

Modelarea compactării cu role în industria farmaceutică folosind tehnici de inteligență computațională



Hossam Mohammed Zawbaa Ismail ELSAYED

Facultatea de Matematică și Informatică

Universitatea Babeș-Bolyai

Rezumatul tezei de doctorat

Doctor în Informatică

Mai 2016

Contents

1	Introducere	8
2	Starea actuală a cercetărilor în domeniu	11
2.1	Învățarea automată	11
2.2	Optimizarea inspirată din biologie	12
3	Preliminarii și fundamente	13
3.1	Învățarea automată	13
3.1.1	Tehnici de clasificare	13
3.1.1.1	Vecinul cel mai K-apropiat	13
3.1.1.2	Pădurea aleatoare	14
3.1.2	Tehnici de regresie	14
3.1.2.1	Rețelele neurale artificiale	14
3.1.2.2	Învățarea automată extremă	14
3.1.3	Modele predictive	15
3.2	Algoritmi de optimizare inspirați din biologie	15
3.2.1	Algoritmi genetici	15
3.2.2	Optimizarea cu roi de particule	15
3.2.3	Colonia artificială de albine	16
3.2.4	Algoritmul licuricilor	16
3.2.5	Căutarea cucului	16
3.2.6	Algoritmul liliacului	16
3.2.7	Algoritmul de polenizare a florilor	16
3.2.8	Optimizarea cu păianjen social	17
3.2.9	Optimizarea cu lupul cenușiu	17

CONTENTS

3.2.10	Algoritmul hibrid maimuță - banc de krill	17
3.2.11	Algoritmul libelulă	17
3.2.12	Optimizarea cu fluture de lampă	18
3.2.13	Optimizarea cu leul-furnicilor	18
4	Sistemul propus și metodologia de cercetare	19
4.1	Distribuții aleatoare	21
4.1.1	Distribuția Gauss	21
4.1.2	Distribuția Lévy flight	21
4.1.3	Distribuția haotică	21
4.2	BIOA cu diverse distribuții aleatoare	22
4.2.1	Versiunea haotică a BIOA	22
4.2.1.1	Optimizarea haotică cu leul furnicilor	22
4.2.1.2	Optimizarea haotică cu lupul cenușiu	23
4.2.2	Versiuni Lévy ale BIOA	23
4.2.2.1	Optimizarea Lévy cu leul furnicilor	24
4.2.2.2	Optimizarea Lévy cu păianjenul social	24
4.3	Versiuni binare ale BIOA	25
4.3.1	Optimizarea binară cu lupul cenușiu	25
4.3.2	Optimizarea binară cu leul furnicilor	25
4.3.3	Optimizarea binară cu fluturele de lampă	25
4.4	Algoritmi folosiți pentru comparare	26
4.5	Metode de inițializare	26
4.6	Metrici de performanță	26
5	Aplicarea optimizării inspirate din biologie în procesele farmaceutice	28
5.1	Analiza farmaceutică și discuții	28
5.1.1	Rezultate de la compactarea cu role	29
5.1.2	Analiza compactării cu ștanțe	29
5.1.3	Acidul poli-lactic-co-glicolic PLGA	30
5.1.3.1	Rezultate și discuții pentru PLGA	31

6	Concluzii și cercetări viitoare	33
6.1	Concluzii	33
6.2	Cercetări viitoare	34
	References	35

Abstract

În această teză am dezvoltat variante noi ale unor algoritmi de optimizare inspirați din natură, precum leul furnicilor binar (binary antlion optimization BALO), lupul cenușiu haotic (chaotic grey wolf optimization CGWO), lupul cenușiu binar (binary grey wolf optimization BGWO), și mulți alții. Toți algoritmi propuși de noi au fost comparați cu tehnici consacrate folosite în selectarea caracteristicilor, precum optimizarea cu roi de particule (particle swarm optimization PSO) și algoritmi genetici (genetic algorithms GA). În industria farmaceutică, pentru dezvoltarea de noi medicamente sau pentru optimizarea proceselor de fabricație se folosesc frecvent cunoștințele referitoare la relațiile de cauzalitate dintre calitatea produsului și atributele (caracteristicile) formulelor chimice. Folosind cantitatea imensă de date colectate de-a lungul practicii de dezvoltare de produse farmaceutice, modelele de inteligență computațională (computational intelligence CI), bazate pe învățarea automată și pe algoritmi de optimizare inspirați din natură pot fi folosite eficient pentru a identifica atributele esențiale de calitate (critical quality attributes CQA) și parametrii esențiali de proces (critical process parameters CPP), asociate formulelor chimice, respectiv proceselor de fabricație. Obiectivul principal al cercetării a fost evaluarea robusteții tehnicilor de învățare automată combinate cu algoritmi de optimizare inspirați din natură și folosite la modelarea proceselor de fabricație a tabletelor. Mai precis, efortul a fost concentrat pe predicția unor proprietăți ale tabletelor precum porozitatea și rezistența la rupere pe baza caracteristicilor pulberilor și panglicilor de cauciuc. În acest scop, au fost efectuate experimente de compactare cu role pentru diverși compuși farmaceutici, din care au rezultat seturi de date ce conțin o gamă largă de caracteristici. Eficiența modelării a fost evaluată pe baza caracteristicilor selectate (reducerea) și folosind rădăcina

pătrată din eroarea medie (root mean square error RMSE). S-a observat că rezultatele prezise sunt într-o bună concordanță cu datele experimentale.

Cuvinte-cheie:

Algoritmi de optimizare inspirați din biologie, Selecția caracteristică, Compactarea cu role în industria farmaceutică, Optimizarea cu leul-furnicilor, Optimizarea cu fluture de lampă, Optimizarea cu lupul cenușiu, Optimizarea cu păianjen social, Algoritmul de polenizare a florilor, Algoritmi genetici, Optimizarea cu roi de particule.

Chapter 1

Introducere

O *caracteristică* de intrare este o proprietate măsurabilă a problemei studiate. În ultimii ani, în aplicațiile de învățare automată și recunoaștere de pattern-uri, domeniul caracteristicilor s-a expandat de la sute la mii de variabile. Volumele mari de date generate astăzi în biologie oferă informație mai detaliată și mai utilă, dar complică analiza acestora, deoarece nu toate informațiile sunt relevante. Selectarea caracteristicilor (atributelor) (feature selection, FS) relevante pentru un set de date este o problemă complexă: *selectarea caracteristicilor* este o tehnică de rezolvare a problemelor de clasificare și regresie, identificând o submulțime a caracteristicilor și eliminând pe cele redundante. Acest mecanism este cu atât mai util cu cât numărul caracteristicilor este mai mare, și nu toate sunt necesare pentru a descrie datele experimentale (1).

Multe lucrări formulează problema selectării caracteristicilor (feature selection problem, FSP) ca fiind o problemă de *optimizare combinatorială*, în care submulțimea rezultată de caracteristici produce cea mai bună potrivire a datelor (2). În aplicațiile din lumea reală, selectarea caracteristicilor este obligatorie, deoarece avem de-a face cu zgomote, cu caracteristici nerelevante sau greșite (3). Acești factori pot avea un impact negativ asupra performanței clasificării. Metodele de selectare a caracteristicilor se clasifică în raport cu două criterii:

1. *strategia de căutare*: metoda folosită pentru a genera submulțimi sau combinații de caracteristici.
2. *calitatea selecției (fitness)*: criteriile folosite pentru a caracteriza calitatea unei submulțimi de caracteristici.

În general, FSP se formulează ca o problemă de optimizare *multi-criterială*: *minimizarea* dimensiunii submulțimii selectate de caracteristici și *maximizarea* performanței predicției. De regulă, aceste obiective sunt contradictorii, iar soluția optimală este un compromis (4).

Mulți dintre noii algoritmi de optimizare sunt *inspirați din natură* (5). Sunt trei mari surse de inspirație: *biologia*, *fizica*, și *chimia*. Toți noii algoritmi de optimizare bazați pe *biologie* se vor numi *bio-inspirați* (bio-inspired optimization algorithms, BIOA) (5).

În general, dimensiunea spațiului de căutare crește exponențial cu numărul de caracteristici ale setului de date studiat (6). Prin urmare, o căutare exhaustivă pentru soluția optimală sau aproape optimală poate fi *impracticabilă*; de asemenea, tehnicile de căutare exhaustivă suferă de *stagnare* în optime locale (2), (7). În toți algoritmi de optimizare bio-inspirați este nevoie de un echilibru convenabil între *explorare* (diversificare, căutare globală) și *exploatare* (intensificare, căutare locală) (8).

Obiectivele cercetării noastre sunt:

- O_1 : utilizarea BIOA la FSP pentru a rezolva problema de clasificare care minimizează numărul de caracteristici selectate și maximizează precizia clasificării.
- O_2 : în problema de regresie, folosirea BIOA pentru a reduce numărul de caracteristici selectate și a minimiza eroarea predicției.
- O_3 : aplicarea BIOA în domeniul farmaceutic, pentru minimizarea numărului de caracteristici selectate și minimizarea erorii predicției. De asemenea, studierea importanței fiecărei variabile de intrare pentru un set de date dat.

