

Modele și tehnici computaționale selectate pentru analiza datelor corelate cu rezistența antimicrobiană

efectuat în cadrul proiectului *Abordarea bioeconomică a
agenților antimicrobieni – utilizare și rezistență*

(cod - PN-III-P1-1.2-PCCDI-2017-0361).

Colectiv de redacție:

Raluca Mureșan, Claudia Zaharia, Radu Moleriu: metode de analiză a asocierilor (sect. 2)

Kristian Miok: modele de învățare automată în procesarea datelor medicale (sect. 3)

Daniela Zaharie: introducere, metode de clasificare cu etichete multiple (sect. 1)

Coordonator: Daniela Zaharie

Membri: Claudia Zaharia, Raluca Mureșan, Radu Moleriu, Kristian Miok

Data finalizării: 15.11.2019

Raport: R.5.1.2.

Versiunea: 1.1

Acknowledgements

Activities under this work were carried out in the *Research Laboratory Complex "Horia Cernescu"* - financed by project *"A bio-economical approach of the antimicrobial agents - use and resistance"*, in the frame of contract PCCDI 7/19.03.2018, code: PN-III P1-1.2-FPRD-2017.

1. Introducere

În general, seturile de date referitoare la bacteriile cu rezistență antimicrobiană pot fi analizate cu modele matematice continue (difuzie în vecinătățile sferice ale celulelor), modele ordinale (nivele izolate care au o concentrație minimă de inhibitori) sau modele dihotomice (de exemplu, folosind doar două categorii: bacterii susceptibile sau rezistente). Analiza asocierii în cazul bacteriilor cu rezistență antimicrobiană este dificilă datorită faptului că setul de date nu verifică ipotezele modelelor clasice din statistică. Modelele continue și ordinale nu au o distribuție normală, iar modelele dihotomice pot conține date rare și dispersate. Pe de altă parte, datele colectate se pot caracteriza prin absența unor valori ceea ce face analiza și mai dificilă.

În acest raport este prezentat contextul în care s-a desfășurat studiul dar și principalele rezultate obținute referitor la:

- Identificarea unor asocieri între tipuri de rezistență a bacteriilor la diferite antibiotice și analiza unor măsuri adecvate de ierarhizare a regulilor de asociere extrase din date. Analiza s-a realizat pe un set de date referitor la Escherichia Coli iar rezultatele obținute sunt în concordanță cu cele obținute folosind rețele de interacțiuni între proteine și un alt set de date referitor la E. Coli ([Cărunta et al., 2019](#)).
- Tratarea valorilor absente folosind ca tehnică de imputare mai multe modele generative de tip auto-encoder ([Miok et al., 2019a](#)) și analiza performanțelor acestora atât în cazul unor date de test tradiționale cât și în cazul unor date referitoare la caracteristicile producției de lapte (inclusiv numărul de celule somatice, interpretate ca potențial indicator al infecției) la Stațiunea de Cercetare-Dezvoltare în Creșterea Bovinelor de la Arad.

Rezultatele prezentate în acest raport au fost publicate în două lucrări (autorii subliniați sunt membri ai echipei de proiect):

- [Kristian Miok](#), Dong Nguyen-Doan, Marko Robnik-Sikonja and [Daniela Zaharie](#), Multiple Imputation for Medical Data using Neural Network-based Autoencoders Combined with Monte Carlo dropout, The 7th IEEE International Conference on E-Health and Bioengineering - EHB 2019, Iasi, Romania, November 21-23, 2019.
- [Claudia Zaharia](#), [Raluca Muresan](#), [Radu Moleriu](#), [Daniela Zaharie](#), Analysis of association measures used to discover antimicrobial resistance patterns, The 7th IEEE International Conference on E-Health and Bioengineering - EHB 2019, Iasi, Romania, November 21-23, 2019.

2. Metode de analiză a asocierilor

2.1. Concepte de bază

Analiza de asociere este o tehnică utilă pentru a descoperi relații interesante între obiecte în seturi de date de mari dimensiuni. Aceste relații se exprimă sub formă de mulțimi de elemente frecvente și reguli de asociere.

În analiza de asociere, datele sunt reprezentate sub forma unei liste de tranzacții. Dacă o mulțime de obiecte $O = \{o_1, o_2, \dots, o_k\}$, vom numi tranzacție o submulțime t a lui O de obiecte înregistrate în baza de date ca urmare a unei anumite acțiuni. Notăm $T = \{t_1, t_2, \dots, t_N\}$ mulțimea tuturor tranzacțiilor. O reprezentare uzuală pentru T este sub formă de matrice binară, cu N linii (corespunzătoare tranzacțiilor) și d coloane (corespunzătoare obiectelor). Elementul de pe poziția (i, j) din matrice va fi 1 dacă obiectul o_j se regăsește în tranzacția t_i , respectiv 0 în caz contrar.

O colecție de obiecte din O poartă numele de *itemset*. O tranzacție t conține un itemset X dacă $X \subseteq t$. Frecvența unui itemset se definește prin numărul tuturor tranzacțiilor care-l conțin:

$$\sigma(X) = \text{card}\{t_i \in T: X \subseteq t_i\}.$$

O regulă de asociere este o implicație de forma $A \rightarrow B$, unde $A, B \subset O$ și $A \cap B = \emptyset$. Mulțimile A și B se numesc antecedentul, respectiv consecventul regulii.

Puterea unei reguli de asociere se măsoară cu ajutorul suportului și încrederii (confidenței). Vom spune că regula $A \rightarrow B$ (și echivalent, itemsetul $A \cup B$) are suportul s dacă o proporție s a tranzacțiilor din T conțin simultan pe A și B . Regula $A \rightarrow B$ are încrederea c dacă B din totalul tranzacțiilor ce îl conțin pe A , proporția celor ce îl conțin și pe B este c . Formal,

$$\text{supp}(A \rightarrow B) = \frac{\sigma(A \cup B)}{\text{card}(T)} = P(A \cap B)$$

și

$$\text{conf}(A \rightarrow B) = \frac{\sigma(A \cup B)}{\sigma(A)} = P(B|A).$$

Dacă o regulă de asociere are suportul mic, atunci este posibil ca ea să fie o coincidență. Aparițiile ei sunt rare, și deci regula poate fi considerată ne semnificativă în aplicații practice. În ceea ce privește încrederea, cu cât aceasta este mai mare, cu atât este mai probabil ca o tranzacție care îl conține pe A să îl conțină și pe B .

Dat un set de tranzacții T , scopul analizei de asociere este de se genera toate regulile de asociere cu suport și încredere mai mari decât anumite valori prag *minsup*, respectiv *minconf*. Întrucât numărul total de reguli ce pot fi formulate pentru un set de k obiecte este extrem de mare ($3^k - 2^{k+1} + 1$), analiza exhaustivă a acestora este prohibitivă din punct de vedere computațional, existând tehnici mult mai avantajoase de filtrare a spațiului de căutare pentru determinarea regulilor ce satisfac cerințele de suport și încredere.

În contextul analizei rezistenței antimicrobiene, în special în contextul rezistenței multiple, este important de identificat antibiotice sau clase de antibiotice care sunt detectate frecvent împreună în seturile de date care cuprind informații privind rezistența sau susceptibilitatea diferitelor tulpini de bacterii la diferite antibiotice.

2.2. Extragerea asocierilor din date

Demersul de extragere a asocierilor din date presupune două etape:

- Generarea tuturor mulțimilor de obiecte (itemset-uri) care au suportul mai mare decât un prag specificat *minsup*. Acestea se vor numi itemset-uri frecvente.
- Generarea tuturor regulilor de încredere mai mare ca *minconf* ce implică itemset-urile frecvente găsite la pasul anterior.

2.2.1. Generarea itemset-urilor frecvente

Un set de date ce conține k obiecte poate genera $2^k - 1$ itemset-uri, excluzând mulțimea vidă. Datorită faptului că în multe aplicații practice numărul de obiecte k implicate în tranzacții este foarte mare, numărul de itemset-uri ce trebuie analizate este potențial enorm.

Se poate defini o structură de tip latice pentru determinarea tuturor submulțimilor unei mulțimi de itemi, un exemplu fiind prezentat în figura de mai jos pentru mulțimea $\{a, b, c, d, e\}$.

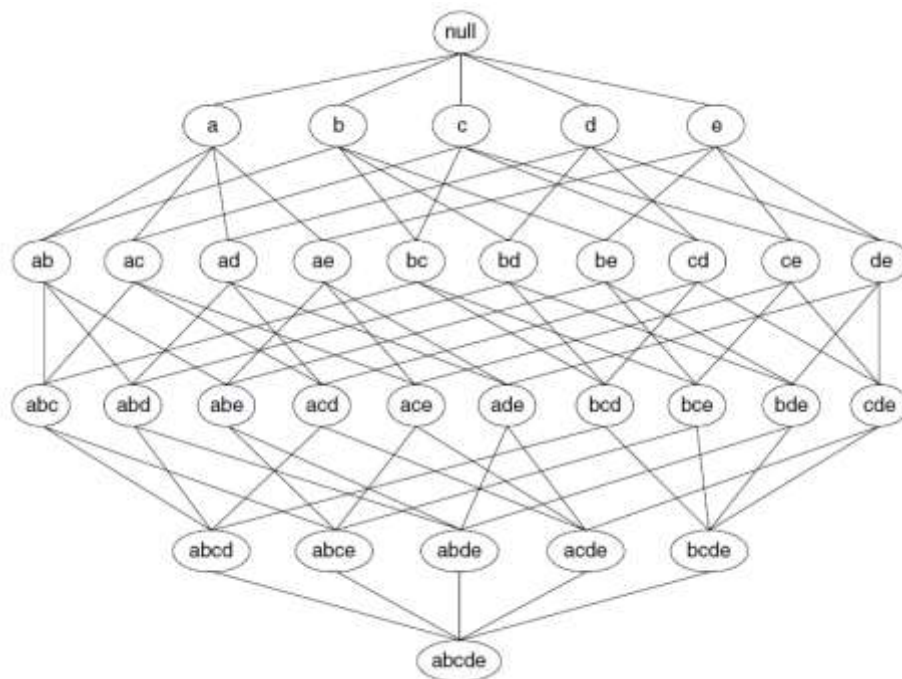


Figura 2.1: Structura de latice a submulțimilor unei mulțimi $\{a, b, c, d, e\}$

O abordare simplistă pentru determinarea itemset-urilor frecvente este calcularea suportului pentru fiecare mulțime candidată din structura latice, comparând-o succesiv cu fiecare tranzacție. Această metodă însă poate fi foarte costisitoare datorită numărului extrem de mare de comparații implicate. O abordare mult mai eficientă folosește algoritmul Apriori (Agrawal et al., 1993, Agrawal & Srikant, 1994).

2.2.2. Algoritmul Apriori

Principiul Apriori. Dacă un itemset este frecvent, atunci toate submulțimile sale sunt de asemenea frecvente.

Invers, dacă itemset-ul nu este frecvent, atunci niciunul dintre itemset-urile care îl conțin nu poate fi frecvent. Acest fapt conduce la o strategie numită filtrare bazată pe suport, care folosește proprietatea de antimonotonie a măsurii suport, anume că suportul unui itemset nu depășește suportul submulțimilor sale.

