

## Problema separării cuvintelor (Separating Words Problem)

### 1. Introducere

Problema separării cuvintelor este un subiect fundamental în domeniul matematicii discrete, al informaticii teoretice și al lingvisticii computaționale. Aceasta se referă la determinarea unei metode eficiente prin care să distingem cuvintele dintr-un limbaj formal sau natural folosind o funcție matematică specifică. Acest subiect are aplicații în diverse domenii, inclusiv recunoașterea automatelor, criptografie, inteligența artificială și analiza lingvistică.

Problema separării cuvintelor este strâns legată de noțiunea de echivalență între cuvinte în cadrul unui limbaj formal. Două cuvinte pot fi considerate echivalente dacă ele nu pot fi distinse prin niciun context din limbaj. Separarea cuvintelor constă în identificarea unui context sau a unui test care poate demonstra diferența dintre două astfel de cuvinte.

Această problemă devine esențială în diverse aplicații practice: de exemplu, în procesarea limbajului natural este necesară pentru a determina limitele în care un cuvânt poate fi interpretat corect. De asemenea, în recunoașterea vocală, sistemele trebuie să fie capabile să deosebească cuvinte similare din punct de vedere fonetic.

### 2. Definiție și formulare matematică a problemei separării cuvintelor

Problema separării cuvintelor este o problemă fundamentală în teoria limbajelor formale și a automatelor, cu implicații directe în clasificarea cuvintelor, recunoașterea modelelor și construirea de algoritmi de decizie. Ea poate fi formulată cu o mare precizie matematică, permițând o analiză formală a relațiilor dintre cuvinte și limbaje.

#### 2.1 Noțiuni preliminare

Să presupunem că avem un alfabet finit  $\Sigma$ , adică o mulțime finită de simboluri (de exemplu,  $\Sigma = \{a, b\}$ ). Din acest alfabet, se pot forma cuvinte finite prin concatenarea simbolurilor. Mulțimea tuturor cuvintelor finite formate din  $\Sigma$  este notată cu  $\Sigma^*$ , numită și monoidul liber generat de  $\Sigma$ .

Un limbaj formal  $L$  este o submulțime a lui  $\Sigma$ , adică  $L \subseteq \Sigma$ . Cuvintele din  $L$  sunt considerate „acceptate” de limbaj, iar cele din  $\Sigma^* \setminus L$  sunt respinse.

#### 2.2 Formularea problemei

Fie două cuvinte  $u, v \in \Sigma$ , cu  $u \neq v$ . Problema separării cuvintelor constă în determinarea unui limbaj  $L \subseteq \Sigma$  care le distinge, adică un limbaj care îl conține pe unul dintre ele și îl exclude pe celălalt. În acest caz, spunem că  $u$  și  $v$  sunt separabile prin limbajul  $L$ .

Formal, spunem că: Două cuvinte  $u$  și  $v$  sunt separabile dacă  $\exists L \subseteq \Sigma^*$  aî.  $u \in L$  și  $v \notin L$  sau  $u \notin L$  și  $v \in L$ .

Această formulare este generală și poate fi particularizată în funcție de clasa de limbaje în care căutăm limbajul separator:

- Dacă  $L$  este un limbaj regulat, atunci spunem că  $u$  și  $v$  sunt regulat separabile.
- Dacă  $L$  este context-free, spunem că ele sunt separabile printr-un limbaj context-free.
- Dacă  $L$  este recursiv (decidabil), avem separabilitate efectivă.

## 2.3 Exemplu simplu de separare

Presupunem  $\Sigma = \{a, b\}$ ,  $u = aab$ ,  $v = aba$ . Dorim să determinăm dacă există un limbaj regulat care îl conține pe  $u$ , dar nu pe  $v$ .

Putem defini următorul limbaj:  $L = \{w \in \Sigma^* \mid w \text{ se termină cu } b\}$

Observăm că  $aab \in L$ , deoarece se termină cu  $b$ , iar  $aba \notin L$ , deoarece se termină cu  $a$ . Astfel,  $L$  separă  $u$  și  $v$ . Mai mult, acest limbaj este regulat, deoarece poate fi recunoscut de un automat finit care verifică ultimul simbol.

## 2.4 Separabilitate față de o mulțime

Problema se poate extinde de la două cuvinte la două mulțimi de cuvinte  $U, V \subseteq \Sigma^*$ . Aici, separarea înseamnă găsirea unui limbaj  $L$  astfel încât:  $U \subseteq L$  și  $V \cap L = \emptyset$ .