În industria farmaceutică, pentru dezvoltarea de produse noi și optimizarea proceselor de fabricație, este necesară o bună înțelegere a relațiilor de cauzalitate dintre calitatea produsului și atributele formulelor chimice. Producția industrială de medicamente are patru procese: amestecare, compactare cu role, măcinare și ștanțare. Compactarea cu role este o metodă de preparare a granulelor de medicament pentru capsule sau tablete cu o densitate adecvată. Cei mai folosiți adjuvanți (excipienți) sunt celuloza microcristalină (MCC), fosfatul de calciu dibasic (DCP) și lactoza. Compactarea cu role este și o tehnică de mărire a dimensiunii particulelor, care granulează materialele sub formă de pulbere pentru a obține subproduse de dimensiuni intermediare, rezultând o granule de calitate înaltă. Este esențială alegerea optimă a parametrilor critici pentru

1. INTRODUCERE

compactarea cu role, precum presiunea constantă de compactare, distanța dintre role, etc.

Teza are 9 capitole, inclusiv cel de față. Introducerea rezumă caracteristicile optimizării inspirate din biologie, importanța selectării caracteristicilor în problemele de clasificare și regresie, precum și impactul selectării caracteristicilor asupra domeniului industriei farmaceutice. Capitolul (2) face o trecere în revistă a cercetărilor existente în învățarea automată, optimizarea inspirată din biologie și aplicațiile acestora. Capitolul (3) prezintă informație de bază referitoare la algoritmi de învățare automată și optimizare inspirată din biologie. Capitolul (4) descrie câteva dimensiuni ale spațiului experimentelor: diferiți generatori de numere aleatoare, diferiți algoritmi propuși de optimizare inspirați din biologie și metrice de evaluare a performanțelor acestora. Următoarele patru capitole descriu contribuțiile originale: cercetări experimentale de folosire a BIOA în clasificare (Capitolul (5)); rezultatele experimentale ale folosirii BIOA în regresie (Capitolul (6)), rezultatele experimentale ale folosirii BIOA în domeniul farmaceutic (Capitolul (7)). Capitolul (8) descrie alte trei studii de caz folosind BIOA. Studiile descrise conțin caracteristicile seturilor de date folosite și analiza rezultatelor obținute, prin compararea cu metode consacrate. În sfârșit, capitolul (9) rezumă concluziile acestui studiu și direcțiile de cercetare viitoare.

Chapter 2

Starea actuală a cercetărilor în domeniu

2.1 Învățarea automată

Tehnicile de învățare automată (machine learning, ML) joacă un rol important în rezolvarea multor probleme complicate de clasificare și regresie. Tehnicile ML se folosesc pentru construirea de modele de predicție din date observaționale. În problemele de clasificare, metodele ML sunt adecvate pentru variabile dependente discrete, având un număr finit de valori neordonate. Pe de altă parte, în cazul problemelor de regresie, variabilele dependente pentru care se folosesc metodele ML au valori continue și ordonate. Tehnicile bazate pe nucleu (kernel-based techniques), precum mașinile cu suport vectorial (support vector machines, SVM), analiza componentei principale, procesele Gaussiene ș.a. reprezintă un important pas înainte în dezvoltarea ML. SVM au fost propuse de Vapnik în anii 1960 pentru clasificare și au devenit recent subiectul unor eforturi intense de cercetare datorită dezvoltării tehnicilor și teoriilor legate de aplicații la probleme de regresie și estimare a densității (9). Rețelele neurale, în special cele cu un singur strat ascuns și feed-forward (single hidden layer feed-forward neural networks, SLFN) sunt considerate unul dintre cele mai utilizate modele de ML în problemele de regresie și clasificare (10). Modelul învățării extreme (extreme learning machine, ELM) a fost propus pentru SLFN, unde conexiunile dintre stratul de intrare și neuronii din stratul ascuns se selectează aleator și rămân nemodificate pe tot parcursul procesului de învățare (11).

2.2 Optimizarea inspirată din biologie

Numeroase tehnici euristice reproduc comportamentul sistemelor din natură, biologice sau fizice, fiind considerate metode robuste de optimizare. Algoritmii genetici (genetic algorithms, GA) au fost primii algoritmi evolutivi introduși în literatură, bazați pe procesul natural de evoluție prin reproducere (12). Un algoritm de FS bazat pe GA și folosind o mulțime fuzzy ca funcție de fitness a fost propus în (13). În PSO, fiecare soluție este o particulă definită de poziție, fitness și un vector de viteză ce reprezintă direcția de mișcare a acesteia (14). Algoritmul de FS folosind optimizarea cu colonia de furnici (ant colony optimization, ACO) a fost aplicat la detectarea intruziunii în rețele (15). ACO folosește rata de discriminare a lui Fisher pentru adoptarea informației euristice și teoria mulțimilor pentru FS (16). Colonia artificială de albine (artificial bee colony, ABC) este un algoritm de optimizare numerică bazat pe comportamentul albinelor la căutarea hranei. În ABC, albinele lucrătoare caută sursa de hrană și le înștiințează pe celelalte albine (17). Algoritmul albinei virtuale (virtual bee algorithm, VBA) se aplică la optimizarea funcțiilor numerice în 2-D folosind un roi de albine virtuale care se mișcă aleator în spațiul caracteristicilor și interacționează pentru detectarea surselor de hrană, conducând astfel la soluțiile posibile ale problemei de optimizare (18). ALO este un algoritm evolutiv recent care reproduce mecanismele de vânatoare ale leului furnicilor (19).

Chapter 3

Preliminarii și fundamente

3.1 Învățarea automată

Clasificarea și regresia sunt tehnici de ML pentru construirea de modele de predicție din date observaționale. Tehnicile de clasificare se folosesc pentru variabile dependente cu valori discrete și neordonate, eroarea de predicție fiind măsurată în funcție de costul clasificărilor greșite. Tehnicile de regresie se aplică la variabile dependente cu valori continue și ordonate, eroarea de predicție (eroarea medie pătratică) calculându-se prin diferența pătratelor dintre valorile observate și cele prezise. Această secțiune prezintă pe scurt tehnicile de clasificare și regresie folosite în teză.

3.1.1 Tehnici de clasificare

Metodele de ML joacă un rol important la rezolvarea problemelor complicate de clasificare. În cele ce urmează facem o prezentare succintă tehnicilor de clasificare folosite.

3.1.1.1 Vecinul cel mai K-apropiat

Vecinul cel mai K-apropiat (K-nearest neighbor, KNN) este un clasificator foarte simplu. În faza de clasificare, un eșantion nou este clasificat pe baza majorității categoriei KNN (K este un întreg predefinit): fiind dat un punct de interogare, algoritmul determină primele K obiecte sau puncte de antrenament cele mai apropiate de el. Se folosește distanța dintre puncte ca măsură de apropiere. După ce sunt găsite clasele KNN, noul eșantion urmează (prezice) clasa majoră a KNN (20).

3. PRELIMINARII ȘI FUNDAMENTE

3.1.1.2 Pădurea aleatoare

Pădurea aleatoare (random forest RF) este considerată una dintre cele mai bune tehnici ML de clasificare și regresie. Ea poate clasifica cu precizie seturi mari de date (21) și constă dintr-o colecție de clasificatori, cu structură arborescentă. Fiecare arbore depinde de valori vectoriale aleatoare eșantionate independent (22). Datele de intrare trec prin rădăcina arborelui și apoi îl traversează. Datele de intrare sunt eșantionate aleator, fiind înlocuite de mulțimi din ce în ce mai mici. Clasa eșantionului se determină folosind arborii RF bazați pe un generator de numere aleatoare (21). Variabila aleatoare specifică cum se efectuează tăierile succesive în timpul construirii arborelui, prin selectarea nodului și a coordonatei care este punctul de diviziune, ca și a poziției acesteia (23).

3.1.2 Tehnici de regresie

3.1.2.1 Rețelele neurale artificiale

Rețelele neurale artificiale (artificial neural networks ANN) au fost create ca generalizări ale modelelor matematice reprezentând sistemele nervoase biologice. În modelul matematic simplificat al neuronului, efectele sinapselor se reprezintă prin ponderi ale conexiunilor care modulează efectul asupra semnalelor de intrare asociate, iar caracteristicile nonlineare ale neuronilor se reprezintă printr-o funcție de transfer. Există multe funcții de transfer pentru prelucrarea intrărilor ponderate și perturbate, dintre care remarcăm patru dintre ele, folosite în mod frecvent la prelucrarea multimedia. Comportamentul rețelei neurale depinde în mare măsură de interacțiunile dintre neuroni (24).