Inițial, fiecare obiect este considerat itemset candidat de dimensiune 1. Dintre acestea, după determinarea suportului fiecăruia, sunt reținute numai itemset-urile de dimensiune 1 frecvente. La iterația următoare, din itemset-urile frecvente de dimensiune 1 sunt alcătuite itemset-uri candidate de dimensiune 2 (niciun itemset frecvent de dimensiune 2 nu poate conține un itemset ne-frecvent), acestora li se determină suportul, și sunt păstrate doar itemset-urile de dimensiune 2 frecvente. Similar, la fiecare pas sunt construite itemset-uri candidate de dimensiune k utilizând itemset-uri frecvente de dimensiune $k - 1$, apoi dintre acestea sunt reținute ca frecvente acelea care satisfac restricția legată de suport. Procedeeul se încheie atunci când pentru o anumită dimensiune k nu mai pot fi generate itemset-uri frecvente.

Complexitatea computațională a algoritmului Apriori depinde de următorii factori:

1. pragul de suport – un prag de suport mai scăzut are ca efect declararea unui număr mai mare de itemset-uri ca frecvente, ceea ce implică necesitatea unui număr mai mare de parcurgeri ale bazei de date.
2. numărul de obiecte din baza de date
3. numărul de tranzacții
4. numărul mediu de obiecte dintr-o tranzacție

2.2.3. Generarea regulilor de asociere

Fiecare itemset frecvent Y de dimensiune k poate produce până la $2^k - 2$ reguli de asociere, fără a include regulile care au antecedent sau consecvent vid ($\emptyset \rightarrow Y$ sau $Y \rightarrow \emptyset$). O regulă de asociere poate fi extrasă prin partiționarea itemset-ului Y în două submulțimi nevide X și $Y - X$ astfel încât $X \rightarrow Y - X$ să satisfacă pragul de încredere. De menționat că toate aceste reguli satisfac pragul pentru măsura suport deoarece sunt generate dintr-un itemset frecvent.

Exemplu: Fie $Y = \{a, b, c\}$ un itemset frecvent. Există șase reguli de asociere candidate care pot fi generate din Y : $\{a, b\} \rightarrow \{c\}$, $\{a, c\} \rightarrow \{b\}$, $\{b, c\} \rightarrow \{a\}$, $\{a\} \rightarrow \{b, c\}$, $\{b\} \rightarrow \{a, c\}$, $\{c\} \rightarrow \{a, b\}$. Suportul fiecărei reguli este egal cu suportul mulțimii Y .

Calcularea încrederii pentru fiecare regulă de asociere nu necesită alte evaluări ale mulțimii tranzacțiilor. Considerând regula $\{a, b\} \rightarrow \{c\}$, încrederea se calculează ca $suport(\{a, b, c\})/suport(\{a, b\})$. Din moment ce suportul pentru fiecare itemset frecvent a fost deja calculat în timpul generării itemset-urilor frecvente, nu mai este necesar să se parcurgă din nou întregul set de date.

Spre deosebire de măsura suport, încrederea nu are nicio proprietate de monotonie. Totuși are loc următorul rezultat, valabil pentru compararea regulilor generate din același itemset frecvent Y :

Dacă o regulă $X \rightarrow Y - X$ nu satisface pragul de încredere, atunci nicio altă regulă $X' \rightarrow Y - X'$, unde X' este o submulțime a lui X , nu va satisface pragul de încredere.

Algoritmul Apriori are o abordare pe nivele pentru generarea regulilor de asociere, unde fiecare nivel corespunde numărului de itemi din consecvent. Inițial toate regulile cu încredere mare care au doar un item în consecvent sunt extrase. Acestea sunt apoi utilizate pentru a genera alte reguli candidate. De exemplu, dacă $\{a, c, d\} \rightarrow \{b\}$ și $\{a, b, d\} \rightarrow \{c\}$ sunt reguli cu încredere mare, atunci regula candidată $\{a, d\} \rightarrow \{b, c\}$ este generată prin reunirea mulțimilor consecvente ale celor două reguli. Figura 2.2 de mai jos ilustrează o structură de tip latice pentru regulile de asociere generate din itemset-ul frecvent $\{a, b, c, d\}$. Dacă un anumit nod din latice are o încredere mică, atunci conform teoremei de mai sus întregul subgraf generat de acest nod poate fi eliminat. Dacă, spre exemplu, valoarea încrederii pentru $\{b, c, d\} \rightarrow \{a\}$ este mică, atunci toate regulile ce conțin item-ul a în consecvent, incluzând $\{c, d\} \rightarrow \{a, b\}$, $\{b, d\} \rightarrow \{a, c\}$, $\{b, c\} \rightarrow \{a, d\}$, $\{d\} \rightarrow \{a, b, c\}$ pot fi eliminate.

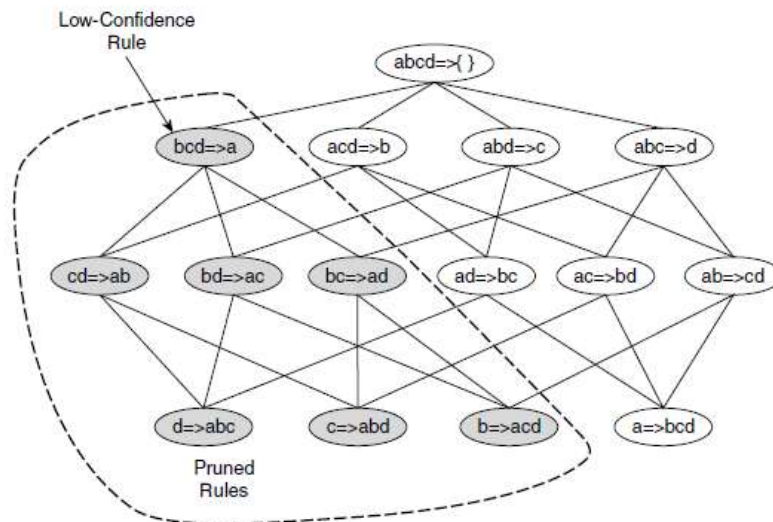


Figura 2.2: Filtrarea regulilor de asociere pe baza încrederii (Tan et al, 2014)

2.3. Măsuri pentru evaluarea calității asocierii

O problemă importantă în analiza regulilor de asociere o constituie faptul că, în general, numărul de reguli ce pot fi extrase este foarte mare, așa încât poate deveni dificil pentru utilizator să le identifice pe cele mai relevante pentru aplicația sa. Din acest motiv au fost introduse o serie de măsuri care să ierarhizeze regulile în funcție de interesul potențial pe care l-ar putea reprezenta pentru utilizator. Aceste măsuri se împart în trei categorii principale (Geng & Hamilton, 2006):

- *măsuri obiective* - bazate strict pe date, provin din domenii precum teoria probabilităților, statistică, teoria informației, învățare automată;
- *măsuri subiective* - depind de date, dar încorporează și cunoștințele sau așteptările utilizatorului. Astfel de măsuri cuantifică, spre exemplu, gradul de noutate al unei reguli pentru utilizator sau gradul de neașteptat al regulii, discordanța față de ceea ce este cunoscut sau anticipat în domeniu;
- *măsuri semantice* - au în vedere contextul și interpretarea regulilor, utilitatea acestora pentru realizarea unui anumit obiectiv.

Modele și tehnici computaționale selectate pentru analiza datelor

În cele ce urmează ne vom concentra asupra măsurilor de interes (MI) de tip obiectiv. În literatură există un număr semnificativ de lucrări științifice dedicate analizei acestui tip de măsuri, a proprietăților lor și a relațiilor dintre acestea - a se vedea spre exemplu (Piatetsky-Shapiro, 1991), (Bayardo & Agrawal, 1999), (Tan et al., 2004), (Vaillant et al., 2004), (Geng & Hamilton, 2006), (Weiß, 2008) sau (Martínez-Ballesteros et al., 2014).

În general, o MI obiectivă pentru o regulă de tipul $A \rightarrow B$ se definește cu ajutorul tabelului de contingență al regulii:

	B	\bar{B}	
A	$n(AB)$	$n(A\bar{B})$	$n(A)$
\bar{A}	$n(\bar{A}B)$	$n(\bar{A}\bar{B})$	$n(\bar{A})$
	$n(B)$	$n(\bar{B})$	N

Tabelul 2.1: Tabelul de contingență al regulii $A \rightarrow B$

Acest tabel conține frecvențele absolute ale tranzacțiilor din baza de date care satisfac anumite condiții cu privire la elementele din A și B , respectiv:

- $n(AB)$ reprezintă numărul de tranzacții care conțin simultan pe A și B ;
- $n(A\bar{B})$ este numărul de tranzacții care conțin pe A , dar nu și pe B ;
- $n(\bar{A}B)$ este numărul de tranzacții care conțin pe B , dar nu și pe A ;
- $n(\bar{A}\bar{B})$ este numărul tranzacțiilor din care lipsesc simultan A și B ;
- $n(A)$ și $n(\bar{A})$ reprezintă numărul de tranzacții care conțin, respectiv nu conțin pe A (idem pentru B)
- N este numărul total de tranzacții din baza de date.

Pe baza acestor valori se definesc probabilitățile corespunzătoare, e.g. $P(A) = \frac{n(A)}{N}$, $P(A \cap B) = \frac{n(AB)}{N}$, $P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{n(AB)}{n(A)}$.

Primele MI obiective utilizate pentru selectarea și ierarhizarea regulilor în analiza de asociere au fost suportul și încrederea, descrise anterior. Alte câteva MI folosite uzual sunt prezentate pe scurt în continuare.

1. *Factorul de interes (liftul)* este definit prin

$$Lift(A \rightarrow B) = \frac{P(A \cap B)}{P(A) \cdot P(B)}$$

Acesta cuantifică în ce măsură apar mai frecvent în tranzacții A și B decât ar fi de așteptat sub ipoteza de independență statistică. Este cunoscut faptul că, dacă A și B sunt independente, atunci $P(A \cap B) = P(A) \cdot P(B)$, prin urmare o valoare mai mare ca 1 a liftului indică o asociere pozitivă între A și B .

2. *Coeficientul de corelație ϕ* este dat de

$$\phi(A \rightarrow B) = \frac{P(A \cap B) - P(A) \cdot P(B)}{\sqrt{P(A)P(B)(1 - P(A))(1 - P(B))}}$$

Acesta provine din statistică și este analogul coeficientului de corelație Pearson folosit în studiul asocierii variabilelor continue. El are valori între -1 (asociere negativă perfectă) și 1 (asociere pozitivă perfectă), valoarea 0 indicând independența statistică între A și B .

Modele și tehnici computaționale selectate pentru analiza datelor

3. Măsura cosinus se definește ca

$$\cos(A \rightarrow B) = \frac{P(A \cap B)}{\sqrt{P(A)P(B)}}$$

și reprezintă cosinusul unghiului între vectorii de incidență a lui A , respectiv B , în tranzațiile din baza de date. O valoare apropiată de 0 indică lipsa de asociere, în timp ce valorile apropiate de 1 indică o asociere puternică între A și B .

