respingerea statistică la unul sau mai multe nivele de

[27] [EasyFit Professional](#) v.50. 2008. Software. MathWave Technologies.

semnificație, și această analiză dată în [Tabelul 1 - Anexa 9](#) este sintetizată în Tabelul 11-1 din Teză. Rezultatele arată că distribuția Fisher-Tippett este singura care nu este respinsă semnificativ statistic (cu riscul de a fi în eroare de 1%) de cel puțin două din cele trei teste statistice utilizate, în timp ce celelalte două distribuții înregistrează un număr semnificativ de respingeri din totalul de 81 de eșantioane: Jonson SB - 20 (~25%), Pert - 18 (~22%).

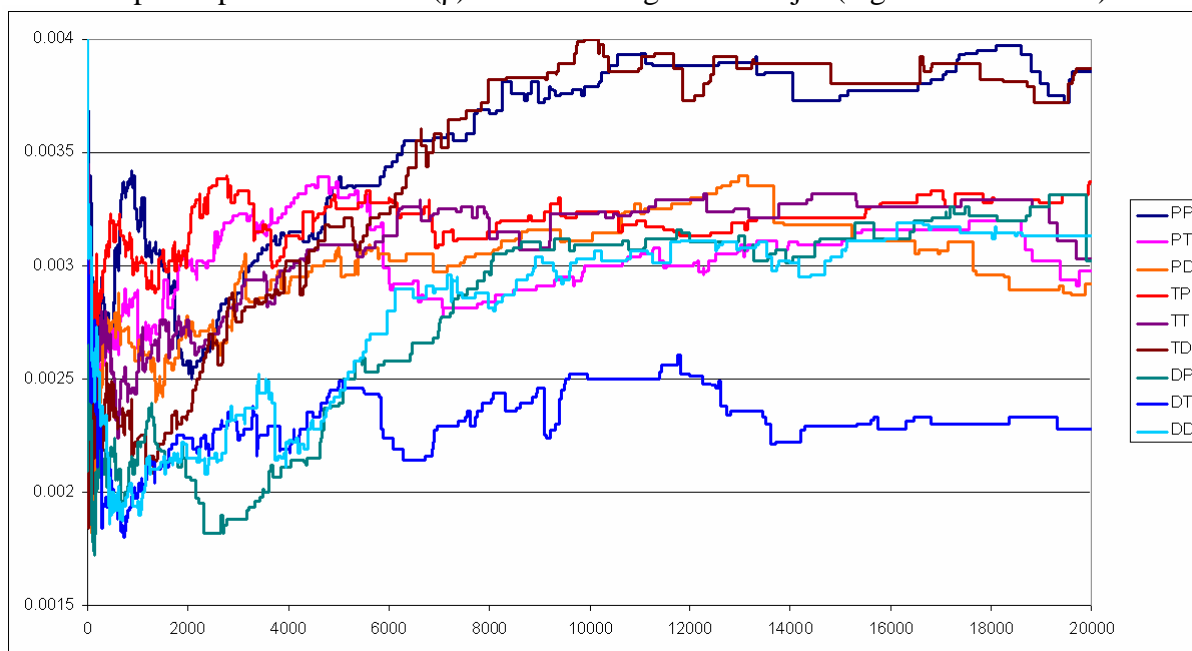
Evoluția locației, scalei și formei distribuției Fisher-Tippett a obiectivului evoluției

S-a folosit macro-ul pentru [Excel](#) al aplicației [EasyFit \(EasyFitXL\)](#) pentru calcularea parametrilor legilor de distribuție Fisher-Tippett pentru fiecare generație și fiecare asociere de selecție și supraviețuire. Estimările parametrilor locație (λ), scală (β) și formă (k) au fost obținute în fiecare caz în parte folosind principiul maximizării ratei șansei dintr-un volum de 46 de observații, și anume cele 46 de execuții independente ale algoritmului genetic. S-a folosit aplicația [Statistica](#) pentru uniformizări exponențiale. Rezultatul pentru parametrul formei (k) este redat în figura de mai jos ([Figura 11-2](#) din Teză).



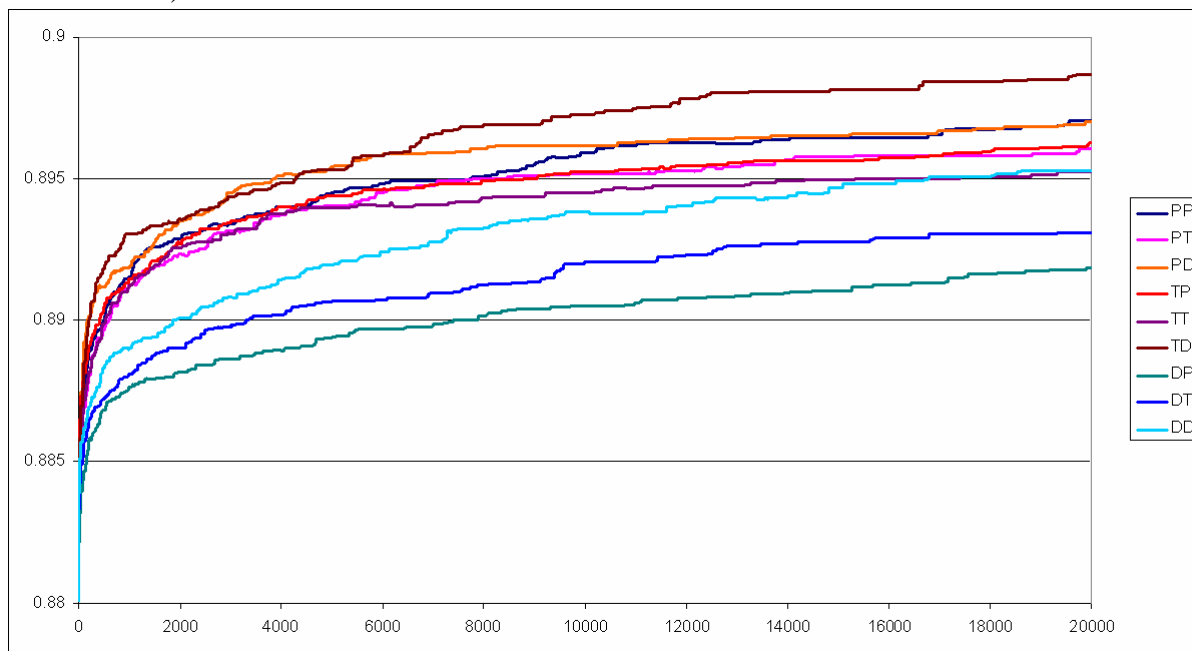
Parametrul formă (k) al distribuției Fisher-Tippett: estimare din observații

Rezultatul pentru parametrul scalei (β) este redat în figura de mai jos ([Figura 11-6](#) din Teză).