3.1.2.2 Învățarea automată extremă

Modelul învățării automate extreme (extreme learning machine ELM) a fost propus pentru SLFN. În acest model, conexiunile dintre stratul de intrare și neuronii ascunși se selectează aleator și rămân nemodificate pe tot parcursul procesului de învățare. Conexiunile de ieșire se ajustează prin minimizarea funcției de cost, folosind un sistem liniar (11). Caracteristicile proprii ELM sunt alegerea aleatorie a ponderilor neuronilor din stratul ascuns și funcția de bias. Ambele operații ale ELM (antrenarea și predicția) sunt mult mai rapide decât cele ale altor tehnici non-liniare. Prin urmare, ELM tinde să ofere o performanță de generalizare bună cu o viteză de învățare ridicată (25).

3.1.3 Modele predictive

1. *Cubist* este un pachet ce implementează arborii de decizie predictivi bazați pe reguli propuși de Quinlan (26). Modelele Cubist introduc ecuații liniare în nodurile lor terminale, fiind capabili să prezică valori numerice.
2. *Pădurea aleatoare Random Forest (RF)* crează o mulțime de arbori de decizie folosind date de intrare aleatoare. Pachetul `randomForest` din mediul `R` a fost folosit în (27).
3. *Perceptronul multistrat monoton Monmlp (monotonic multilayer perceptron)* (28, 29) a fost folosit pentru a exploata avantajele învățării fără propagare înapoi.
4. *Rețelele neurale cu învățare profundă Deep learning neural networks* sunt folosite pentru rezolvarea unor probleme complexe, pe baza combinării unor soluții mai simple. Aceste sisteme pot fi exploatare în medii reale (30).
5. *FugeR* este destinat antrenării sistemelor fuzzy bazate pe algoritmi evolutivi.

3.2 Algoritmi de optimizare inspirați din biologie

3.2.1 Algoritmi genetici

Algoritmii genetici (genetic algorithms, GA) au fost primii algoritmi evolutivi cunoscuți în literatură, fiind dezvoltați de Holland în deceniile 7 și 8 ale secolului trecut. GA pot rezolva probleme de optimizare complexe și neliniare.

3.2.2 Optimizarea cu roi de particule

Optimizarea cu roi de particule (particle swarm optimization, PSO) este o metodă euristică de optimizare globală dezvoltată de Kennedy and Eberhart în 1995 (14). Algoritmul PSO face parte din categoria algoritmilor de inteligență a roiurilor și se bazează pe comportamentul de mișcare al păsărilor (31). El este utilizat pe scară largă pentru rezolvarea problemelor de optimizare și de FS (32).

3. PRELIMINARII ȘI FUNDAMENTE

3.2.3 Colonia artificială de albine

Algoritmul de optimizare bazat pe comportamentul albinelor pentru căutarea hranei se numește colonia artificială de albine (artificial bee colony, ABC) și a fost propus de Karaboga în 2007 (17). În ABC, albinele lucrătoare caută o sursă de hrană și apoi le înștiințează pe celelalte. Albinele culegătoare le urmează pe cele lucrătoare iar albinele cercetaș zboară spontan pentru a descoperi surse de hrană mai bune (8).

3.2.4 Algoritmul licuricilor

Algoritmul licuricilor (firefly algorithm, FFA) este o metodă stohastică de optimizare globală dezvoltată de Yang în 2008 (33). FFA imită mecanismele de împerechere și de schimb de informație ale licuricilor care folosesc scipiri luminoase. În FFA, mișcarea unui licurici este determinată în principal de atractivitatea altor licurici.

3.2.5 Căutarea cucului

Algoritmul căutarea cucului (cuckoo search, CS) este un algoritm euristic de căutare propus de Yang în 2009 (34) și folosit la rezolvarea problemelor de optimizare continuă. Cucii au o strategie agresivă de reproducere și-și depun ouăle în cuiburile altor păsări gazdă, care pot fi de altă specie. Pasărea gazdă poate descoperi că oul nu este al ei, caz în care fie că aruncă oul din cuib, fie abandonează cuibul și-și construiește altul în alt loc.

3.2.6 Algoritmul liliacului

Yang a dezvoltat algoritmul liliacului (bat algorithm, BA) în 2010 (35). BA este o tehnică a meta-heuristică care folosește comportamentul de eclocație pentru căutarea prăzii și detectarea sau evitarea obstacolelor. Pentru navigare, liliecii emit pulsuri sonore puternice și ascultă ecoul reflectat de obiectele înconjurătoare.

3.2.7 Algoritmul de polenizare a florilor

Algoritmul de polenizare a florilor (flower pollination algorithm, FPA) este o metodă meta-heuristică de optimizare bazat pe operația de polenizare a florilor plantelor și a fost dezvoltat de Yang în 2012 (36). Obiectivul principal al florii este reproducerea prin polenizare. Polenizarea florilor este asociată de obicei cu transferul de polen, realizat de

3.2 Algoritmi de optimizare inspirați din biologie

polenizatori precum insectele, păsările, lilieci și a. (37). Polenizarea este de două tipuri: *auto-polenizare (abiotică)* și *polenizare încrucișată (biotică)*, care corespund căutării locale, respectiv căutării globale în algoritm.

3.2.8 Optimizarea cu păianjen social

Optimizarea cu păianjen social (social spider optimization, SSO) este unul dintre algoritmii recenți din categoria inteligenței roiuilor și a fost propus de Cuevas în 2013 (38). Algoritmul SSO mimează comportamentul coloniei de păianjeni din natură și are două componente: membrii sociali și pânza comună. Păianjenul cooperează interactiv cu ceilalți membri ai coloniei (39).

3.2.9 Optimizarea cu lupul cenușiu

Algoritmul de optimizare cu lupul cenușiu (grey wolf optimization, GWO) a fost dezvoltat recent de Mirjalili în 2014. GWO imită modul în care haita de lupi caută hrana și supraviețuiește prin evitarea inamicilor (40). Fiecare lup cenușiu din haită își alege poziția, mișcându-se continuu spre un loc mai bun și fiind atent la amenințările potențiale. GWO are un parametru ce semnifică probabilitatea de amenințare, imitând incidentele când lupii se ciocnesc de inamicii lor.

3.2.10 Algoritmul hibrid maimuță - banc de krill

Algoritmul hibrid maimuță - banc de krill (monkey and krill herd algorithm, MAKHA) a fost dezvoltat de Khalil în 2015 (41). În general, algoritmii hibridi folosesc operatori dintr-un algoritm combinați cu operatori din alt algoritm în ideea de a folosi ce e mai bun din fiecare algoritm pentru a îmbunătăți performanța. Algoritmul hibrid MAKHA folosește cei mai performanți operatori din algoritmul maimuță (monkey algorithm, MA) și algoritmul bancului de krill (krill herd algorithm, KHA), evitând folosirea operatorilor mai puțin performanți sau care necesită putere mare de calcul din respectivii algoritmi (41).

3.2.11 Algoritmul libelulă

Mirjalili a propus anul trecut algoritmul libelulă (dragonfly algorithm, DA) (42). Libelulele sunt insecte uimitoare, existând peste 3000 de specii diferite (43). Ciclul de

3. PRELIMINARII ȘI FUNDAMENTE

viață al libelulei are două faze principale: nimfă și adult. Cea mai mare parte a vieții lor, libelulele sunt în faza nimfă. La terminarea acesteia, ele devin adulte și se comportă ca niște prădători mici, vânând aproape toate celelalte insecte mici din mediul lor. DA imită comportamentele de roi static (vânătoare) și dinamic (migrare) ale libelulelor (44).

3.2.12 Optimizarea cu fluture de lampă

Tot anul trecut, Mirjalili a dezvoltat și algorithmul de optimizare cu fluture de lampă (moth-flame optimization, MFO). Prin evoluție, fluturii au ajuns să zboare noaptea orientându-se după lună și folosind pentru navigare o metodă numită orientare transversală: direcția de zbor a fluturelui menține un unghi constant în raport cu luna (adică sursa de lumină) (45). Această metodă este considerată o tehnică foarte eficientă pentru zborul în linie dreaptă pe distanțe lungi (46).

3.2.13 Optimizarea cu leul-furnicilor

Optimizarea cu leul-furnicilor (antlion optimization, ALO) este un algoritm propus tot anul trecut și tot de Mirjalili (19). ALO imită mecanismele de vânătoare folosite de leul furnicilor în natură. Leii-furnicilor (antlions, doodlebugs) aparțin familiei Myrmeleonidae din ordinul Neuroptera. Ei vânează în stadiul de larvă, iar în stadiul de adult se reproduc (19).

Chapter 4

Sistemul propus și metodologia de cercetare

O caracteristică comună a algoritmilor de inteligență roiurilor este că informația este partajată între mai mulți agenți. La fiecare iterație a procesului de optimizare, toți agenții/o parte dintre agenți își actualizează/modifică poziția pe baza informație de poziție a celorlalți agenți sau a lor.