4. Coeficientul Jaccard

$$\zeta(A \rightarrow B) = \frac{P(A \cap B)}{P(A \cup B)}$$

reprezintă proporția tranzațiilor care conțin simultan pe A și B din numărul aceluia care conțin cel puțin una dintre A și B . Domeniul de valori al coeficientului Jaccard este $[0,1]$, unde 0 indică lipsa de asociere, iar 1, asociere pozitivă perfectă.

5. Raportul de șanse ("odds ratio")

$$\alpha(A \rightarrow B) = \frac{P(A \cap B)P(\bar{A} \cap \bar{B})}{P(\bar{A} \cap B)P(A \cap \bar{B})}$$

este o măsură de asociere provenită din statistică. Interpretarea acestuia este următoarea: dacă B este prezent într-o tranzație, șansele să fie prezent și A sunt $\frac{P(A \cap B)}{P(\bar{A} \cap B)}$, în timp ce dacă B este absent, șansele ca A să fie prezent în tranzația respectivă sunt date de $\frac{P(A \cap \bar{B})}{P(\bar{A} \cap \bar{B})}$. Dacă raportul acestor două cantități este supraunitar, prezența lui B sporește șansele ca A să fie prezent, deci A și B sunt asociate pozitiv. Dacă raportul de șanse este subunitar, asocierea este negativă. În sfârșit, dacă raportul de șanse este 1, A și B nu sunt asociate.

Coeficienții Y și Q ai lui Yule sunt variante normalizate ale raportului de șanse, definite respectiv prin

$$Q = \frac{\alpha - 1}{\alpha + 1}, Y = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$$

O serie de alte MI obiective folosite în analiza de asociere sunt incluse în Tabelul 2.1 de mai jos. Acestea sunt descrise pe larg în (Hahsler, 2015).

Măsura de interes	Definiția
Coeficientul Goodman-Kruskal (λ)	$\frac{\sum_j \max_k P(A_j \cap B_k) + \sum_j P(A_j \cap B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
Cohen κ	$\frac{P(A \cap B) + P(\bar{A} \cap \bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
Informație mutuală	$\frac{\sum_i \sum_j P(A_i \cap B_j) \log \frac{P(A_i \cap B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$

Modele și tehnici computaționale selectate pentru analiza datelor

Măsura J	$\max(P(A \cap B) \log \frac{P(B A)}{P(B)} + P(A \cap \bar{B}) \log \frac{P(\bar{B} A)}{P(\bar{B})}, P(A \cap B) \log \frac{P(A B)}{P(A)} + P(\bar{A} \cap B) \log \frac{P(\bar{A} B)}{P(\bar{A})})$
Indicele Gini	$\max(P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] - P(B)^2 - P(\bar{B})^2, P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] - P(A)^2 - P(\bar{A})^2)$
Laplace	$\max\left(\frac{NP(A \cap B) + 1}{NP(A) + 2}, \frac{NP(A \cap \bar{B}) + 1}{NP(B) + 2}\right)$
Convingere	$\max\left(\frac{P(A)P(\bar{B})}{P(A\bar{B})}, \frac{P(B)P(\bar{A})}{P(B\bar{A})}\right)$
Piatetsky-Shapiro	$\frac{P(A \cap B) - P(A)P(B)}{1 - P(A)P(B)}$
Factor de certitudine	$\max\left(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)}\right)$
Valoare adăugată	$\max(P(B A) - P(B), P(A B) - P(A))$
Putere colectivă	$\frac{P(A \cap B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \cdot \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A \cap B) - P(\bar{A}\bar{B})}$
Klosgen	$\sqrt{P(A \cap B)} \max(P(B A) - P(B), P(A B) - P(A))$

Tabelul 2.2: Măsuri obiective pentru reguli de asociere

În funcție de proprietățile sale, o MI poate fi adecvată pentru o anumită clasă de aplicații și nepotrivită pentru altele. Principalele proprietăți teoretice pe care le poate avea o MI obiectivă m , propuse și studiate în (Piatetski – Shapiro, 1991) și (Tan et al., 2004), sunt enumerate în continuare.

- (P1) $m(A \rightarrow B) = 0$ atunci când A și B sunt independente
- (P2) m crește cu $P(A \cap B)$ atunci când $P(A)$ și $P(B)$ rămân constante
- (P3) m scade cu $P(A)$ (respectiv $P(B)$) atunci când $P(A \cap B)$ și $P(B)$ (respectiv $P(A)$) rămân constante
- (P4) $m(A \rightarrow B) = m(B \rightarrow A)$ (simetrie la permutarea variabilelor)
- (P5) m este invariantă la scalarea liniilor sau coloanelor tabelului său de contingență cu un factor pozitiv
- (P6) $m(A \rightarrow B) = -m(A \rightarrow \bar{B}) = -m(\bar{A} \rightarrow B)$
- (P7) $m(A \rightarrow B) = m(\bar{A} \rightarrow \bar{B})$
- (P8) valoarea lui $m(A \rightarrow B)$ nu este influențată de tranzacțiile care nu le conțin pe A și B

Măsura de interes	P1	P2	P3	P4	P5	P6	P7	P8
Φ	da	da	da	da	nu	da	da	nu
Goodman-Kruskal	da	nu	nu	da	nu	nu*	da	nu
Raport de șanse	da*	da	da	da	da	da*	da	nu
Yule Q	da	da	da	da	da	da	da	nu

Modele și tehnici computaționale selectate pentru analiza datelor

Yule Y	da	da	da	da	da	da	da	nu
Cohen κ	da	da	da	da	nu	nu	da	nu
Informație mutuală	da	da	da	nu	nu	nu*	da	nu
Măsura J	da	nu	nu	nu	nu	nu	nu	nu
Indicele Gini	da	nu	nu	nu	nu	nu*	da	nu
Suport	nu	da	nu	da	nu	nu	nu	nu
Încredere	nu	da	nu	nu	nu	nu	nu	da
Laplace	nu	da	nu	nu	nu	nu	nu	nu
Convingere	nu	da	nu	nu	nu	nu	da	nu
Lift	da*	da	da	da	nu	nu	nu	nu
Cosinus	nu	da	da	da	nu	nu	nu	da
Piatetsky-Shapiro	da	da	da	da	da	da	da	nu
Factor de certitudine	da	da	da	nu	nu	nu	da	nu
Valoare adăugată	da	da	da	nu	nu	nu	nu	nu
Putere colectivă	nu	da	da	da	nu	da*	da	nu
Jaccard	nu	da	da	da	nu	nu	nu	da
Klosgen	da	da	da	nu	nu	nu	nu	nu

da*: da, dacă măsura este normalizată

nu*: simetrie la permutarea de linii sau coloane

Tabelul 2.3: Proprietățile principalelor măsuri de interes obiective (cf. Tan et al., 2004)

Dat fiind faptul că nu toate MI existente satisfac toate aceste proprietăți (a se vedea de exemplu Tabelul 2.3 de mai sus), este important de stabilit care dintre (P1) – (P8) sunt relevante în contextul analizei de asociere a rezistenței antimicrobiene, în care variabilele sunt dihotomice (prezență/absență a rezistenței la un anumit antibiotic) și asimetrice (prezența rezistenței este considerată în investigație a fi mai importantă decât absența).

Proprietatea (P1) afirmă că o asociere care apare strict aleator ar trebui să aibă interes 0. În practică această cerință este prea restrictivă. Ea poate fi relaxată impunând o valoare constantă pentru situațiile de independență (spre exemplu, în cazul liftului și al raportului de șanse, această valoare este 1).

Proprietatea (P2) spune că, presupunând suporturile lui A și B fixate, regula va fi cu atât mai interesantă cu cât cele două apar mai frecvent împreună.

Proprietatea (P3) spune că dacă suportul comun al lui A și B , precum și cel al uneia dintre acestea, sunt fixate, regula va fi cu atât mai interesantă cu cât cealaltă mulțime are

Modele și tehnici computaționale selectate pentru analiza datelor

suport mai mic (i.e., există mai puține tranzacții care o conțin doar pe ea). Atât (P2), cât și (P3) sunt proprietăți de dorit pentru o MI utilizată în analiza de asociere a rezistenței la antibiotice.

În ceea ce privește (P4), în contextul analizei rezistenței antimicrobiene, existența relației de asociere între A și B se consideră a fi mai importantă decât direcția implicației între antecedent și consecvent, așadar putem considera că proprietatea nu este esențială în alegerea unei MI adecvate.

Proprietatea (P5) este utilă în general în studiul asocierilor pentru variabile nominale. Proprietatea (P6) spune că m poate identifica atât asocierile pozitive, cât și pe cele negative, și este un caz particular al lui (P7). În (Tan et al., 2004) se arată că măsurile cu proprietatea (P7), numite simetrice, nu sunt potrivite pentru aplicații ce necesită un tratament diferențiat pentru prezență vs. absență – prin urmare, MI folosite în analiza de asociere a rezistenței la antibiotice trebuie să fie asimetrice.

Măsurile de interes cu proprietatea (P8) țin seama doar de tranzacțiile în care se regăsesc cel puțin una dintre A și B . Aceasta implică faptul că valoarea lui $m(A \rightarrow B)$ nu se modifică atunci când în baza de date se adaugă tranzacții care nu implică pe A sau B (un exemplu de astfel de MI este încrederea). Acest lucru poate fi util când se analizează asocieri între antibiotice la care rezistența individuală este foarte rară – astfel de asocieri, deși potențial relevante pentru cercetător, nu vor putea fi detectate dacă în baza de date există prea multe probe simultan sensibile la respectivele antibiotice.

2.3.1. Consistența diverselor măsuri de interes pentru regulile de asociere

Este de notat faptul că diferite MI pot furniza ierarhizări diferite, uneori contradictorii, ale unui set de reguli, o regulă putând fi considerată interesantă după un anumit criteriu dar nu și după un altul. Prin urmare, alegerea celor mai potrivite MI este esențială pentru descoperirea celor mai relevante asocieri într-o aplicație dată.

O pereche de MI se consideră ca fiind consistente dacă ele ierarhizează similar un set de reguli. Gradul de asemănare între ordonările regulilor date de două MI poate fi evaluat, spre exemplu, folosind coeficienții de corelație Pearson sau Spearman, măsura cosinus sau inversa normei L2. În (Tan et al., 2004) se arată faptul că metodele enumerate mai sus furnizează rezultate similare în aplicații, și mai mult, dacă valorile rangurilor regulilor din fiecare ordonare sunt unice, toți acești coeficienți sunt monoton dependenți unul de celălalt.

2.3.2. Decizii statistice în alegerea valorilor prag utilizate în filtrarea regulilor

În analiza regulilor de asociere trebuie ținut seama de faptul că nu întotdeauna regulile descoperite reflectă asocieri reale, existând posibilitatea ca unele așa-zise asocieri să fie descoperiri false sau coincidențe. Extragerea regulilor de asociere se face pe baza unui set de tranzacții care reprezintă doar un eșantion din populația de interes, așa încât este de așteptat ca relațiile descoperite pe acest eșantion să nu corespundă întocmai celor prezente la nivelul întregii populații. Se pune deci problema unei abordări statistice pentru a identifica asocieri reale, care reflectă proprietățile existente la nivel de populație, și care vor fi valabile și pentru instanțe viitoare (a se vedea (Hämäläinen & Nykänen, 2008), (Riondato & Vandin, 2014), (Hämäläinen & Webb, 2019)).