Parametrul scală (β) al distribuției Fisher-Tippett: estimare din observații

Rezultatul pentru estimarea parametrului locație (λ) este redat în figura de mai jos (Figura 11-9 din Teză).



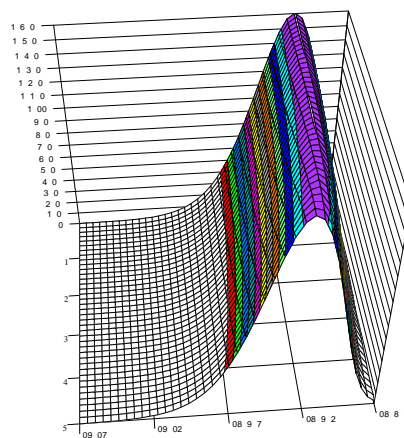
În valorile parametrilor distribuțiilor Fisher-Tippett s-au găsit ecuații de tendință, care sunt redate în tabelul de mai jos (Tabelul 11-7 din Teză).

Ecuțiile de tendință pentru formă (k), scală (β) și locație (λ) ale distribuțiilor Fisher-Tippett

SS	$k(G) = a_0 + a_1 \cdot G$		$\beta(G) = a_0 + a_1 \cdot G$		Tendință $\lambda(G)$	a_0	a_1	a_2
PP	-0.1912	$-1.47 \cdot 10^{-6}$	3.541E-3	5.5E-9	$\lambda(G) = a_0 + a_1 \cdot \ln(G + a_2)$	0.89357	$1.82 \cdot 10^{-4}$	0.867
PD	-0.0961	$3.12 \cdot 10^{-7}$	2.983E-3	1.9E-9		0.89422	$1.55 \cdot 10^{-4}$	-0.344
TP	-0.0833	$1.24 \cdot 10^{-7}$	3.192E-3	8.9E-10		0.89333	$1.54 \cdot 10^{-4}$	-0.213
TT	-0.1476	$5.58 \cdot 10^{-7}$	3.072E-3	2.9E-9		0.89286	$1.40 \cdot 10^{-4}$	-0.348
PT	-0.2108	$1.08 \cdot 10^{-6}$	2.996E-3	8.2E-10	$\lambda(G) = a_0 + a_1 \cdot \ln(G)$	0.89309	$1.69 \cdot 10^{-4}$	
TD	-0.1352	$-1.47 \cdot 10^{-6}$	3.419E-3	7.9E-9	$\lambda(G) = a_0 + a_1 \cdot G^{a_2}$	0.89465	$6.84 \cdot 10^{-4}$	0.117
DP	-0.0193	$-1.32 \cdot 10^{-6}$	2.730E-3	7.1E-9		0.88916	$2.02 \cdot 10^{-4}$	0.171
DT	-0.0797	$-1.35 \cdot 10^{-7}$	2.296E-3	6.1E-10		0.89016	$3.19 \cdot 10^{-4}$	0.151
DD	-0.0207	$-9.52 \cdot 10^{-7}$	2.745E-3	5.6E-9		0.89173	$2.93 \cdot 10^{-4}$	0.172

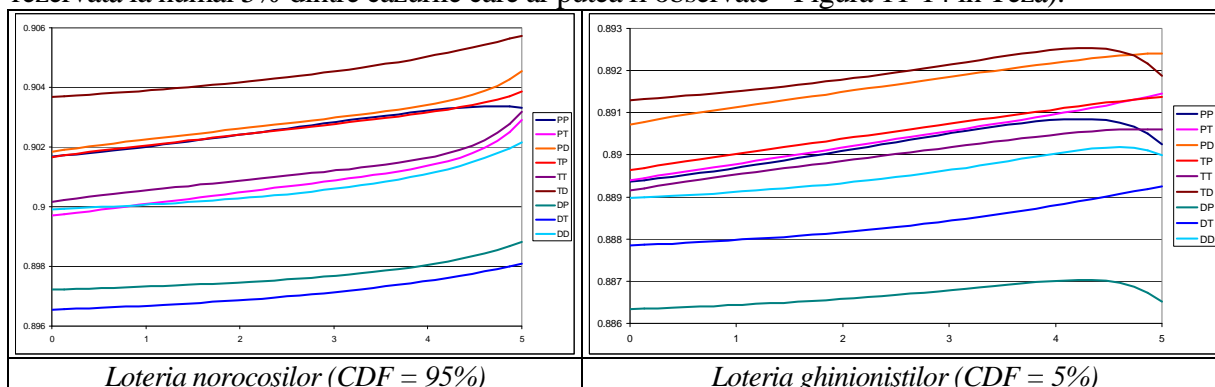
Similar cu ecuațiile de mai sus se pot obține și expresiile care dau tendința funcției de densitate de probabilitate (PDF). Acestea sunt însă mult mai complicate pentru a fi redate sub formă de expresii matematice, însă mult mai sugestivă este reprezentarea tridimensională a acestora. [Figura 1 - Anexa 10](#) redă reprezentările tridimensionale ale funcției de densitate de probabilitate în care în locul variabilei generație s-a folosit logaritmul în baza 10 al acesteia (scară logaritmică).

Figura alăturată este extrasă din [Figura 1 - Anexa 10](#) și reprezintă tendința în densitatea de probabilitate Fisher-Tippett a coeficientului de determinare pentru strategia de selecție în turnir și strategia de supraviețuire deterministă - DTFT_{PDF}($r^2, \log_{10}G$).