Explorarea sau *căutarea globală* se poate defini ca achiziție de informație nouă prin căutare (33). Explorarea este preocuparea principală a tuturor algoritmilor de optimizare deoarece ea permite determinarea de noi regiuni de căutare care ar putea conține soluții mai bune. *Exploatarea* sau *căutarea locală* se poate defini ca folosire a informației cunoscute. Regiunile bune se exploatează prin căutare locală. Procesul de selecție trebuie să fie echilibrat între selectarea aleatoare și selectarea greedy, pentru a dirija căutarea către soluții candidat mai potrivite (exploatare) și în același timp prin promovarea diversității utile în populație (explorare) (33).

Algoritmii BIOA propuși sunt folosiți pentru a determina un număr *minim* de caracteristici care *maximizează* performanța de predicție. Spațiul de căutare reprezintă fiecare caracteristică ca dimensiune individuală, amplitudinea fiecărei dimensiuni fiind între 0 și 1; este nevoie de o metodă inteligentă de căutare pentru determinarea mulțimii optime de caracteristici în spațiul de căutare imens care maximizează o funcție de potrivire dată. În cazul clasificării, funcția generală de potrivire pentru BIOA propuși este maximizarea performanței de clasificare (4.1) peste mulțimea de validare, fiind

4. SISTEMUL PROPUȘ ȘI METODOLOGIA DE CERCETARE

cunoscută mulțimea de antrenare, cu condiția menținerii unui număr minim de caracteristici selectate:

$$\downarrow Fitness = \alpha(1 - P) + \beta \frac{|R|}{|C|}, \quad (4.1)$$

unde R este lungimea submulțimii de caracteristici selectate, C este numărul total de caracteristici în setul de date, α și β sunt doi parametri caracterizând importanța performanței clasificatorului și lungimea submulțimii, $\alpha \in [0, 1]$ and $\beta = 1 - \alpha$, iar P este performanța clasificării măsurată ca în relația (4.2):

$$P = \frac{N_c}{N}, \quad (4.2)$$

unde N_c este numărul instanțelor de date clasificate corect, iar N este numărul total de date din setul de date.

În cazul regresiei, funcția generală de potrivire pentru BIOA propuși este minimizarea erorii de predicție (4.3) peste mulțimea de validare, fiind cunoscută mulțimea de antrenare, cu condiția menținerii unui număr minim de caracteristici selectate:

$$\downarrow Fitness = \alpha * E + \beta \frac{\sum_i \theta_i}{N}, \quad (4.3)$$

unde E este eroarea predicției, θ reprezintă un vector N -dimensional cu elemente 0 și 1 reprezentând caracteristici neselectate/selectate, iar N este numărul total de caracteristici din setul de date.

Toate valorile asociate caracteristicilor sunt în intervalul $[0, 1]$, acolo unde valoarea asociată unei caracteristici se apropie de 1; respectiva caracteristică este candidată a fi selectată în predicție. În calculele individuale de potrivire, există un prag în raport cu care se decide dacă o caracteristică se selectează în etapa de evaluare.

Termenul aleator de ponderare α este folosit cu o valoare mare astfel ca el să fie adecvat unui spațiu de caracteristici cu mai multe minime locale. El este folosit pentru a echilibra compromisul între explorare și exploatare și trebuie tratat cu atenție.

În timpul fazei de antrenare, poziția fiecărui agent (fluture, leul-furnicilor, lup cenușiu, furnică, albina, etc.) reprezintă o submulțime de caracteristici. Mulțimea de antrenare se folosește pentru a evalua modelele de clasificare (KNN) și de regresie (ELM sau ANN) pe mulțimea de validare în timpul optimizării pentru a ghida procesul

de selectare a caracteristicilor. Fiecare set de date se împarte în trei părți egale pentru *antrenare*, *validare* și *testare*. Mulțimea de *antrenare* se folosește pentru a antrena modelul de predicție prin optimizare și la evaluarea finală. Mulțimea de *validare* se folosește pentru evaluarea performanței modelului de predicție în timpul procesului de optimizare. Mulțimea de *testare* se folosește pentru evaluarea caracteristicilor selectate la evaluarea finală. Modelele de clasificare și de regresie (KNN sau ELM sau ANN) se folosesc pentru a garanta calitatea caracteristicilor selectate și sunt evaluate pe mulțimea de validare în interiorul funcției de validare, pe parcursul procesului de optimizare (31).

4.1 Distribuții aleatoare

În această secțiune prezentăm succint diverse modele de distribuții aleatoare (distribuția Gauss, *levy flight* și haotică) precum și modul în care se aplică acestea în diferiți BIOA. O *variabilă aleatoare* se poate considera ca o expresie a cărei valoare este realizarea sau rezultatul unor evenimente asociate unui proces aleator (33). Variabila aleatoare este o funcție reprezentată printr-o funcție care transformă evenimentele în numere reale. Domeniul acestei transformări se numește spațiu de eșantionare. Fiecare variabilă aleatoare se reprezintă printr-o funcție densitate de probabilitate care exprimă distribuția sa de probabilitate. În studiul nostru folosim trei distribuții diferite; detalii despre acestea se găsesc în subsecțiunile următoare și în (47).

4.1.1 Distribuția Gauss

Distribuția Gaussiană (normală) este cea mai populară și se aplică la multe variabile fizice, precum intensitatea luminii, eroarea/incertitudinea măsurărilor, ș.a.

4.1.2 Distribuția Lévy flight

Mulți cercetători au studiat comportamentul de zbor al păsărilor și insectelor care are caracteristicile tipice zborurilor *levy* (48). Distribuția *levy flight* a fost aplicată la problemele de optimizare și rezultatele preliminare sunt promițătoare (49).

4.1.3 Distribuția haotică

Haosul înseamnă o condiție sau un loc cu o mare dezordine sau confuzie (50). Sistemele haotice sunt sisteme deterministe care expun comportament neregulat (sau chiar

4. SISTEMUL PROPUȘ ȘI METODOLOGIA DE CERCETARE

aleator) și o dependență majoră de condițiile inițiale. Haosul este unul dintre cele mai populare fenomene care se petrec în sistemele neliniare, a căror acțiune este complexă și prezintă caracteristici de imprevizibilitate (51). Teoria haosului studiază comportamentul sistemelor care respectă legi deterministe dar care este aleator și imprevizibil, adică sistemele dinamice. Variabilele haotice pot trece prin toate stările în anumite intervale, în funcție de regularitatea lor, fără repetiție (51). O *transformare haotică* este o transformare care prezintă un anumit tip de comportament haotic (50). În lucrarea de față, am folosit trei transformări haotice: logistică, Tent, și Singer.

4.2 BIOA cu diverse distribuții aleatoare

4.2.1 Versiunea haotică a BIOA

Datorită unor proprietăți interesante precum mixingul topologic, orbitele periodice dense, ergodicitatea și stochasticitatea intrinsecă, sistemele haotice se pot folosi pentru a modela parametrul care caracterizează comutarea între explorare și exploatare. În FS, căutarea haotică este mai aptă să evadeze din optimele locale decât căutarea aleatoare.

4.2.1.1 Optimizarea haotică cu leul furnicilor

În optimizarea haotică cu leul furnicilor (chaotic antlion optimization, CALO) ALO selectează iterativ o furnică pentru vânare într-o manieră ruletă și efectuează doi pași aleatori: (a) pasul aleator al furnicilor în jurul celei selectate și (b) pasul aleator în jurul leului furnicilor elită/cel mai bun. Furnica selectată își adaptează poziția sa din cei doi pași aleatori anteriori. Parametrul I controlează raportul dintre explorare și exploatare în ALO, descrescând liniar pentru a permite mai multă exploatare la începutul procesului de optimizare, explorarea devenind mai importantă spre sfârșitul optimizării. Prin urmare, în explorare se folosesc jumătate din resursele de optimizare, iar restul de timp este dedicat exploatării. Sistemele haotice se pot folosi pentru a adapta acest parametru, ce permite comutarea între explorare și exploatare. Algoritmul CALO propus este prezentat schematic în figura (4.1). Strategia de căutare a abordării bazate pe wrapper explorează spațiul caracteristicilor pentru a găsi o submulțime de caracteristici ghidat de performanța individuală a submulțimilor de caracteristici.