Rata descoperirilor false poate fi minimizată prin filtrarea setului de reguli după o anumită MI, alegând un prag convenabil. După cum s-a văzut anterior, filtrarea în algoritmul Apriori se bazează pe măsura suport. În (Tan et al., 2004) se arată faptul că filtrarea regulilor

Modele și tehnici computaționale selectate pentru analiza datelor

prin impunerea unui prag inferior de suport are ca și consecință eliminarea cu precădere a asocierilor foarte slabe sau negative. Acesta este un aspect benefic în analiza de asociere privind rezistența la antibiotice, întrucât în acest caz de interes sunt asocierile pozitive.

Pe baza considerațiilor din (Megiddo & Srikant, 1998), în lucrarea (Cazer et al., 2019) se propune o metodologie de simulare pentru estimarea ratei de descoperiri false corespunzătoare unei anumite valori prag a măsurii de interes utilizate în filtrarea regulilor. Pornind de la această idee, în lucrarea (Zaharia et al., 2019) se abordează problema complementară, aceea de a determina pragul inferior pentru o anumită MI astfel încât filtrarea să asigure o rată de descoperiri false mai mică decât o valoare impusă (e.g. 5%).

Procedeeul se bazează pe generarea unui număr mare de seturi de date artificiale ("nule") conținând același număr de tranzacții ca și setul de date original, în care apariția fiecărui item într-o tranzacție este simulată ca o variabilă aleatoare Bernoulli de parametru p =frecvența relativă de apariție a itemului în setul de date original, independent de ceilalți itemi. Prin urmare, într-un set de date nul, itemii vor apărea în tranzacții cu aceeași probabilitate ca în datele real observate, dar fără a fi asociați între ei. Ulterior pe fiecare set de date se aplică algoritmul Apriori generându-se toate regulile de asociere posibile și se determină cuantila de ordin 0.95 a seriei de valori a MI. Ideea este aceea că, folosind valoarea cuantilei de ordin 0.95 ca prag inferior în filtrarea după acea MI, cel mult 5% din asocierile rezultate strict aleator (coincidențe) pot fi descoperite ca reguli. Estimarea pragului inferior de filtrare în setul de date original care asigură rata de descoperiri false dorită (sub 5%) se obține apoi ca media cuantilelor de ordin 0.95 calculate pentru toate seturile de date nule.

2.4. Studiu de caz în analiza rezistenței antimicrobiene

Rezultatele prezentate în această secțiune se bazează pe lucrarea (Zaharia et al., 2019).

2.4.1. Descrierea setului de date

Pentru exemplificarea tehnicilor descrise mai sus s-a considerat un set de date conținând specimene de *Escherichia Coli* provenite de la subiecți umani și informații privind susceptibilitatea sau rezistența acestora la mai multe antibiotice. Datele au fost extrase din *National Antimicrobial Resistance Monitoring System for Enteric Bacteria (NARMS)* pentru perioada 1996-2016 (NARMS, 2019). Această bază de date a fost înființată în 1996 în Statele Unite ale Americii ca un sistem național de monitorizare a rezistenței antimicrobiene pentru bacterii precum *Salmonella enterica*, *Campylobacter* și *Escherichia coli* regăsite la subiecți umani, animale domestice, carne și produse din carne.

Setul de date extras din NARMS conține 3583 de specimene, fiecare fiind clasificat ca rezistent sau susceptibil, pe baza pragurilor de concentrație minimă inhibitoare, la următoarele antibiotice: *ampicilină* (AMP), *amoxicilină – acid clavulanic* (AUG), *ceftriaxonă* (AXO), *cloramfenicol* (CHL), *ciprofloxacina* (CIP), *acid nalidixic* (NAL), *gentamicină* (GEN), *streptomycină* (STR), *trimetoprim – sulfametoxazol* (COT) și *tetraciclină* (TET).

Din cele 3583 de specimene considerate inițial, 1 (0.03%) a fost rezistent la 8 din cele 10 antibiotice, 13 (0.36%) la cel puțin 6 antibiotice, 19 (0.53%) la cel puțin 5, 49 (1.37%) la cel puțin 4, 97 (2.71%) la cel puțin 3, 164 (4.58%) la cel puțin 2 și 329 (9.18%) la cel puțin un antibiotic. Prevalența rezistenței la fiecare antibiotic poate fi observată în Tabelul 2.3 de mai jos.

Modele și tehnici computaționale selectate pentru analiza datelor

	AMP	AUG	AXO	CHL	CIP	NAL	GEN	STR	COT	TET
N	106	14	13	60	10	60	15	147	48	200
%	2.96	0.39	0.36	1.67	0.28	1.67	0.42	4.10	1.34	5.58

Tabelul 2.4: Prevalența rezistenței la cele 10 antibiotice în setul de date inițial (număr total și procent)

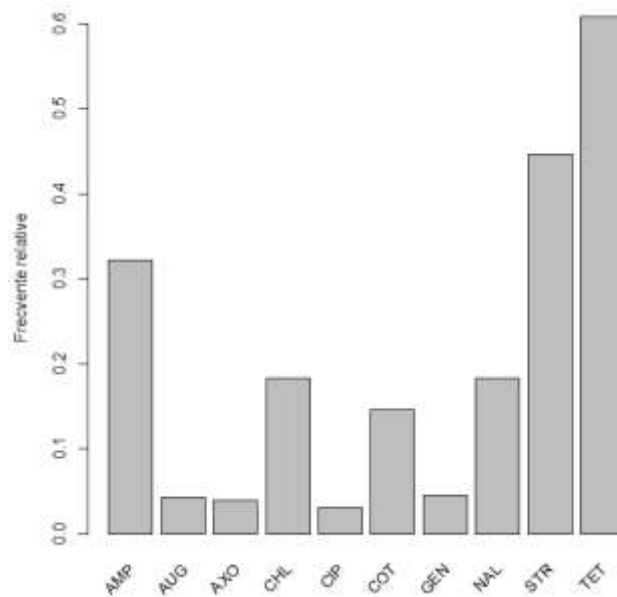


Figura 2.2: Prevalența rezistenței la cele 10 antibiotice în setul de date utilizat în analiza de asociere

Cele 329 de specimene găsite ca fiind rezistente la cel puțin un antibiotic au fost folosite în continuare pentru analiza asocierilor. Extragerea regulilor de asociere s-a realizat utilizând pachetul R “arules” (Hahsler et al., 2019). Pragurile de suport și încredere au fost stabilite ca $1/(\text{număr de specimene})$ pentru a genera toate regulile posibile de lungime cel puțin 2. S-au obținut în total 1329 de reguli de asociere. Acestea au fost evaluate folosind cele 46 de MI implementate în pachetul “arules” (a se vedea (Hahsler et al., 2019) pentru mai multe detalii). O primă selecție a fost efectuată, eliminând următoarele MI: frecvența absolută (coliniară cu suportul, deci redundantă), χ^2 (o variantă scalată a lui Φ^2), și acele MI care nu au putut fi calculate pentru toate regulile (convingere, raport de șanse, informație mutuală, măsura J, Sebag, coeficienții Yule Q și Y) sau care au obținut valoarea ∞ (coeficientul de îmbunătățire).

A fost evaluată apoi consistența MI rămase, calculând coeficientul de corelație Pearson pentru vectorii de ranguri ale regulilor generate din setul de date. Distanța dintre două MI a fost definită ca $1 - |r|$, unde r este coeficientul de corelație Pearson corespunzător. Distanțele dintre măsuri au avut valori între 0.001 și 0.999. Se poate observa deci că există perechi de măsuri care ierarhizează similar regulile, cât și perechi care oferă rezultate discordante. Aplicând, apoi, un algoritm de clustering ierarhic cu metoda “average linkage” au fost identificate cinci grupuri de măsuri cu un comportament similar (a se vedea dendrograma din Figura 2.3).

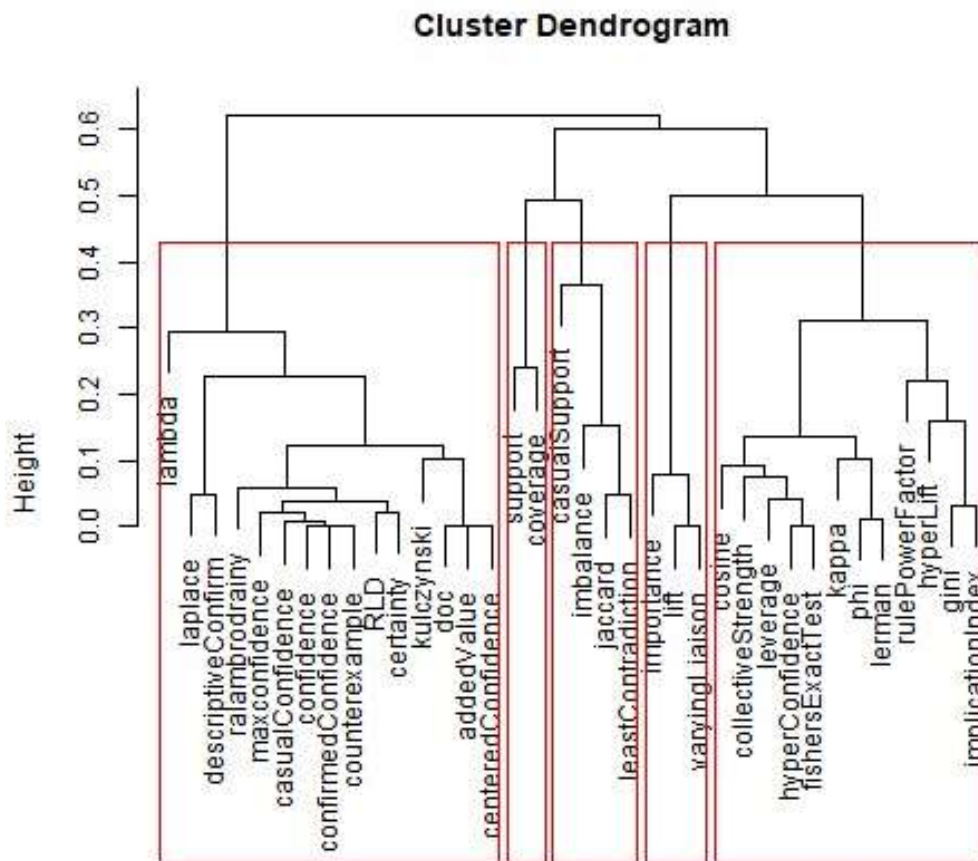


Figura 2.3: Clustering ierarhic al măsurilor de interes

Fiecare cluster conține măsuri care sunt consistente din punctul de vedere al rangurilor atribuite regulilor extrase din setul de date AMR. Astfel, este suficient să selectăm o măsură reprezentativă din fiecare cluster pentru a determina un clasament agregat al regulilor de asociere.

