

**UNIVERSITATEA "BABEȘ-BOLYAI" CLUJ-NAPOCA
FACULTATEA DE CHIMIE ȘI INGINERI CHIMICĂ**

**REZUMAT
TEZĂ DE DOCTORAT**

STUDII DE SIMILARITATE TOPOLOGICĂ

**CONDUCĂTOR ȘTIINȚIFIC
Prof.Dr. MIRCEA V. DIUDEA**

**DOCTORAND
CRISTINA DORINA MOLDOVAN**

2012

CUPRINS

INTRODUCERE

CAPITOLUL 1. SIMILARITATEA TOPOLOGICĂ

1.1. Descrierea cantitativă a similarității

1.1.1. Similaritatea structurilor moleculare 2D și 3D

1.1.2. Similaritatea structurilor moleculare descrise prin câmpuri de forțe

1.1.3. Similaritatea structurilor moleculare descrise cuanto-chimic

1.1.4. Concluzii privitoare la modalitățile de descriere a similarității

CAPITOLUL 2. MATRICI TOPOLOGICE

2.1. Matricea de adiacență

2.2. Matricea de conectivitate

2.3 Matricea distanțelor

2.4. Matrici strat

2.4.1. Matrici SeM (sequence matrices)

2.4.2. Matrici LeM (layer matrices)

CAPITOLUL 3. INDICI TOPOLOGICI

3.1. Construcția indicilor topologici

3.2. Principalii indici topologici

3.2.1. Indici bazați pe matricea de adiacență

3.2.1.1. Indicele adiacenței totale:

3.2.1.2. Indicele RANDIC

3.2.2. Indici bazați pe matricea distanțelor

3.2.2.1 Indicele WIENER

3.2.2.4. Indicele BALABAN

3.2.3. Indici bazați pe matrici pătrate dense

3.2.3.1. Indicele hyper-Wiener, R

3.2.4. Indici bazați pe matrici strat

3.2.4.1. Indici de centricitate

3.2.4.2. Indici de centrocomplexitate

3.2.5 Indici de similaritate moleculară

3.2.5.1. Indicele de similaritate Cosine

3.3.5.2. Indicele de similaritate Dice

3.3.5.3. Indicele de similaritate Richards

3.3.7. Metode de ponderare

CAPITOLUL 4. RELAȚII CANTITATIVE STRUCTURĂ- ACTIVITATE BIOLOGICĂ

4.1. Metoda Hansch clasică. Parametri structurali pentru QSAR

4.1.1. Introducere

4.1.2. Ecuația Hansch

4.2. Modele avansate în QSAR

4.2.1. Analiza substructurilor

CAPITOLUL 5. ANALIZA DE DATE

5.1. Noțiuni introductive

5.2. Regresii liniare

5.3 Analize de date multivariate

5.3.1. Analiza clusterilor

5.3.1.1. Obiectivele analizelor de clusteri

5.3.1.2. Unități asemănătoare și preprocesarea informațiilor

- 5.3.1.3. Algoritmul clusterării
- 5.3.2. Analiza componentelor principale
 - 5.3.2.1. Considerente teoretice
- 5.3.3. Analiza factorială
 - 5.3.3.1. Ecuația de bază a analizei factoriale
 - 5.3.3.2. Valori proprii, vectori proprii
 - 5.3.3.4. Numărul factorilor de extras

CONTRIBUTII PERSONALE

CAPITOLUL 6. METODE SI MODELE DE CLUSTERARE

- 6.1. Reprezentarea și căutarea structurilor chimice
- 6.2. Metode de clusterare a sistemelor informatice chimice
- 6.3. Compuși steroidici
 - 6.3.1. Calcularea descriptorilor moleculari
 - 6.3.2. Prelucrarea datelor
- 6.4. Compuși cu acțiune antibacteriană
 - 6.4.1. Obținerea descriptorilor moleculari
 - 6.4.2. Analiza și prelucrarea datelor
- 6.5. Noi modele QSAR pentru prezicerea activității biologice a derivaților de benzoxazol/benzimidazol
 - 6.5.1. Obținerea descriptorilor moleculari
 - 6.5.2. Prelucrarea și analiza datelor

CONCLUZII

BIBLIOGRAFIE

Cuvinte cheie: matrice, indice, descriptori topologici, activitate biologică, QSAR, analize de date, APC, statistica, coeficient de corelare, similaritate, androstan, benzimidazol, benzoxazol, SIMIL, TOPOCLUJ, HyperChem, DRAGON

INTRODUCERE

Matematizarea chimiei are o lungă și diversificată istorie care începe acum aproape două secole. În 1874, Alexander Crum Brown (1838-1922), unul dintre fondatorii teoriei structurii chimice a prezis următoarele: "...chimia va deveni o ramură a matematicii aplicate; dar aceasta nu va înceta a fi o știință experimentală. Matematica ne va permite o privire retrospectivă pentru a justifica rezultatele obținute de experiment, ce va fi utilă în cercetări și câteodată chiar la prezicerea în întregime a noilor descoperiri. Noi nu știm când schimbarea va avea loc sau dacă aceasta va fi treptată sau bruscă...". Această prezicere pare a fi în curând împlinită.

Proiectarea de structuri moleculare cu proprietăți fizico-chimice ori biologice dorite reprezintă una din principalele obiective ale diferitelor ramuri ale industriei (chimice, farmaceutice, etc.), în special și a cercetătorului în particular. Realizarea unui produs nou cu o anumită proprietate chimică sau activitate biologică, presupun mari cheltuieli atât material cât și umane. Design-ul de medicamente este un proces iterativ care începe cu un compus care afișează un profil interesant biologic și se termină cu optimizarea profilului de activitate pentru molecula de sinteză. Procesul este inițiat atunci când chimistul concepe o ipoteză care leagă caracteristicile chimice ale moleculei (sau serie de molecule), de activitatea biologică. Fără o înțelegere detaliată a procesului biochimic, responsabil pentru activitatea biologică, în general, ipoteza este rafinată prin examinarea asemănărilor și diferențelor structurale pentru moleculele active și inactiv.

În acest sens, prezenta teză dezvoltă conceptele utilizate pentru a determina diferite modele QSPR/QSAR precum și proiectarea unor noi structuri pe baza acestora. Aceasta este structurată în două părți distincte. Prima parte se referă la prezentarea domeniului topologiei moleculare, a metodelor și tehnicilor folosite, iar a doua parte se referă la contribuțiile personale în modelarea relațiilor structură-proprietate și structură-activitate biologică:

Capitolul I "Similaritate topologica" prezintă conceptul și metodele de similaritate moleculară.

Capitolul II "Matrici topologice", prezintă pe lângă matricile clasice, noi matrici de tip Szedged și Cluj.

Capitolul III "Indici topologici" tratează descriptorii topologici obținuți pe baza matricilor topologice. Sunt prezentați principalii indici topologici utilizați și noi indici topologici propuși, SP, indici de tip Szedged și Cluj.

Capitolul IV "Relații cantitative structură- activitate biologică" tratează diferitele metode și modele utilizate pentru descrierea proprietăților fizico-chimice și biologice ale compușilor.

Capitolului V "Analiza de date " se referă la metodele de prelucrare statistică utilizate în proiectarea modelelor QSPR/QSAR. De asemenea sunt prezentate câteva aplicații ale acestora pe diferite seturi de structuri.

În partea de contribuții proprii sunt prezentate rezultatele obținute ca urmare a analizei unor seturi de structuri chimice, cu ajutorul metodelor de similaritate, precum și prezicerea proprietăților unor compuși noi cu potențială activitate biologică.

CAPITOLUL 1. SIMILARITATEA TOPOLOGICĂ

Similaritatea structurilor moleculare exprimă existența unor trăsături comune într-un set de molecule. Similaritatea este definită pe baza unor variate criterii și/sau proceduri și ea generează clase de echivalență intermoleculară, în cadrul setului de molecule.

Similaritatea moleculară, la fel ca ramificarea, este o noțiune intuitivă, astfel că nu se poate defini o măsură unică și non-ambiguă de similaritate. Multe dintre ariile științei cum sunt: chimia organică sintetică, biologia structurală, farmacologia și toxicologia sunt necesare în procesul de dezvoltare a medicamentelor. Costurile foarte înalte și timpul îndelungat, caracteristice acestui proces, accentuează necesitatea pentru investiții în tehnologii care să accelereze procesul de proiectare a noilor structuri chimice, cu efectele biologice dorite, și să scurteze timpul până când acestea sunt puse pe piață (aceiași lucru este valabil și pentru pesticide, fungicide, etc., chimicale pentru agricultură). Aceste investiții au dus la dezvoltarea sistemelor foarte sofisticate de stocare, căutare și procesare a unei variate game de informații chimice.

Descrierea moleculară, utilizată în analiza similarității moleculare, se realizează cu ajutorul descriptorilor moleculari. Orice descriere moleculară induce o partiționare în clase de echivalență a setului de molecule. Este necesar să definim aici relația de echivalență:

Fie S un set de structuri moleculare și R o relație binară pe S legând perechi de molecule. Dacă $x, y \in S$ sunt astfel legate, atunci se poate scrie xRy . Relația R este o relație de echivalență dacă satisface următoarele proprietăți:

$$xRx, \text{ pentru oricare } x \in S \quad (\text{reflexivitate}) \quad (1.1)$$

$$\text{dacă } xRy, \text{ atunci } yRx \quad (\text{simetrie}) \quad (1.2)$$

$$\text{dacă } xRy \text{ și } yRz, \text{ atunci } xRz \quad (\text{tranzitivitate}) \quad (1.3)$$

Subsetul elementelor $y \in S$, aflate în relația xRy , reprezintă clasa de echivalență a lui x . Impunând relația de echivalență R pe setul S rezultă o partiționare a S în subseturi disjuncte, numite clase de echivalență sub R . Un astfel de subset, notat S/R , se mai scrie și S modulo R .¹

Fie f o funcție *mapping* (de acoperire, potrivire, asociere, etc.) a elementelor setului S peste elementele unui set oarecare Y . Adică, pentru oricare $x \in S$, f atribuie o valoare corespunzătoare $y = f(x)$ în Y . Această corespondență poate fi scrisă ca $f: S \rightarrow Y$. Dacă Y este setul descrierilor, funcția de suprapunere asociază o descriere moleculară fiecărei molecule în S . Acele molecule din S sunt echivalente care au aceeași descriere moleculară. O astfel de funcție f poate fi o numerotare, un cod, ori un simplu proces de măsurare. Se poate demonstra că diferitele descrieri moleculare împreună cu reprezentarea lor algebrică formează un grup.^{1,2}

O potrivire (*matching*) se poate realiza prin suprapunerea a două molecule. O astfel de operație poate indica aspecte comune ale celor două molecule sau ale descrierilor acestora.

O *ordonare parțială* se referă la o anumită ordonare locală, indusă de o acoperire (potrivire) parțială (i.e. substructure matching) între moleculele unui set. Matematic, *relația de ordonare* implică proprietatea de antisimetrie (1.4)

$$\text{dacă } xRy \text{ și } yRx \text{ atunci } x = y, \quad (1.4)$$

în locul celei de simetrie (1.2) (vezi mai sus). Randić³ a raportat o ordonare parțială a isomerilor alcanilor indusă de numărul căilor (path numbers) p_2 și p_3 . Oricare alți descriptori graf-teoretici (e.g., indici topologici, secvențe ale distanțelor, etc.) pot fi utilizați în scopul caracterizării și ordonării parțiale ca și al grupării (*clustering*) structurilor moleculare.

Un întreg volum al revistei *MATCH (Communications in Mathematical and in Computer Chemistry, 2000, 42)* este dedicat ordonării parțiale în chimie.

Compușii poziționați mai aproape, în secvență, este de așteptat să prezinte valori mai apropiate ale unor proprietăți (i.e., proprietăți similare).

Proximitatea exprimă în esență două categorii: *similaritatea* și *disimilaritatea*.

Similaritatea exprimă asemănarea a două molecule, printr-un număr mare, dacă descrierile lor moleculare sunt apropiate și respectiv printr-un număr tinzând spre zero în cazul în care aceste descrieri nu au nimic comun.⁴ De exemplu, la suprapunerea a două molecule, raportul numărului de atomi și legături care se suprapun și numărul total al acestora în molecula întregă, multiplicat cu același raport pentru molecula cu care se compară⁵ a fost propus ca măsură a similarității între două molecule. O astfel de măsură are proprietatea corelației (*zero* pentru non-corelație și *unu* pentru corelație deplină).

Disimilaritatea exprimă asemănarea a două molecule, cu un număr apropiat de zero când descrierile lor moleculare sunt apropiate și respectiv printr-un număr mare în cazul în care aceste descrieri sunt diferite.⁴ În exemplul de suprapunere de mai sus, numărul atomilor și legăturilor care nu se suprapun poate fi luat ca măsură a disimilarității între două molecule. Acest caz particular de disimilaritate⁶ este cunoscut în literatură ca *distanță chimică*.^{7,8,9,10}

Cercetări recente efectuate în domeniul proiectării substanțelor bioactive ("*drug design*"), în special medicamente dar și insecticide, erbicide, etc., au relevat importanța deosebită a *similarității*¹ compușilor implicați, pe care-i vom denumi în continuare *efectori* (E), pentru interacția lor cu un anumit *receptor biologic* (R), oricare ar fi natura lui. Deoarece în domeniul *drug design* se urmărește obținerea unei relații cantitative între structura chimică a efectorilor E și activitatea lor biologică - măsurată de obicei ca $\log(1/C)$, unde C reprezintă concentrația molară (doza) care determină un răspuns biologic constant a devenit cuantificarea similarității, respectiv disimilarității compușilor chimici. Aceasta presupune găsirea unor descriptori cantitativi adecvați care să poată fi utilizați în relații liniare (sau neliniare) structură chimică - activitate biologică (**QSAR = Quantitative Structure - Activity Relationship**) și/sau pentru proiectarea seriilor QSAR.

Deși ideea intuitivă de similaritate, așa cum a fost ea utilizată de chimiștii organicieni - fiind "măsurată" prin numărul de caracteristici structurale comune pentru doi compuși și prin aranjamentul lor reciproc - pare simplă, în realitate există mai multe abordări datorate multiplelor reprezentări posibile ale compușilor chimici organici (E): prin formula moleculară, topologic - cu ajutorul grafurilor moleculare constituționale plane (2D), sau, mai exact, prin introducerea unei metrici euclidiene peste structurile moleculare reprezentate în spațiul cartezian (3D) - formulele configuraționale precum și cele conformaționale, sau pe baza formei moleculelor așa cum este ea descrisă prin suprafețele van der Waals (vdW) sau cu ajutorul unor câmpuri de forțe moleculare.

