



(12) CERERE DE BREVET DE INVENȚIE

(21) Nr. cerere: a 2022 00333

(22) Data de depozit: 15/06/2022

(41) Data publicării cererii:
29/12/2023 BOPI nr. 12/2023

(71) Solicitant:
• ARTIFICIAL INTELLIGENCE EXPERT
S.R.L., STR.ALEXANDRU VLAHUȚĂ,
BLOC LAMA C, NR.45, CLUJ-NAPOCA, CJ,
RO

(72) Inventatori:
• FLOARES ALEXANDRU,
STR. ALEXANDRU VLAHUȚĂ, BL. LAMA C,
AP. 45, CLUJ-NAPOCA, CJ, RO;
• ZETY ADRIAN, STR.AUREL VLAICU,
BL.17, AP.11, CLUJ-NAPOCA, CJ, RO

(54) SISTEM INTELIGENT DE ANALIZĂ AUTOMATĂ A DATELOR
MICRORNA MICROARRAY

(57) Rezumat:

Invenția se referă la un sistem și la un algoritm de analiză a datelor microRNA brute, obținute prin tehnologiile Microarray, cu ajutorul unui flux de lucru ce combină etape bioinformatică și de inteligență artificială pentru a obține biomarkeri relevanți pentru problema investigată și modele predictive performante pentru soluționarea acestora, cu aplicare în domenii precum zootehnie, agricultură, criminalistică și medicină. Algoritmul de analiză, conform invenției, cuprinde următoarele etape:

-preprocesarea datelor brute miRNA microarray și controlul calității, etapă ce include și normalizarea datelor,

-obținerea profilului de expresie al microRNA,
-adnotarea miRNA care constă în atașarea denumirii fiecărui miRNA la ID-urile din tabelul cu expresia fiecărei molecule, obținându-se datele de intrare pentru un expert AI care aplică în mod automat algoritmi de inteligență artificială care efectuează:

-analiza exploratorie și curățarea datelor prin verificarea formei datelor, identificarea eventualelor date lipsă și procentul acestora, eventualele anomalii prezente și efectuarea unei analize statistice sumare precum: distribuția datelor, valori minime, maxime, medii, etc., și ulterior verifică corelațiile dintre variabilele prezente, putându-se observa în acest fel potențialul ca unele variabile să fie mai importante în analiza și construcția modelelor predictive,

-preprocesarea datelor înainte de antrenarea modelelor,
-dezvoltarea modelelor predictive și
-optimizarea parametrilor algoritmilor de modele.

Revendicări: 1
Figuri: 2



Fig. 2

Cu începere de la data publicării cererii de brevet, cererea asigură, în mod provizoriu, solicitantului, protecția conferită potrivit dispozițiilor art.32 din Legea nr.64/1991, cu excepția cazurilor în care cererea de brevet de invenție a fost respinsă, retrasă sau considerată ca fiind retrasă. Întinderea protecției conferite de cererea de brevet de invenție este determinată de revendicările conținute în cererea publicată în conformitate cu art.23 alin.(1) - (3).



Descriere

Sistem Inteligent de Analiza Automată a Datelor microRNA Microarray

Invenția se referă la un sistem și un algoritm de analiza a datelor microRNA brute, obținute prin tehnologiile Microarray, cu ajutorul unui flux de lucru ce combina etape Bioinformatică și de Inteligența Artificială, pentru a obține biomarkeri relevanți pentru problema investigată și modele predictive performante, pentru a o soluționa.

Dat fiind rolul important al miRNA de reglare posttranslatională a activității genelor în întreaga lume vie a pluricelularelor, descoperirea biomarkerilor și dezvoltarea de modele predictive este importantă în domenii precum zootehnia, agricultura, criminalistica, și medicina.

MicroRNA (miRNA) sunt molecule mici de ARN necodificatoare de proteine care joacă un rol reglator important în traducerea genelor prin degradarea sau blocarea unor molecule de ARN mesager. Acestea influențează numeroase procese biologice majore, incluzând diferențierea, proliferare, apoptoza, inflamația și metabolismul¹. Datorită rolului proeminent pe care miRNA îl joacă în expresia genică și funcționalitatea normală a organismelor, nu e surprinzător faptul că expresia lor aberantă poate duce la o multitudine de boli incluzând cancerul, bolile neurodegenerative, diabetul, condițiile cardiace, disfuncții ale rinichilor și ficatului, alterări ale sistemului imun. În plus, pe lângă contribuția la cauza a numeroase boli, microRNA pot fi utilizate și în terapiile țintite și în generarea de biomarkeri pentru diagnosticarea precoce a multor afecțiuni. Dintre acestea, un rol important îl joacă moleculele de miRNA circulante care pot prelua rolul de biomarkeri² fiind detectabile din probe de sânge, stand la baza metodelor neinvazive de detecție a numeroase boli³.