Aceasta este o problemă mai generală și este adesea investigată în contextul învățării automate (machine learning), unde o clasă pozitivă și una negativă trebuie separate de un clasificator formal.

## 2.5 Caracterizarea în termeni de automate

Problema poate fi interpretată și în termeni de automate finite deterministe (DFA). Dacă există un DFA care acceptă  $u$  și respinge  $v$ , atunci cuvintele sunt regulat separabile. Însă, dacă orice DFA care le acceptă/respinge nu reușește să le distingă, ele nu sunt regulat separabile.

Se poate dovedi că pentru orice două cuvinte distincte există un limbaj (nu neapărat regulat) care le separă. Însă, nu toate perechile de cuvinte sunt regulat separabile, adică nu întotdeauna putem găsi un limbaj regulat care să facă această distincție. Un exemplu clasic este dat de cuvintele cu număr egal de „a”-uri și „b”-uri în poziții variabile - acest tip de limbaj nu este regulat și nu poate separa întotdeauna astfel de cazuri.

## 2.6 Separabilitate și complexitate

O întrebare naturală este: care este complexitatea de a determina dacă două cuvinte sunt separabile printr-un limbaj regulat? Aceasta depinde de modul în care sunt reprezentate cuvintele și limbajul. Dacă folosim automate finite, problema poate fi rezolvată în timp polinomial. În schimb, dacă ne extindem la limbaje context-free sau limbaje recursive, complexitatea poate crește semnificativ — ba chiar poate deveni indecidabilă în anumite cazuri.

## 2.7 Alte exemple și construcții

Să luăm încă un exemplu:  $u = abab$  și  $v = abba$ . Un limbaj care separă aceste două cuvinte ar putea fi:  $L = \{w \in \Sigma^* \mid \text{al doilea caracter este } b\}$ .

Acest limbaj poate fi descris regulat și este capabil să distingă între cele două, dacă poziția caracterelor este criteriul relevant.

În schimb, dacă diferența dintre cele două cuvinte este una subtilă, de exemplu o literă în mijlocul unui cuvânt lung, este posibil ca separarea să necesite limbaje mai complexe, cum ar fi limbajele context-free, unde stivele pot ajuta la „ținut minte” informație structurală.

## 2.8 Relația cu conceptul de limbaj universal

Există și o legătură între separabilitate și existența unui limbaj universal care „codifică” toate regulile de separare. În acest sens, teoria separării se leagă de noțiuni precum diagrama Myhill-Nerode, care caracterizează limbajele regulate în funcție de numărul de clase de echivalență ale cuvintelor. Două cuvinte sunt inseparabile regulat dacă aparțin aceleiași clasă de echivalență Myhill-Nerode.

## 2.9 Diagrama Myhill-Nerode și legătura cu separabilitatea

O abordare profundă și elegantă a problemei separării cuvintelor provine din teoria Myhill-Nerode, care oferă o caracterizare formală a limbajelor regulate și un criteriu riguros pentru a determina când două cuvinte sunt indistincte în raport cu un limbaj regulat.

### 2.9.1 Teorema Myhill-Nerode

Această teoremă fundamentală afirmă că un limbaj  $L \subseteq \Sigma^*$  este regulat dacă și numai dacă există un număr finit de clase de echivalență induse de relația:

$$u \equiv_L v \Leftrightarrow \forall z \in \Sigma^*, uz \in L \Leftrightarrow vz \in L.$$

Aceasta înseamnă că două cuvinte  $u$  și  $v$  sunt echivalente în sensul lui Myhill-Nerode dacă extensia lor cu orice sufix  $z$  nu le poate distinge în raport cu apartenența la  $L$ . Cu alte cuvinte, niciun context nu poate face distincția între  $u$  și  $v$  în ceea ce privește acceptarea în limbajul  $L$ .

### 2.9.2 Aplicație la problema separării

Această noțiune este esențială în înțelegerea separabilității regulate a cuvintelor. Dacă două cuvinte  $u$  și  $v$  sunt în aceeași clasă de echivalență Myhill-Nerode, atunci nu poate exista niciun limbaj regulat care să le distingă - adică nu sunt regulat separabile. Dimpotrivă, dacă ele sunt în clase diferite, atunci există cel puțin un limbaj regulat care le separă.

Această relație este esențială în proiectarea automatelor finite: fiecare clasă de echivalență corespunde unei stări a unui DFA minim care recunoaște  $L$ . Astfel, dacă  $u$  și  $v$  duc automatul în aceeași stare, ele sunt echivalente și inseparabile regulat.