În consecință, noțiunea de similaritate, deși intuitiv este simplă, depinde în mod esențial de punctul de vedere din care este abordată. Din aceste motive, au fost propuse pentru fiecare din modalitățile de descriere a structurii moleculare prezentate mai sus *măsuri diferite de similaritate*, capabile să evalueze cantitativ această noțiune vagă (în sens matematic). Trebuie subliniat faptul că acești *indicatori cantitativi de similaritate* au fost utilizați, aproape în totalitate în domeniul proiectării substanțelor bioactive. În afara similarității, dar în strânsă legătură cu ea, există alte două noțiuni, totuși distincte - *disimilaritatea* și *complementaritatea* - de asemenea frecvent utilizate în QSAR. Fiecare din aceste noțiuni trebuie definită în mod univoc din punct de vedere matematic pentru a evita orice ambiguitate. Noțiunea de complementaritate a fost propusă la sfârșitul secolului trecut de către E. Fischer, care a pus bazele chimiei proteinelor și carbohidraților, pentru a explica specificitatea acțiunii biologice a unor molecule E, prin teoria sa "*cheie în broască*". De atunci noțiunea de complementaritate, ca bază a interacțiunilor R-E, a fost continuu dezvoltată și rafinată, pe baza ei fiind propuse o serie de metode QSAR (MSA², MSD³, MTD³, etc.) și de modele moleculare (CoMFA⁴, etc).

Cele trei concepte menționate mai sus - *similaritatea, disimilaritatea, și complementaritatea* - ar putea fi încadrate pe o scară similară celei de corelare liniară simplă în care un coeficient de corelare $r = +1$ indică o corelare liniară perfectă $y=f(x)$, toate punctele fiind dispuse strict pe o dreaptă cu panta pozitivă (similaritate totală între moleculele Y descrisă prin $y_i, i=1,N$ și X caracterizată prin $x_i, i=1,N, x_i$ și y_i fiind descriptorii structurali), $r = -1$ corespunde unei corelări perfecte inverse (Y și X sunt perfect complementare), dreapta având panta negativă, iar $r = 0$ semnifică absența totală a oricărei corelări (Y și X sunt în totalitate disimilare). Luând în considerare scara de mai sus, apare ca foarte probabil ca două molecule să fie disimilare în ceea ce privește, de exemplu, repartiția atomilor în spațiul 3D sau, topologic, în 2D, având un coeficient de corelare $r \rightarrow 0$ definit în cadrul unei metrici adecvate și să prezinte o similaritate avansată, cu r apropiat de 1 ($r \rightarrow 1$), din punctul de vedere al formei, volumului sau suprafeței, etc.

Similaritatea structurilor moleculare 2D și 3D

Graful molecular este o reprezentare (în general 2D, dar acest fapt nu diminuează generalitatea metodei) a modului de legare a atomilor între ei (conectivitatea moleculară). Atomii unui sistem molecular dat formează nodurile (vârfurile) grafului molecular respectiv iar legăturile chimice dintre ei sunt asimilate muchiilor. Acest graf molecular poate fi descompus, succesiv, în fragmente din ce în ce mai mici, adică în subgrafuri (elemente de teorie legate de această problemă se găsesc în capitolul 1 al primei părți). Odată definite subgrafurile de un anumit tip, compararea a două structuri moleculare pentru a stabili gradul de similaritate sau de disimilaritate poate fi realizată prin găsirea subgrafului comun de mărime maximă (**MCS**) sau prin luarea în considerare a tipului atomilor și a conectivităților lor. Acest ultim mod de cuantificare poate asigura valorilor obținute o mai mare semnificație fizică, fapt neglijat de multe ori în studiile QSAR.

Metoda **MCS** poate fi extinsă la spațiul 3D prin luarea în considerare a matricei de distanță. Astfel, dacă fiecărui fragment considerat i se asociază matricea de distanțe în limitele de încredere adecvate, a fost definită următoarea măsură de similaritate, r_{AB} , între moleculele A și B :

$$r_{AB} = \frac{MCS_{A,B}}{N_A + N_B - MCS_{A,B}} \quad (1)$$

În relația (1) N_A și N_B reprezintă numărul de atomi din, respectiv, moleculele A și B . Principala utilizare a acestui indicator cantitativ este legată de procesul de căutare în bazele mari de date care conțin molecule bioactive, pentru a extrage moleculele similare din punct de vedere structural.

Concluzii privitoare la modalitățile de descriere a similarității

Găsirea unei metode de evaluare cantitativă a similarității este doar una dintre problemele care trebuie rezolvate de către "*proiectantul*" de substanțe bioactive. Odată stabilit modul de calcul este necesar un algoritm care să permită compararea celor două molecule, A și B , a căror similaritate sau disimilaritate urmează a fi evaluată cantitativ. De obicei, una dintre molecule este menținută fixă iar cealaltă este rotită pentru a suprapune anumite caracteristici moleculare cu scopul de a maximiza similaritatea structurilor A și B .

O problemă dificilă este construirea unui algoritm de suprapunere. Una dintre metode se bazează pe rotația unei molecule, B , în jurul axelor Euler ale celeilalte molecule, A . Metoda a fost dezvoltată de Oxford Molecular Ltd. În cadrul programului "*ANACONDA*". Și în acest

caz, metoda nu dă rezultate bune dacă formele structurilor moleculare sunt foarte disimilare. Principalele dificultăți sunt legate de modul în care este fixat centrul de rotație al moleculei B. Astfel, centrul de inerție este adecvat pentru o moleculă sferică dar nu și pentru o moleculă cu structură elipsoidală. În plus, dacă molecula este flexibilă, dificultățile cresc enorm deoarece centrul de inerție trebuie continuu deplasat.

În concluzie, cea mai bună metodă de suprapunere a două molecule disimilare constă în dezvoltarea unor metode invariante la rotație și translație. În acest sens, metodele care utilizează matricea distanțelor (topologice sau cu diferite metrice) sunt cele mai adecvate pentru a compara puncte distribuite în spațiu deoarece, în acest caz, nu este necesar ca structurile moleculare care se compară să fie centrate identic și, de asemenea, nu sunt necesare transformări ale punctelor în cursul procesului de superpoziție.

CAPITOLUL 2. MATRICI TOPOLOGICE

Un graf molecular poate fi reprezentat prin: un număr, o secvență de numere, o matrice sau un polinom¹¹. Aceste reprezentări se doresc a fi unice, pentru o structură dată. Randić consideră¹² că matricile topologice pot fi acceptate ca bază rațională pentru dezvoltarea de indici topologici, utili în studii corelaționale ori de similaritate.

Matricea de adiacență

În 1874 Sylvester³ a arătat că o moleculă organică adecvat numerotată, poate fi reprezentată printr-o matrice de adiacență, $\mathbf{A}(\mathbf{G})$. Aceasta este o tabelă pătratică de dimensiuni $\mathbf{N} \times \mathbf{N}$, ale cărei elemente $[\mathbf{A}]_{ij}$ se definesc astfel:

$$[\mathbf{A}]_{ij} = \begin{cases} 1 & \text{dacă } i \neq j \text{ și } (i, j) \in E(G) \\ 0 & \text{dacă } i = j \text{ sau } (i, j) \notin E(G) \end{cases} \quad (2)$$

iar matricea $\mathbf{A}(\mathbf{G})$:

$$\mathbf{A}(\mathbf{G}) = \{[\mathbf{A}]_{ij}; i, j \in \mathbf{V}(\mathbf{G})\} \quad (3)$$

$\mathbf{A}(\mathbf{G})$ caracterizează graful până la izomorfism, din ea putându-se reconstitui \mathbf{G} . Matricea $\mathbf{A}(\mathbf{G})$ este simetrică față de diagonala principală, astfel că transpusa ei, $\mathbf{A}^T(\mathbf{G})$, lasă matricea de adiacență neschimbată¹³:

$$\mathbf{A}^T(\mathbf{G}) = \mathbf{A}(\mathbf{G}) \quad (4)$$

Matricea de conectivitate

Matricea \mathbf{A} nu ia în considerare caracterul de multigraf (legătura multiplă). Pentru a indica tipul legăturii se utilizează matricea de conectivitate $\mathbf{C}(\mathbf{G})$, definită prin relațiile:

$$[\mathbf{C}]_{ij} = \begin{cases} b_{ij} & \text{dacă } i \neq j \text{ și } (i, j) \in E(G) \\ 0 & \text{dacă } i = j \text{ sau } (i, j) \notin E(G) \end{cases} \quad (5)$$

unde b_{ij} reprezintă ordinul convențional de legătură: 0; 1; 2; 3; 1.5 pentru nelegătură, legătură simplă, dublă, triplă și respectiv aromatică. Matricea $\mathbf{C}(\mathbf{G})$ va fi:

$$\mathbf{C}(\mathbf{G}) = \{[\mathbf{C}]_{ij}; i, j \in \mathbf{V}(\mathbf{G})\} \quad (6)$$

Matricea distanțelor

Matricea distanțelor, $\mathbf{D}(G)$, este o tabelă pătratică de dimensiuni $N \times N$, ale cărei elemente, $[D]_{ij}$, se definesc astfel:

$$[D]_{ij} = \begin{cases} \text{numarul de arce pe drumul} \\ \text{cel mai scurt } (i, j), \text{ dac} \acute{a} \ i \neq j \\ 0 \quad \text{dac} \acute{a} \ i = j \end{cases} \quad (7)$$

iar matricea $\mathbf{D}(G)$ va fi:

$$\mathbf{D}(G) = \{[D]_{ij}; i, j \in V(G)\} \quad (8)$$

CAPITOLUL 3. INDICI TOPOLOGICI

Un număr care reprezintă o structură chimică, în termeni graf-teoretici, se cheamă descriptor topologic. Fiind un invariant structural, el nu depinde de numerotarea atomilor ori de reprezentarea pictorială a grafului molecular. Cu toată pierderea considerabilă de informație prin "proiectarea" ca singur număr a structurii, astfel de invarianți și-au găsit largi aplicații în corelarea și prezicerea a numeroase proprietăți moleculare^{14,15} și de asemenea în testele de similaritate și izomorfism^{16,12}. Când un descriptor topologic se corelează cu o proprietate moleculară el poate fi numit indice molecular ori indice topologic (*IT*).

Indici bazați pe matricea de adiacență

Indicele RANDIC

Indicii care operează după catene (în particular arce) se denumesc indici de conectivitate. Indicele χ (de conectivitate, similar cu M_2) a fost introdus de Randic pentru caracterizarea legăturilor în graf¹⁹:

$$\chi = \sum_{(ij) \in E(G)} (k_i * k_j)^{-1/2} \quad (9)$$

Diudea⁶ a definit indicele χ per vârf:

$$\chi_i = \sum_{j:(ij) \in E(G)} (k_i * k_j)^{-1/2} \quad \text{si} \quad \chi = \frac{1}{2} \sum_i \chi_i \quad (10)$$

Indici bazați pe matricea distanțelor

Indicele WIENER

Wiener³⁰ a definit numărul W ca "suma distanțelor (ca număr de legături carbon-carbon) dintre oricare doi atomi de carbon în moleculă". În structuri aciclice, autorul a calculat W ca sumă a contribuțiilor "per legătură" ("bond contributions", care se corelează cu proprietățile termodinamice ale hidrocarburilor aciclice):

$$W = W(G) = \sum_e W_e = \sum_e N_{L,e} * N_{R,e} \quad (11)$$

unde:

$$N_{L,e} + N_{R,e} = N(G) \quad (12)$$

N_L , și N_R fiind numărul de vârfuri la stânga și la dreapta arcului e , însumarea făcându-se după toate vârfurile în G .

Indici de similaritate moleculară

Structura chimică a fiecărei molecule A este decodată într-un set de n descriptori structurali (SD) colectați într-un vector de tip $X=X(A)$,

$$X(A) = \{SD_1, SD_2, SD_3, \dots, SD_n\} \quad (13)$$

Pentru un set M de molecule, $M=\{A, B, C, \dots\}$ toți descriptorii structurali sunt colectați într-o matrice de dimensiuni $m \times n$ unde fiecare rând corespunde unei molecule și fiecare coloană corespunde unui descriptor structural particular. Pentru calcularea indicilor de similaritate, descriptorii structurali pot fi standardizați prin metoda Z_Score (autoscaling) care dă valorile ce au media zero și sunt scalate la varianță unitate.

Indicele de similaritate Cosine

Coefficientul Cosine, C_s pentru similaritatea între două molecule A și B este dat de relația:

$$C_s(A, B) = \frac{\sum_{i=1}^n X(A)_i X(B)_i}{\left[\sum_{i=1}^n X(A)_i^2 \sum_{i=1}^n X(B)_i^2 \right]^{1/2}} \quad (14)$$

cu proprietatea ca C_s este cuprins în intervalul $[-1, 1]$. Carbo¹⁷ a utilizat o formă a indicelui de similaritate Cosine definit pe integrala de densitate electronică pe tot spațiul.

Indicele de similaritate Dice

Coefficientul Dice, D_s pentru similaritatea dintre doi vectori a doi descriptori structurali $X(A)$ și $X(B)$ este dat de:

$$D_s(A, B) = \frac{2 \sum_{i=1}^n X(A)_i X(B)_i}{\sum_{i=1}^n X(A)_i^2 + \sum_{i=1}^n X(B)_i^2} \quad (15)$$

cu proprietatea că $-1 \leq D_s \leq 1$.

Indicele de similaritate Richards

Indicele de similaritate Richards este definit prin ecuația:

$$R_s(A, B) = \frac{1}{n} \sum_{i=1}^n \left(1 - \frac{|X(A)_i - X(B)_i|}{\max(|X(A)_i|, |X(B)_i|)} \right) \quad (16)$$

CAPITOLUL 4. RELAȚII CANTITATIVE STRUCTURĂ- ACTIVITATE BIOLOGICĂ

Metodele moderne utilizate în scopul "proiectării" unor structuri moleculare cu activitate biologică specificată, de exemplu medicamente, insecticide, ierbicide și fungicide, se bazează pe

cuantificarea bioactivității ca o funcție de structura moleculară¹⁸. Acest mod de abordare își are originea în lucrările lui Meyer și Overton¹⁹ de la sfârșitul secolului trecut și începutul secolului nostru. Astfel, ei au demonstrat cu succes pentru prima oară o dependență a bioactivității de un parametru fizico-chimic, coeficientul de partiție, care este o funcție de structura moleculară.

Introducerea conceptului de situs al receptorului biologic a constituit un element vital, de o importanță excepțională pentru acest domeniu al cercetării științifice; el a fost intuit de Langley,²⁰ dar fundamentat și dezvoltat de către Ehrlich,²¹ părintele chimioterapiei. Conform acestui model, activitatea biologică depinde de recunoașterea substratului bioactiv (efector) de către situsul receptorului; această etapă este urmată de legarea efectorului în situs.

Descoperirea dependenței dintre bioactivitate și configurație²² a condus la recunoașterea faptului că efectele sterice, indiferent care este tipul sau natura lor, joacă un rol esențial în cadrul interacțiunilor receptor-efector, condiționând și modelând potența biologică a efectorului.