Problema fundamentală întâlnită în toate domeniile de aplicabilitate a analizei miRNA ce este departe de a fi satisfăcător rezolvată, este aceeași. Utilizarea statisticii convenționale conduce la o listă de miRNA ce au o expresie diferită în două sau mai multe situații de interes distincte (expresie diferențială), de exemplu Cancer sau Normal. O astfel de listă este de foarte mică utilitate practică. Inteligența Artificială (utilizată în sistemul nostru) face posibilă obținerea unor modele predictive, cu valoare practică (de exemplu teste de diagnostic al cancerului cu acuratețe >95%). Simpla aplicare a AI, chiar a celor mai avansate tehnici, nu conduce automat la astfel de performanțe întrucât sistemele vii au o proprietate aparte - redundanță funcțională. Acesta este de regulă neglijată (nu și de către sistemul nostru) conducând la modele ce își

¹ O'Brien J, Hayder H, Zayed Y, Peng C. Overview of microRNA biogenesis, mechanisms of actions, and circulation. *Front. Endocrinol.* 9, 402 (2018).

² Fabbri M (2010) miRNAs as molecular biomarkers of cancer. *Expert Rev* 10(4):435–444

³ Chen X et al. (2008) Characterization of microRNAs in serum: a novel class of biomarkers for diagnosis of cancer and other diseases. *Cell Res* 18(10):997–1006

degradează performanța când sunt aplicate la cazuri noi, diferite de cele folosite la învățare. Există numeroase domenii ce pot beneficia de pe urma analizei miRNA precum zootehnia, agricultura, medicina.

Asadar, printre principalele avantaje ale sistemului nostru se enumera și faptul ca iau în considerare redundanță funcțională conducând la modele robuste ce generalizează bine la noi cazuri.

În oricare din domeniile de aplicabilitate de mai sus, se doresc soluții pragmatice și cu o precizie cât mai înaltă. Dezavantajul major al soluțiilor tehnice cunoscute, implementate în pachete software comerciale (precum Agilent miRNA Microarray Software, Partek Flow, sau open source (cum ar fi diverse pachete prezente în biblioteca R Bioconductor, de exemplu, AgiMicroRna, sau pachete individuale precum Chipster, este o consecință a utilizării statisticii convenționale, ce conduce doar la o listă de miRNA ce au o expresie cantitativă diferită în două sau mai multe situații de interes distincte (expresie diferențială), de exemplu Cancer sau Normal. O astfel de soluție tehnică este de foarte mică utilitate practică întrucât:

- Metoda statistică de obținere a subsetului de miRNA relevanți pentru problema dată, din întregul set de miRNA cunoscuți și determinați cantitativ de echipamentul microarray de laborator, nu primește feedback de la acuratețea cu care este rezolvată problema (de exemplu, diagnosticarea cancerului). Aceasta are impact asupra acurateții soluției.
- Conține un număr foarte mare de miRNA exprimați diferit. Determinarea cantitativă a tuturor acestora, pentru aplicarea la situații concrete, este atât costisitoare cât și inutilă.
- Pentru reducerea costurilor se alege un prag arbitrar de expresie diferențială pentru a scurta lista.
- Expresia biomarkerilor miRNA, obținuți prin trunchierea listei complete trebuie apoi comparată, mai mult sau mai puțin manual, cu expresia aceluiași miRNA la un nou caz la care se aplică.
- Drept rezultat, soluțiile tehnice cunoscute au adesea acuratețe scăzută, sunt costisitoare și dificil de aplicat.

Deși etapele de procesare bioinformatică a datelor miRNA microarray sunt similare la toate soluțiile (inclusiv a noastră), fiind standard, invenția noastră înlătură toate aceste dezavantaje ale soluțiilor tehnice cunoscute, înlocuind analiza expresiei diferențiale (statistica convențională) cu un întreg flux de lucru bazat pe Inteligența Artificială (AI). Mai mult, pentru a putea fi folosită cu ușurință și de către cei fără cunoștințe de programare sau AI, face apel la metode de automatizare a Machine Learning (AutoML). Datorită integrării unui flux AI automatizat ce face invenția utilizabilă și de către utilizatorii fără cunoștințe de AI sau programare am denumit pipeline-ul AI "Expert AI".