Caracterizarea statistică a distribuției evoluției

Folosind ecuațiile de tendință pentru formă (k), scală (β) și locație (λ) se pot obține frontierele de probabilitate 95% (în care șansa de a fi atinsă o valoare superioară acestei frontiere îi este rezervată la numai 5% dintre cazurile care ar putea fi observate - Figura 11-13 în Teză) și respectiv de probabilitate 5% (în care șansa de a fi atinsă o valoare inferioară acestei frontiere îi este rezervată la numai 5% dintre cazurile care ar putea fi observate - Figura 11-14 în Teză).



Se remarcă în figurile de mai sus comportarea selecției deterministe însoțite de supraviețuirea în turnir (DT) care este singura asociere de selecție și supraviețuire al cărui interval de normalitate (fără șansă/noroc și fără ghinion) rămâne aproximativ același pe parcursul evoluției. Tendința generală observată este că norocul crește ceva mai accentuat decât în tendință liniară odată cu evoluția - doar selecția proporțională însoțită de supraviețuire proporțională (PP) se abate de la această regulă în timp ca ghinionul crește ceva mai puțin accentuat decât în tendință liniară odată cu evoluția (doar PT și DT se abat de la această regulă).

Estimarea parametrilor distribuției Fisher-Tippett pornind de la observațiile coeficientului de determinare în cele 46 de execuții independente face posibilă urmărirea valorilor parametrului statistic k (forma distribuției) pentru fiecare generație observată (0..20000) și pentru fiecare asociere de selecție și supraviețuire. În tabelul următor sunt redată frecvențele observate în raport cu cele 3 categorii definite (Tabelul 11-8 din Teză).

Tipul formei distribuției Fisher-Tippett

Valori extreme de tip	PP	PT	PD	TP	TT	TD	DP	DT	DD
I ($ k < 10^{-2}$) \approx Gumbel	778 (3.9%)	0 (0%)	317 (1.6%)	63 (0.3%)	23 (0.1%)	992 (5%)	3237 (16.2%)	1091 (5.5%)	292 (1.5%)
II ($k > 10^{-2}$) = Fréchet	324 (1.6%)	0 (0%)	299 (1.5%)	0 (0%)	36 (0.2%)	2158 (10.8%)	9012 (45.1%)	1619 (8.1%)	0 (0%)
III ($k < -10^{-2}$) = Weibull	18899 (94.5%)	20001 (100%)	19385 (96.9%)	19938 (99.7%)	19942 (99.7%)	16851 (84.3%)	7752 (38.8%)	17291 (86.5%)	19709 (98.5%)

Ce lege urmează momentele de apariție a evoluției?

Familia de curbe log-Pearson de tipul III

În această secțiune s-a căutat răspuns le următoarea întrebare: *Ce lege urmează momentele de apariție a evoluției?*

În acest scop, s-a aplicat o procedură de transformare asupra datelor primare, procedură exemplificată în tabelul de mai jos (Tabelul 12-1 în Teză, exemplu din run 1, DP).

Transformarea momentelor evoluției la durate relative

Numărul generației	0	15	136	188	246	528	5423	11887
Momentele evoluției	1	16	137	189	247	529	5424	11888
Durata până la viitoarea evoluție	15	121	52	58	282	4895	6464	
Durata relativă la momentul evoluției	1.5E+1	7.6E+0	3.8E-1	3.1E-1	1.1E+0	9.3E+0	1.2E+0	1.7E-4

Un prim răspuns la întrebarea “*Ce lege urmează momentele de apariție a evoluției?*” se poate obține completând întrebarea astfel “*Ce lege urmează momentele de apariție a evoluției independent de strategia de evoluție urmată?*”.

Răspunsul la această întrebare a fost obținut astfel: s-au obținut momentele relative așa cum descrie tabelul de mai sus pentru fiecare run (1..46) și fiecare strategie (PP, PT, PD, TP, TT, TD, DP, DT, DD); fiecare șir de tipul celui din linia “Durata relativă la momentul evoluției” a tabelului de mai sus reprezintă o observație a evoluției; două astfel de linii reprezintă două observații ale aceluiași fenomen, dar în același timp poate fi privită ca o singură observație din punctul de vedere al distribuției valorilor, una venind să o completeze pe cealaltă în ceea ce privește observarea momentelor; prin extensie, toate la un loc (în număr de 414) sunt cea mai largă observație a momentelor evoluției care poate fi constituită din datele înregistrate din execuția algoritmului genetic; [fișierul rezultat](#) descris din procedura de mai sus cumulează 11347 momente de evoluție; acestea au intrat în analiza de distribuție; analiza de distribuție a fost realizată cu EasyFit având ca alternative un număr de 65 de legi de distribuție continue; au fost folosite pentru? măsurarea agrementului între observație și model un număr de trei statistici (C-S, A-D și K-S); [rezultatele obținute](#) au fost edificatoare cu privire la ipoteza de distribuție care poate fi formulată asupra momentelor evoluției; prima porțiune a acestei analize este redată în tabelul de mai jos (Tabelul 12-2 în Teză, cuprinzând primele trei legi de distribuție în ordinea agrementului între observație și model).

Cele mai probabile legi de distribuție pentru momentele relative ale evoluției

Dist\Stat	K-S	p(K-S)	Rang	A-D	p(A-D)	Rang	C-S(df)	p(C-S)	Rang
Log-P-3	0.01197	0.07683	1	2.4264	0.05617	1	41.731(13)	7.3E-05	1
Burr	0.01635	4.57E-03	3	6.7901	3.23E-04	3	46.345(13)	1.25E-05	2
Burr-4P	0.01592	6.27E-03	2	6.0813	7.48E-04	2	51.408(13)	1.71E-06	3

Dist: Lege de distribuție; Stat: Statistică; Rang: Rangul statisticii în lista celor 65 de alternative
Log-P-3: log-Pearson de tipul 3

Rezultatele din tabelul de mai sus arată următoarele:

1. O singură lege de distribuție din cele 65 de alternative se califică ca ipoteză de distribuție pentru momentele relative ale evoluției; este singura pentru care se obțin riscuri de a fi în eroare la respingerea distribuției mai mari de 1%, care sunt de fapt în cazul distribuției Log-Pearson 3 mai mari de 5% (7.68% K-S și 5.62% A-D).
2. Din modalitatea de calcul se poate evidenția că statistica K-S măsoară agrementul între rangurile observațiilor, în timp ce statistica C-S măsoară agrementul între valorile observațiilor, iar statistica A-D chiar dacă folosește ranguri (ca și K-S) este totuși o măsură care o apropie de C-S;
3. Statistica C-S pentru Log-Pearson 3 are valoarea $7.3 \cdot 10^{-5}$, adică observații mai defavorabile agrementului cu modelul dat de Log-Pearson 3 se obțin în mai puțin de 0.08% din cazuri. valoarea probabilității de observare scade de la 7.7% pentru K-S, la 5.6% pentru A-D ca să

- ajungă la 0.08% pentru C-S.
- Având în vedere ce măsoară statisticile C-S, A-D și K-S (remarca 2 de mai sus) și valorile obținute pentru aceste statistici (remarca 3 de mai sus) se desprinde concluzia că agrementul între rangurile observate ale momentelor relative ale evoluției și rangurile din distribuția teoretică Log-Pearson 3 este mai probabil decât agrementul între valorile observate ale momentelor relative ale evoluției și valorile din distribuția teoretică Log-Pearson 3.
 - Există o explicație pentru remarca 4, și anume eșantionul de evoluții a cuprins observații din toate cele nouă asocieri de strategii; este astfel așteptat ca agrementul să fie mult mai mare pentru ranguri (statisticile K-S și A-D) decât pentru valori (statistica C-S).
 - Având în vedere remarcile 4 și 5 de mai sus, concluzia care se obține în urma analizei este că se poate accepta cu riscul de a fi în eroare de 5% (superior probabilităților de a greși din Tabelul 12-2 pentru Log-Pearson 3 folosind statisticile K-S & A-D) că Log-Pearson 3 este legea de distribuție a momentelor relative ale evoluției obținute prin procedura menționată.

Concluzia că “*Log-Pearson 3 este legea de distribuție a momentelor relative ale evoluției, independent de strategia de selecție și de supraviețuire*” se poate verifica. În acest sens, pentru fiecare pereche de strategie de selecție și supraviețuire din cele 46 de execuții independente au fost cumulate momentele relative ale evoluției. Rezultatele sunt disponibile online:

http://l.academicdirect.org/Horticulture/GAs/MLR_MDF_selection_vs_survival/evolut_SS.txt

Rezultatele au intrat în aceeași analiză a agrementului între observație și model, și fișierele rezultat se află la aceeași adresă de mai sus, denumirile fișierelor fiind date în tabelul de mai jos (Tabelul 12-3 în Teză), împreună cu semnificațiile statistice ale agrementului între observații și modelul dat de distribuția Log-Pearson 3.

Agrementul între observație și modelul Log-Pearson 3 pentru distribuția momentelor relative ale evoluției

Stra\Stat	nr.Obs	K-S	p _{K-S}	A-D	p _{A-D}	C-S(df)	p _{C-S}
DD	TT	1379	0.02284	0.46	0.63251	0.47	12.3(10)
DP	TD	1429	0.01224	0.98	0.23477	0.75	3.3064(10)
DT	TP	1318	0.02691	0.29	1.2118	0.24	14.35(10)
PD	DT	996	0.02845	0.39	0.73496	0.41	10.628(9)
PP	DD	1084	0.01919	0.81	0.34184	0.66	8.1401(10)
PT	DP	851	0.02416	0.69	0.6234	0.47	6.8598(9)
TD	PT	1463	0.0203	0.58	0.70531	0.43	12.512(10)
TP	PD	1474	0.03055	0.13	0.93998	0.33	8.6564(10)
TT	PP	1353	0.01212	0.99	0.23201	0.75	3.5574(10)
Stra (DD, DP, DT, PD, PP, PT, TD, TP, TT): strategie							
Stat (nr.Obs, K-S, p _{K-S} , A-D, p _{A-D} , C-S(df), p _{C-S}): statistică							

Rezultatele din tabelul de mai sus demonstrează că ipoteza “*Log-Pearson 3 este legea de distribuție a momentelor relative ale evoluției, independent de strategia de selecție și de supraviețuire*” formulată pe baza analizei cuprinzând cumulat toate observațiile și se verifică pentru fiecare strategie de selecție și supraviețuire în parte.

Valorile din tabelul de mai sus evidențiază că: nu există nici o respingere a ipotezei formulate la un risc de a fi în eroare de 10% sau mai mic; există două respingeri (din 27 de cazuri) la un risc de a fi în eroare de 20% sau mai mic (PD pentru K-S & TP pentru C-S) care este o situație așteptată (fixându-se un risc de a fi în eroare de 20% s-a făcut în fapt o eroare inferioară lui 20% de $2/27 = 7.4\%$); dacă statistica C-S respingea ipoteza cu un risc de a fi în eroare mai mic de 0.08% pentru eșantionul reunind observațiile de la toate strategiile, atunci când se analizează separat fiecare strategie în parte valoarea riscului de a fi în eroare în respingerea ipotezei urcă dramatic, cea mai mică valoare a acestuia fiind 16% și având o valoare medie de 53%; explicația pentru dezagrementul măsurat de statistica C-S între observație și modelul Log-Pearson 3 remarcat în experimentul cu observațiile cumulate și explicația de mai sus pentru agrementul măsurat de statistica C-S între observație și modelul Log-Pearson 3 remarcat în experimentul cu observațiile

separate pe strategii aduce pe cale de consecință că:

- ÷ ipoteza “Log-Pearson 3 este legea de distribuție a momentelor relative ale evoluției *independent de strategia de selecție și de supraviețuire*” este susținută de observațiile asupra fiecărei strategii în parte;
- ÷ parametrii legii de distribuție Log-Pearson 3 sunt însă dependenți de strategia de selecție și de supraviețuire (dezagrement la X^2 la cumularea observațiilor, agrement la X^2 individual);
- ÷ nu există nici un motiv pentru a presupune că parametrii legii de distribuție Log-Pearson 3 sunt dependenți de altceva decât de strategia de selecție și de supraviețuire (pentru experimentul desfășurat) - motivul - agrementul la observațiile ce provin dintr-o anumită strategie înregistrat în toate cele nouă asocieri de strategii de selecție și supraviețuire;

Sintetizând, o concluzie importantă a fost extrasă din studiul momentelor relative ale evoluției, și anume că “Log-Pearson 3 este legea de distribuție a momentelor relative ale evoluției, iar parametrii distribuției depind de strategia de selecție și de supraviețuire”.