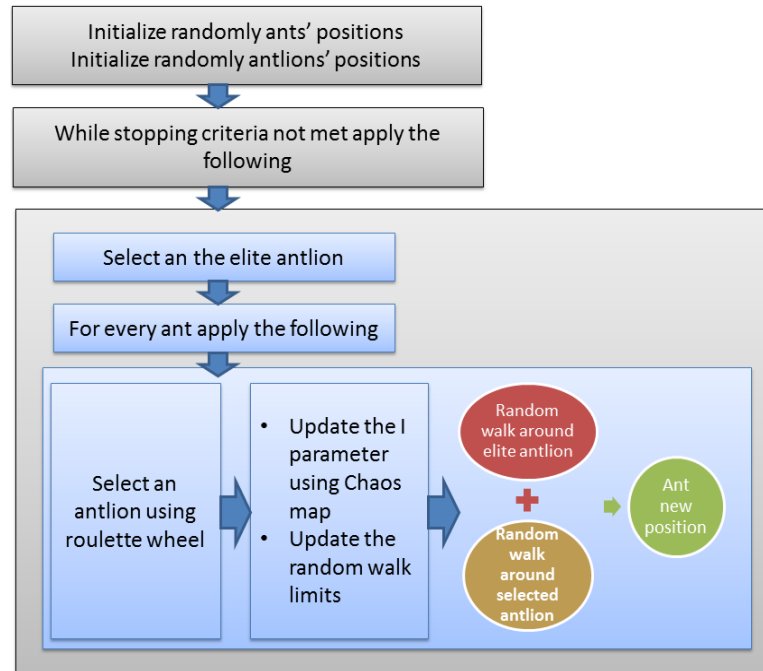


Figure 4.1: Optimizarea haotică cu leul furnicilor propusă (CALO)

4.2.1.2 Optimizarea haotică cu lupul cenușiu

În GWO (40), un singur parametru, \vec{a} , a fost propus pentru a controla comutarea dintre explorare și exploatare. Acest parametru a fost prevăzut să descrească liniar pentru a permite mai multă explorare la începutul optimizării, pe când exploatarea devine mai importantă la sfârșitul optimizării. Propunerea permite consumarea a jumătate din timpul de optimizare pentru explorare, în timp ce exploatarea ocupă cealaltă jumătate. În optimizarea haotică cu lupul cenușiu (Chaotic Grey Wolf Optimization CGWO), parametrul \vec{a} este astfel adaptat ca să permită perioade succesive de explorare să fie urmate de exploatare, prin distribuirea explorării pe toată perioada de optimizare, urmată, de fiecare dată, de exploatare.

4.2.2 Versiuni Lèvy ale BIOA

Lèvy flight este o modalitate de parcurgere eficientă în explorarea spațiului de căutare imens, fiind caracterizat de salturi mari și abrupte. O parcurgere bazată pe lèvy flight

4. SISTEMUL PROPUȘ ȘI METODOLOGIA DE CERCETARE

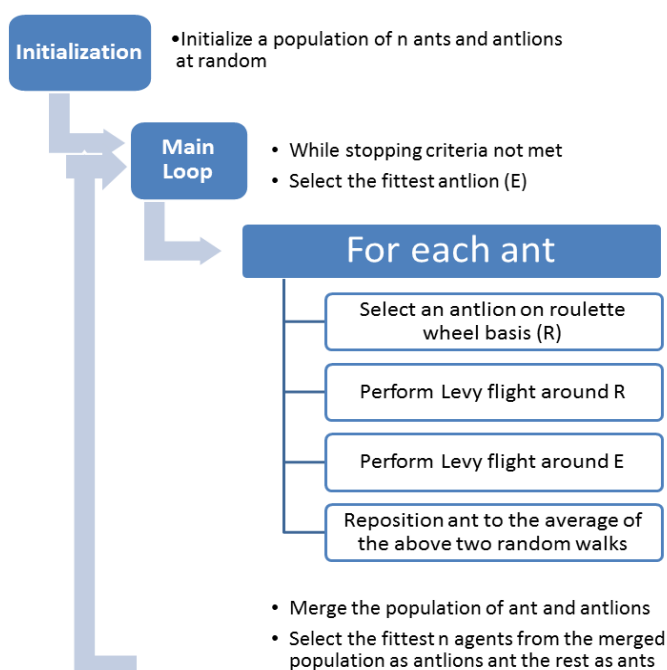


Figure 4.2: Algoritm de optimizare lèvy cu leul furnicilor (LALO)

este preferată distribuției uniforme atunci când se dorește creșterea vitezei de convergență și convergența la optime globale.

4.2.2.1 Optimizarea Lèvy cu leul furnicilor

Optimizarea Lèvy cu leul furnicilor (lèvy antlion optimization LALO) este descrisă în figura (4.2). Randomizarea joacă un rol important atât în *explorare*, cât și în *exploatare*, esența acesteia fiind parcurgerea aleatoare (33) - proces aleator constând din efectuarea unei secvențe de pași aleatori consecutivi.

4.2.2.2 Optimizarea Lèvy cu păianjenul social

Optimizarea Lèvy cu păianjenul social (Lèvy social spider optimization LSSO) este o variantă a algoritmului SSO în care lungimea pasului urmează distribuția lèvy; o astfel de parcurgere aleatoare se numește zbor lèvy (lèvy flight) sau parcurgere lèvy (lèvy walk). Din punct de vedere matematic, versiunile simple de zboruri lèvy sunt mai

eficiente decât parcurgerile aleatoare Browniene la explorarea necunoscutului, căutarea în spațiile mari fiind caracterizată de salturi abrupte și mari .

4.3 Versiuni binare ale BIOA

În unele probleme specifice, precum FS, soluțiile au doar valori binare $\{0, 1\}$, fapt ce sugerează folosirea unor versiuni binare ale BIOA. Aceste versiuni prezintă avantajul reducerii dimensiunii spațiului de căutare, simplificând sarcina determinării soluțiilor optimale.

4.3.1 Optimizarea binară cu lupul cenușiu

În optimizarea continuă GWO, lupii își modifică continuu pozițiile în spațiu. Versiunea binară a optimizării cu lupul cenușiu (BGWO) a fost propusă pentru realizarea FS (52). Formula de actualizare a poziției lupilor este funcție de trei vectori de poziție, $x_\alpha, x_\beta, x_\delta$, care atrag fiecare lup spre primele trei cele mai bune soluții. În algoritmul BGWO, mulțimea soluțiilor este în formă binară la orice moment: toate soluțiile sunt vârfuri ale unui hipercub.

4.3.2 Optimizarea binară cu leul furnicilor

Algoritmul ALO are numeroase avantaje precum *explorarea îmbunătățită*, *evitarea optinelor locale*, *exploatarea*, și *convergența* (19), fapt ce sugerează folosirea lui și în alte aplicații. Algoritmii de optimizare continuă se folosesc pentru determinarea de combinații de caracteristici care maximizează performanța clasicatorului iar agenții de căutare sunt poziționați într-un spațiu d -dimensional de căutare pe poziții în intervalul $[0, 1]$. În cazul algoritmilor binari de optimizare, spațiul de căutare este limitat la doar două valori $\{0, 1\}$ pe dimensiune, deci algoritmul va funcționa mai eficient. De asemenea, operatorul binar este mult mai simplu decât cel continuu.

4.3.3 Optimizarea binară cu fluturele de lampă

În algoritmul MFO, fluturii își modifică continuu poziția în spațiu pe o spirală (mișcare în spirală, spiral moving SM). SM este componenta computațională principală a algoritmului, deoarece ea decide cum se repositionează fluturii în jurul surselor de lumină,

4. SISTEMUL PROPUȘ ȘI METODOLOGIA DE CERCETARE

fapt ce permite fluturelul să zboare în jurul sursei sale de lumină și nu în spațiul dintre surse.

4.4 Algoritmi folosiți pentru comparare

Algoritmii folosiți pentru comparare în acest studiu sunt de dată recentă: ALO, MFO, GWO, SSO, DA, MAKHA, FPA, BAT, CS și FFA. Toți algoritmii de optimizare au fost comparați cu doi algoritmi consacrați pentru FS: PSO and GA, ca și cu variantele fiecăruia dintre algoritmii propuși.

4.5 Metode de inițializare

În studiile noastre, am folosit patru metode de inițializare a parametrilor, capabile să asigure convergența din diverse poziții inițiale: *inițializarea mică*, *inițializarea uniformă*, *inițializarea mare* și *inițializarea MRMR*. *Inițializarea mică* se folosește pentru a testa capabilitatea de căutare globală a unui algoritm de optimizare dat atunci când pozițiile inițiale ale agenților de căutare sunt departe de cele optime. *Inițializarea mare* evaluează capabilitatea de căutare locală a unui algoritm de optimizare atunci când pozițiile agenților de căutare sunt în vecinătatea soluției optime și este necesară doar căutarea locală pentru a ajunge la soluția optimală. Modelele de *inițializare uniformă* se folosesc atunci când pozițiile agenților de căutare se setează aleator în spațiul de căutare folosind distribuția uniformă pentru fiecare dimensiune a spațiului. În sfârșit, *inițializarea MRMR* este o metodă de inițializare bazată pe filtre care exploatează datele indiferent ce clasificator se folosește pentru a selecta o submulțime de caracteristici.