Etapele din cadrul Expertului AI având proprietatea de a se desfășura automat prin aplicarea AutoML. Toate acestea fac posibilă obținerea unor biomarkeri relevanți și a unor modele predictive cu valoare practică (de exemplu teste de diagnostic al cancerului cu acuratețe >95%).

Problema tehnică pe care o rezolvă invenția este de a descoperi cât mai ușor, din datele brute miRNA microarray, biomarkeri și modele predictive cu o înaltă acuratețe, la costuri pe cât posibil mai reduse, ce pot soluționa probleme complexe din domeniile amintite (de exemplu, diagnosticarea cancerului). Aceste avantaje sunt posibile datorită înlocuirii metodei statistice de analiză a expresiei diferențiale cu un flux de lucru bazat pe Inteligența Artificială Automată (AutoML).

Pentru descoperirea biomarkerilor miRNA relevanți și transformarea acestora în instrumente pragmatice, de soluționare a problemelor în domeniile menționate mai sus, prin intermediul unor modele predictive, introducem metodologii de Inteligența Artificială adaptate specificului sistemelor vii, caracterizate prin redundanța funcțională la nivel molecular (miRNA aici), în locul analizei statistice convenționale - expresia diferențială a miRNA.

Se dă, în continuare, un exemplu de realizare a invenției, în legătură cu figurile 1, 2 și 3, care reprezintă:

- Figura 1. Sistem - vedere de ansamblu
- Figura 2. Workflow Bioinformatic specific datelor Microarray și Expert AI

Se dă, în continuare un exemplu de realizare a invenției pe un set de date publice miRNA circulante microarray GSE168227 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE168227>) a 48 pacienți cu 30 Carcinom celular scuamos oral și 18 Normal. Este importantă discriminarea cancerului de normal, motiv pentru care grupăm cazurile de cancer în clasa Cancer (30 pacienți cu cancer) și clasa normal (18). Trebuie subliniat foarte clar că pașii necesari obținerii datelor - colectarea sangelui/plasmei, extracția ARN, pregătirea ADN complementar (cDNA), alegerea și utilizarea platformei de Microarray - NU fac obiectul invenției. Invenția noastră se aplică DOAR datelor nu și probelor pacienților. Etapele realizării invenției sunt următoarele:

Preprocesarea datelor și controlul calitatii. Acest pas include normalizarea datelor brute ceea ce verifică apariția variațiilor tehnice între array-urile din cadrul experimentului. Scopul normalizării este de a îndepărta variația cauzată de motive tehnice păstrând variația biologică a datelor. Procesul începe cu vizualizarea datelor brute în vederea controlului calitatii, există numeroase pachete software disponibile pentru realizarea acestui pas. Tot în cadrul acestui pas se vor identifica și înlătura coloanele cu valori nule ("NA"). O logaritmare în baza 2 a datelor poate fi realizată pentru a normaliza valorile de expresie pentru a urma o distribuție Gaussiană.

Obținerea profilului de expresie al microRNA. Pentru acest pas exista numeroase pachete comerciale si open-source disponibile, dupa cum au fost mentionate si anterior. Datorita naturii pachetelor open-source, acestea sunt mai des reinoite de catre comunitatile care le intretin fata de pachetele comerciale, astfel o strategie tipica pentru analiza bioinformatica a datelor obtinute cu tehnologii mature precum microarray este utilizarea unor pachete software open-source (ex: cele prezente in R Bioconductor) in construirea unui pipeline customizat.

Adnotarea miRNA. Tabelul cu profilul de expresie al miRNA necesita legarea ID-urilor fiecărei molecule determinate de numele miRNA specific platformei si bazei de date utilizata de echipament (ex: miRBase). Astfel, e necesar un pas de atasare al denumirii fiecarui miRNA la ID-urile din tabelul cu expresia fiecărei molecule. Dupa acest pas se obtine input-ul ce va fi introdus in expertul AI.

Analiza Exploratorie și Curatarea Datelor. Primul pas al workflow-ului Expertului AI este analiza exploratorie a datelor. Aceasta începe cu verificarea formei datelor în care se observa numărul de linii și coloane. Una din aceste coloane va deveni ținta pentru care va fi realizata predicția utilizand algoritmi AI. Datele de input și output trebuie sa aiba formatul corect - coloanele miRNA conțin numere reale, iar ținta - diagnosticul aici - este variabila categoriala.

Se vor identifica eventualele date lipsa și procentul acestora, eventualele anomalii prezente în date și se va proceda la o analiza statistica sumara (distributia datelor, valori minime, maxime, medii, etc.). Ulterior se vor verifica corelatiile dintre variabilele prezente, astfel se poate observa potențialul ca unele variabile sa fie mai importante în analiza și construcția modelelor predictive.