Asocieri în familia de curbe de evoluție

Există asocieri între parametrii formă, scală și locație ai distribuțiilor Log-Pearson 3 obținuți pentru momentele relative ale evoluției urmând diferite strategii (Tabelul 12-4). Analiza de corelație susține această afirmație: $r(\alpha, \beta) = 0.857$; $r(\alpha, \gamma) = 0.994$; $r(\beta, \gamma) = 0.885$. $r(\alpha, \gamma) = 0.994$ arată că există o asociere liniară între formă (α) și locație (γ). S-a evidențiat de asemenea existența unei dependențe după o funcție putere ($r(\beta, \text{pow}(\gamma)) > 0.999$) între scală (β) și locație (γ). Astfel familia de curbe log-Pearson de tipul III este degenerată în caracterizarea evoluției:

$$LP3(x; \alpha_j, \beta_j, \gamma_j) \text{ degenerază în } LP3(x; 8.77 \cdot \gamma_j - 68.3, -0.14 - 144 \cdot \gamma_j^{-2.57}, \gamma_j)$$

unde $j \in \{TT, TD, TP, DT, DD, DP, PT, PD, PP\}$.

Pentru a verifica ipoteza de degenerare sugerată de regresiiile parametrilor, s-a estimat din nou parametrul locație pentru expresiile degenerate ale distribuției log-Pearson de tipul III cu un singur parametru independent, și tabelul următor (Tabelul 12-10 în Teză) prezintă pentru comparație valorile probabilităților de observare ale distribuțiilor din eșantioane în ipoteza că ele provin din legea de distribuție log-Pearson de tipul III.

Agrementul observație - model Log-Pearson 3 pentru distribuția momentelor relative ale evoluției în ipoteza de asociere liniară între formă și locație și neliniară între scală și locație

SS	nr.Obs	p _{K-S}		p _{A-D}		p _{C-S}	
TT	1379	0.46	0.09	0.47	0.17	0.27	0.12
TD	1429	0.98	0.98	0.75	0.74	0.97	0.77
TP	1318	0.29	0.30	0.24	0.19	0.16	0.10
DT	996	0.39	0.47	0.41	0.52	0.3	0.55
DD	1084	0.81	0.88	0.66	0.66	0.62	0.47
DP	851	0.69	0.14	0.47	0.15	0.65	0.21
PT	1463	0.58	0.68	0.43	0.46	0.25	0.36
PD	1474	0.13	0.08	0.33	0.24	0.56	0.44
PP	1353	0.99	0.90	0.75	0.64	0.97	0.80

÷ p_{K-S}, p_{A-D}, p_{C-S}:
probabilități de observare întâmplătoare

÷ prima valoare p:
din MLE parametrii independenți

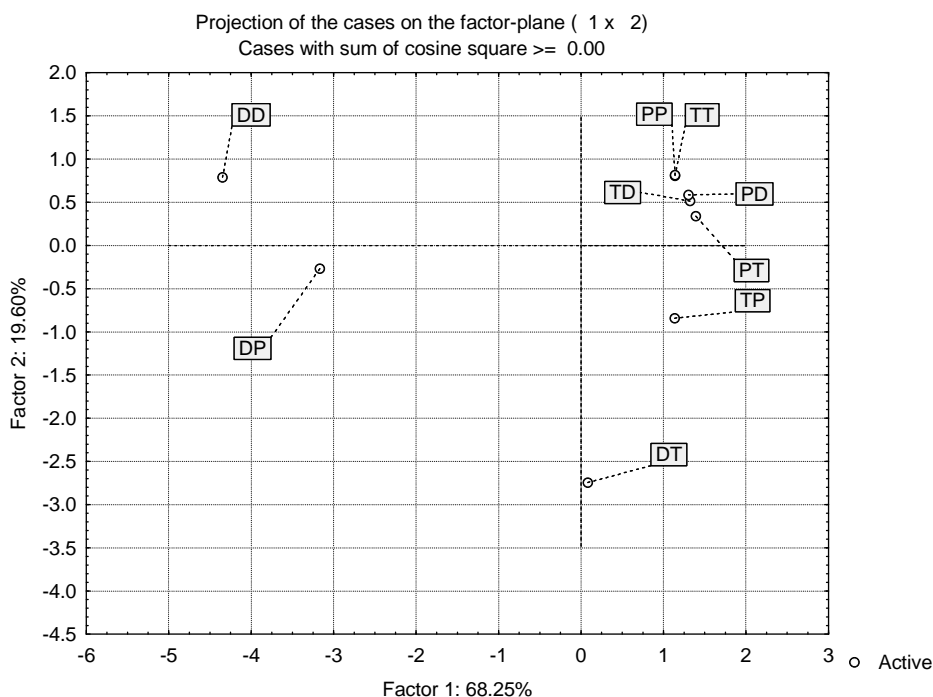
÷ a doua valoare p:
din MLE cu un parametru independent (γ)

Analiza sintetizată în tabelul de mai sus nu oferă nici un motiv de a respinge ipoteza formulată de asociere liniară între parametrul formă (α) și parametrul locație (γ) și de asociere neliniară între parametrul scală (β) și parametrul locație (γ) ai distribuțiilor log-Pearson de tipul III pentru momentele relative ale evoluției în seria de strategii de evoluție investigată (DD, DP, DT,

PD, PP, PT, TD, TP, TT).

Diferențe în familia de curbe de evoluție

S-au pus în evidență o serie de asocieri între valorile statisticilor distribuțiilor degenerate. Astfel, valori pereche au: DP și DD pentru locație, medie, abatere standard, oblicitate și exces de boltire, DT și TP pentru locație și modă, și așa mai departe. Pentru a evidenția aceste asocieri de valori între statisticile distribuțiilor (locație, medie, modă, mediană, deviație standard, asimetrie, exces de boltire), s-a construit o analiză de componente principale (figura de mai jos, în Figura 12-3 în Teză, obținută folosind aplicația [Statistica](#)).