Modele avansate în QSAR

Un compus bioactiv, introdus într-un organism viu, induce un răspuns biologic, o reacție specifică din partea organismului. Răspunsul este condiționat de structura și identitatea chimică a compusului bioactiv.

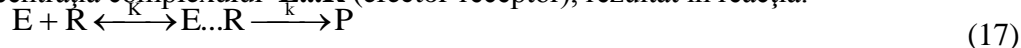
Interacția compusului bioactiv cu organismul se face la nivelul molecular, în așa-numiții receptori biologici. Aceștia sunt situs-uri active localizate în macromolecule proteice, în interiorul celulei vii sau pe membranele celulare.

Receptorii biologici^{23,24,25,26} au următoarele caracteristici :

- (i) specificitate : receptorii recunosc enantiomerii activi (“eutomeri” spre deosebire de “distomeri” care nu sunt activi biologic) sau diastereoizomerii agoniștilor sau antagoniștilor;
- (ii) saturabilitate : numărul situs-urilor active într-o formațiune celulară este finit;
- (iii) în general, receptorii se află în celula care generează răspunsul biologic.

E. FISCHER (1894) a formulat prima teorie a interacției compusului bioactiv (efector) cu receptorul. Conform acesteia, receptorul este privit ca o cavitate rigidă, în care efectorul trebuie să se “potrivească” asemănător “cheii în broască” (“kee in lock”). Ulterior s-a admis că receptorii sunt semi-rigizi, ei căutând optimizarea mutuală (deși limitată) cu efectorii săi. Ambii parteneri ai complexului efector-receptor se abat de la starea conformațională de minim energetic, pentru a realiza complexul cu cea mai mare stabilitate (ex. care induce răspunsul biologic). Despre alosterism și alte aspecte legate de interacțiunea efector-receptor vezi^{148,149,150,151}.

Răspunsul biologic provocat de un compus chimic pătruns în organism este proporțional cu concentrația complexului **E...R** (efector-receptor), rezultat în reacția:



Cât timp complexul **E...R** există, se manifestă o reacție specifică a organismului, numită răspuns biologic. Complexul poate disocia în componente (echilibrul caracterizat prin constanta **K**) sau poate forma (cu viteza caracterizată prin constanta **k**) produsul **P**. Concentrația acestuia variază în timp cumform relației:

$$\frac{d[P]}{dt} = k[E...R] \quad (18)$$

concentrația complexului putând fi aproximată din relația de echilibru:

$$[E...R] = [E][R]K = [E][R] \exp(-\Delta G / RT) \quad (19)$$

CAPITOLUL 5. ANALIZA DE DATE

Coeficientul de corelație. Coeficientul de corelație teoretic ρ_{xy} a două variabile aleatoare x și y este covarianța variabilelor normate corespunzătoare:

$$\rho_{xy} = \text{cov}\left(\frac{x-\bar{x}}{\sigma_x}, \frac{y-\bar{y}}{\sigma_y}\right) = \text{Med}\left[\left(\frac{x-\bar{x}}{\sigma_x}\right)\left(\frac{y-\bar{y}}{\sigma_y}\right)\right] = \frac{\text{Med}[(x-\bar{x})(y-\bar{y})]}{\sigma_x\sigma_y} = \frac{\text{cov}(x, y)}{\sigma_x\sigma_y} \quad (20)$$

Valorile lui fiind cuprinse în domeniul: $-1 \leq \rho_{xy} \leq 1$. În cazul unei selecții, relația (20) se va scrie:

$$\rho_{xy} = \frac{\text{cov}(x, y)}{S_x S_y} = r_{xy} \quad (21)$$

r_{xy} fiind coeficientul de corelație empiric. Când variabilele x și y sunt independente, $\text{cov}(x, y) = 0$ și $r_{xy} = 0$. *Reciproca nu este adevărată.* Subliniem că ρ (coeficientul de corelație teoretic) se referă la întreaga populație iar r (coeficientul de corelație empiric) se referă la o selecție. În cazul în care între x_i și y_i este definită o *funcție de regresie* $y(x_i) = f(x_i)$, se calculează un *coeficient de corelație (regresie) parțial*, r_{xy} , al variabilei x_i .

Coeficientul de corelație (regresie) global, R^2 , este dat de raportul dispersiei valorilor calculate de la medie cu dispersia valorilor empirice de la medie:

$$R^2 = 1 - \frac{\sum_i (y_i - y_{i,\text{calc}})^2}{\sum_i (y_i - \bar{y})^2} = \frac{\sum_i (y_{i,\text{calc}} - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} \quad (22)$$

R^2 ia valori în domeniul $[0, 1]$. Valori apropiate de 1 arată că dependența liniară este potrivită pentru descrierea relației dintre y și x .

Regresii liniare

Fie setul de structuri chimice C_1, C_2, \dots, C_n și valorile observate ale unei proprietăți moleculare y_1, y_2, \dots, y_n (y_i corespunde compusului C_i). Estimarea valorilor y_i cu ajutorul variabilelor independente x_{ij} (numite și variabile predictor sau explicative) se poate face conform relației:

$$y_{i,\text{calc}} = b_0 + \sum_{j=1}^m b_j x_{ij} \quad (23)$$

Variabilele x_{ij} codifică numeric caracteristica structurală (topologică) sau fizico-chimică j - prezentă în compusul C_i .

Statistica Fischer F, indică nivelul de semnificație al ecuației de regresie. Estimatorul F se calculează cu relația:

$$F = \frac{SS_{\text{reg}}}{SS_e / (n - k)} \quad (24)$$

unde SS_{reg} este suma pătratelor erorilor atribuite regresiei:

$$SS_{\text{reg}} = \sum_i (y_{i,\text{calc}} - \bar{y})^2 \quad (25)$$

Statistica t - Student. Estimatorul t indică nivelul de semnificație al coeficienților b_j ; el se calculează cu formula:

$$t_j = \frac{|b_j|}{\sigma_{b_j}} \quad (26)$$

unde σ_{b_j} este eroarea standard a coeficientului de regresie b_j . Din compararea valorii t_j calculate cu valori tabelate (pentru un prag de semnificație impus, care este funcție de gradele de libertate ale regresiei - vezi și estimatorul F) se validează sau nu aportul unei variabile x_{ij} la corelația globală.^{27,28}

Analiza clusterilor

Obiectivele analizelor de clusteri

Noțiunea **analiză de clusteri (CA)** se încadrează într-o familie de metode care este folosită în principal pentru găsirea și scoaterea în evidență a structurilor din interiorul datelor. Din acest punct de vedere este mai degrabă o optimizare a datelor decât o manipulare a lor. În această direcție analiza clusterilor poate fi văzută ca și o metodă numită **model de cunoaștere**. Este folosită pentru aceasta o denumire sinonimă cu numere taxonomice și clasificare automată. Dacă scopul analizelor de date sau întrebările ce se pun despre ele sunt fixate clar și duc la obținerea unor rezultate, următoarea întrebare va fi: care este cauza structurii găsite?

Să presupunem că avem un număr rezonabil de obiecte, n , și proprietățile lor, ce au fost selectate și aranjate în matricea X ca și în exemplul de mai jos:

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix} \quad (27)$$

și să începem extragerea de informații din această matrice.

Unități asemănătoare și preprocesarea informațiilor

În ordinea de găsire a structurilor în grupul de date sau descoperirea similarităților probelor, organismelor,...care mai jos vor fi denumite **obiecte**, mai înainte de toate avem nevoie de **unități asemănătoare**. Cele mai simple unități asemănătoare pot deriva din geometrie. Fără a demonstra că conceptele intuitive de similaritate și distanță sunt complementare în natură și amintindu-ne legea lui PITAGORA, distanța d dintre două puncte O_1 și O_2 într-un sistem rectangular cu două axe x și y este:

$$d(O_1, O_2) = \sqrt{(y_1 - y_2)^2 + (x_1 - x_2)^2} \quad (28)$$

Aceasta este prezentată în figura 5.1.

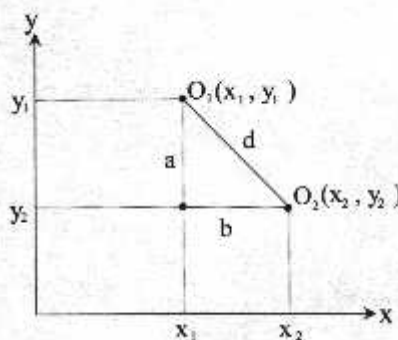


Figura 1. Distanța dintre două obiecte în spațiul plan (legea lui Pitagora)

Întinderea acestei legi peste mai mult de două dimensiuni a “spațiului PITAGORA” conduce la distanța euclidiană a oricăror două obiecte O_i și O_k care mai jos vor fi simplu scrise ca și $d(i,k)$:

$$d(i,k) = \sqrt{\sum_{j=1}^m (x_{ij} - x_{kj})^2} \quad (29)$$

unde: m -numărul trăsăturilor.

Analiza factorială

Analiza factorială a apărut ca o metodă de reducere a efectivului de variabile propuse pentru descrierea unui domeniu, prin construcția unor noi variabile (numite *factori*), în număr mult mai mic, și prin determinarea unor relații matematice care să precizeze legătura dintre variabilele inițiale și factori, astfel încât prin aceste noi variabile să se reproducă, în cea mai mare măsură, informația dată de variabilele inițiale. Metoda datează încă de la începutul secolului trecut și este datorată lui K. Pearson (1901) și C. Spearman (1904).

Tehnicile statistice cunoscute sub această denumire au ca obiectiv comun reducerea numărului de variabile ce caracterizează o mulțime de obiecte la un număr mai mic de variabile, de obicei diferite de cele inițiale.

Ecuția de bază a analizei factoriale

Din matricea de variabilelor standardizate poate fi dedusă matricea corelațiilor dintre variabilele inițiale:

$$R = \frac{1}{n} Z \cdot Z^T, \quad r_{ij} = \frac{1}{n} (z_{i1}z_{j1} + z_{i2}z_{j2} + \dots + z_{in}z_{jn}) \quad (30)$$

Numărul factorilor de extras

O modalitate empirică a modului suficient de descompunere cu k factori a fost definită de Malinowki (1977) prin funcția IND. Valoarea minimă IND indică numărul probabil al factorilor relevanți.

$$IND = \frac{RE}{(m-k)^2} \quad (31)$$

$$RE^2 = \frac{1}{n(m-k)} \sum_{j=k+1}^m \lambda_j \quad (32)$$

unde m – numărul variabilelor
 n – numărul obiectelor
 k – numărul extras de factori
 RE – „eroarea reală”

CONTRIBUTII PERSONALE

Tehnicile relațiilor cantitative structură-activitate (QSAR) devin indispensabile în toate aspectele cercetării privind interpretarea moleculară a proprietăților biologice.²⁹ Este de înțeles ca proprietățile fizice, chimice sau biologice ale compușilor depind de aranjamentul 3D (tri-dimensional) al atomilor din moleculă. Abilitatea de a produce corelații cantitative între structura moleculară 3D și activitatea biologică este importantă în alegerea căilor de sinteza a substanțelor biologice active.³⁰

Activitatea biologică a steroizilor variază considerabil cu modificarea foarte ușoară a structurii. Aceste familii importante de molecule prezintă caracteristici foarte schimbătoare

pentru orice metodă de predicție, datorită flexibilității relative scăzute a scheletului de cyclopentanoperhidrofenantrenă. Din acest motiv, foarte multe modele QSAR bazate pe proprietățile 2D, cum ar fi descriptorii topologici, au o calitate comparabilă cu modelele provenite de la metodele complexe 3D.^{31,32}

În acesta parte sunt identificate și prezentate aspectele legate de structura moleculară care sunt relevante în particular pentru activitatea biologică (afinitatea de legare cu receptorul) pentru diferite clase de substanțe cu activitate biologică.

Compuși steroidici

Setul de 31 de structuri steroidice (androstan) AS (Figura 2, Tabel 1), cu afinitatea de legare a corticosteroidilor de globulină (CBG) a fost luat din publicațiile lui Dunn *et al.*³³ și Tuppurainen, *et al.*³⁴ Acest set de structuri a fost de multe ori utilizat pentru evaluarea performanțelor noilor metode de analiză QSAR. Oricum, menționând alți autori putem spune că foarte multe publicații ce au utilizat acest set au inclus erori în structura steroidilor. Structurile utilizate în acest studiu au fost verificate foarte atent cu scopul de a evita o viitoare propagare de erori. Calitativ, moleculele cu substituenți cum ar fi oxigen sau hidroxil în poziția 17 pe structura steroidică, au activitate CBG crescută, în timp ce prezența unei catene voluminoase cum ar fi $-COCH_2OH$, conduce la scăderea acestei activități.

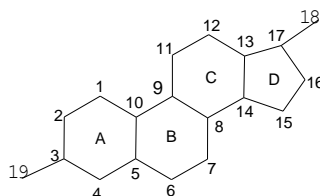


Figura 2. Structura Androstanului

Tabelul 1. Setul de structuri AS:

Compuși	Activitate	Compuși	Activitate
1 aldosterone	6.279	17 pregnenolone	5.255
2 androstanediol	5.000	18 17-hydroxypregnenolone	5.000
3 androstenediol	5.000	19 progesterone	7.380
4 androstenedione	5.763	20 17-hydroxyprogesterone	7.740
5 androsterone	5.613	21 testosterone	6.724
6 corticosterone	7.881	22 prednisolone	7.512
7 cortisol	7.881	23 cortisol 21-acetate	7.553
8 cortisone	6.892	24 4-pregnene-3,11,20-trione	6.779
9 dehydroepiandrosterone	5.000	25 epicorticosterone	7.200
10 deoxycorticosterone	7.653	26 19-nortestosterone	6.144
11 deoxycortisol	7.881	27 16R,17-dihydroxy-4-pregnene-3,20-dione	6.247
12 dihydrotestosterone	5.919	28 16-methyl-4-pregnene-3,20-dione	7.120
13 estradiol	5.000	29 19-norprogesterone	6.817
14 estriol	5.000	30 11β,17,21-trihydroxy-2R-methyl-4-pregnene-3,20-dione	7.688
15 estrone	5.000	31 11β,17,21-trihydroxy-2R-methyl-9R-fluoro-4-pregnene-3,20-dione	5.797
16 etiocholanolone	5.255		

Subsetul parametrilor electronici include descriptorii moleculari derivați din sarcinile atomice parțiale. Cu programul TOPOCLUJ, sarcinile parțiale Ch_i au fost calculate astfel:

$$Ch_{i,j} = \log(S_j / S_i)^{1/(d_{i,j})^2} \quad (33)$$

$$Ch_i = \sum_j ch_{i,j} \quad (34)$$

În ambele cazuri S_i , S_j reprezintă electronegativitățile Sanderson de grup calculate pentru grupuri hibride (de exemplu pentru atomi grei care sunt înconjurați de atomi de hidrogen) în molecule, în timp ce d_{ij} este distanța Euclidiană ce separă atomii i și j din structura chimică optimizată cu energie minimă (HyperChem).¹³⁷ $Ch_{i,j}$ este perturbația electronegativității atomului i produsă de orice atom j din moleculă în timp ce Ch_i este rezultatul acestor perturbații pe atomul i . Pentru alte calcule topologice ale sarcinii parțiale se pot consulta referințele.^{35,36}

Orice compus steroidic poate fi descris prin aceste sarcini parțiale care caracterizează pozițiile substituie sau nesubstituie și heteroatomii. Bazându-ne pe acest lucru am definit un nou descriptor global flexibil (**CD**) care poate fi definit ca o funcție aditivă de autocorelare ponderată cu sarcinile parțiale corespunzătoare ale atomului j considerat:

$$CD = \sum_j c_j \cdot Ch_j \quad (35)$$

unde c_j este coeficientul de regresie dat de regresia multivariată $\log(A_i \text{obs}) = f(Ch_j)$. Această “ad-hoc” ponderare depinde de setul de molecule luate în considerare și de asemenea de descriptorii locali utilizați. Sarcinile parțiale (Ch_j) corespund următoarelor poziții de pe structura de bază: 3, 10, 11, 13, 17, 18, 19 (Figura 2).