### 2.9.3 Exemplu ilustrativ

Fie  $\Sigma = \{a\}$ , și să considerăm limbajul:  $L = \{a^n \mid n \text{ este par}\}$ .

Conform Myhill-Nerode, două cuvinte  $a^m$  și  $a^n$  sunt echivalente dacă  $m \bmod 2 = n \bmod 2$ . Astfel, există doar două clase de echivalență: una pentru lungimi pare, alta pentru lungimi impare. Dacă luăm  $u = a^2$  și  $v = a^3$ , observăm că  $u$  nu este echivalent cu  $v$  în raport cu  $L$ , deci sunt regulat separabile (de exemplu, de limbajul  $L$  dat mai sus).

Dacă însă luăm  $u = a^2$  și  $v = a^4$ , atunci ele sunt în aceeași clasă și nu pot fi separate de niciun limbaj care recunoaște  $L$ .

### 2.9.4 Implicații teoretice și practice

Această caracterizare are consecințe importante:

- Oferă un criteriu decidabil pentru separabilitate regulată: dacă putem determina că două cuvinte sunt în clase Myhill-Nerode diferite, știm că sunt separabile.
- Permite minimizarea automatelor: automatul minim are un număr de stări egal cu numărul claselor de echivalență.
- Ghidează proiectarea de filtre sau clasificatori în procesarea automată a textului, unde clasele echivalente corespund unor modele lingvistice recurente.

### 2.9.5 Limitele teoremei

Este important de menționat că teorema Myhill-Nerode este aplicabilă exclusiv limbajelor regulate. În cazul limbajelor context-free sau recursive, alte tehnici sunt necesare pentru a analiza separabilitatea. De exemplu, în cazul limbajelor context-free, relația de echivalență nu este, în general, finit generată, ceea ce face problema mult mai complexă sau chiar indecidabilă.

### 3. Aplicații în teoria limbajelor formale

Limbajele formale reprezintă un pilon central în informatică teoretică, logică matematică și știința computațională, oferind un cadru riguros pentru definirea și analiza sintaxei și semnificației în diferite contexte: limbaje de programare, procesare automată a limbajului natural, inteligență artificială, dar și în biologie computațională sau lingvistică computațională. În această ecuație, problema separării cuvintelor se dovedește a fi un instrument metodologic fundamental, având multiple aplicații ce vizează atât aspecte teoretice, cât și practice ale limbajelor formale.

#### 3.1 Recunoașterea și clasificarea limbajelor

Una dintre cele mai directe aplicații ale problemei separării este în clasificarea cuvintelor în cadrul unui limbaj formal. Dacă avem o mulțime de cuvinte și dorim să determinăm dacă ele aparțin unui limbaj regulat (sau unei alte clase), putem utiliza metode de separare pentru a construi un automat care să accepte doar acele cuvinte.

Prin această metodă, se determină dacă există un automat finit care poate distinge între două cuvinte - adică poate accepta unul și respinge pe celălalt. Această proprietate este relevantă pentru clasificarea limbajelor în clasele Chomsky: limbaje regulate, limbaje independente de context, limbaje sensibile la context și limbaje recursiv enumerabile.

**Exemplu:** Dacă putem construi un DFA care acceptă "abba" dar respinge "abab", atunci înseamnă că aceste două cuvinte sunt regulat separabile, și implicit aparțin unor trasee distincte în automat.

Acest proces ajută și în generarea automată de reguli de acceptare — în special în contexte în care limbajul nu este cunoscut în avans, dar trebuie aproximat sau învățat din exemple (learning from examples).

#### 3.2 Minimizarea automatelor și complexitatea de stat

În teoria automatelor, minimizarea DFA-urilor presupune reducerea numărului de stări fără a schimba limbajul acceptat. Această tehnică este esențială pentru:

- Reducerea complexității de implementare
- Creșterea eficienței în recunoaștere
- Scăderea consumului de memorie și resurse

Pentru a realiza minimizarea, trebuie să identificăm stările echivalente - adică acele stări care nu pot fi distinse de niciun cuvânt. Dacă nu există niciun cuvânt separator care, aplicat din cele două stări, produce comportamente diferite (una acceptă, cealaltă nu), atunci stările sunt echivalente și pot fi comasate.

Această idee se regăsește în algoritmul Hopcroft pentru minimizarea DFA-urilor, unul dintre cele mai eficiente din punct de vedere computațional.