Proiecțiile primilor doi factori (principali) în valorile $(\gamma, \mu, \hat{\mu}, \tilde{\mu}, \sigma, \gamma_1, \gamma_2)$

Figura de mai sus (Figura 12-3 în Teză) evidențiază următoarele: strategiile DD și DP definesc și se află aproape în lungul unuia din cei doi factori principali - cel cu 68.25% explicare cantitativă - datorită motivelor evidențiate în analiza valorilor din - Tabelul 12-11 din Teză - și anume valorile mari pe care aceste strategii le produc pentru majoritatea statisticilor; strategia DT se află aproape în lungul unuia din cei doi factori principali - cel cu 19.6% explicare cantitativă - în valorile din Tabelul 12-11 din Teză - evidențiindu-se doar cu locație, medie și abatere standard relativ mari în raport cu celelalte strategii, dar cu valori mici - de fapt cele mai mici - pentru asimetrie și exces de boltire; un grup de strategii - PP, TT, TD, PD și PT este poziționat în planul primilor doi factori principali compact cu diferențe extrem de mici; strategia TP este situată în lungul celui de-al doilea factor principal la distanță relativ egală de grupul compact de strategii (format din PP, TT, TD, PD și PT) și strategia DT, situându-se în cadranul acesteia din urmă.

Ce distribuție urmează numărul de evoluții?

S-au derulat 10 analize, 9 dintre ele cu datele obținute în fiecare strategie în parte din fiecare execuție (46 observații în eșantionul numărului de evoluții) și a 10-a cu toate la un loc (414 observații în eșantionul numărului de evoluții). A fost folosită aplicația [EasyFit](#) pentru a obține estimările parametrilor distribuțiilor și a da măsurile agrementelor. S-a întocmit un clasament pe ranguri după statistici (C-S, A-D, K-S) care a dat cea mai probabilă distribuție ca fiind Fisher-Tippett (rang 284 față de peste 420 restul alternativelor în număr de peste 60).

Tabelul următor (Tabelul 12-14 în Teză) conține agrementul între observații și modelul de distribuție Fisher-Tippett pentru fiecare din cele 10 eșantioane care au fost supuse analizei.

Agrementul pe care îl realizează distribuția Fisher-Tippett cu observațiile experimentale ale numărului de evoluții este de-a dreptul remarcabil. Nu numai că nu se semnalează respingeri semnificativ statistice la nici un risc de a fi în eroare uzual între 1% și 20%, ci mai mult, agrementul

între model și date este în medie de 86.51% în acord cu statistica Kolmogorov-Smirnov, de 72.33% în acord cu statistica Anderson-Darling și de 71.73% în acord cu statistica Chi-Square.

Agrementul între Fisher-Tippett și numărul de evoluții dintr-o execuție a algoritmului genetic

Strategie	Obs	K-S	P _{K-S}	A-D	P _{A-D}	C-S/df	P _{C-S}
TT	46	0.0924	0.7931	0.4183	0.6028	5.17/5	0.3956
TD	46	0.1199	0.4859	0.5976	0.4877	3.57/4	0.4671
TP	46	0.0454	0.9999	0.0818	0.8972	0.96/5	0.9661
DT	46	0.0632	0.9873	0.2303	0.7527	1.27/5	0.9381
DD	46	0.0615	0.9906	0.215	0.7665	0.72/5	0.9816
DP	46	0.0954	0.7612	0.2766	0.7127	3.76/4	0.4389
PT	46	0.0712	0.9608	0.2052	0.7754	4.23/5	0.5171
PD	46	0.0634	0.9869	0.1693	0.8090	0.99/5	0.9632
PP	46	0.0665	0.9787	0.2428	0.7417	0.69/5	0.9835
Toate	414	0.0342	0.7066	0.307	0.6875	7.14/8	0.5218

În cea mai mică măsură agrementul a fost înregistrat de statistica Chi-Square (39.56%) pentru strategia TT și de statisticile Kolmogorov-Smirnov (48.59%) și Anderson-Darling (48.77%) pentru strategia TD. Tabelul următor (Tabelul 12-15 în Teză) redă statisticile comune cu privire la legile de distribuție obținute.

Statistici ale distribuțiilor Fisher-Tippett ale numărului de evoluții către optim

Strategia	Distribuția F-T(α ; β ; γ)	μ	$\hat{\mu}$	$\tilde{\mu}$	σ	γ_1	γ_2
TT	F-T(-0.0771; 8.0028; 26.929)	31.0	28	29.8	9.38	0.739	0.849
TD	F-T(-0.19367; 8.9378; 28.367)	32.1	30	31.5	9.44	0.276	-0.095
TP	F-T(0.04267; 8.7648; 24.208)	29.7	24	27.4	11.93	-1.420	3.975
DT	F-T(-0.0309; 7.0811; 18.775)	22.7	19	21.4	8.74	0.966	1.635
DD	F-T(-0.30349; 9.3813; 21.38)	24.6	25	24.6	9.26	-0.079	-0.289
DP	F-T(-0.27344; 8.0192; 16.622)	19.5	19	19.4	8.05	0.013	-0.280
PT	F-T(-0.15998; 8.6245; 29.02)	32.8	31	32.1	9.35	0.398	0.074
PD	F-T(-0.12837; 9.3279; 28.721)	33.0	30	32.1	10.39	0.520	0.299
PP	F-T(-0.24824; 9.8865; 26.7)	30.4	29	30.2	10.07	0.093	-0.249
Toate	F-T(-0.16044; 9.6882; 24.161)	28.4	26	27.6	10.50	0.396	0.072

μ : Media; $\hat{\mu}$: Moda; $\tilde{\mu}$: Mediana; σ : Deviația standard; γ_1 : Asimetria; γ_2 : Excesul de boltire