Softul *Dragon 2.1*¹³⁶ a fost folosit pentru a calcula 1600 de descriptori moleculari pentru compușii studiați. Cei mai relevanți dintre acești descriptori utilizați aici sunt cei cu funcții de distribuție radială (RDF), indici de autocorelare și descriptori geometrici.

Descriptorii aparținând clasei cu funcție de distribuție radială sunt bazați pe distribuția distanței în reprezentarea geometrică a moleculei. În plus față de distanțele interatomice în întreaga moleculă, RDS aduc informații despre distanțele legăturii, tipul inelului, sisteme planare sau neplanare, tipuri de atomi și alte modificări structurale importante. Prin utilizarea diferitelor scheme de ponderare, care includ tipuri de atomi, electronegativități, masa atomică (*RDF090m*) sau distanțe van der Walls, RDF poate fi ajustat să dea evoluții importante ale descriptorilor în studiile QSAR.

Al doilea grup de descriptori este obținut aplicând funcții de autocorelare bidimensionale pe graful molecular. Astfel de descriptori exprimă corelația dintre valorile numerice, care pot fi ponderate statistic utilizând proprietăți atomice, la intervale egale de valoare dată.³⁷ De exemplu, *MATS1p*- auto-corelație Moran-lag 1/ponderată cu polarizabilitățile atomice; *MATS4e*- auto-corelație Moran-lag 4/ponderată cu electronegativitățile Sanderson atomice. Aplicațiile electronegativităților Sanderson ca și coeficienți de ponderare, luate în acest context în unele cazuri schimbă distribuția în interiorul moleculei.

Descriptorii geometrici indică mărimea moleculei, ei sunt derivați din coordonatele tridimensionale ale nucleelor atomice și a maselor atomice și / sau distanța atomică din moleculă.

Prelucrarea datelor

Ținând cont de complexitatea interacțiilor dintre molecula receptor și moleculele cu potențial inhibitor, este destul de dificil să modelezi setul de structuri folosind doar modele simple de regresie liniară. În vederea analizei datelor obținute s-a propus următorul model:

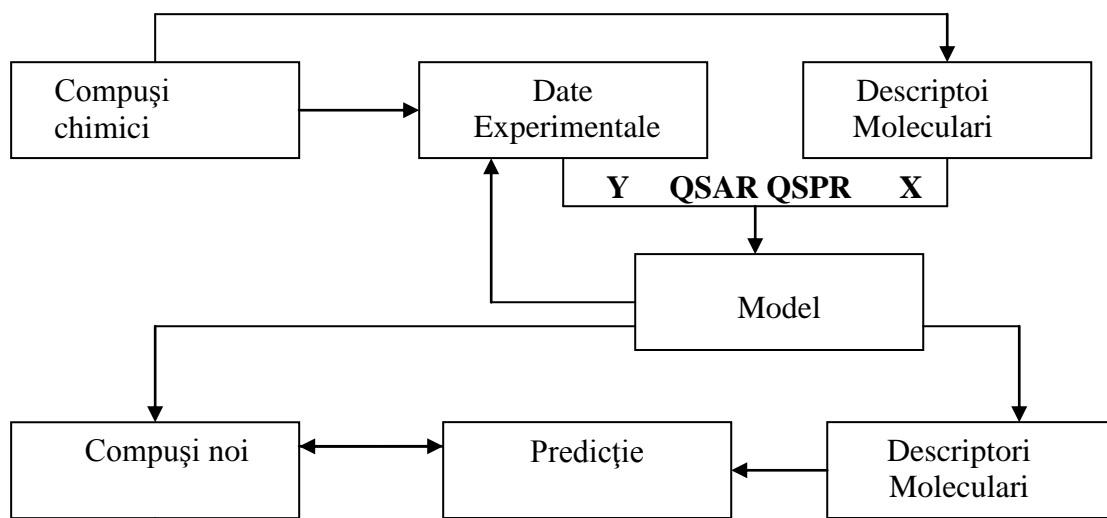


Figura 3. Reprezentarea schematică a construcției modelului.

Analiza QSAR constă în următorii pași:

- (i) optimizarea structurii folosind metoda semiempirică PM3;
- (ii) calcularea descriptorilor moleculari;
- (iii) regresie multivariabilă pentru găsirea coeficienților de autocorelare;
- (iv) împărțirea setului de date în unul școală (calibrarea regresiei) și unul de predicție (pentru validarea modelului);
- (v) testarea capacității de a prezice modelul;
- (vi) analiza componentelor principale (PCA);
- (vii) găsirea unei funcții de regresie pentru model;
- (viii) testarea capacității predictive a acestui model;
- (ix) interpretarea modelului.

În ambele regresii monovariate și bivariate (tabelul 2) compusul 13 apare ca fiind *outlier*. Acest compus nu a fost inclus în analizele viitoare. Cu acest *outlier* exclus, am observat o îmbunătățire a corelației cu descriptorii prezentați mai jos.

Tabelul 2. Modelele utilizate și rezultatele procedurii de *cross-validare*

Setul de date	Numărul observațiilor (n)	Modelul	Coeficientul de corelație (r^2) (înainte de LOO)	Structura Outlier	Coeficientul de corelație (r^2) (după LOO)
AS	31	$\log(P_{i\text{calc}}) = f(CD)$	0.891	13	0.920
		$\log(P_{i\text{calc}}) = f(CD, L/B_w)$	0.931	13	0.939

Regresia trivariabilă a dat rezultate asemănătoare dar nu a adus îmbunătățiri esențiale. Rezultatele obținute vor fi prezentate mai departe.

Setul AS de structuri

Am împărțit setul de structuri AS în doua seturi: setul școală ($n = 20$) și setul de predicție (validare) ($n = 11$) așa cum este arătat în Tabelul 3.

(a) Setul Școală ($n = 20$)

Descriptorul electronic CD , actualul CDP , a fost calculat *de novo* pe setul școală ținând cont de ecuația 6.7, aceasta datorită greutatea de corelare ale lui c_j potrivit doar cu

proprietatea selectată și cu setul dat (în acest caz, setul școală având 20 de structuri). În Tabelul 3 sunt prezentați cei mai relevanți descriptori pentru setul de structuri date. Cele mai bune modele pentru setul de structuri AS sunt prezentate mai jos.

Regresie monovariată

$$\log P_{i\text{calc}} = 7.236 + 1.033 \cdot \text{CDP}_i \quad (36)$$

$$n = 20 \quad R^2 = 0.903 \quad s = 0.10 \quad F = 159.99$$

Regresie bivariată

$$\log P_{i\text{calc}} = 5.737 + 0.914 \cdot \text{CDP}_i + 0.194 \cdot \text{L/Bw}_i \quad (37)$$

$$n = 20 \quad R^2 = 0.931 \quad s = 0.31 \quad F = 114.12$$

Regresie multiplă

$$\log P_{i\text{calc}} = 6.268 + 0.17 \cdot \text{L/Bw}_i - 0.166 \cdot \text{RDF090m}_i + 0.904 \cdot \text{CDP}_i \quad (38)$$

$$n = 20 \quad R^2 = 0.962 \quad s = 0.24 \quad F = 129.73$$

Tabelul 3. Descriptorii topologici și coeficientul de partiție observat $\log P$ pentru setul de structuri AS

Structura	L/Bw	RDF090m	CDP	$\log P$ obs.
<i>Setul Școala</i>				
1	7.3	3.17	-1.722	5
2	6.5	3.127	-2.101	5
4	6.5	3.505	-2.062	5
5	6.9	1.755	-2.456	5
7	6.3	2.803	-2.188	5
8	4.3	0.191	-1.950	5.255
9	7.1	1.91	-2.087	5.255
12	7	4.01	-0.880	5.797
13	6.1	3.497	-0.187	5.919
14	7.2	0.61	-1.212	6.144
17	6.9	2.729	-0.336	6.724
18	6	3.288	0.492	6.779
20	7.1	3.976	-0.243	6.892
22	6.6	2.525	-0.046	7.2
23	7.1	0.942	0.085	7.38
25	9.5	1.704	0.205	7.553
26	9.1	0.673	0.006	7.653
28	7.8	1.122	0.317	7.74
29	8.7	2.31	0.040	7.881
30	9	1.077	0.309	7.881
<i>Setul de predicție (validare)</i>				
3	6.7	2.669	-2.350	5
6	6.6	0.811	-2.117	5
10	6.3	1.338	-2.184	5.613
11	6.2	1.272	-0.479	5.763
15	6.9	1.986	0.400	6.247
16	6.4	1.759	-0.716	6.279
19	9.1	0.406	-0.810	6.817
21	6.9	2.014	0.046	7.12
24	5.9	2.777	0.110	7.512
27	6.9	3.512	0.675	7.688
31	8.7	1.423	0.225	7.881

(b) Setul de predicție (validare) ($n = 11$)

Orice model QSAR trebuie să fie validat cu un set de predicție extern. Calcularea descriptorului CDP în setul de predicție s-a realizat pe baza sarcinilor parțiale și a parametrilor c_j generați pentru setul școală (vezi tabelul 4.). Noi presupunem că activitatea biologică a setului de predicție este necunoscută.

Tabelul 4. Setul de predicție pentru structurile AS ($n = 11$)

Structuri	$\log P_{i,obs}$	$\log P_{i,calc}$ (ecuația 6)	$\log P_{i,calc}$ (ecuația 7)	$\log P_{i,calc}$ (ecuația 8)
3	5.000	4.808	4.890	4.836
6	5.000	5.049	5.083	5.339
10	5.613	4.980	4.964	5.140
11	5.763	6.741	6.504	6.676
15	6.247	7.650	7.444	7.472
16	6.279	6.496	6.326	6.415
19	6.817	6.400	6.765	7.015
21	7.120	7.284	7.121	7.147
24	7.512	7.350	6.984	6.907
27	7.688	7.934	7.696	7.466
31	7.881	7.468	7.633	7.713
R²		0.716	0.766	0.728
CV%		7.089	5.010	6.586

Pentru setul de structuri AS abilitatea de prezicere pare a fi mult mai bună în cazul regresiei bivariate (tabelul 4).

Figura 4. a-c prezintă valorile experimentale vs valorile calculate pentru afinitatea de legare a receptorului pentru structurile din AS : (a) valorile calculate conform ecuației 36; (b) valorile calculate conform ecuației 37, (c) valorile calculate conform ecuației 38.

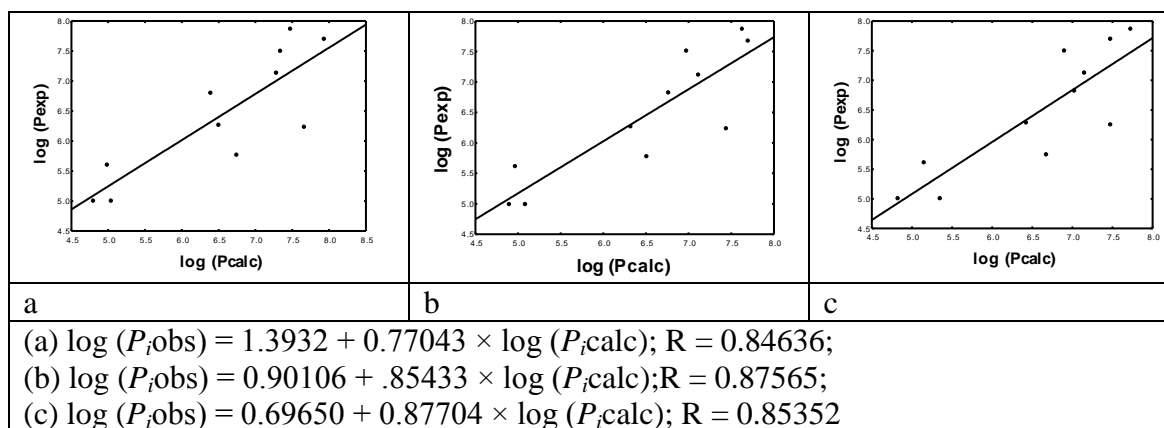


Figura 4. Reprezentarea grafică a valorilor experimentale versus valorile calculate ale afinității de legare a receptorului pentru setul de structuri AS

În scopul explicării contribuției fiecărui substituent al grafului din setul de structuri AS la afinitatea de legare a receptorului (CBG), am generat un descriptor electronic simplu în funcție de sarcinile atomice parțiale, corelat parțial cu proprietatea studiată.

Modelele QSAR descrise în acest studiu indică faptul că acest descriptor global este unul dintre cei mai semnificativi în predicția activității compușilor noștri. El poate indica pozițiile celor mai importanți substituenți. Astfel, CD calculate pentru atomii din pozițiile mai sus

menționate, fără substituenți în poziția 17 pe structura steroidică, duc la o varianța de 25% a activității CBG în acest set de structuri **AS** în timp ce incluzând poziția 17 crește la 89%. Modelul a fost validat cu un set extern de predicție. Modelul derivat pentru activitatea moleculară utilizând descriptorul **CD** și descriptorii obținuți prin *factor loading* al **PCA**-ului, este comparabil cu celelalte modele descrise în literatura de specialitate, având o abilitate bună de predicție. De notat faptul că modelul simplu 2D, la fel ca acesta pe care noi l-am dezvoltat, este comparabil cu rezultatele obținute cu ajutorul modelelor complexe 3D (CoMFA, COMSA, GRIND, EEVA etc)^{27,38,39}, care cer mult mai multe resurse computaționale.