Preprocesare datelor inaintea de antrenarea modelelor:

Identificare valorilor lipsa și imputarea acestora cu metode avansate.

Identificarea de Outliers prin multiple metode și îndepărtarea acestora.

Identificare de eventuale denumiri multiple pentru aceeași coloana și corectare.

Logaritmare în baza 2 a datelor. Se considera ca distribuția datelor se apropie de cea Gausiana, prin aceasta transformare.

Normalizarea. Se alege una din următoarele metode: scalarea datelor între 0 și 1, sau între -1 si 1 sau, mai frecvent, standardizarea – media datelor devine 0 iar deviatia standard devine 1.

Dezvoltarea modelelor predictive

Acest pas începe cu inițializarea interfeței Expertului AI care va rula în cadrul unei instanțe de tip notebook (ex: Jupyter Notebooks). Prin utilizarea metodelor AutoML, modelele predictive vor fi dezvoltate automat. Totusi, vor fi cerute anumite setari din partea utilizatorului.

Se vor incarca librăriile și pachetele necesare mediului de lucru pentru Expertul AI.

Se vor defini perechile input - valorile de expresie a miRNA si output - coloana tinta (diagnostic) cu valorile "Cancer" si "Non-Cancer").

Se definesc seturile de antrenare si testare. Daca numărul cazurilor permite, se defineste si cel de al 3 lea set (de validare). Cel mai adesea, numărul cazurilor este relativ mic, motiv pentru care se definesc doar un set de antrenare și testare și se utilizează cross validation pentru antrenare. De exemplu, se împart cazurile de antrenare în 10 subseturi, învățarea are loc pe 9 dintre ele iar testarea pe al 10 lea, subseturile fiind permutate.

In acest exemplu, performanta clasificarii binare Cancer - Non-Cancer poate fi reprezentată de mai multor metrici precum Acuratețea, AUC (Area Under the Curve), Precizia, Kappa si MCC (coeficient de corelatie Matthews).

Se antrenează toate modelele disponibile in librariile incarcate si se evaluează performanța acestora, rezultand o ierarhie de modele.

Optimizarea parametrilor algoritmilor de modelare:

Orice algoritm de modelare are el însuși o serie de parametri (de exemplu, adancimea arborilor decizionali). Valorile obișnuite (default) ale acestora pot sa nu conducă la performanta maxima posibila, dar pot fi optimizate. În acest scop se utilizează subsetul de validare care, cand sunt putine cazuri, poate fi un subset al setului de antrenare. Daca performanța obținută este superioară se retin valorile corespunzătoare ale parametrilor algoritmului

Prin aplicarea invenției se obțin următoarele avantaje:

- Se descoperă biomarkeri cu adevărat relevanți, intrucat acestia sunt doar o componenta a unor modelelor predictive a caror performanta este masurata. Daca modelul bazat pe un set de biomarkeri are o acuratete de peste 95%, in mod clar markerii sunt relevanti. Daca, inasa, acuratețea este doar de 60% de exemplu, valoare biomarkerilor descoperiți este cvasi nula, intrucat alegerea clasei la întâmplare are o acuratețe de 50%, destul de apropiata.
- Se descoperă modele predictive de mare performanta, intrucat acestea extrag toata informația conținută în inputuri, relativ la outputuri, fiind și direct aplicabile la noi cazuri, spre deosebire de simplele liste.
- Îmbunătățiri semnificative ale eficienței însoțite de scăderea costurilor în domeniile de aplicabilitate.

1

Revendicari

Inventia permite descoperirea de biomarkeri și modele predictive cu o înaltă acuratețe. Tehnicile AI din inventie elimina variabilele ce nu diferă cantitativ între situațiile de discriminat păstrand doar variabilele informative pentru predicție, selectarea acestora se desfășoară utilizând un mecanism de feedback de la acuratețea predicției. În plus prin utilizarea abordărilor de tip AutoML se testează în mod automat numeroși algoritmi AI fapt care permite experimentarea și alegerea celui mai bun algoritm față de abordările de analiză statistică convențională.