Conform cu aspectele discutate în secțiunea 2.3, din fiecare cluster a fost selectată o MI reprezentativă, cu proprietăți adecvate (care să nu satisfacă P7, dar să satisfacă P2 și P3, asigurându-ne că includem atât MI cu proprietatea P8, cât și MI care nu satisfac această proprietate). Măsurile selectate sunt: *încredere*, *support*, *coeficientul Jaccard*, *lift* și *cosinus*.

În etapa următoare de analiză, a fost determinat pragul minim de suport care să asigure o rată de descoperiri false sub 5%, urmând procedura detaliată în secțiunea 2.3. Au fost considerate 10000 seturi de date nule de câte 329 de tranzacții, în care fiecare rezistență la un anumit antibiotic a fost simulată independent ca variabilă aleatoare Bernoulli de parametru egal cu prevalența rezistenței la antibioticul respectiv în cele 329 tranzacții inițiale, și a fost determinată cuantila de ordin 0.95 a suportului pentru fiecare din cele 10000 seturi de reguli obținute.

Distribuția valorilor acestor cuantile ale suportului poate fi vizualizată în Figura 2.4.

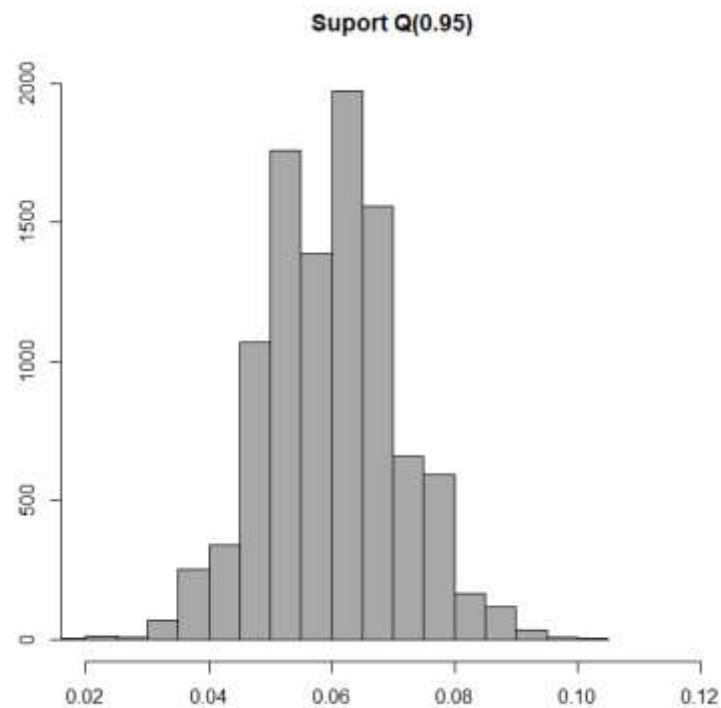


Figura 2.4: Histograma cuantilelor de ordin 0.95 pentru suport în cele 10000 seturi de date simulate

Pragul minim al suportului ce asigură o rată de descoperiri false mai mică decât 5% este 0.06, cu o abatere standard de 0.01. Filtrând setul inițial de reguli folosind această valoare a suportului, numărul de reguli relevante a fost redus la 42. Acestor reguli le-au fost atribuite ranguri folosind fiecare dintre cele cinci măsuri reprezentative menționate anterior. Pentru fiecare regulă s-a obținut un rang agregat, calculând media celor cinci ranguri corespunzătoare acestor măsuri.

Regulile ordonate după valoarea rangului agregat (R) sunt prezentate în Tabelul 2.5 de mai jos. Regula cu cel mai mare rang agregat conține antibioticele CHL, STR și TET, pentru care cazuri de rezistență încrucișată au fost documentate clinic (a se vedea de exemplu [\(Tadesse et al., 2012\)](#)).

A	B	suport	încredere	lift	cosinus	Jaccard	R
{CHL,TET}	{STR}	0.128	0.857	1.918	0.495	0.273	9.9
{STR}	{TET}	0.322	0.721	1.186	0.618	0.440	11
{STR,TET}	{CHL}	0.128	0.396	2.173	0.527	0.339	11.8
{CHL}	{STR}	0.140	0.767	1.716	0.490	0.286	11.9
{AMP,STR}	{COT}	0.082	0.500	3.427	0.530	0.360	12
{TET}	{STR}	0.322	0.530	1.186	0.618	0.440	12.4
{AMP,STR,TET}	{COT}	0.067	0.595	4.075	0.522	0.349	13.3
{AMP,TET}	{COT}	0.076	0.500	3.427	0.510	0.342	13.8
{CHL}	{TET}	0.149	0.817	1.343	0.447	0.232	14.8

Modele și tehnici computaționale selectate pentru analiza datelor

{COT}	{STR}	0.116	0.792	1.772	0.452	0.242	15.1
{CHL,STR}	{TET}	0.128	0.913	1.502	0.438	0.206	15.8
{STR}	{CHL}	0.140	0.313	1.716	0.490	0.286	16.5
{COT,TET}	{STR}	0.100	0.846	1.894	0.436	0.216	17
{STR,TET}	{COT}	0.100	0.311	2.134	0.463	0.273	17.1
{COT}	{AMP}	0.094	0.646	2.005	0.435	0.252	17.2
{COT,STR}	{AMP}	0.082	0.711	2.205	0.425	0.231	18.2
{AMP}	{STR}	0.164	0.509	1.140	0.433	0.271	19
{AMP,TET}	{STR}	0.112	0.740	1.656	0.432	0.231	19.2
{STR}	{COT}	0.116	0.259	1.772	0.452	0.242	20.1
{AMP,COT}	{STR}	0.082	0.871	1.949	0.400	0.179	20.2
{AMP}	{COT}	0.094	0.292	2.005	0.435	0.252	20.8
{STR}	{AMP}	0.164	0.367	1.140	0.433	0.271	20.8
{TET}	{CHL}	0.149	0.245	1.343	0.447	0.232	21.2
{COT}	{TET}	0.119	0.813	1.337	0.398	0.187	22.2
{COT,STR}	{TET}	0.100	0.868	1.429	0.379	0.161	23
{COT,TET}	{AMP}	0.076	0.641	1.990	0.389	0.208	23.2
{AMP,COT,TET}	{STR}	0.067	0.880	1.970	0.363	0.147	24.5
{COT}	{CHL}	0.061	0.417	2.285	0.373	0.227	24.6
{COT,STR,TET}	{AMP}	0.067	0.667	2.069	0.372	0.188	24.7
{CHL}	{COT}	0.061	0.333	2.285	0.373	0.227	26.3
{AMP,STR}	{CHL}	0.064	0.389	2.132	0.369	0.226	26.6
{AMP}	{TET}	0.152	0.472	0.776	0.343	0.195	27.5
{AMP,CHL}	{STR}	0.064	0.840	1.880	0.346	0.139	27.8
{TET}	{COT}	0.119	0.195	1.337	0.398	0.187	28.6
{AMP,STR}	{TET}	0.112	0.685	1.127	0.356	0.171	28.8
{STR,TET}	{AMP}	0.112	0.349	1.083	0.349	0.211	30.2
{TET}	{AMP}	0.152	0.250	0.776	0.343	0.195	30.3
{AMP,COT}	{TET}	0.076	0.806	1.327	0.318	0.121	30.6
{AMP,COT,STR}	{TET}	0.067	0.815	1.340	0.299	0.107	31.5
{CHL}	{AMP}	0.076	0.417	1.293	0.313	0.177	33.4
{CHL,STR}	{AMP}	0.064	0.457	1.417	0.301	0.160	34.2
{AMP}	{CHL}	0.076	0.236	1.293	0.313	0.177	35.9

Tabelul 2.5: Regulile finale extrase din setul de date ordonate după rangul agregat

În (Tan et al., 2004), Tan et al. investighează efectul folosirii unui prag inferior de suport asupra distribuției coeficientului Φ al regulilor obținute. Autorii au observat că procedura elimină cel mai adesea reguli cu Φ negativ sau apropiat de 0 (modele necorelate sau slab corelate). În Figura 2.5 se poate observa că fixând un prag minim pentru suport se obțin rezultate similare pentru măsurile Jaccard și cosinus. Anume, cele mai multe reguli eliminate au valori mici ale celor doi indici, iar cele rămase reprezintă asocieri mai puternice.

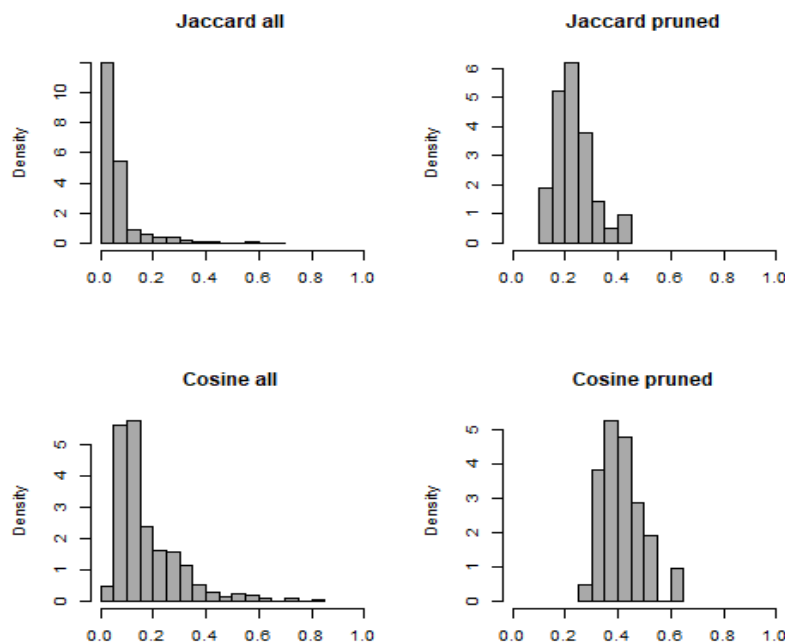


Figura 2.5: Distribuția măsurilor Jaccard și cosinus aplicate regulilor înainte și după filtrarea bazată pe suport

Rezultatele obținute ilustrează faptul că extragerea regulilor de asociere poate fi un instrument util în determinarea tiparelor de rezistență antimicrobiană încrucișată atâta timp cât sunt considerate măsuri de interes relevante.

3. Modele de învățare automată în procesarea datelor medicale

Rezultatele prezentate în această secțiune se referă la utilizarea unor modele bazate pe învățare automată în procesarea seturilor de date caracterizate prin valori absente și se bazează în principal pe lucrările (Miok et al, 2019a), (Miok et al, 2019b).

3.1. Tehnici de imputare a valorilor absente

3.1.1. Problema valorilor absente

Datele incomplete reprezintă una dintre problemele frecvente în analiza datelor și poate cauza inferența din date a unor modele de calitate scăzută și/sau obținerea unor concluzii eronate (Lall., 2016). În contextul datelor medicale, prevalența valorilor absente poate fi semnificativă din cauza dificultății înregistrării datelor (disfuncționalități ale dispozitivelor de măsură, necompletarea unor câmpuri etc.) iar strategia simplă de a exclude înregistrările cu valori absente nu este cea mai adecvată. Tehnicile de tratare a valorilor absente depind de categoria în care se încadrează (Camino et al., 2019):

- Valorile absente sunt complet aleatoare (*missing completely at random* - MCAR) – în acest caz absența unei valori poate fi interpretată ca fiind un eveniment aleator

Modele și tehnici computaționale selectate pentru analiza datelor

independent atât în raport cu valorile posibile ale atributului pentru care lipsește valoarea cât și în raport cu valorile celorlalte atribute.