Rezultatele din tabelul de mai sus evidențiază că legile de distribuție sunt foarte apropiate una de cealaltă. În fapt, nu a fost respinsă nici ipoteza că toate împreună au o singură lege de distribuție și strategia nu creează populații distincte în numărul de evoluții. Se poate remarca însă că valori la modă sub moda globală (26) au strategiile de selecție deterministă (DP: 19; DT: 19; DD: 25) și selecția în turnir cuplată cu supraviețuirea proporțională (TP: 24) în timp ce valori la modă superioare modei globale (26) au strategiile de selecție proporțională (PT: 31; PD: 30; PP:29) și în turnir (TD: 30; TT: 28) - cu o excepție, cea de mai sus (TP: 24). Tot rezultatele din tabelul de mai sus evidențiază că o singură lege de distribuție este a valorilor extreme de tipul II (Fréchet) - TP ($\alpha > 0$), toate celelalte, inclusiv distribuția globală fiind a valorilor extreme de tipul III (Weibull). Aceeași strategie (TP) produce și cel mai mare exces de boltire (3.975) fiind aproape dublu decât în cazul oricărei alte strategii. Dacă se admite că valoarea asimetriei pentru strategia DD (-0.079) este aproape nulă, atunci strategia TP rămâne singura cu asimetrie negativă (-1.42), fiind în același timp distribuția cu cea mai mare deviație standard (11.93). Un calcul simplu asupra deviației standard date în tabelul de mai sus în ipoteza de independență între eșantioanele extrase din strategii diferite ($\sigma_{\Sigma}^2 = (\sigma_{TT}^2 + \dots + \sigma_{PP}^2)/9$) ne permite să separăm varianța totală ($\sigma_T^2 = 10.5^2$) în varianță în interiorul strategiilor ($\sigma_{\Sigma}^2 = 9.68^2$) și varianță între strategii (4.07^2). Acest din urmă rezultat pune în evidență o măsură în care alegerea unei strategii influențează evoluția, contribuția cantității varianței între strategii fiind importantă în valoarea varianței totale.

Concluzii și recomandări

Prezenta lucrare demonstrează că algoritmi genetici, ca tehnici adaptive de căutare euristică, bazate pe principiile geneticii și selecției naturale, pot fi eficient utilizați în simularea procesului biologic al evoluției și în cel de ameliorare a plantelor.

Modelele informatice elaborate prin intermediul algoritmilor genetici, emulează modelele biologice evoluționiste, asigurând rezolvarea unor probleme concrete de optimizare sau căutare în experiențele de genetica și ameliorarea plantelor. Prin intermediul elementelor individuale, reprezentate sub forma șirurilor binare, și a operatorilor de natură biologică definiți asupra populației și a modelului molecular, algoritmi genetici manipulează cele mai promițătoare șiruri, evaluate conform unei funcții obiectiv, căutând soluții mai bune, tinzând în esență spre cea “optimă”, dar acceptând în final una apropiată de optim.

(Concluzie)

În lucrare, eficiența aplicării algoritmilor genetici în optimizarea relației structură-activitate în seria de compuși PCB, care prezintă potență biologică distructivă asupra mediului vegetal și animal, a fost probată într-un experiment de evoluție, folosind diferite strategii pentru estimarea efectului procesului de selecție și supraviețuire asupra indivizilor în cadrul populațiilor.

(Concluzie)

S-a constatat că eșantionul de genotipuri supus evoluției tinde relativ rapid către optim, probabilitățile din funcțiile cumulative de distribuție obținute asigurând obținerea a 99% din optim în 1000 de generații (analizându-se mai puțin de $2 \cdot 10^{-11}$ din numărul de regresii posibile în întreaga populație) pentru strategia TD în 55% din cazuri, PD - 67%, PP - 68%, TP - 73%, PT - 78%, TT - 80%, DD - 87%, DP - 95% și DT în 97% din cazuri.

(Recomandare)

Asigurarea condițiilor optime pentru plantele cultivate, cu privire la ansamblul cerințelor tehnologice, incluzând necesarul de nutrienți și, în particular, a necesarului de apă, este esențială pentru reușita culturilor horticoale. Acestea pot contribui semnificativ la productivitatea și eficiența economică a culturilor horticoale, precum și la calitatea materialului biologic și, în mod deosebit, a calității fructelor. Controlul optimal al proceselor fiziologice de creștere, dezvoltare, fructificare etc. în asociere cu genotipul (cultivarul), necesarul de apă și nutrienți pentru culturile de câmp sau seră (spații protejate), în asociere cu factorii tehnologici, de cultură (ex. suplimentul de apă furnizat prin sistemele de irigare) poate fi asistat de calculator folosind algoritmi genetici, care, după cum s-a demonstrat în teză, sunt capabili să ofere soluții optimale la problemele complexe de evoluție în condiții specifice de mediu, în timp scurt.

(Concluzie)

Definirea unui design de experiment corect este esențială în obținerea, prelucrarea și interpretarea datelor experimentale; designul de experiment realizat în prezenta lucrare, a permis studiul evoluției mai multor parametri ce definesc materialul genetic/cultivarul, cu ajutorul algoritmului genetic; analiza evoluției a cuprins atât evoluția în ansamblu, cât și evoluția bazată pe diferite strategii de selecție și supraviețuire. Astfel, s-a constatat că obiectivul fixat al evoluției (coeficientul de determinare) se distribuie după o lege de distribuție generalizată: distribuția Fisher-Tippett. Numărul de evoluții către optim se distribuie după aceeași lege de distribuție. Momentele relative ale evoluției se distribuie după o lege degenerată (punându-se în evidență două dependențe între cei trei parametri) din familia de curbe log-Pearson de tipul III.

(Recomandare)

Procesele de ameliorare care, prin selecții repetate, au ca scop amplificarea unui caracter, cum este de exemplu mărimea fructului la o specie horticolă oarecare, trebuie realizate într-un design experimental urmărit cu atenție pe întreg parcursul proceselor de încrucișare repetată. Conform datelor obținute, asemenea caractere nu urmează pe parcursul

selecției o lege de distribuție normală; astfel, media care se obține din observațiile pe întreg eșantionul unei generații supuse înmulțirii repetate reprezintă o statistică suficientă pentru ilustrarea evoluției. În explicarea unui asemenea fenomen, selecția naturală sau cea artificială, ultima cu scop de ameliorare, este necesară o analiză mai detaliată, complexă, în care variabilitatea și heritabilitatea caracteristicilor specifice trebuie abordată diferențiat.