4.6 Metrici de performanță

Fiecare algoritm a fost rulat de $K * M$ ori cu o poziționare aleatoare a agenților de căutare, cu excepția soluției care includea toate caracteristicile, care a fost forțată să fie pe poziția unuia dintre agenții de căutare. Forțarea soluției complete garantează că toate submulțimile ulterioare de caracteristici, dacă sunt selectate ca soluție globală, sunt mai potrivite decât aceasta. Execuțiile repetate ale algoritmilor de optimizare au fost folosite pentru testarea convergenței. Indicatorii (măsurile) folosite pentru a compara algoritmii, precum media valorilor de potrivire, cea mai bună (cea mai slabă) valoare

4.6 Metrice de performanță

de potrivire, abaterea standard a valorilor de potrivire, eroarea predicției și numărul (media) caracteristicilor selectate.

Chapter 5

Aplicarea optimizării inspirate din biologie în procesele farmaceutice

În acest capitol discutăm aplicarea BIOA în diverse procese din industria farmaceutică. ELM și ANN se folosesc ca funcții de potrivire în modelele de regresie, ca în formula (4.3). S-au prelucrat trei seturi de date din industria farmaceutică: compactarea cu role, compactarea cu ștanțe și acidul poli-lactic-coglicolic (poly-lactic-co-glycolic acid PLGA). Rezultatele prezentate în acest capitol sunt preluate din lucrările publicate de autor (53), (54) și (55).

Seturile de date.

Am folosit trei seturi de date, preluate din cele trei procese farmaceutice menționate mai sus în cadrul proiectului IPROCUM. De asemenea, pentru PLGA am folosit și alte seturi de date din literatură, deoarece PLGA are o importanță deosebită în multe aplicații farmaceutice.

5.1 Analiza farmaceutică și discuții

În datele farmaceutice, pentru realizarea FS folosind modelarea CI, ANN se folosește ca model de regresie pentru evaluarea performanței predicției fiecărui algoritm folosit. Scopul acestei secțiuni este folosirea BIOA la FS pentru domeniul farmaceutic, pentru minimizarea numărului caracteristicilor selectate și reducerea erorii predicției. De asemenea, vom evidenția caracteristicile selectate și importanța fiecăreia dintre ele în modelul de predicție.

5.1 Analiza farmaceutică și discuții

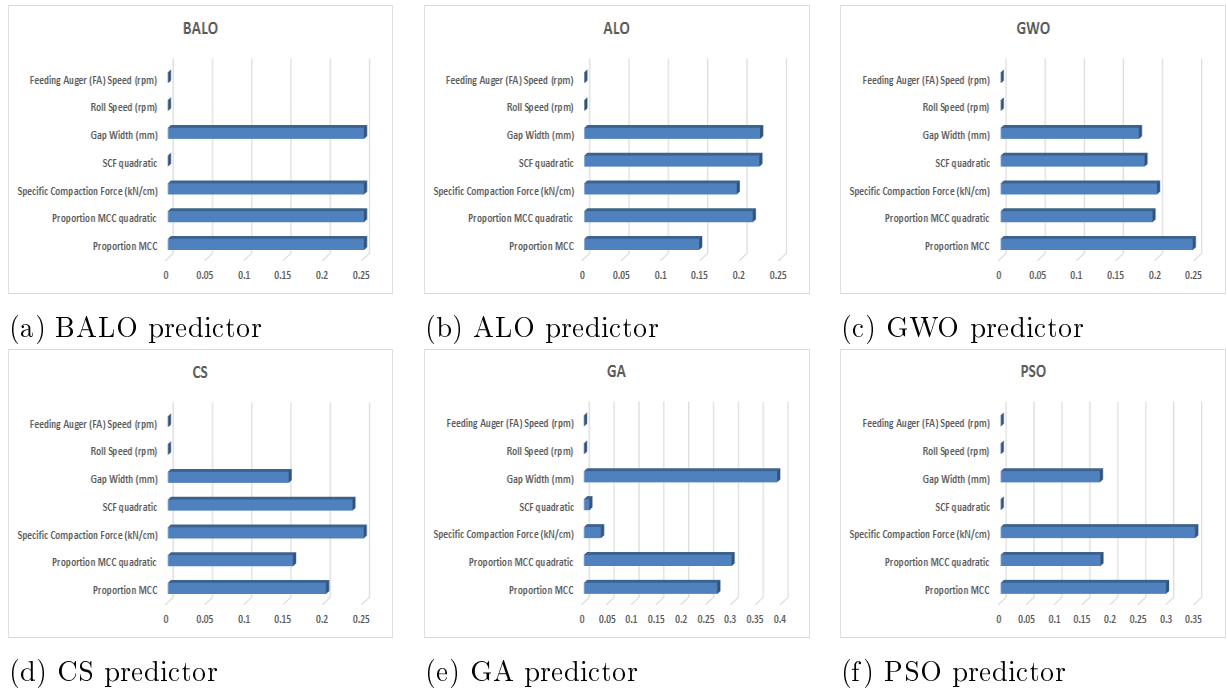


Figure 5.1: Exemplu de importanță a percentilelor pentru ieșirea granule X10

5.1.1 Rezultate de la compactarea cu role

În acest experiment, s-au folosit toate datele (caracteristicile) de intrare (7) pentru a face predicția a 4 date (caracteristici) de ieșire (Granule X10, Granule X50, Granule X90 și granule reziduale). S-au aplicat 6 BIOA pentru FS, combinați cu ELM pentru regresie (predicția celor 4 ieșiri). S-a observat că cele mai importante caracteristici sunt proporția de MCC, proporția de MCC pătratică, forța specifică de compactare (kN/cm) și distanța dintre role (mm), după cum se poate observa în figura (5.1). Mai mult, se poate conchide că GWO este cel mai bun dintre algoritmi de optimizare folosiți în acest experiment, pentru două ieșiri (predicția de percentile pentru Granule X50 și granule reziduale), realizând un bun compromis între cele două obiective contradictorii, RMSE și reducere.

5.1.2 Analiza compactării cu ștanțe

Algoritmi BIOA s-au folosit pentru FS, pentru a obține date de calitate pentru ANN. Pe urmă, ANN este folosită pentru predicția celor două ieșiri continue (porozitatea și

5. APLICAREA OPTIMIZĂRII INSPIRATE DIN BIOLOGIE ÎN PROCESELE FARMACEUTICE

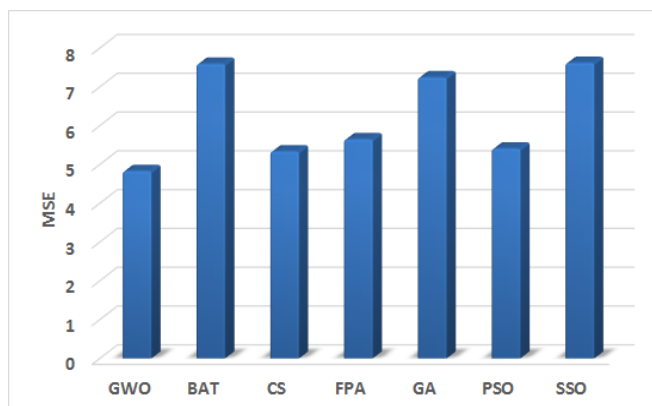


Figure 5.2: Pătratul mediu al erorii (Mean Square Error MSE) pentru porozitate

rezistența la trațiune). Rezultatele prezentate în figura (5.2) conțin valorile MSE pentru fiecare algoritm (20 de rulări diferite), iar în figura (5.3) este reprezentată reducerea medie de caracteristici pentru cele două ieșiri amintite anterior. S-a constatat că algoritmul GWO este cel mai precis pentru predicția porozității, iar SSO pentru predicția rezistenței la trațiune. Cea mai mare reducere a numărului de caracteristici a fost de 60%, cu un pătrat mediu al erorii (MSE) de 7.2 pentru predicția porozității, respectiv 5.1 pentru predicția rezistenței la trațiune. Cumulativ, algoritmul GWO a obținut cel mai bun compromis între MSE și reducere pentru ambele ieșiri. Cele mai importante intrări au fost "presiunea de compactare", urmată de "material" și de "limita superioară a dimensiunii granulei" (53).