Compuși cu acțiune antibacteriană

Setul de 38 de derivați de 2-furiletilenă, cu activitate antibacteriană a fost luat din publicațiile lui Miguel Angel Cabrera Pérez⁴⁰ precum și Yovani Marrero Ponce⁴¹. Derivații de 2-furiletilenă sunt substanțe biologic active cu spectru larg antimicrobial, antispasmodic, citotoxic dar în anumite cazuri având activități carcinogenice și mutagenetice (Yahagi *et al*, 1974; Dore și Viel, 1975⁴²; Miyaji, 1976; McCalla, 1979; McCalla, 1983; Kelloval *et al*, 1984; Estrada, 1998). Acest interes în studierea derivaților de 2-furiletilenă a crescut mult în ultimii ani ca și consecință a descoperirilor a unor noi compuși cu potențial microcidal cu această structură chimică (McCoy and Thornburgh, 1992; Castañedo *et al*, 1994⁴³; Blondeau *et al*, 1999⁴⁴). Modelul utilizat în acest caz este prezentat schematic în figura 5.

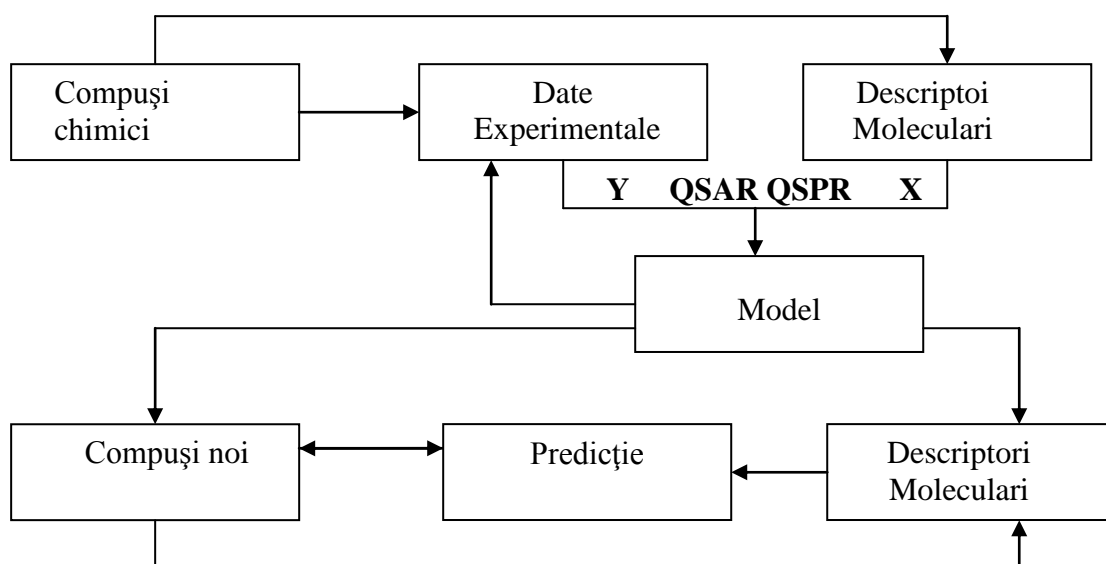


Figura 5. Reprezentarea schematică a construcției modelului.

Softul *Dragon 2.1* a fost folosit pentru a calcula pentru compușii studiați a 1600 de descriptori moleculari. Ce mai relevanți dintre acești descriptori utilizați aici sunt cei constituționali (*MW* – masă moleculară), de conectivitate (*X4v* – indicele de conectivitate al valenței *chi-4*) și geometrici (*G(N..O)* – suma distanțelor geometrice între N și O).

Descriptorii topologici calculați de programul TOPOCLUJ sunt derivați din matricele topologice sau polinoamele care descriu grafurile moleculare. Algoritmii de calcul includ câteva scheme de ponderare a grafurilor moleculare cu sarcini atomice parțiale, electronegativități și mase ale fragmentelor moleculare. Matricele topologice calculate includ pe cele de bază: adiacență, distanță, detour, conectivitate, Wiener, etc. precum și matrici dezvoltate de grupul TOPO Cluj care includ: Cluj și Cluj Fragmental, matricea W, LM, SM, operatorul matricial $W_{(M1,M2,M3)}$. Cei mai relevanți dintre acești descriptori sunt :

- IE[CfMax[Density]],
- VAAI,
- VADI,
- CS[LM[Electronegativity]],
- CS[Sh[W4[Charge_Adjacency]]].

Tabelul 5. Setul de structuri a derivaților de 2-furiletlenă.

	R₁	R₂	R₃	log P
1	H	NO ₂	COOCH ₃	1.879
2	CH ₃	NO ₂	COOCH ₃	2.439
3	Br	NO ₂	COOCH ₃	2.739
4	COOCH ₃	NO ₂	COOCH ₃	1.869
5	NO ₂	NO ₂	COOCH ₃	1.599
6	NO ₂	COOC ₂ H ₅	COOC ₂ H ₅	2.504
7	NO ₂	H	NO ₂	1.303
8	H	H	NO ₂	1.583
9	NO ₂	H	CONHC ₂ H ₅	1.386
10	NO ₂	H	CONH(CH ₂) ₂ CH ₃	1.86
11	NO ₂	H	CONHCH(CH ₃) ₂	1.803
12	NO ₂	H	CONH(CH ₂) ₃ CH ₃	2.356
13	NO ₂	H	CONHCH ₂ CH(CH ₃) ₂	2.225
14	NO ₂	H	CONHCH(CH ₃)C ₂ H ₅	2.284
15	NO ₂	H	CONHC(CH ₃) ₃	2.333
16	NO ₂	H	CONHCH ₂ C(CH ₃) ₃	2.605
17	NO ₂	H	COOCH ₃	1.652
18	NO ₂	H	COOC ₂ H ₅	2.098
19	NO ₂	H	COO(CH ₂) ₂ CH ₃	2.673
20	NO ₂	H	COOCH(CH ₃) ₂	2.641
21	NO ₂	H	COO(CH ₂) ₃ CH ₃	2.827
22	NO ₂	H	COOCH ₂ CH(CH ₃) ₂	3.135
23	NO ₂	H	COOCH(CH ₃)C ₂ H ₅	3.091
24	NO ₂	H	COOC(CH ₃) ₃	3.06
25	NO ₂	H	COO(CH ₂) ₄ CH ₃	3.404
26	NO ₂	H	Br	2.447
27	NO ₂	H	CN	1.05
28	NO ₂	H	OCH ₃	1.591
29	NO ₂	H	H	1.611
30	NO ₂	CN	COOCH ₃	1.488
31	I	NO ₂	COOCH ₃	2.999
32	NO ₂	H	CONH ₂	0.649
33	NO ₂	H	CONHCH ₃	0.984
34	NO ₂	H	CON(CH ₃) ₂	0.819
35	Br	NO ₂	Br	2.820
36	Br	NO ₂	CH ₃	2.730
37	H	NO ₂	H	1.290
38	H	NO ₂	CH ₃	1.940

Analiza și prelucrarea datelor

Coeficientul de partiție *n*-octanol/apă (log P) joacă un rol important în înțelegerea comportamentului biologic al acestor derivați de 2-furiletlenă.

Ținând cont de complexitatea interacțiilor dintre molecula receptor și moleculele cu potențial inhibitor, este destul de dificil să modelezi setul de structuri folosind doar simple modele de regresie liniară.

În pachetul software SIMIL sunt implementate proceduri de filtrare folosind secvențele valențelor vârfurilor utilizate la clusterarea sistemelor informatice chimice.

Analiza QSAR constă în următorii pași:

- (x) optimizarea structurii folosind metoda semiempirică PM3;
- (xi) calcularea descriptorilor moleculari;
- (xii) împărțirea setului de date în unul școală (calibrarea regresiei) și unul de predicție (pentru validarea modelului) pe baza similarității moleculare;
- (xiii) analiza componentelor principale (PCA);
- (xiv) testarea capacității de a prezice modelul;
- (xv) găsirea unei funcții de regresie pentru model;
- (xvi) testarea capacității predictive a acestui model;
- (xvii) interpretarea modelului.

Tabelul 6. Similaritatea pentru setul de derivați de 2-furiletlenă față de structurile lider alese în setul de predicție și coeficientul lor de partiție *n*-octanol/apă (log P).

Nr. structurii	35	36	37	38	log P
1	0.59524	0.72024	0.71429	0.78571	1.879
2	0.55556	0.67222	0.66667	0.73333	2.439
3	0.67222	0.8	0.66667	0.73333	2.739
4	0.462963	0.56019	0.55556	0.61111	1.869
5	0.490196	0.59314	0.58824	0.64706	1.599
6	0.459375	0.48151	0.45125	0.50114	2.504
7	0.64103	0.64103	0.76923	0.6993	1.303
8	0.83333	0.83333	1	0.90909	1.583
9	0.50139	0.6125	0.60167	0.66818	1.386
10	0.47005	0.57422	0.56406	0.62642	1.86
11	0.47005	0.57422	0.56406	0.62642	1.803
12	0.490196	0.54044	0.53088	0.58957	2.356
13	0.490196	0.54044	0.53088	0.58957	2.225
14	0.490196	0.54044	0.53088	0.58957	2.284
15	0.490196	0.54044	0.53088	0.58957	2.333
16	0.462963	0.51042	0.50139	0.55682	2.605
17	0.5372	0.65625	0.64464	0.71591	1.652
18	0.50139	0.6125	0.60167	0.66818	2.098
19	0.47005	0.57422	0.56406	0.62642	2.673
20	0.47005	0.57422	0.56406	0.62642	2.641
21	0.490196	0.54044	0.53088	0.58957	2.827
22	0.490196	0.54044	0.53088	0.58957	3.135
23	0.490196	0.54044	0.53088	0.58957	3.091
24	0.490196	0.54044	0.53088	0.58957	3.06

25	0.462963	0.51042	0.50139	0.55682	3.404
26	0.83523	0.68371	0.82045	0.74587	2.447
27	0.62674	0.76563	0.75208	0.83523	1.05
28	0.62674	0.62674	0.75208	0.68371	1.591
29	0.75208	0.75208	0.9025	0.82045	1.611
30	0.47005	0.57422	0.56406	0.62642	1.488
31	0.55556	0.67222	0.66667	0.73333	2.999
32	0.57853	0.70673	0.69423	0.77098	0.649
33	0.5372	0.65625	0.64464	0.71591	0.984
34	0.50139	0.6125	0.60167	0.66818	0.819
35	1	0.84028	0.83333	0.75758	2.820
36	0.84028	1	0.83333	0.91667	2.730
37	0.83333	0.83333	1	0.90909	1.290
38	0.75758	0.91667	0.90909	1	1.940

Deoarece coeficientul de corelare este supus fluctuațiilor de selecție, o valoare r mare trebuie privită cu circumspecție dacă numărul observațiilor este mic și, în plus, el nu poate fi utilizat ca termen de comparație pentru ecuații cu număr de date diferit.

Din tabelul 6.6 se observă similaritatea amprentelor moleculare pentru moleculele alese în setul de predicție față de celelalte structuri din setul de bază. Putem observa că, cele patru structuri pot face parte dintr-un cluster deoarece similaritatea între ele este foarte ridicată. De asemenea am ales în setul școală acele structuri care au coeficientul de similaritate mai mare de 0.70 pentru cel puțin una dintre aceste structuri.

Setul școală și setul de predicție

Din setul școală fac parte structurile prezentate în tabelul 7, structuri cu coeficienți de similaritate ridicați față de structurile lider (tabelul 8) din setul de predicție, fiind prezentați aici și cei mai semnificativi descriptori calculați cu programul Dragon alături de coeficientul de partiție n-octanol/apă $\log P$.

Tabelul 7. Descriptorii topologici și coeficientul de partiție observat $\log P$ pentru setul școală.

Nr. Structură	MW	X4v	G(N..O)	$\log P$
1	197.16	0.799	9.619	1.879
2	211.19	0.956	9.613	2.439
3	276.05	1.15	9.626	2.739
8	139.12	0.52	4.386	1.583
17	197.16	0.748	14.97	1.652
26	218.01	0.853	0	2.447
27	164.13	0.649	22.162	1.05
28	169.15	0.665	6.598	1.591
29	139.12	0.55	0	1.611
31	323.05	1.265	9.641	2.999
32	182.15	0.694	29.088	0.649
33	196.18	0.771	29.063	0.984

Tabelul 8. Descriptorii topologici și coeficientul de partiție observat logP pentru setul de predicție.

Nr. Structură	MW	X4v	G(N..O)	logP
35	296.9	1.121	4.405	2.49
36	232.04	0.984	4.516	2.37
37	139.12	0.52	4.386	1.56
38	153.15	0.632	4.517	1.92

Ecuțiile utilizate pentru predicție sunt cele calibrate în clusterul școală (în forme lor normalizată de Matlab). Aceste ecuații prezintă un coeficient de corelație bun ($R^2 = 0.9843$ în cazul regresiei bivariante, ec. 39 și $R^2 = 0.98653$ în cazul regresiei multiple, cu 3 descriptorii, ec. 40). De precizat este faptul că în studiile anterioare la cel mai bun model obținut în predicția lui log P s-au folosit ecuații de regresie multiple cu 7 descriptorii (ecuația 17 din publicația lui Yovani Marrero Ponce et all.) cu $R^2 = 0.968$ și cu abilitatea de precizie $R^2 = 0.938$, valori mult mai mici decât rezultatele obținute de noi.

Regresie bivariabilă

$$\log P = 0.086403 + 3.442395 \cdot X_{4v} + 2.163165 \cdot G(N..O) \quad (39)$$

n = 12 $R^2 = 0.9843$ s = 0.091837 F = 282.3465

Regresie multi-variabilă

$$\log P = 0.391261 - 0.003500 \cdot MW + 3.287415 \cdot X_{4v} - 0.043587 \cdot G(N..O) \quad (40)$$

n = 12 $R^2 = 0.98653$ s = 0.07885 F = 195.3105

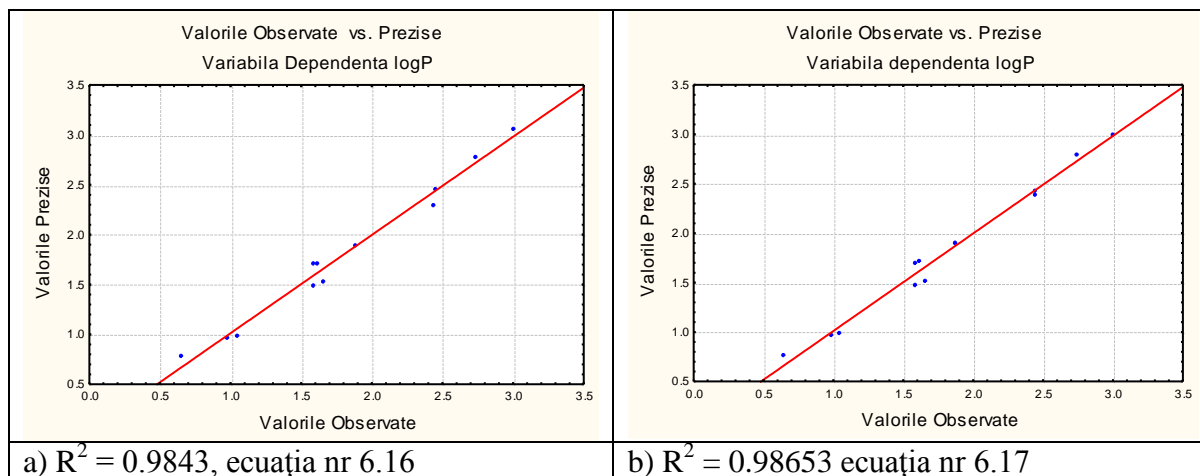


Figura 6. Reprezentarea grafică a proprietății observate vs. prezise pentru ec. 6.16 și 6.17 pentru structurile din setul școală.