(Concluzie)

Analiza compoziției eșantionului de material genetic pe parcursul evoluției a dus la următoarele concluzii:

- ÷ Selecția deterministă, indiferent de strategia de supraviețuire a indivizilor în populație, are drept consecință o scădere semnificativă din punct de vedere statistic a numărului total de genotipuri distincte reprezentate în cultivar;
- ÷ Supraviețuirea deterministă, asociată cu selecția turnir sau proporțională, favorizează o creștere semnificativă din punct de vedere statistic (fapt evidențiat de către toate cele trei statistici folosite: Pearson-Fisher, Anderson-Darling și Kolmogorov-Smirnov) a numărului total de genotipuri reprezentate în cultivar;
- ÷ Supraviețuirea deterministă favorizează o creștere semnificativă din punct de vedere statistic a numărului de genotipuri în grupul celor mai frecvente 23 din 46 de execuții independente în generațiile ce produc evoluție, în timp ce selecția deterministă determină scăderea semnificativ statistică a acestora.

(Recomandare)

În procesele de ameliorare a plantelor horticole, care au ca scop general îmbunătățirea productivității, a calității producției precum și creșterea rezistenței plantelor la atacul diferitelor boli și dăunători, utilizarea algoritmilor genetici poate reduce considerabil durata procesului de obținere a unor genotipuri superioare.

Algoritmii genetici, considerați de specialiști o aplicație a inteligenței artificiale, au conform studiului realizat perspective de a deveni și apanajul geneticienilor și amelioratorilor, respectându-se condiția ca, în procesul de selecție al celor mai bune genotipuri destinate înmulțirii, folosirea unei anumite strategii de selecție (proporțională, în turnir și deterministă fiind cele analizate în lucrare) trebuie asociată cu obiectivele generale de ameliorare urmărite.

Un asemenea deziderat se suprapune aplicațiilor practice prin care algoritmii genetici contribuie la rezolvarea problemelor de optimizare, planificare ori căutare în chimie, informatică, matematică, modelare moleculară etc., dar și în genetica și ameliorarea plantelor. În acest sens, prezenta lucrare deschide noi perspective aplicative în domeniu, și se constituie totodată într-o cercetare fundamentală originală, inedită în horticultură, deschizătoare de noi direcții de cunoaștere a fenomenelor și proceselor biologice, capabilă să permită formularea și verificarea de noi ipoteze, modele conceptuale și teorii.

Lucrări reprezentative publicate

- ÷ On about what Can Be Done and what Cannot Be Done with Genetic Algorithms in Phylogenetic Tree and Gene Sequence Analyses. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Radu E. SESTRĂȘ. *Bulletin UASVM, Horticulture* 65(1):63-70, 2008. ([Jäntschi & alții, 2008](#))
- ÷ Hard Problems in Gene Sequence Analysis: Classical Approaches and Suitability of Genetic Algorithms. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Radu E. SESTRĂȘ. *Biotechnology & Biotechnological Equipment* 23(2):1275-1280, 2009. ([Jäntschi & alții, 2009](#))
- ÷ Classical Approaches of Genetic Algorithms and their Suitability. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Radu E. SESTRĂȘ. *Asian Journal of Chemistry* 22(3):2275-2284, 2010. ([Jäntschi & alții, 2010](#))
- ÷ Distribution Fitting 1. Parameters Estimation under Assumption of Agreement between Observation and Model. Lorentz JÄNTSCHI, *Bulletin UASVM, Horticulture* 66(2):684-690, 2009. ArXiv electronic library permanent link (July 16, 2009): <http://arxiv.org/abs/0907.2829> (Subject: Statistics - Methodology).
- ÷ Distribution Fitting 2. Pearson-Fisher, Kolmogorov-Smirnov, Anderson-Darling, Wilks-Shapiro, Kramer-von-Misses and Jarque-Bera statistics. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ. *Bulletin UASVM, Horticulture* 66(2):691-697, 2009. ArXiv electronic library permanent link (July 16, 2009): <http://arxiv.org/abs/0907.2832> (Subject: Statistics - Methodology).
- ÷ Distribution Fitting 3. Analysis under Normality Assumption. Sorana D. BOLBOACĂ, Lorentz JÄNTSCHI. *Bulletin UASVM, Horticulture* 66(2):698-705, 2009. ([Bolboacă & alții, 2009](#))
- ÷ Distribution Fitting 4. Benford test on a sample of observed genotypes number from running of a genetic algorithm. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Carmen E. STOENOIU, Mihaela IANCU, Monica M. MARTA, Elena M. PICĂ, Monica ȘTEFU, Adriana F. SESTRĂȘ, Marcel M. DUDA, Radu E. SESTRĂȘ, Ștefan ȚIGAN, Ioan ABRUDAN, Mugur C. BĂLAN. *Bulletin UASVM, Agriculture* 66(1):82-88, 2009. ([Jäntschi & alții, 2009](#))
- ÷ Recording Evolution Supervised by a Genetic Algorithm for Quantitative Structure-Activity Relationship Optimization. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Radu E. SESTRĂȘ. *Applied Medical Informatics* 26(2):89-100, 2010. ([Jäntschi & alții, 2010](#))
- ÷ Meta-heuristics on quantitative structure-activity relationships: study on polychlorinated biphenyls. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Radu E. SESTRĂȘ. *Journal of Molecular Modeling* 16(2):377-386, 2010, DOI: [10.1007/s00894-009-0540-z](https://doi.org/10.1007/s00894-009-0540-z).
- ÷ A Study of Genetic Algorithm Evolution on the Lipophilicity of Polychlorinated Biphenyls. Lorentz JÄNTSCHI, Sorana D. BOLBOACĂ, Radu E. SESTRĂȘ. *Chemistry and Biodiversity*, 2010, DOI: [10.1002/cbdv.200900356](https://doi.org/10.1002/cbdv.200900356).
- ÷ A genetic algorithm for structure-activity relationships: software implementation. Lorentz JÄNTSCHI. ArXiv electronic library permanent link (June 26, 2009): <http://arxiv.org/abs/0906.4846> (Subject: Neural and Evolutionary Computing).