5.1.3 Acidul poli-lactic-co-glicolic PLGA

Pentru consistența comparării rezultatelor ambelor abordări, datele de intrare au fost preluate din literatură, iar structura lor a fost menținută la fel ca în Szłek et al. (54). Pe scurt datele s-au colectat din aproximativ 200 de lucrări științifice. Datele extrase constă din ratele de livrare a 68 formule de PLGA din 24 publicații. Vectorul de intrare original avea 320 variabile (descriptori moleculari ai proteinelor, excipienți, caracteristici ale formulilor și condiții experimentale) și 745 observații (înregistrări). Selectarea variabilelor esențiale a fost realizată în ipoteza că predictibilitatea și simplitatea modelului sunt la fel de importante pentru rezultatul final. La început, s-au folosit 4 BIOA – ALO, BALO, GWO and SSO – ca instrument pentru FS. Pentru fiecare mulțime

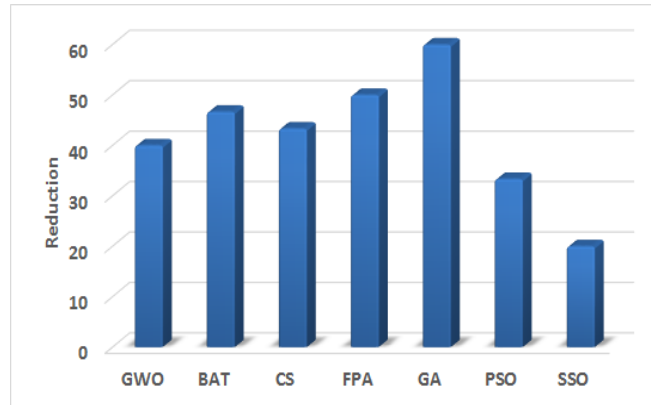


Figure 5.3: Reducerea medie pentru porozitate și rezistența la tracțiune

de caracteristici propusă de instrumentul de FS, s-a folosit o procedură de screening pentru determinarea erorii de generalizare minime peste diversele modele predictive și setările/arhitecturile acestora. Pentru modelarea predictivă s-au folosit Cubist, RF, ANN (MLP monoton, MLP deep learning) și FugeR.

5.1.3.1 Rezultate și discuții pentru PLGA

Măsurile descrise în secțiunile următoare se folosesc pentru cuantificarea calității rezultatelor obținute de modelele computaționale. Pentru măsurarea preciziei clasificării datelor, s-a folosit RMSE (root mean squared error, radical din eroarea medie pătratică) atât la FS cât și la modelele predictive. Dimensiunea reducerii, indicator de dimensiune, este folosit la metodele de FS. În total, considerând toate metodele și mulțimile 10-cv, aproape 18,000 de modele au fost antrenate și testate. RMSE normalizată (NRMSE, normalized RMSE) a variat de la 31.1 la 15.9%. Algoritmul RF (Random Forest RF) a produs cea mai mică eroare; în consecință, a fost folosit pentru selectarea vectorilor de intrare optimali, vezi tabelul (5.1). Modelul RF dezvoltat pe 9 seturi de date de intrare, 9in(2), selectate cu algoritmul BALO, a produs una dintre cele mai mici NRMSE, vezi tabela (5.2). S-au obținut rezultate comparabile cu cele ale lui Szlęk et al. (54), adică 15.97% față de 15.4%, dar vectorul intrărilor a fost mai mic, adică 9 (în loc de 11) (55).

5. APLICAREA OPTIMIZĂRII INSPIRATE DIN BIOLOGIE ÎN PROCESELE FARMACEUTICE

Table 5.1: NRMSE pentru vectori de intrare selecta cti de BIOA

Metoda FS	Nr. intrări	Cubist	Mon-mlp	RF	FugeR
ALO	8	22.45	24.55	21.95	-
	12	25.95	25.15	22.19	-
	20	18.73	20.20	16.33	20.15
BALO	9(1)	21.20	20.63	18.81	-
	9(2)	18.26	17.31	15.97	18.09
	9(3)	22.60	21.88	19.79	-
	11	19.40	19.35	18.70	-
	12	17.26	18.17	16.56	18.73
GWO	15	19.30	18.88	16.73	19.10
	18	20.65	18.58	17.63	-
	24	20.30	22.30	17.90	-
	25	20.04	19.29	15.86	19.10
	26	17.32	22.22	16.22	-
SSO	8	30.49	31.12	28.89	-
	13	27.09	25.82	24.86	-

Table 5.2: Resultate pentru 9in (2), antrenat și testat pe 10cv seturi de date.

Algoritmul	NRMSE	R2
Cubist	18.26	0.611
Monmlp	17.31	0.652
ANN cu deep learning	16.87	0.655
FugeR	18.09	0.612
RF	15.97	0.692

Chapter 6

Concluzii și cercetări viitoare

6.1 Concluzii

În lucrarea de față, au fost propuși și aplicați BIOA pentru FS în modul wrapper (învelitoare). Au fost studiați cei mai recentți BIOA, iar rezultatele obținute au fost comparate cu metode consacrate de FS, precum PSO și GA. Evaluarea s-a efectuat folosind o mulțime de criterii de evaluare, în scopul evidențierii unor aspecte multiple. Selectarea caracteristicilor (feature selection FS) este formulată ca problemă de optimizare multicriterială, cu funcția de potrivire reflectând performanța predicției și gradul de reducere a numărului de caracteristici. S-a folosit o mulțime de indicatori de evaluare pentru caracterizarea performanțelor algoritmilor de optimizare, aplicați pentru peste 21 de seturi de date în problemele de clasificare și 10 seturi de date în problemele de regresie, date preluate din repozitoriul UCI.

De asemenea, au fost evaluați algoritmi de FS inspirați din biologie în modelarea proceselor de fabricație a tabletelor din industria farmaceutică, în particular pentru predicția unor proprietăți precum porozitatea și rezistența la trațiune a pulberilor. Eficiența modelării a fost evaluată în raport cu reducerea medie a numărului de caracteristici și RMSE. Cel mai precis în predicția porozității s-a dovedit a fi algoritmul GWO, iar pentru predicția rezistenței la trațiune campionul a fost SSO. În final, putem conchide că BIOA sunt algoritmi eficienți de căutare, fiind adecvați pentru rezolvarea problemelor de FS.

6.2 Cercetări viitoare

Pe baza cercetărilor de până acum, rezultă cel puțin următoarele șase noi direcții de cercetare:

1. Evaluarea algoritmilor BIOA propuși pe seturi de date de complexitate mare.
2. Prelucrarea avansată (din punct de vedere statistic) a rezultatelor experimentale.
3. Folosirea BIOA la rezolvarea unor probleme interesante și din alte domenii.
4. Folosirea și a altor tehnici de ML, precum SVM, SVR și RF, pentru evaluarea potrivirii bazată pe învelitoare (wrapper-based).
5. Propunerea unei funcții de potrivire multi-obiectiv care folosește BIOA pentru determinarea unui subset optimal de caracteristici.
6. Propunerea unor combinații ale tehnicilor de BIOA recente la rezolvarea problemelor de FS.

References

- [1] B. CHIZI, L. ROKACH, AND O. MAIMON. **A Survey of Feature Selection Techniques.** *IGI Global*, pages 1888–1895, 2009. 8
- [2] R.O. DUDA, P.E. HART, AND D.G. STORK. *Pattern Classification, 2nd Edition.* Wiley-Interscience, 2000. 8, 9
- [3] Y. CHEN, D. MIAO, AND R. WANG. **A rough set approach to feature selection based on ant colony optimization.** *Pattern Recognition Letters*, **31**(3):226–233, 2010. 8
- [4] S. SHOGHIAN AND M. KOUZEHGAR. **A Comparison among Wolf Pack Search and Four other Optimization Algorithms.** *Computer, Electrical, Automation, Control and Information Engineering*, **6**(12):1619–1624, 2012. 9
- [5] I.F. JR., X.S. YANG, I. FISTER, J. BREST, AND D. FISTER. **A Brief Review of Nature-Inspired Algorithms for Optimization.** *Elektrotehniski Vestnik*, **80**(3):116–122, 2013. 9
- [6] I. GUYON AND A. ELISSEEFF. **An introduction to variable and attribute selection.** *Machine Learning Research*, **3**:1157–1182, 2003. 9
- [7] B. XUE, M. ZHANG, AND W.N. BROWNE. **Particle swarm optimization for feature selection in classification: a multi-objective approach.** *IEEE transactions on cybernetics*, **43**(6):1656–1671, 2013. 9
- [8] R. AKBARI, A. MOHAMMADI, AND K. ZIARATI. **A novel bee swarm optimization algorithm for numerical function optimization.** *Communications in Nonlinear Science and Numerical Simulation*, **15**:3142–3155, 2010. 9, 16