Orice model QSAR trebuie să fie validat cu un set de predicție extern. În acest caz setul de predicție este format din cei 4 derivați de 2-furiletină la care prezicem valorile coeficientului de partiție n-octanol/apă logP (Tabelul 9).

Tabelul 9. Proprietatea log P observată și prezisă pentru structurile din setul de validare.

Nr. structură	log P observat	log P prezis ec 45	log P prezis ec. 46
35	2.820	2.924846	2.845186
36	2.730	2.580922	2.617008
37	1.290	1.438395	1.422570
38	1.940	1.709787	1.735940
R²		0.9333	0.9585

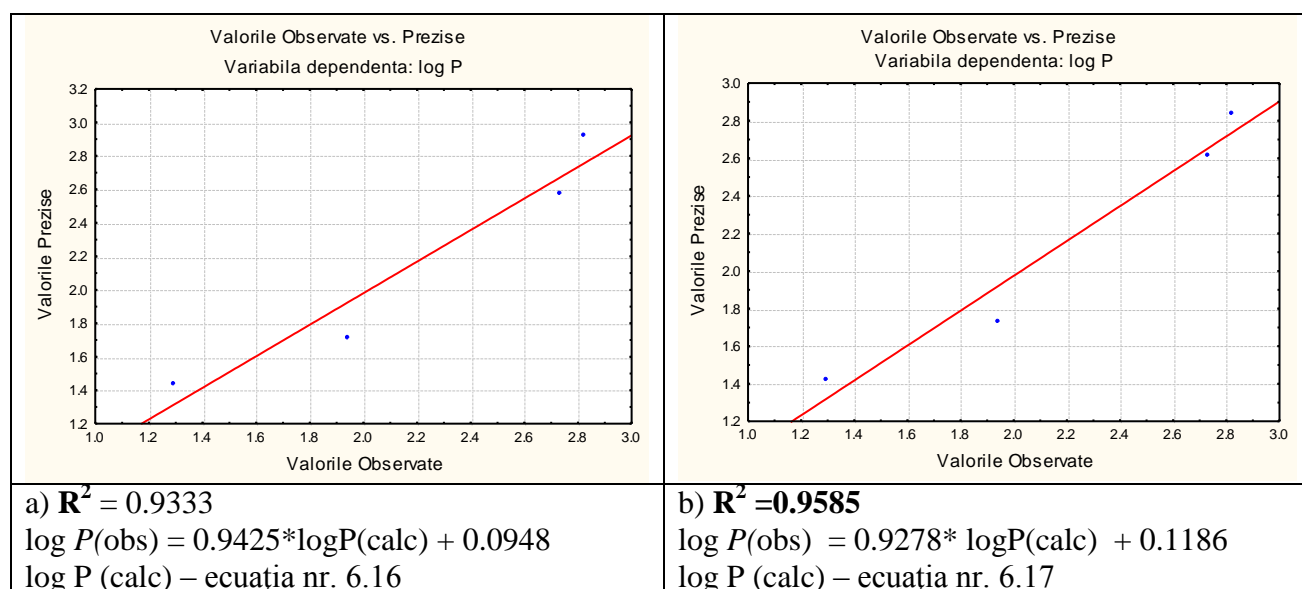


Figura 7. Reprezentarea grafică a proprietăților observate vs. prezise pe setul extern de validare.

Ecuția explică 95,85% din varianța lui log P, această valoare arată abilitatea crescută a modelului de predicție (ecuația 39). În figura nr. 7.b este reprezentată relația dintre valorile observate și cele prezise ale lui log P. În acest sens, ecuația obținută cu descriptorii MW, X4v, G(N..O) este cea mai bună pentru predicția proprietății log P.

Parcurgând aceleași etape am realizat prezicerea proprietății pe setul școală și validarea lui pe setul de predicție a coeficientului de partiție n-octanol/apă log P, dar de această dată cu ajutorul descriptorilor calculați cu programul TOPOCLUJ. Rezultatele obținute sunt prezentate în continuare. Cei mai semnificativi descriptorii calculați cu programul TOPOCLUJ sunt prezentați în tabelul 10 alături de proprietatea log P pentru structurile din setul școală determinat pe baza similarității moleculare cu liderii din setul de predicție extern.

Tabelul 10. Descriptorii topologici și coeficientul de partiție observat logP pentru setul de predicție.

Nr. Structură	CS[LM [Electronegativity]]	CS[Sh[W4 [Charge_Adjacency]]]	IE[CfMax [Density]]	VAA1	VAD1	Pobs
35	410	22	190	7.1	36	2.49
36	410	19	170	7.1	36	2.37
37	290	11	62	6.1	27	1.56
38	340	17	120	6.6	31	1.92

Ecuțiile utilizate pentru predicție sunt cele calibrate în clusterul școală (în forma lor normalizată de Matlab). Aceste ecuații prezintă un coeficient de corelație bun ($R^2 = 0.986337$ - ec. 6.20 , $R^2 = 0.960186$ -ec. 6.19, $R^2 = 0.911236$ - ec.41) (Figura 8).

Regresie bivariabilă

$$\log P = 3.562116 + 0.015455 \cdot \text{IE}[\text{CfMax}[\text{Density}]] - 0.109763 \cdot \text{VAD1} \quad (41)$$

n = 12 $R^2 = 0.911236$ s = 0.519627 F = 46.19599

Regresie multi-variabilă

$$\log P = 7.99004 + 0.05992 \cdot \text{CS}[\text{Sh}[\text{W4}[\text{Charge_Adjacency}]]] + 0.01853 \cdot \text{IE}[\text{CfMax}[\text{Density}]] - 1.45456 \cdot \text{VAA1} \quad (42)$$

n = 12 $R^2 = 0.960186$ s = 0.233074 F = 64.31076

$$\log P = 12.43950 + 0.01138 \cdot \text{CS}[\text{LM}[\text{Electronegativity}]] + 0.06207 \cdot \text{CS}[\text{Sh}[\text{W4}[\text{Charge_Adjacency}]]] + 0.01658 \cdot \text{IE}[\text{CfMax}[\text{Density}]] - 2.70391 \cdot \text{VAA1} \quad (43)$$

n = 12 $R^2 = 0.986337$ s = 0.079982 F = 126.3344

Modelele QSAR trebuie să fie validate pe un set de predicție extern. Setul de predicție este format din cei 4 derivați de 2-furiletină la care prezicem valorile coeficientului de partiție n-octanol/apă logP (Tabelul 12).

Tabelul 12. Proprietatea log P observată și prezisă pentru structurile din setul de validare

Nr. structură	log P observat	log P prezis ec 47	log P prezis ec. 48	log P prezis ec. 49
35	2.820	2.547183	2.501652	2.424993
36	2.730	2.238073	1.951290	1.907205
37	1.290	1.556749	0.925237	0.957875
38	1.940	2.014115	1.632225	1.509150
R^2		0.924	0.9037	0.9143

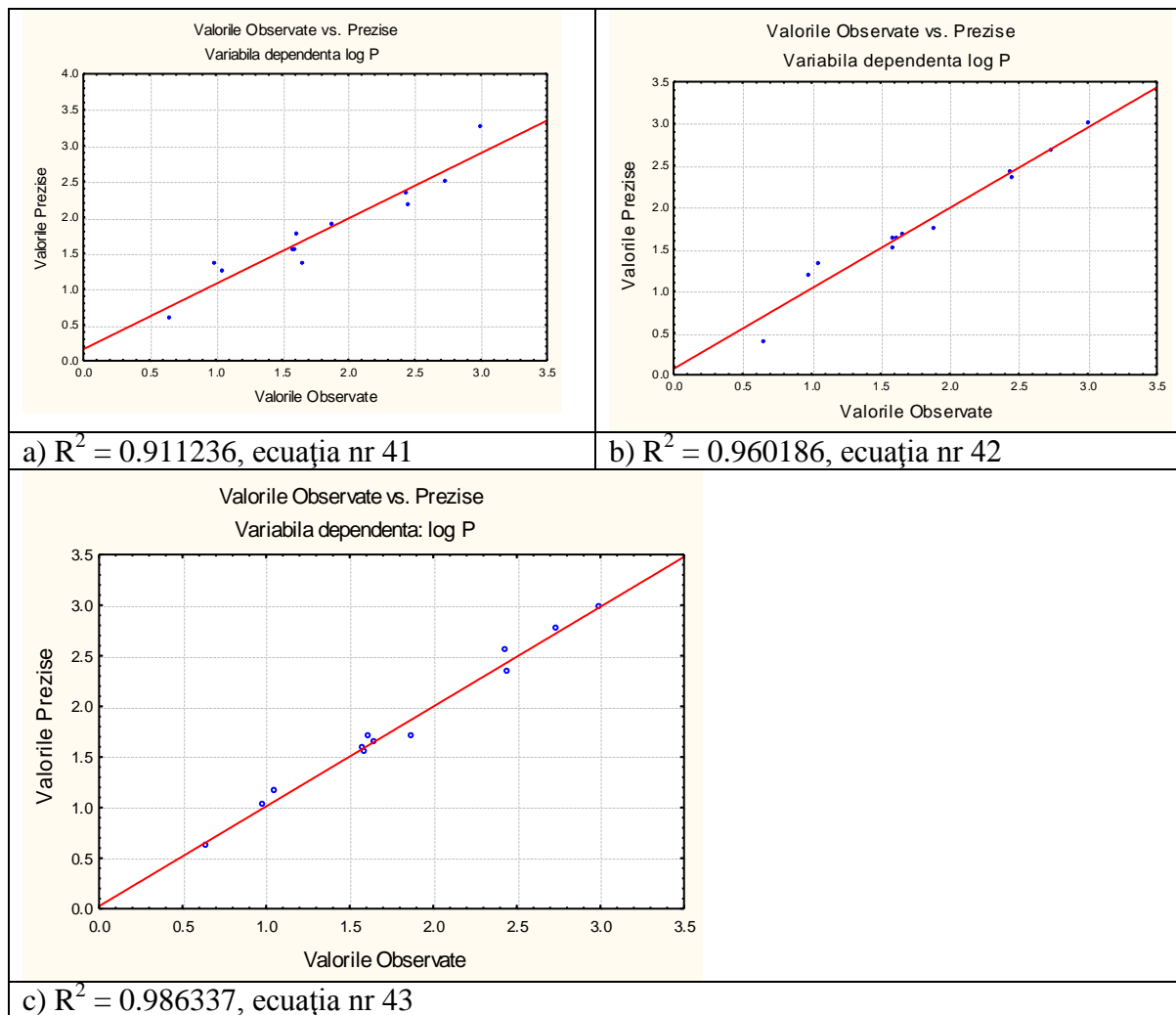


Figura 8. Reprezentarea grafică a proprietății observate vs. prezise pentru ec. 6.18, 6.19 și 6.20 pentru structurile din setul școală.

Descriptorii topologici implicați în această ecuație sunt $IE[CfMax[Density]]$ și $VAD1$. Reprezentarea grafică a acestora este prezentat in figura 9.

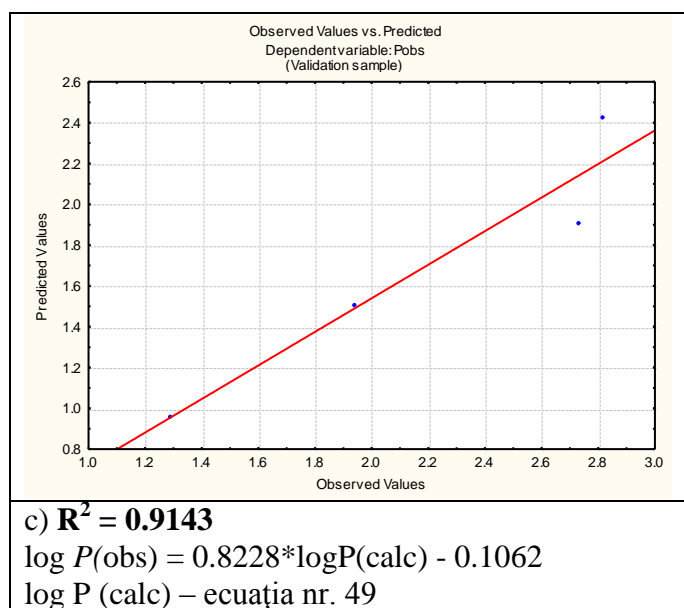
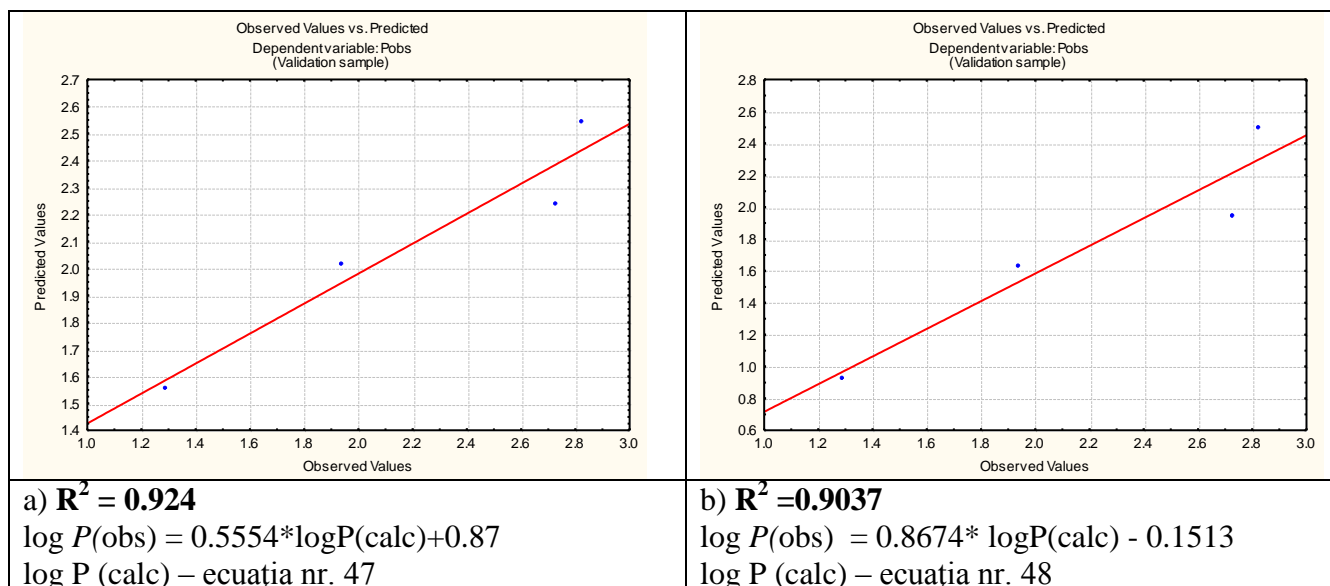


Figura 9. Reprezentarea grafică a proprietăților observate vs. prezise pe setul extern de validare.