1

21

Desene

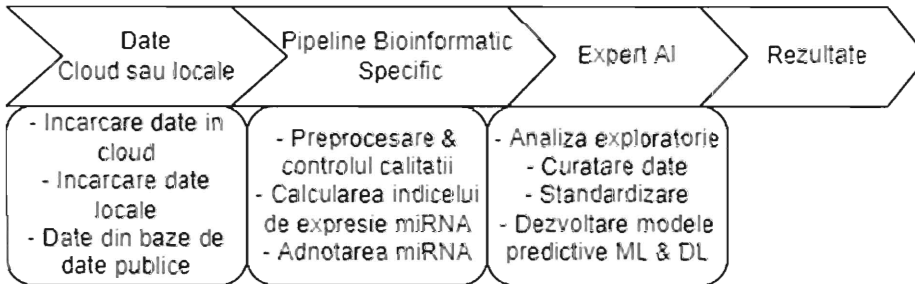


Figura 1. Sistem - vedere de ansamblu

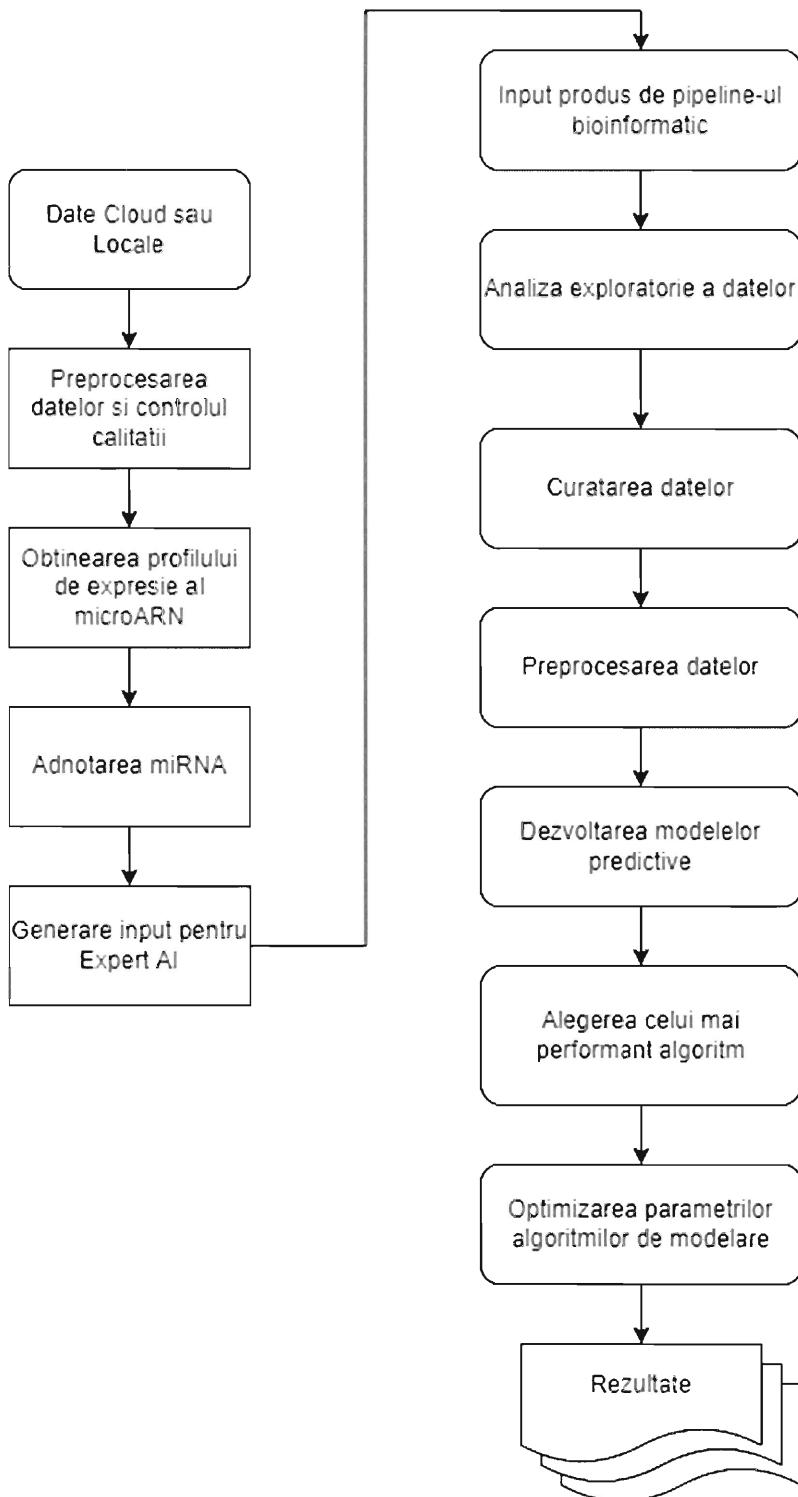


Figura 2. Workflow Bioinformatic specific datelor Microarray si Expert AI