REFERENCES

- [9] V. VAPNIK. **Statistical Learning Theory**. *IEEE Transactions on Neural Networks*, **10**(5):988–999, 1999. 11
- [10] Y. MICHE, A. SORJAMAA, P. BAS, O. SIMULA, C. JUTTEN, AND A. LENDASSE. **OP-ELM: Optimally Pruned Extreme Learning Machine**. *IEEE Transactions Neural Networks*, **21**(1):158–162, 2010. 11
- [11] C. JIUWEN AND L. ZHIPING. **Extreme Learning Machines on High Dimensional and Large Data Applications: A Survey**. *Mathematical Problems in Engineering*, **2015**(1):1–13, 2015. 11, 14
- [12] J.H. HOLLAND. *Adaptation in natural and artificial systems*. MIT Press, Cambridge, MA, USA, 1992. 12
- [13] B. CHAKRABORTY. **Genetic algorithm with fuzzy fitness function for feature selection**. In *International Symposium on Industrial Electronics*, pages 315–319. IEEE, 2002. 12
- [14] R. EBERHART AND J. KENNEDY. **A New Optimizer Using Particle Swarm Theory**. In *International Symposium on Micro Machine and Human Science*, pages 39–43. IEEE, 1995. 12, 15
- [15] H.H. GAO, H.H. YANG, AND X.Y. WANG. **Ant colony optimization based network intrusion feature selection and detection**. In *International Conference on Machine Learning and Cybernetics*, pages 3871–3875. IEEE, 2005. 12
- [16] H. MING. **A rough set based hybrid method to feature selection**. In *International Symposium on Knowledge Acquisition and Modeling*, pages 585–588. IEEE, 2008. 12
- [17] D. KARABOGA AND B. BASTURK. **A powerful and efficient algorithm for numerical function optimization: artificial bee colony (ABC) algorithm**. *Journal of Global Optimization*, **39**:459–471, 2007. 12, 16
- [18] X.S. YANG. **Engineering optimizations via nature-inspired virtual bee algorithms**. *Artificial Intelligence and Knowledge Engineering Applications: A Bioinspired Approach*, **3562**:317–323, 2005. 12

-
- [19] S. MIRJALILI. **The Ant Lion Optimizer**. *Advances in Engineering Software*, **83**:80–98, 2015. 12, 18, 25
- [20] A. KIBRIYA AND E. FRANK. **An empirical comparison of exact nearest neighbour algorithms**. *European Conference on Principles and Practice of Knowledge Discovery in Databases*, **4702**:140–151, 2007. 13
- [21] V.Y. KULKARNI AND P.K. SINHA. **Efficient Learning of Random Forest Classifier using Disjoint Partitioning Approach**. *World Congress on Engineering*, **2**, 2013. 14
- [22] L. BREIMAN. **Random forests**. *Machine learning*, **45**(1):5–32, 2001. 14
- [23] G. BIAU, L. DEVROYE, AND G. LUGOSI. **Consistency of Random Forests and Other Averaging Classifiers**. *Machine Learning Research*, **9**:2015–2033, 2008. 14
- [24] D.S. HUANG. **A constructive approach for finding arbitrary roots of polynomials by neural networks**. *IEEE Transaction on Neural Networks*, **15**(2):477–491, 2004. 14
- [25] G.B. HUANG, Q.Y. ZHU, AND C.K. SIEW. **Extreme Learning Machine: theory and applications**. *Neurocomputing*, **70**(1):489–501, 2006. 14
- [26] J.R. QUINLAN. **Learning with continuous classes**. *Proceedings of the 5th Australian joint Conference on Artificial Intelligence*, **92**:343–348, 1992. 15
- [27] A. LIAW AND M. WIENER. **Classification and Regression by randomForest**. *R News*, **2**(3):18–22, 2002. 15
- [28] A.J. CANNON. **monmlp: Monotone Multi-Layer Perceptron Neural Network** [online]. 2015 [cited 2016 Jan 15]. 15
- [29] H. ZHANG AND Z. ZHANG. **Feedforward networks with monotone constraints**. *International Joint Conference on Neural Networks*, **3**:1820–1823, 1999. 15
- [30] Y. BENGIO. **Learning deep architectures for AI**. *Foundations and Trends in Machine Learning*, **2**(1):1–127, 2009. 15

REFERENCES

- [31] B. XUE, M. ZHANG, AND W.N. BROWNE. **Particle swarm optimisation for feature selection in classification: Novel initialisation and updating mechanisms.** *Applied Soft Computing*, **18**:261–276, 2014. 15, 21
- [32] B. CHAKRABORTY. **Feature subset selection by particle swarm optimization with fuzzy fitness function.** In *Third International Conference on Intelligent System and Knowledge Engineering*, pages 1038–1042. IEEE, 2008. 15
- [33] X.S. YANG. **Nature-Inspired Metaheuristic Algorithms.** *Luniver Press, UK*, **2nd Edition**, 2010. 16, 19, 21, 24
- [34] X.S. YANG AND S. DEB. **Cuckoo Search via Levy Flights.** *World Congress on Nature and Biologically Inspired Computing*, 2009. 16
- [35] X.S. YANG. **A New Metaheuristic Bat-Inspired Algorithm.** *Nature Inspired Cooperative Strategies for Optimization*, **284**:65–74, 2010. 16
- [36] X.S. YANG. **Flower pollination algorithm for global optimization.** *Unconventional Computation and Natural Computation, Lecture Notes in Computer Science*, **7445**:240–249, 2012. 16
- [37] X.S. YANG, M. KARAMANOGLU, AND X. HE. **Multi-objective Flower Algorithm for Optimization.** *International Conference on Computational Science, Procedia Computer Science*, **18**:861–868, 2013. 17
- [38] E. CUEVAS, M. CIENFUEGOS, D. ZALDIVAR, AND M. PEREZ-CISNEROS. **A swarm optimization algorithm inspired in the behavior of the social-spider.** *Expert Systems with Applications*, **40**(16):6374–6384, 2013. 17
- [39] T. JONES AND S. RIECHERT. **Patterns of reproductive success associated with social structure and microclimate in a spider system.** *Animal Behaviour*, **76**(6):2011–2019, 2008. 17
- [40] S. MIRJALILI, S.M. MIRJALILI, AND A. LEWIS. **Grey Wolf Optimizer.** *Advances in Engineering Software*, **69**:46–61, 2014. 17, 23
- [41] A.M.E. KHALIL, S.K. FATEEN, AND A. BONILLA-PETRICIOLET. **MAKHAÛA New Hybrid Swarm Intelligence Global Optimization Algorithm.** *Algorithms*, **8**(2):336–365, 2015. 17

-
- [42] S. MIRJALILI. **Dragonfly algorithm: a new meta-heuristic optimization technique for solving single-objective, discrete, and multi-objective problems.** *Neural Computing and Applications*, **1**(1):1–21, 2015. 17
- [43] J.H. THORP AND D.C. ROGERS. *Thorp and Covich's freshwater invertebrates, 4th Edition.* Elsevier, 2014. 17
- [44] M. WIKELSKI, D. MOSKOWITZ, J.S. ADELMAN, J. COCHRAN, D.S. WILCOVE, AND M.L. MAY. **Simple rules guide dragonfly migration.** *Biology letters*, **2**(1):325–329, 2006. 18
- [45] S. MIRJALILI. **Moth-Flame Optimization Algorithm: A Novel Nature-inspired Heuristic Paradigm.** *Knowledge-Based Systems*, **89**:228–249, 2015. 18
- [46] K.J. GASTON, J. BENNIE, T.W. DAVIES, AND J. HOPKINS. **The ecological impacts of nighttime light pollution: a mechanistic appraisal.** *Biological reviews*, **88**:912–927, 2013. 18
- [47] A.E. HASSANIEN AND E. EMARY. **Swarm Intelligence: Principles, Advances, and Applications.** *CRC Press, Taylor & Francis Group*, 2015. 21
- [48] I. PAVLYUKEVICH. **Levy flights, non-local search and simulated annealing.** *Computational Physics*, **226**(2):1830–1844, 2007. 21
- [49] A.M. REYNOLDS AND M.A. FRYE. **Free-flight odor tracking in *Drosophila* is consistent with an optimal intermittent scale-free search.** *PLoS One*, **2**(4), 2007. 21
- [50] R. VOHRA AND B. PATEL. **An Efficient Chaos-Based Optimization Algorithm Approach for Cryptography.** *Communication Network Security*, **1**(4):75–79, 2012. 21, 22
- [51] B. REN AND W. ZHONG. **Multi-objective optimization using chaos based PSO.** *Information Technology*, **10**(10):1908–1916, 2011. 22
- [52] E. EMARY, H.M. ZAWBAA, AND A.E. HASSANIEN. **Binary Grey Wolf Optimization Approaches for Feature Selection.** *Neurocomputing, Journal indexed*

REFERENCES

- in 'Journal Citation Reports' (Thomson Reuters), Impact Factor (2014): 2.083, 172:371–381, 2016. 25*
- [53] H.M. ZAWBAA, S. SCHIANO, L. PEREZ-GANDARILLAS, C. GROSAN, C.Y. WU, AND A. MICHRAFY. **An Evaluation of Bio-inspired Feature Selection Techniques for Computational Intelligence Modeling of Die Compaction.** *International congress on Particle Technology (PARTEC)*, 2016. 28, 30
- [54] J. SZŁĘK, A. PACLAWSKI, R. LAU, R. JACHOWICZ, AND A. MENDYK. **Heuristic modeling of macromolecule release from PLGA microspheres.** *Nanomedicine*, 8:4601–4611, 2013. 28, 30, 31
- [55] H.M. ZAWBAA, J. SZLEK, C. GROSAN, A. MENDYK, AND R. JACHOWICZ. **Computational modelling and optimization of the macromolecule release from PLGA microspheres.** *PLOS ONE*, 2016. 28, 31