***Noi modele QSAR pentru prezicerea activității biologice a derivaților de benzoxazol/
benzimidazol***

Bolile micotice sunt produse de microorganisme încadrate sistematic în regnul Fungi. Ele se găsesc în număr mare în mediul înconjurător. Majoritatea au adoptat mediul de viață saprobită, dar unele dintre ele s-au adaptat la viața parazitară. Se citează că peste 300 de specii au fost găsite a fi patogene pentru animale. Ciupercile obligatoriu parazite aparțin, în cea mai mare parte, categoriei dermatomicetelor. Ele nu se dezvoltă în mediu ci doar pot supraviețui și se pot transmite la organismele vii prin contagiune directă sau indirectă.

Obținerea descriptorilor moleculari

Recent s-au raportat sinteza și activitățile a câtorva derivați benzoxazolici/ benzimidazolici împotriva speciei *C. Albicans*. Analiza relațiilor cantitative structură activitate (QSAR) a fost mai mult răspândită și eficient folosită în studiile teoretice ale designului medicamentelor. Modelele QSAR propun cercetătorilor în domeniu sinteze orientate, permițând realizarea unor economii de timp, bani și energie, ameliorând cercetarea farmaceutică (Figura 10, Tabelul 13).

R ₁ = H, Cl, F, NO ₂ , CH ₃ , OCH ₃	Y= —, CH ₂ , CH ₂ O, CH ₂ S
R ₂ = H, Cl, Br, NH ₂ , CH ₃ , OCH ₃	X= O, NH
R ₅ =H, CH ₃	
R ₆ = H, CH ₃	

Figura 10. Activitatea antimicrobică a derivaților de benzoxazol și benzimidazol împotriva speciei *C.Albicans*

Tabelul 13. Valorile log 1/C observat pentru setul de derivați benzoxazolici/benzimidazolici.

	R ₁	R ₂	R ₅	R ₆	X	Y	log 1/C Obs.
1	Cl	H	CH ₃	H	O	-	3.989
2	OCH ₃	H	CH ₃	H	O	-	3.980
3	NO ₂	H	CH ₃	H	O	-	4.007
4	Cl	Cl	CH ₃	H	O	-	4.046
5	CH ₃	CH ₃	CH ₃	H	O	-	3.977
6	OCH ₃	OCH ₃	CH ₃	H	O	-	4.032
7	Cl	H	H	CH ₃	O	-	3.989
8	OCH ₃	H	H	CH ₃	O	-	3.980
9	F	H	H	CH ₃	O	-	3.958
10	NO ₂	H	H	CH ₃	O	-	4.007
11	Cl	Cl	H	CH ₃	O	-	4.046
12	CH ₃	CH ₃	H	CH ₃	O	-	3.977
13	OCH ₃	OCH ₃	H	CH ₃	O	-	4.032
14	H	H	CH ₃	H	O	CH ₂	4.251
15	H	Br	CH ₃	H	O	CH ₂	4.383
16	H	NH ₂	CH ₃	H	O	CH ₂	4.280
17	H	H	H	CH ₃	O	CH ₂	4.251
18	H	H	CH ₃	H	NH	CH ₂	4.249
19	H	Cl	CH ₃	H	NH	CH ₂	4.312
20	H	Br	CH ₃	H	NH	CH ₂	4.382
21	H	NH ₂	CH ₃	H	NH	CH ₂	4.278
22	H	H	CH ₃	H	O	CH ₂ O	3.980
23	H	H	CH ₃	H	O	CH ₂ S	4.009
24	H	H	CH ₃	H	NH	CH ₂ S	4.007
25	H	Cl	CH ₃	H	NH	CH ₂ O	4.037

Descriptorii topologici utilizați în studiile QSAR sunt accesibili și pot fi ușor calculați cu ajutorul programelor software. Setul descriptorilor moleculari utilizați în acest studiu a fost calculat cu pachetul software DRAGON.⁴⁵ Structurile au fost optimizate folosind metoda Hamiltoniană semi empirică PM3 disponibilă în softul HyperChem.¹³⁷

Prelucrarea și analiza datelor

Scopul acestui studiu este de a dezvolta un model QSAR nou și mult mai eficient decât celelalte obținute anterior pentru prezicerea activității antifungice a bezoxazolilor /benzimidazolilor substituiți în poziția 2 și cu grupare metil în poziția 5 sau 6.

Descriptorii topologici obținuți cu programul Dragon și activitatea biologică (antifungică) log 1/C pentru acest set de structuri sunt prezentați în tabelul 14.

Tabelul 14. Descriptorii topologici obținuți cu programul Dragon, valorile experimentale și valorile precise log 1/C ale derivaților de benzoxazoli/ benzimidazoli.

	MW	nCs	nHDon	Obs. log 1/C	Calc. log 1/C
1	243.7	0	0	3.989	3.988
2	239.29	0	0	3.980	3.981
3	254.26	0	0	4.007	4.006
4	278.14	0	0	4.046	4.046
5	237.32	0	0	3.977	3.977
6	269.32	0	0	4.032	4.031
7	243.7	0	0	3.989	3.988
8	239.29	0	0	3.980	3.981
9	227.25	0	0	3.958	3.960
10	254.26	0	0	4.007	4.006
11	278.14	0	0	4.046	4.046
12	237.32	0	0	3.977	3.977
13	269.32	0	0	4.032	4.031
14	223.29	1	0	4.251	4.251
15	302.18	1	0	4.383	4.384
16	238.31	1	2	4.280	4.279
17	223.29	1	0	4.251	4.251
18	222.31	1	1	4.249	4.251
19	256.75	1	1	4.312	4.309
20	301.2	1	1	4.382	4.384
21	237.33	1	3	4.278	4.278
22	239.29	0	0	3.980	3.981
23	255.36	0	0	4.009	4.008
24	254.38	0	1	4.007	4.007
25	272.75	0	1	4.037	4.038

Descriptorii moleculari semnificativii calculați cu programul Dragon sunt: MW- masa moleculară, nCs- numărul atomilor de C secundar (sp^3) și nHDon- numărul atomilor donori ai legăturii de H (N și O).

S-a găsit ecuația QSAR utilizând analiza de regresie multiliniară. Calitatea modelului obținut este dată de pătratul coeficientului de regresie (R^2), raportul Fischer, eroarea standard estimată (s) și de procedeul leave-one-out (LOO) ca și procedură de cross-validare.

Ecuația prezintă o valoare bună a coeficientului de corelație $R^2=0.999$ în regresia multivariabilă cu 3 descriptori. Aceasta corelație este semnificativ mai bună decât cele raportate anterior ($R^2=0.94$). Valorile experimentale și cele calculate ale lui $\log 1/C$ pentru acest set sunt reprezentate grafic în Figura 11.

$$\log 1/C = 3.577225 + 0.001686 * MW + 0.297347 * nCs + 0.00122 * nHDon \quad (44)$$

n=25 $R^2=0.999932$ $s=0.0029$ $F=103092.3$

Pobs Std.Dev.= 0.145464

Pobs Std.Err = 0.029093

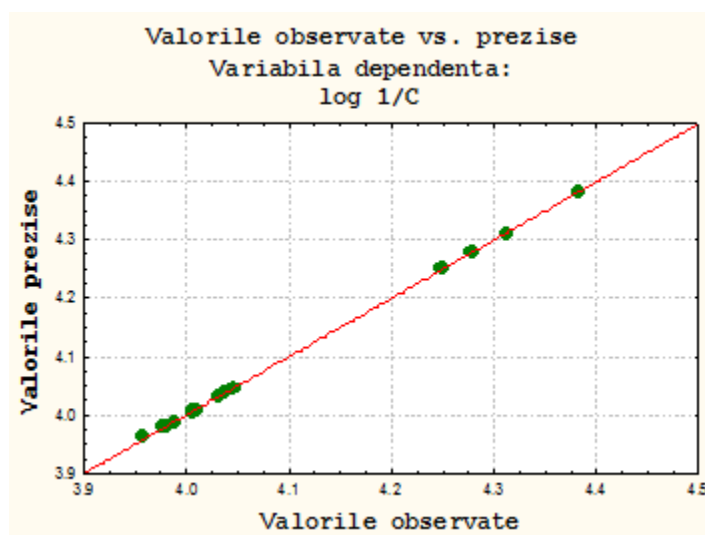


Figura 11. Reprezentarea grafica a valorilor observate vs. prezise pentru ecuația 6.21.

Pe baza datelor obținute s-a propus un nou model dezvoltat pentru predicția directă a activității biologice. S-au construit clusteri de similaritate pentru fiecare moleculă din setul scoală și s-au realizat predicții în subseturile congenerice obținute. Descriptorii moleculari semnificativi sunt MW și nCs. S-au găsit ecuațiile QSAR utilizând analiza regresiei multivariate (Tabelul 15).

Scopul acestei metode este acela de obține cel mai bun model pentru prezicerea viitoare a activității antifungice pentru alte structuri noi. Pentru fiecare structură, luată ca și lider în testul de similaritate, s-a calculat ecuația de predicție bivariată pentru fiecare subset de similaritate. Fiecare din cele 25 de ecuații de învățare au $R^2=0.9999$ iar această cifră variază doar începând de la a 5-a zecimală. Ecuația de predicție obținută prezintă un coeficient de corelație foarte bun, $R^2=0.999$, fiind o regresie bivariabilă (Figura 12).

Corelația este de asemenea mult mai bună decât cea raportată anterior ($R^2=0.94$) astfel încât ea este excelentă.

Tabelul 15. Valorile parametrilor regresiei precum și valorile observate și calculate pentru log 1/C, ec. $Y_{calc}=a + b*MW + c*nCs$.

Nr.	MW	nCs	log 1/C Obs.	Parametrii ecuațiilor de predicție			log 1/C Calc.
				a	b	c	
1	243.70	0	3.989	3.529574	0.001885	0.300421	3.989006
2	239.29	0	3.980	3.578017	0.001684	0.298368	3.980871
3	254.26	0	4.007	3.577699	0.001684	0.298482	4.005963
4	278.14	0	4.046	3.577586	0.001685	0.298405	4.046282
5	237.32	0	3.977	3.577923	0.001684	0.298388	3.97753
6	269.32	0	4.032	3.586488	0.001650	0.299659	4.030916
7	243.70	0	3.989	3.577964	0.001683	0.299659	3.988105
8	239.29	0	3.980	3.578017	0.001684	0.298368	3.980871
9	227.25	0	3.958	3.579723	0.001677	0.298309	3.960813
10	254.26	0	4.007	3.577699	0.001684	0.298482	4.005963
11	278.14	0	4.046	3.586385	0.001651	0.299609	4.045527
12	237.32	0	3.977	3.586968	0.001649	0.299502	3.978302
13	269.32	0	4.032	3.586488	0.001650	0.299659	4.030916
14	223.29	1	4.251	3.529649	0.001885	0.300392	4.250924
15	302.18	1	4.383	3.582466	0.001666	0.299986	4.385951
16	238.31	1	4.280	3.585095	0.001656	0.299324	4.279045
17	223.29	1	4.251	3.529649	0.001885	0.30032	4.250924
18	222.31	1	4.249	3.529748	0.001884	0.300463	4.249153
19	256.75	1	4.312	3.585044	0.001656	0.298955	4.309212
20	301.20	1	4.382	3.583753	0.001661	0.299843	4.383946
21	237.33	1	4.278	3.585631	0.001654	0.299458	4.277595
22	239.29	0	3.980	3.587052	0.001649	0.299486	3.981574
23	255.36	0	4.009	3.585960	0.001652	0.299638	4.007876
24	254.38	0	4.007	3.585968	0.001652	0.299608	4.006286
25	272.75	0	4.037	3.586233	0.001651	0.299593	4.036649

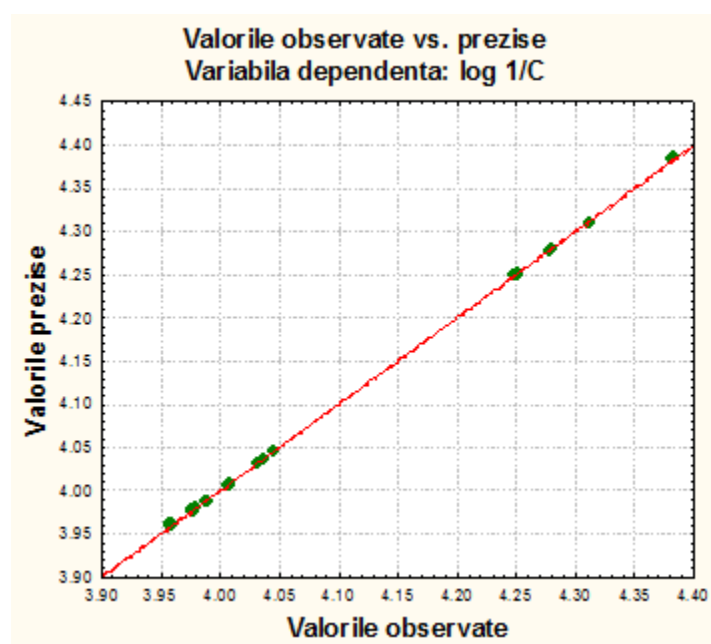


Figura 12. Reprezentarea grafica a log 1/C observat fata de log 1/C prezis.

Setul de derivați de benzoxazoli/benzimidazoli substituiți în poziția 2 și având gruparea metil în poziția 5 sau 6, testat anterior pentru activitatea antibacteriană împotriva lui *C. Albicans*, a fost analizat prin relația cantitativă structură-activitate biologică și contribuțiile la activitate pentru efectele structurale și funcționale au fost determinate utilizând procedeul de regresie multiplă. Rezultatele QSAR obținute relevă faptul că substituentul din poziția Y (descriptorul nCs) joacă un rol important și își aduce o contribuție importantă la activitatea antibacteriană.

CONCLUZII

Sinteza și implementarea pe piață a unor medicamente performante care să răspundă într-un grad cât mai ridicat la soluționarea unor probleme de sănătate acute societății reprezintă un deziderat important al industriei farmaceutice în special, și al cercetătorului în particular.