- Valorile absente sunt aleatoare (*missing at random* - MAR) – în acest caz absența unei valori poate fi interpretată ca fiind un eveniment aleator care depinde de valorile celorlalte atribute și prin urmare probabilitatea ca o valoare să fie absentă poate fi estimată în funcție de valorile celorlalte atribute.
- Valorile absente nu sunt aleatoare (*missing not at random* – MNAR) – în acest caz absența unei valori depinde atât de valorile posibile ale atributului analizat cât și de valorile corespunzătoare celorlalte atribute.

3.1.2. Metode de tratare a valorilor absente

Problema valorilor absente poate fi abordată, în funcție de tiparul specific de absență a valorilor, folosind diferite strategii (Van Buuren., 2018):

- *Eliminarea atributului.* În cazul în care pentru unul dintre atribute lipsesc valori corespunzătoare multor instanțe, acel atribut nu conține suficientă informație pentru a fi inclus în analiză, și ca atare este eliminat.
- *Eliminarea instanței.* În cazul în care există instanțe pentru care lipsesc valori ale unuia sau mai multor atribute atunci aceste instanțe pot fi excluse din analiză. Este varianta standard de tratare a valorilor absente, fiind varianta implicită în multe pachete de analiză statistică, cum sunt SPSS, SAS și Stata. Funcția `na.omit()` din S-PLUS și R realizează aceeași prelucrare. Avantajele acestei abordări sunt: (i) ușor de implementat; (ii) în cazul în care valorile absente sunt complet aleatoare, prin eliminarea instanțelor cu valori absente nu sunt alterate caracteristicile statistice ale datelor (medie, abatere standard sau estimările coeficienților modelelor de regresie). Principalul dezavantaj este faptul că prin eliminarea datelor se pierde informație ceea ce conduce creșterea erorilor modelelor statistice construite pornind de la setul redus de date.
- *Eliminare la nivel de perechi.* Ideea acestei metode este de a nu exclude instanțe în întregime ci doar de a nu lua în considerare instanța dacă analiza implică un atribut pentru care nu este completată valoarea în instanța respectivă. Spre exemplu în analiza corelației între atribute, pentru fiecare pereche de atribute se vor lua în considerare toate valorile prezente (în felul acesta o instanță va contribui la valoarea coeficientului de corelație în care intervin atributele care au valori completate dar nu va contribui în cazul în care sunt implicate atribute care nu au valori completate). Abordarea aceasta este adecvată doar dacă valorile absente sunt aleatoare.
- *Imputare pe baza mediei.* Tehnicile de imputare au ca scop completarea valorilor absente, pentru a se evita pierderea de informație. Cea mai simplă strategie este înlocuirea valorilor absente cu media valorilor existente (în cazul atributelor numerice) respective cu moda valorilor existente (în cazul atributelor categoriale). Principalul avantaj este simplitatea, și faptul că nu alterează media atributului (în cazul în care valorile absente sunt aleatoare), iar principalul dezavantaj este că va conduce la o subestimare a varianței. Se recomandă utilizarea acestei strategii doar dacă alte strategii nu pot fi aplicate.

- *Imputare pe baza unui model de regresie.* Această tehnică exploatează relațiile existente între atribute cu scopul de a implementa strategii inteligente de imputare. Primul pas al acestei strategii constă în construirea unui model de regresie pornind de la datele existente. În al doilea pas se calculează, folosind modelul de regresie construit la primul pas, estimări ale valorilor absente. Modelul de regresie poate fi unul simplu sau unul bazat pe tehnici specifice învățării automate. Principalul avantaj al abordării este că se ține cont de corelațiile între atribute, iar principalul dezavantaj este că este o tehnică mai costisitoare decât celelalte. În plus, poate conduce la o sub-estimare a variabilității datelor. Din acest motiv este recomandată imputarea multiplă, când unei valori absente i se asociază mai multe valori.

3.1.3. Particularitățile imputării multiple

Tehnicile clasice de imputare generează o singură valoare pentru o observație care lipsește. Dezavantajul acestei abordări este faptul că în procesul de imputare nu se ia în considerare incertitudinea, presupunându-se implicit că valorile generate sunt corecte, adică în procesul de imputare nu intervin erori (ceea ce nu este întotdeauna adevărat). Pentru a elimina acest dezavantaj, pentru fiecare valoare absentă ar putea fi generate mai multe valori, variabilitatea setului de valori putând fi interpretată ca o măsură a incertitudinii.

Această tehnică este denumită *imputare multiplă* (Little & Rubin., 2019). Problema principală în implementarea unei tehnici de imputare multiplă este alegerea adecvată a unui model de imputare (Morris et al., 2014). Unele modele de imputare multiplă nu permit tratarea datelor mixte (catorogorice și continue), altele impun restricții privind distribuțiile utilizate (de exemplu normalitate) iar altele nu pot gestiona tipare arbitrare de valori absente. Alte dificultăți pot să apară din cauza dependențelor/ interacțiunilor neliniare între atribute sau a necesității de a procesa un volum mare de date.

3.2. Modele de tip auto-encoders

3.2.1. Invățare automată și modele cu structură profundă

Invățarea automată (Machine Learning) este un domeniu al inteligenței artificiale având ca scop extragerea de modele pornind de la date. Aceste modele permit descrierea relațiilor dintre date și sunt construite pornind de la date folosind algoritmi de învățare (antrenare). În funcție de specificul modelului (model de clasificare în clase predefinite sau model de grupare în grupuri de date similare), algoritmi de antrenare pot fi *supervizați* (pentru fiecare dintre datele din setul de antrenare se cunoaște clasa în care se încadrează) sau *nesupervizați* (se cunosc doar datele, grupurile nefiind prestabilite). Problemele de clasificare (de exemplu, asocierea unei bacterii cu una sau mai multe familii de rezistență antimicrobiană) se încadrează în categoria celor care se rezolvă folosind modele a căror parametri se determină prin algoritmi de antrenare supervizată. Modelele diferă între ele prin structura lor, adică prin modul în care sunt combinate variabilele de intrare și modul în care intervin parametrii antrenabili. Cele mai simple modele sunt cele liniare, în care variabila de ieșire depinde liniar atât de variabilele de intrare cât și de parametri.

Modele și tehnici computaționale selectate pentru analiza datelor

Modele liniare sunt cu un singur nivel, în sensul că variabila de ieșire poate fi calculată simplu, prin $Y=WX$, unde X este vectorul cu valorile de intrare iar W este o matrice ce conține valorile parametrilor antrenabili. Aceste modele au avantajul simplității dar dezavantajul că nu pot capta interacțiuni complexe între variabile. Modelele neliniare ierarhice se caracterizează prin faptul că se compun nivele de transformări liniare și neliniare: $Y=F_n(W_n F_{n-1}(W_{n-1} F_{n-2}(\dots F_2(W_2 F_1(W_1 X))\dots)))$, unde F_i reprezintă o funcție vectorială neliniară. În această familie de modele intră și rețelele neuronale cu structură profundă (deep neural networks) care au devenit, în ultimii ani, modele populare în captarea interacțiunilor complexe dintre variabile.

3.2.2. Structura generală a unui auto-encoder

Modelele ierarhice cu structură profundă permit extragerea unor caracteristici ale datelor precum și ale unor reprezentări mai compacte ale acestora (de exemplu, transformarea unor date caracterizate prin N atribute în date caracterizate prin $M < N$ atribute). Noile atribute pot codifica caracteristici ascunse în datele de intrare și extragerea lor poate fi benefică în construirea modelelor de procesare ulterioară a datelor (clasificare, regresie sau reconstruire a datelor inițiale).

Modelele de tip *auto-encoder* (AE) au ca scop transformarea datelor inițiale în reprezentări mai compacte (care captează cât mai mult din informația purtată de datele inițiale) care ulterior sunt retransformate în date de dimensiunea celor inițiale, cu scopul de a elimina tot ceea ce este irelevant în datele inițiale și a obține o reprezentare nouă (latentă) a informației stocate în date. Auto-encoder-ele sunt constituite din două componente (Rumelhart et al., 2014):

- un modul de codificare (*encoder*), al cărui scop este de a comprima date de dimensiuni mari în reprezentări de dimensiuni mai mici;
- un modul de decodificare (*decoder*), al cărui scop este să reconstruiască datele inițiale.

Intrucât în etapa de decodificare, datele sunt reconstruite pornind de la reprezentarea în spațiul de variabile latente, modelele de tip auto-encoder pot fi utilizate atât pentru a reduce dimensionalitatea datelor (în acest caz se folosesc valorile generate de modulul de codificare) cât și pentru a genera noi date pornind de la distribuția datelor inițiale (în acest caz se folosește modulul de decodificare).

În 2013, Kingma și Welling (Kingma & Welling, 2013) au propus modelele generative de tip variațional (*Variational Autoencoder* - VAE) care se particularizează în raport cu modelele de tip auto-encoder deterministe prin faptul că în procesul de antrenare se estimează parametrii distribuției care generează valorile variabilelor latente (Doersch, 2016). Prin urmare în VAE variabilele latente sunt stocate implicit prin intermediul parametrilor unei distribuții probabiliste, ceea ce permite generarea unui număr nelimitat de valori având un efect similar procesului de interpolare din cazul determinist.

Arhitectura generală a unui model de tip VAE bazat pe o rețea neuronală multi-nivel este ilustrată în Figura 3.1.

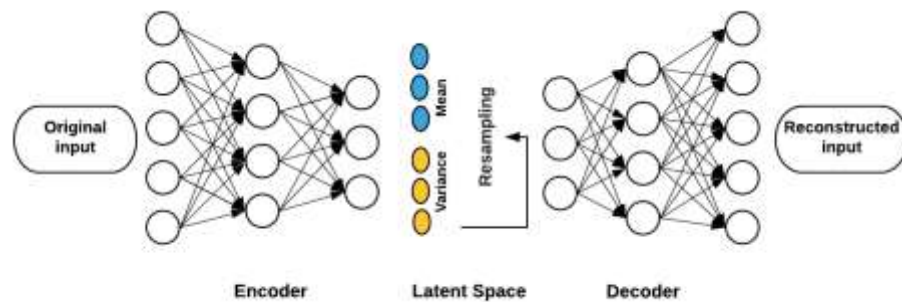


Figura 3.1. Arhitectura unei rețele neuronale generative de tip VAE

3.2.3. Exemple de aplicații unde pot fi utilizate modelele de tip auto-encoder

Deși, inițial modelele de tip AE și VAE au fost propuse în contextul prelucrării imaginilor, cu scopul generării de exemple care să extindă seturile de antrenare, la ora actuală sunt utilizate în diferite alte contexte. Printre cele mai populare aplicații sunt: reducerea dimensiunii datelor, extragerea caracteristicilor, generarea de noi date. Principiul de bază în utilizarea AE pentru reducerea dimensiunii datelor este că informația purtată de datele inițiale ar putea fi încorporată într-un număr mai mic de variabile, iar aceste variabile sunt cele obținute la ieșirea modulului de compresie a unui AE. Principala dificultate este identificarea arhitecturii potrivite, adică a numărului și dimensiunii nivelelor din structura ierarhică (a rețelei neuronale).

3.2.4. Utilizarea în contextul imputării valorilor absente

Pentru a se realiza imputarea valorilor absente se generează, pornind de la reprezentarea latentă, valori noi folosind tehnica MCMC (Markov Chain Monte Carlo). S-a demonstrat (McCoy et al., 2013) că distribuția generată folosind MCMC converge către distribuția marginală a valorilor absente, prin urmare valorile selectate folosind această distribuție sunt în concordanță cu distribuția datelor inițiale.

Principalele etape ale procesului de imputare bazat pe modele de tip VAE sunt:

- se înlocuiesc valorile absente cu valori stabilite folosind o tehnică simplă de imputare (de exemplu valoarea medie sau moda);
- se transmit datele modulului de codificare din VAE și se determină parametrii distribuției corespunzătoare reprezentării latente a datelor;
- se generează valori folosind distribuția estimate folosind modulul de codificare
- valorile generate care corespund unor componente absente în datele inițiale sunt utilizate pentru imputare

Pentru evaluarea calității unui model de tip VAE se utilizează ca referință un set de date complete și se elimină din acestea, în manieră aleatoare, o parte dintre valori. Datele de cu valori eliminate sunt utilizate ca date de intrare asupra cărora se aplică etapele enumerate mai sus, iar rezultatul se compară cu datele inițiale.

Modele și tehnici computaționale selectate pentru analiza datelor

În etapa de antrenare, parametrii modelului de tip VAE sunt ajustați atât timp cât diferența dintre datele ce conțin valori imputate și datele inițiale este mai mare decât un prag dat. Procesul de antrenare este influențat de o serie de parametri de control (hiperparametri): arhitectura și dimensiunea modelului (numărul de unități funcționale din rețeaua neuronală), parametrii algoritmului de antrenare (algoritmul de antrenare, rata de învățare, numărul maxim de iterații, dimensiunea subsetului de date care constituie un pachet de antrenare, toleranța de eroare), iar alegerea acestora are impact asupra calității modelului, deci implicit asupra calității modelului de imputare.

3.3. Instrumente software

Există o serie de instrumente software, majoritatea dezvoltate pentru a fi utilizate din aplicații implementate în Python, care permit specificarea și antrenarea unor modele bazate pe rețele neuronale cu structură adâncă. Unul dintre cele mai populare instrumente este TensorFlow, dezvoltat de Google Brain (<https://www.tensorflow.org/>) care este o bibliotecă open-source pentru a efectua eficient calculele numerice implicate în modelele specifice învățării automate, în particular variate arhitecturi de rețele neuronale cu multe nivele. Accesul la funcțiile din TensorFlow este facilitat de biblioteca Keras (<https://keras.io/>) care oferă o interfață intuitivă și modalități simple de descriere a structurii modelelor pornind de la o serie de componente pre-definite.

Fig. 3.2. prezintă un exemplu de descriere în Keras a unui model de tip auto-encoder cu trei nivele ascunse având 2000, 500 respectiv 500 de unități funcționale aplicat pentru imagini de 28*28 pixeli liniarizate ca vectori de 784 de elemente.

```
#defining input placeholder for autoencoder model
input_img = Input(shape=(784,))
# "enc_rep" is the encoded representation of the input
enc_rep = Dense(2000, activation='relu')(input_img)
enc_rep = Dense(500, activation='relu')(enc_rep)
enc_rep = Dense(500, activation='relu')(enc_rep)
enc_rep = Dense(10, activation='sigmoid')(enc_rep)

# "decoded" is the lossy reconstruction of the input from encoded representation
decoded = Dense(500, activation='relu')(enc_rep)
decoded = Dense(500, activation='relu')(decoded)
decoded = Dense(2000, activation='relu')(decoded)
decoded = Dense(784)(decoded)

# this model maps an input to its reconstruction
autoencoder = Model(input_img, decoded)
```

Fig. 3.2. Exemplu de descriere în Keras a unui model de tip auto-encoder

3.4. Studiu de caz în tratarea valorilor absente în seturi de date

3.4.1. Descrierea setului de date

Pentru a ilustra modul în care modelele de tip auto-encoder pot fi utilizate pentru imputarea valorilor absente în cazul datelor bio-medicale au fost utilizate două categorii de date:

- Seturi de date de la UCI (University of California Irvine) repository ([Asuncion & Newman, 2007](#))
- Date colectate la Stațiunea de cercetare-dezvoltare pentru creșterea bovinelor Arad care conțin informații privind calitatea laptelui precum și informații relevante pentru analiza rezistenței antimicrobiene (numărul de celule somatice). Datele au fost colectate de la 264 de vaci și conțin informații cum sunt: cantitate lapte, cazeină, grăsime, lactoză, nivel Ph, proteine, uree și număr de celule somatice.

Caracteristicile datelor utilizate în analiză sunt ilustrate în Tabelul 3.1.

Table 3.1: Caracteristici ale setului de date: N – număr de instanțe, a – număr de atribute, Num – număr de atribute numerice, Missing – procentul valorilor absente, disc – numărul de atribute discrete, C – numărul de clase.

Datasets	N	a	Num	disc	Missing %	C
Study-MILK	610	11	11	1	3.9	2
Brest-WISC	699	9	9	1	0	2
PIMA-diabetes	768	8	8	1	0	2

3.4.2. Descrierea arhitecturii

Modelul utilizat este cel de tip VAE combinat cu o strategie de Monte Carlo Dropout (MCD) ([Gal & Ghahramani, 2016a](#)). Tehnica de tip dropout constă în inactivarea temporară a unor unități funcționale și a fost aplicată în contextul a diferite arhitecturi de rețele neuronale : feed-forward, convoluționale, recurente ([Gal & Ghahramani, 2015](#)), ([Gal & Ghahramani, 2016b](#)). În contextul analizei derulate în cadrul proiectului BioAMR, proprietățile strategiei VAE+MCD ([Helmbold & Philip, 2017](#)) au fost utilizate în scopul estimării valorilor absente. Aplicând tehnica dropout în componenta de decodificare se pot genera mai multe valori la ieșire, iar prin medierea acestora se generează valorile pentru imputare. Combinarea MCD cu modelele de tip AE și VAE a fost inițial propusă în ([Miok et al., 2019a](#)) cu scopul de a construi modele generative. Avantajul combinării MCD cu VAE este că în etapa de imputare este captată distribuția specifică datelor și interacțiunile dintre atribute și utilizează mecanismul reducerii dimensiunii datelor cu scopul reducerii efectului zgomotului asupra datelor și asupra aleatorității valorilor absente. Arhitectura propusă este descrisă în Fig. 3.3.

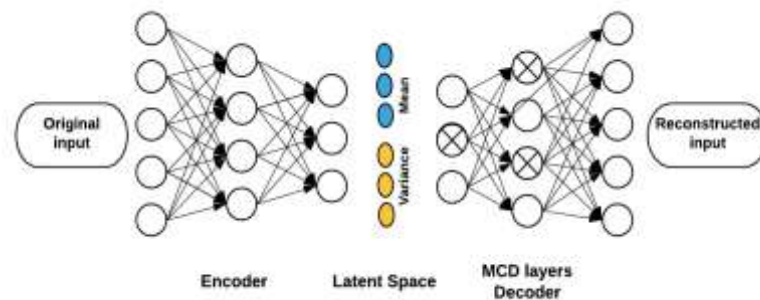


Figure 3.3. Arhitectura modelului VAE+MCD.

3.4.3. Metodologia de testare

Pentru a analiza performanța modelului de imputare, procesul de testare trebuie efectuat într-un mediu controlat, ceea ce înseamnă că se pornește de la un set de date complete (fără valori absente) din care se elimină în mod aleator, dar controlat, valori pentru a se genera date de test. În acest context, datele de test au fost generate eliminând, în manieră aleatoare, 10%, 30% respectiv 50% din valori (în implementare, valorile selectate au fost inactivate prin înlocuirea lor cu -1). Fiecare set de date este divizat într-un subset de antrenare (90%) și un subset de testare (10%) utilizat pentru analiza performanței modelului.

Pentru analiza performanței modelului sunt folosite două măsuri:

- Eroarea de imputare, măsurată folosind eroarea medie pătratică dintre datele inițiale (ce conțin valorile reale ale variabilelor absente) și datele obținute în urma imputării.
- Acuratețea clasificării (estimată folosind un model de clasificare construit pe baza datelor inițiale, respectiv pe baza datelor imputate).

Abilitatea de a conserva performanța unui model de clasificare construit folosind setul complet de date a fost evaluată folosind abordarea propusă în (Robnik-Šikonja, 2019). Ideea de bază este compararea dintre acuratețea modelului construit folosind setul complet de date (dataset1) și acuratețea modelului construit folosind setul de date cu valori imputate (dataset2): $\Delta_{acc} = accuracy(\text{dataset1}) - accuracy(\text{dataset2})$. O valoare mică a diferenței sugerează că prin imputarea valorilor absente performanța predictivă a setului de date poate fi conservată.

3.4.4. Rezultate

Arhitectura modelului conține două nivele ascunse constituite din 80 respectiv 20 unități funcționale. Rata de dropout utilizată este 0.2 (o valoare utilizată frecvent în aplicații) iar numărul de epoci de antrenare a fost setat la 300. Întrucât folosind MCD se generează mai multe valori precise rezultatul este stabilit folosind media valorilor generate. Implementarea este accesibilă la https://github.com/KristianMiok/MI_MCD_VAE.

Rezultatele obținute folosind validare încrucișată cu 5 felii sunt prezentate în Tabelele 3.2, 3.3. și 3.4 (Miok et al., 2019c). Cazurile corespunzătoare celei mai mici valori (eroare sau variație în acuratețe) sunt marcate.

Modele și tehnici computaționale selectate pentru analiza datelor

Table 3.2: Eroarea de imputare (RMSE) pentru cazul în care procentul valorilor absente este 10%:

valoare medie (abatere standard).