Găsirea unor metode teoretice prin care se reduc costurile și timpul necesar găsirii și sintezei compușilor biologic activi cu aplicabilitate practică, reprezintă țelul oricărui cercetător. Una dintre aceste metode îl oferă chimia matematică prin folosirea diferitelor tehnici și modele, prin care structurile sunt evaluate și cuantificate într-un număr.

Teza de față urmărește tratarea și obținerea diferitelor modele teoretice cu aplicabilitate în design-ul de molecule bioactive, în special prin prisma similarității structurale. Obținerea unor metode de evaluare cantitativă a similarității este doar una dintre problemele care trebuie rezolvate de către "proiectantul" de substanțe bioactive. Noțiunile centrale ale oricărui sistem de căutare a similarității intermoleculare sunt criteriul de similaritate și măsura utilizată pentru cuantificarea similarității.

Astfel, este tratată problema similarității moleculare, urmărind și dezvoltând diferite aspecte ale acestei teorii. În acest sens au fost prezentate pe larg:

- descrierea cantitativă a similarității
- indici de similaritate moleculară
- descriptorii graf-moleculari
- programul software TOPOCLUJ-SIMIL

Au fost analizate diferite tehnici și metode de suprapunere moleculară, similaritate, aplicate pe clase de structuri cu activitate biologică.

Aceste metode includ:

- analiza de clusteri
- analiza componentelor principale
- analiza de factori

S-au propus și realizat diferite modele de clusterare pe baza similarității, utile pentru modelarea și prezicerea proprietăților moleculare (biologice sau fizico-chimice), având o contribuție majoră la studiile cantitative ale relației structură - activitate biologică QSAR pentru diferite tipuri de molecule.

Rezultatele modelelor QSAR propuse urmăresc într-un mod excelent proprietățile biologice ale claselor de structuri propuse. Astfel, modelele propuse cu 3 variabile independente, MW, X4v și G(N...O), având $R^2=0.98653$, prezintă o capacitate avansată de predicție a proprietății log P în setul de validare având $R^2=0.9585$.

Rezultatele obținute devansează rezultatele similare din literatură, având o abilitate mai bună de predicție și în același timp un număr redus de descriptori utilizați în ecuația de regresie, fiind utile în estimarea proprietăților fizico-chimice și biologice a unor noi derivați de 2-furiletană.

Modelele QSAR de succes obținute prezintă o corelație statistică semnificativă între caracteristicile chimice ale compușilor (descriptorilor) și activitatea biologică. În urma analizelor ecuațiilor obținute în urma modelării seturilor de compuși, s-a ajuns la concluzia, că, este necesar în toate cazurile să se evite predicția activității biologice a compușilor care au structură foarte diferită de a compușilor din setul școală. Oricum, oricare dintre aceste proceduri poate fi folosită separat ca parte a studiilor computaționale cu multiple scopuri de atins. Sunt necesare aceste tehnici analitice din ce în ce mai mult pentru generarea,

interpretarea și redarea mai multor informații despre molecule cu potențial biologic și nu numai.

De asemenea, s-a observat faptul că, numai un subset de descriptori ai structurilor moleculare, care sunt cei mai importanți și semnificativi din punct de vedere statistic, sunt selectați pentru a descrie o activitate biologică aleasă.

Rezultatele obținute au fost publicate în următoarele reviste:

1. Costescu, A., Moldovan, C.D., Diudea, M.V., QSAR modeling of steroid hormones, *Match* 55 (2), pp. 315-329, **2006**.
2. Moldovan, C.D., Costescu, A., Katona, G., Diudea, M.V., A novel QSAR approach in modeling antifungal activity of some 5-or 6-methyl-2-substituted benzoxazoles/benzimidazoles against *C. albicans* using molecular descriptors, *Match* 60 (3), pp. 977-984, **2008**.
3. Costescu, A., Moldovan, C.D., Katona, G., Diudea, M.V., QSAR modeling of human catechol O-methyltransferase enzyme kinetics, *Journal of Mathematical Chemistry* 45 (2), pp. 287-294, **2009**.
4. Moldovan, C.D., Costescu, A., Katona, G., Diudea, M.V., Application to QSAR studies of 2-furylethylene derivatives, *Journal of Mathematical Chemistry* 45 (2), pp. 442-451, **2009**.

INDEX

- A**
Activitatea biologică, 77, 86, 96, 100
Analiza clusterilor, 45
Analiza componentelor principale, 48, 67, 82
Analiza factorială, 51
Afinitatea de legare, 62,76
- C**
Covarianța, 42, 49
Cross-validare, 73, 100, 103
- D**
Descriptorii topologici, 5, 14, 88, 89, 92, 95
Dragon 2.1, 63, 66, 71, 80, 88, 97
Disimilaritatea, 7, 9, 12, 15, 59, 84
Distanța Euclidiană, 14, 46, 64
- E**
Electronegativitățile Sanderson de grup, 64
- G**
Geometry optimization, 69
Graf molecular, 7, 10, 16
- H**
HyperChem, 63, 69, 80, 85
Hyper-Wiener, 28, 33
- L**
Layer matrices, 19
Leave-one-out, 73, 103
Lipofilicitate, 6, 37, 56, 77
- I**
Indici topologici, 24
Invariant 17, 22, 23, 24, 32
- M**
Matricea de adiacență, 16
Matricea de conectivitate, 16
Matricea distanțelor, 17
Matricea Wiener, 21
Matrici strat, 18
Matrici topologice, 16
Molecular mechanics, 69
- O**
Outliers, 73
- P**
PM3, 68
- Q**
QSAR, 14, 34, 66, 68, 69, 75, 76, 79, 86, 95, 96, 97, 103, 116
QSPR, 14, 34, 69, 86
- R**
Regresii liniare, 43
- S**
Setul de validare (predicție), 73, 90, 94, 119
Setul școală, 73,
SIMIL, 84
Similaritatea topologica, 7, 9, 10, 14
Single point, 69
STATISTICA, 68, 99
- T**
TOPOCLUJ, 64, 69, 80

Bibliografie selectiva

- 1 Rosen, R. in: Johnson, M. A.; Maggiora, G. M. Eds. Concepts and Applications of Molecular Similarity, Wiley, New York, **1990**, Chap. 12, 369-382.
- 2 Mezey, P. G. Three-Dimensional Topological Aspects of Molecular Similarity. In: Johnson, M.A.; Maggiora, G. M. Eds. Concepts and Applications of Molecular Similarity, Wiley, New York, **1990**, Chap. 11, 321-368.
- 3 Randić, M. Design of Molecules with Desired Properties. A Molecular Similarity Approach to Property Optimization. In: Johnson, M. A.; Maggiora, G. M. Eds. Concepts and Applications of Molecular Similarity, Wiley, New York, **1990**, Chap. 5, 77-145.
- 4 Maggiora, G. M.; Johnson, M. A. Introduction to Similarity in Chemistry. In: Johnson, M. A.; Maggiora, G. M. Eds. Concepts and Applications of Molecular Similarity, Wiley, New York, **1990**, Chap. 1, 1-13.
- 5 Tsai, C. -c.; Johnson, M. A.; Nicholson, V.; Naim, M. Eds., Graph Theory and Topology in Chemistry, Elsevier, Amsterdam, **1987**, 231.
- 6 Balaban, A. T.; Chiriac, A.; Motoc, I.; Simon, Z. Steric Fit in QSAR (Lecture Notes in Chemistry, Vol. 15), Springer, Berlin, **1980**, Chap. 6.
- 7 Kvasnička, V.; Pospichal, J. Fast Evaluation of Chemical Distance by Tabu Search Algorithm. *J. Chem. Inf. Comput. Sci.* **1994**, 34, 1109-1112.
- 8 Diudea, M.V. Layer Matrices in Molecular Graphs, *J. Chem. Inf. Comput. Sci.* **1994**, 34, 1064-1071.
- 9 Ugi, I.; Wochner, M.A.; Fontain, E.; Bauer, J.; Gruber, B.; Karl, R. Chemical Similarity, Chemical Distance, and Computer Assisted Formalized Reasoning by Analogy, in: Maggiora, G. M. Eds. Concepts and Applications of Molecular Similarity, Wiley, New York, **1990**, Chap. 9, 239-288.
- 10 Basak, S.C.; Magnusson, V.R.; Niemi, G.J.; Regal, R.R., Determining Structural Similarity of Chemicals Using Graph-Theoretic Indices, *Discr. Appl. Math.* **1988**, 19, 17-44.
- 11 M. Randić, W.L. Woodworth, A. Graovac, Unusual random Walks, *Int. J. Quantum Chem.*, **1983**, 24, 435-452.
- 12 M. Randić, Generalized molecular descriptors, *J. Math. Chem.*, **1991**, 7, 155-168.
- 13 N. Trinajstić, Chemical Graph Theory, CRC Press. Inc., Boca Raton Florida, **1983**.
- 14 A.T. Balaban, I. Moțoc, D. Bonchev, O. Mekenyan, Topological Indices for Structure - Activity Correlations, *Top. Curr. Chem.*, **1993**, 114, 21-55.
- 15 D.H. Rouvray, The challenge of characterizing branching in molecular species, *Discr. Appl. Math.*, **1988**, 19, 317-338.
- 16 M. Randić, Design of molecules with desired properties. A molecular similarity approach to property optimization, in "Concepts and Applications of Molecular Similarity, M.A. Johnson and G.M. Maggiora, Eds., John Wiley & Sons, Inc., **1990**.
- 17 R. Carbó, L. Leyda, and M. Arnau, *Int. J. Quantum Chem.*, **1980**, 17, 1185-1189.
- 18 D. Ciubotariu, S. Mureșan, V. Gogonea, M. Medeleanu, D. Dragoș, Relații Cantitative Structură Chimică-Activitate Biologică (QSAR), Ed.Mirton, Timișoara, **1996**
- 19 E. Overton, Studien über die Narkose, Fischer, Jena, **1901**; A. Meyer, *Arch. Exptl. Pathol. Pharmacol.*, **1899**, 42, 110.
- 20 J.N. Langley, *J. Physiol.*(London), **1908**, 1, 339.
- 21 P. Ehrlich, *Ber. Dtsch. Chem. Ges.*, **1909**, 42, 17.
- 22 A. R. Cushny, Biological Relations of Optically Isomeric Substances, Balliere, Tindall and Cox, London, **1926**.
- 23 I. Moțoc, Structura moleculelor și activitatea biologică, Ed. Facla, Timișoara, 1980.
- 24 E.J. Ariëns, Stereochemistry: A Source of Problems in Medicinal Chemistry, *Med.Res.Rev.*, **1986**, 6, 451-466.

-
- 25 E.J. Ariëns, Stereochemistry in the Analysis of Drug-Action, Part II. *Med. Res. Rev.*, **1987**, 7, 367-387.
 - 26 E.J. Ariëns, Stereochemical Implications of Hybrid and Pseudohybrid Drugs, Part III. *Med. Res. Rev.*, **1988**, 8, 309-320.
 - 27 I.D. Resa, S. Petrescu, M. Precupas, A. Căra, Probleme de statistica rezolvate pe calculator, Ed. Facla, Timisoara, **1984**.
 - 28 D. McCormick, A. Roach, Measurement, Statistics and Computation, John Wiley & Sons, London, **1987**.
 - 29 Katritzky, A. R.; Maran, U.; Lobanov, V. S.; Karelson, M. Perspective: Structurally diverse quantitative structure-property relationship correlations of technologically relevant physical properties. *J. Chem. Inf. Comput. Sci.* **2000**, 40, 1-18.
 - 30 Robinson, D. D.; Winn, P. J.; Lyne, P. D.; Richards, W. G. Self-Organizing Molecular Field Analysis: A Tool for Structure-Activity Studies. *J. Med. Chem.* **1999**, 42, 573-583.
 - 31 Cramer, R. D., I.; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis (CoMFA). 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, 110, 5959-5967.
 - 32 Coats, E. The CoMFA Steroids as a Benchmark Dataset for Development of 3D QSAR Methods. *Perspect. Drug DiscoV. Design.* **1998**, 12/13/14, 199-213.
 - 33 Dunn, J. F.; Nisula, B. C.; Rodbard, D. Transport of Steroid Hormones: Binding of 21 Endogeneous Steroids to Both Testosterone- Binding Globulin and Corticosteroid-Binding Globulin in Human Plasma. *J. Clin. Endocrin. Metab.* **1981**, 53, 58-68.
 - 34 Tuppurainen K, Viisas M, Laatikainen R, Perakyla M.: Evaluation of a Novel Electronic Eigenvalue (EEVA) Molecular Descriptor for QSAR/QSPR Studies: Validation Using a Benchmark Steroid Data Set. *J. Chem. Inf. Comput. Sci.* **2002**, 42, 607-613
 - 35 Rios-Santamaria I, Garcia-Domenech R, Cortijo J, Santamaria P, Morcillo E J and Galvez J, *Internet Electronic Journal of Molecular Design*, **2002**, 1, 70.
 - 36 Galvez J, Garcia-Domenech R, Salabert M T and Soler R, *J Chem Inf Comput Sci*, **1994**, 34, 520.
 - 37 Todeschini, R.; Consonni, V. Handbook of molecular descriptors. Wiley-VCH: Weinheim, Germany, **2000**.
 - 38 Polanski, J.; Walczak, B. The Comparative Molecular Surface Analysis (COMSA): a novel Tool for Molecular Design. *Comput. Chem.* **2000**, 24, 615-625.
 - 39 M. Pastor, G. Cruciani, I. McLay, S. Pickett and Sergio Clementi: Grid Independent Descriptors (GRIND): A Novel Class of Alignment- Independent Three-Dimensional Molecular Descriptors. *J. Med. Chem.* **2000**, 43, 3233-3243.
 - 40 M.A. Cabrea Pérez et al.: Experimental and theoretical determination of physicochemical properties in a novel family of microcidal compounds. *European Bulletin of Drug Research*, **2001**, 9, 1.
 - 41 Yovani Marrero Ponce et al.: Atom, Atom-Type, and Total Linear Indices of the Molecular Pseudograph's Atom Adjacency Matrix: Application to QSPR/QSAR Studies of Organic Compounds, *Molecules* **2004**, 9, 1100-1123.
 - 42 Dore, J. Ch.; Viel, C. Antitumoral Chemoterapy. X. Cytotoxic and Antitumoral Activity of β -Nitrostyrenes and Nitrovinyl Derivatives. *Farmaco.* **1975**, 30, 81-109.
 - 43 Castañedo, N.; Goizueta, R.; Perez, J.; Gonzalez, J.; Silveira, E. Cuesta, M.; Martinez, A.; Lugo, E.; Estrada, E.; Carta, A.; Navia, O.; Delgado, M. Cuban Pat. 22446, 1994; Can. Pat. 2,147,594, **1999**.
 - 44 Blondeau, J. M.; Castañedo, N.; Gonzalez, O.; Medina, R.; Silveira, E. In Vitro Evaluation of G-1: A Novel Antimicrobial Compound. *Antimicrob. Agents Chemother.* **1999**, 11, 1663-1669.