10% valori absente (RMSE)			
Model	MILK	WISC	PIMA
AE	0.05364[0.0011]	0.07649[0.0093]	0.06565[0.0059]
VAE	0.05027[0.0015]	0.06014[0.0038]	0.06909[0.0083]
MCD-AE	0.0479 [0.0015]	0.06048[0.0053]	0.06649[0.0082]
MCD-VAE	0.0465 [0.0012]	0.05939[0.0029]	0.06462[0.0088]

Table 3.3: Eroarea de imputare (RMSE) pentru cazul în care procentul valorilor absente este 30%:

valoare medie (abatere standard)

30% valori absente (RMSE)			
Model	MILK	WISC	PIMA
AE	0.09755[0.0080]	0.12444[0.0170]	0.11103[0.0046]
VAE	0.08049[0.0049]	0.11229[0.0091]	0.11666[0.0054]
MCD-AE	0.08685[0.0049]	0.1129 [0.0105]	0.11410[0.0048]
MCD-VAE	0.07827[0.0051]	0.1059 [0.0080]	0.11221[0.0051]

Table 3.4: Eroarea de imputare (RMSE) pentru cazul în care procentul valorilor absente este 50%

valoare medie (abatere standard)

50% valori absente (RMSE)			
Model	MILK	WISC	PIMA
AE	0.12491[0.0104]	0.14901[0.0216]	0.14132[0.0108]
VAE	0.10002[0.0050]	0.13753[0.0092]	0.14057[0.0072]
MCD-AE	0.10559[0.0052]	0.12488[0.0091]	0.13829[0.0074]
MCD-VAE	0.09764 [0.0053]	0.12706[0.0128]	0.13815[0.0072]

Valorile raportate în tabelele 3.2-2.4 sugerează că eroarea de imputare este redusă când se utilizează MCD pentru toate seturile de date și toate procentele de valori absente, cu excepția setului PIMA (cu 30% valori absente).

Table 3.5: Compararea proprietăților predictive ale modelelor construite folosind datele inițiale respectiv datele imputate (cazul a 10% valori absente).

Model	MILK	WISC	PIMA
AE	-0.0131	0.01428	0.1365
VAE	-0.0172	-0.00714	0.1429
MCD-AE	-0.009	0	0.1230
MCD-VAE	-0.010	0	0.1492

Rezultatele arată ca variantele de imputare bazate pe modele care includ MCD au o mai bună abilitate în conservarea proprietăților predictive decât cele care nu includ dropout.

3.4.5. Discuție

Valorile absente reprezintă o provocare în analiza datelor și există diferite strategii de imputare, fiecare dintre ele cu avantaje și dezavantaje. În contextul proiectului BioAMR a fost investigată performanța unor tehnici de imputare bazate pe modele de tip auto-encoder combinate cu strategii de inactivare a unor unități funcționale (dropout) cu scopul de a realiza imputare multiplă. Rezultatele obținute indică superioritatea modelelor care încorporează dropout, atât din perspectiva erorii de imputare cât și din perspectiva proprietăților predictive ale setului de date completat cu valorile obținute în urma imputării. Principalul dezavantaj al abordării propuse este acela că presupune că valorile absente sunt complet aleatoare, ceea ce nu se întâmplă întotdeauna în seturile de date reale.

Bibliografie:

1. (Agrawal et al., 1993) R. Agrawal, T. Imielinski, A. Swami, Mining association rules between sets of items in large databases, in: Proceedings of 1993 ACM-SIGMOD International Conference on Management of Data, Washington, DC, May 1993, pp. 207-216.
2. (Agrawal & Srikant, 1994) R. Agrawal, R. Srikant, Fast algorithms for mining association rules, in: Proceedings of the 20th International Conference on Very Large Databases, Santiago, Chile, 1994, pp. 487-499.
3. (Asuncion & Newman, 2007) Asuncion, Arthur, and David Newman. "UCI machine learning repository." (2007).
4. (Bayardo & Agrawal, 1999) R. Bayardo, R. Agrawal, Mining the most interesting rules, in: Proceedings of the Fifth International Conference on Knowledge Discovery and Data Mining, San Diego, CA, August 1999, pp. 145-154.
5. (Camino et al., 2019) Camino, Ramiro D., Christian A. Hammerschmidt, and Radu State. "Improving Missing Data Imputation with Deep Generative Models." *arXiv preprint arXiv:1902.10666* (2019).
6. (Cazer et al., 2019) C. L. Cazer, M. A. Al-Mamun, K. Kaniyamattam, W. J. Love, J. G. Booth, C. Lanzas, Y. T. Gröhn, Shared multidrug resistance patterns in chicken-associated *Escherichia coli* identified by association rule mining, *Frontiers in Microbiology*, 10, 2019, 687.
7. (Cărunta et al., 2019) Alina Cărunta, Mihai Pleșu, Mircea Marin, Antimicrobial resistance patterns detection using gene - EHB 2019, Iasi, Romania, November 21-23, 2019.
8. (Doersch., 2016) Doersch, Carl. "Tutorial on variational autoencoders." *arXiv preprint arXiv:1606.05908* (2016).
9. (Gal & Ghahramani, 2015) Gal, Yarin, and Zoubin Ghahramani. "Bayesian convolutional neural networks with Bernoulli approximate variational inference." *arXiv preprint arXiv:1506.02158* (2015).
10. (Gal & Ghahramani, 2016a) Gal, Yarin, and Zoubin Ghahramani. "Dropout as a bayesian approximation: Representing model uncertainty in deep learning." *international conference on machine learning*. 2016.
11. (Gal & Ghahramani, 2016b) Gal, Yarin, and Zoubin Ghahramani. "A theoretically grounded application of dropout in recurrent neural networks." *Advances in neural information processing systems*. 2016.
12. (Geng & Hamilton, 2006) L. Geng, H. J. Hamilton, Interestingness measures for data mining: a survey, *ACM Computing Surveys*, 38, no. 3, 2006, article 9.
13. (Hahsler, 2015) M. Hahsler, A Probabilistic Comparison of Commonly Used Interest Measures for Association Rules, 2015, URL: http://michael.hahsler.net/research/association_rules/measures.html
14. (Hahsler et al., 2019) M. Hahsler, C. Buchta, B. Gruen, K. Hornik, arules: Mining association rules and frequent itemsets. R package version 1.6-4, 2019, <https://CRAN.R-project.org/package=arules>.
15. (Hämäläinen & Nykänen, 2008) W. Hämäläinen, M. Nykänen, Efficient discovery of statistically significant association rules, Eighth IEEE International Conference on Data Mining, 2008, pp. 203-212.

Modele și tehnici computaționale selectate pentru analiza datelor

16. (Hämäläinen & Webb, 2019) W. Hämäläinen, G. I. Webb, A tutorial on statistically sound pattern discovery, *Data Mining and Knowledge Discovery* 33, 2019, pp.325-377.
17. (Helmbold & Philip, 2017) Helmbold, David P., and Philip M. Long. "Surprising properties of dropout in deep networks." *The Journal of Machine Learning Research* 18.1 (2017): 7284-7311.
18. (Kingma & Welling, 2013) Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." *arXiv preprint arXiv:1312.6114* (2013).
19. (Lall., 2016) Lall, Ranjit. "How multiple imputation makes a difference." *Political Analysis* 24.4 (2016): 414-433.
20. (Little & Rubin., 2019) Little, Roderick JA, and Donald B. Rubin. *Statistical analysis with missing data*. Vol. 793. John Wiley & Sons, 2019.
21. (Martinez – Ballesteros et al., 2014) M. Martínez-Ballesteros, F. Martínez-Álvarez, A. Troncoso, J. C. Riquelme, Selecting the best measures to discover quantitative association rules, *Neurocomputing*, 126, 2014, pp. 3-14.
22. (McCoy et al., 2013) McCoy, John T., Steve Kroon, and Lidia Auret. "Variational Autoencoders for Missing Data Imputation with Application to a Simulated Milling Circuit." *IFAC-PapersOnLine* 51.21 (2018): 141-146.
23. (Megiddo & Srikant, 1998) N. Megiddo, R. Srikant, Discovering predictive association rules, in: *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining*, New York, NY: AAAI Press, 1998.
24. (Miok et al., 2019a) Miok, Kristian, Dong Nguyen-Doan, Daniela Zaharie, and Marko Robnik-Šikonja. "Generating Data using Monte Carlo Dropout." *arXiv preprint arXiv:1909.05755* (2019).
25. (Miok et al., 2019b) Kristian Miok, Dong Nguyen-Doan, Marko Robnik-Sikonja and Daniela Zaharie Multiple Imputation for Medical Data using Neural Network-based Autoencoders Combined with Monte Carlo dropout, *The 7th IEEE International Conference on E-Health and Bioengineering - EHB 2019*, Iasi, Romania, November 21-23, 2019.
26. (Morris et al., 2014) Morris, Tim P., Ian R. White, and Patrick Royston. "Tuning multiple imputation by predictive mean matching and local residual draws." *BMC medical research methodology* 14.1 (2014): 75.
27. (NARMS, 2019) National Antimicrobial Resistance Monitoring System (NARMS) Now: Human Data. Atlanta, Georgia: U.S. Department of Health and Human Services, CDC. 09/18/2019. <https://www.cdc.gov/narmsnow> Accessed 7/19/2019.
28. (Piatetski – Shapiro, 1991) G. Piatetsky-Shapiro, Discovery, analysis and presentation of strong rules, in: G. Piatetsky-Shapiro, W. Frawley (Eds.), *Knowledge Discovery in Databases*, MIT Press, Cambridge, MA, 1991, pp. 229-248.
29. (Riondato & Vandin, 2014) M. Riondato, F. Vandin, Finding the true frequent itemsets, *Proceedings of the 2014 SIAM International Conference on Data Mining*, Philadelphia, USA, April 24-26, 2014, pp. 497-505.
30. (Robnik-Šikonja., 2019) Robnik-Šikonja, Marko. "Dataset comparison workflows." *International Journal of Data Science* 3, no. 2 (2018): 126-145.
31. (Rumelhart et al., 2014) Rumelhart, David E., Geoffrey E. Hinton, and Ronald J. Williams. *Learning internal representations by error propagation*. No. ICS-8506. California Univ San Diego La Jolla Inst for Cognitive Science, 1985.

32. (Tadesse et al., 2012) D.A. Tadesse, S. Zhao, E. Tong, S. Ayers, A. Singh, M.J. Bartholomew et.al., “Antimicrobial drug resistance in Escherichia coli from humans and food animals”, United States,1950–2002, Emerging Infectious Diseases, 18 (5), 2012, pp. 741-749.
33. (Tan et al., 2004) P. -N. Tan, V. Kumar, J. Srivastava, Selecting the right objective measure for association analysis, Information Systems, 29, 2004, pp. 293-313.
34. (Tan et al., 2014) P. -N. Tan, M. Steinbach, V. Kumar, Introduction to Data Mining, Pearson Education Ltd, 2014.
35. (Vaillant et al., 2004) B. Vaillant, P. Lenca, S. Lallich, A clustering of interestingness measures, in: Suzuki E., Arikawa S. (eds) Discovery Science. DS 2004. Lecture Notes in Computer Science, vol 3245, Springer, Berlin, Heidelberg, 2004.
36. (Van Buuren., 2018) Van Buuren, Stef. *Flexible imputation of missing data*. Chapman and Hall/CRC, 2018.
37. (Weiß, 2008) C. H. Weiß, Statistical mining of interesting association rules, Statistics and Computing 18, 2008, pp. 185-194.
38. (Zaharia et al., 2019) Claudia Zaharia, Raluca Muresan, Radu Moleriu, Daniela Zaharie, Analysis of association measures used to discover antimicrobial resistance patterns, The 7th IEEE International Conference on E-Health and Bioengineering - EHB 2019, Iasi, Romania, November 21-23, 2019.