**Exemplu:** Considerăm un DFA cu stările  $q_1$  și  $q_2$ . Dacă niciun sufix nu produce un rezultat diferit (acceptare vs. respingere), atunci  $q_1 \equiv q_2$ .

Separarea cuvintelor este astfel nu doar o unealtă teoretică, ci și un criteriu algoritmic concret în optimizarea sistemelor de recunoaștere.

#### 3.3 Verificarea echivalenței și incluziunii între limbaje

În cadrul analizei formale, o problemă frecventă este verificarea dacă două limbaje sunt egale sau dacă unul este inclus în celălalt. Pentru aceasta, identificarea unui cuvânt separator care aparține unuia dintre limbaje, dar nu și celuilalt, este esențială.

Această abordare este utilizată în:

- Verificarea formală a software-ului (model checking)

- Validarea de compilatoare
- Analiza diferențială între specificații

Un algoritm obișnuit ar presupune construirea automatului diferenței  $L_1 \setminus L_2$ . Dacă acesta nu este gol, înseamnă că există cel puțin un cuvânt separator.

**Exemplu practic:** Într-un compilator, putem verifica dacă versiunea nouă a unei reguli gramaticale include toate cuvintele din versiunea anterioară. Un cuvânt separat poate indica o regresie semantică.

### 3.4 Construcția gramaticilor și învățarea limbajelor

Separarea cuvintelor este o tehnică fundamentală în gramatici inductive, unde scopul este de a genera o gramatică pe baza unor exemple pozitive și negative. Prin compararea cuvintelor și identificarea trăsăturilor lor distinctive, se pot deduce:

- Reguli de producție
- Constrângeri sintactice
- Structuri ierarhice (ex. arbori de derivare)

#### Aplicații relevante:

- Lingvistică computațională – extragerea de reguli morfologice și sintactice
- Învățarea automată a limbajelor (grammar induction) – învățarea unui limbaj formal din date
- Compilatoare educaționale – generarea automată de reguli pe baza intrărilor elevilor

**Exemplu:** Avem cuvintele: “aaabbb”, “aabb”, “aaaabbbb”. Putem deduce o regulă generală:  $a^n b^n$ , ceea ce corespunde unui limbaj context-free. Cuvântul “aabb” poate fi folosit ca separator pentru a testa dacă se respectă proporția 1:1 între a și b.

### 3.5 Aplicații în verificarea și testarea formală

În testarea formală a sistemelor, separarea cuvintelor poate genera exemple concrete de comportamente incorecte sau incomplete. Dacă un sistem are un model de referință formal și implementarea sa se abate, cuvântul care declanșează abaterea este un separator valoros.

**Exemplu:** Într-un protocol de comunicație, secvența “ACK SYN RST” poate fi acceptată de modelul formal dar respinsă de implementare. Aceasta înseamnă că implementarea este subspecificată, iar cuvântul constituie un contra-exemplu util.

Această metodă este utilizată în:

- Testarea automată a aplicațiilor critice
- Verificarea de modele finite (FSM) în industria hardware
- Detectarea bug-urilor de tip corner-case

### 3.6 Aplicații interdisciplinare: biocomputing și inteligență artificială

În afara informaticii teoretice, separarea cuvintelor este relevantă și în:

- Biocomputing: analiza secvențelor de ADN pentru a detecta motifuri sau pentru a separa genele funcționale de zgomotul genomic
- Inteligență artificială: în NLP, separarea servește la tokenizarea corectă, la detectarea entităților și la parsing-ul semantic

**Exemplu:** În recunoașterea vocală, sistemul trebuie să distingă între there și their, care pot suna similar fonetic. Un separator contextual poate fi o propoziție ca: “They left their car outside.”

### Concluzii parțiale

Aplicabilitatea problemei separării cuvintelor este profundă și extinsă, depășind granițele teoriei pentru a pătrunde în domenii aplicate ale informaticii moderne. De la analiza sintactică și

semantică, până la optimizare, testare și învățare automată, separarea cuvintelor devine o unealtă versatilă pentru cercetători, ingineri și dezvoltatori de sisteme.

#### 4. Complexitatea problemei și aspecte computaționale

Problema separării cuvintelor nu este doar interesantă din punct de vedere teoretic, ci ridică și o serie de provocări semnificative din perspectiva complexității computaționale și a realizabilității algoritmice. Încercarea de a determina dacă două cuvinte pot fi separate printr-un anumit tip de limbaj implică analiza resurselor necesare (timp, memorie, număr de tranziții) pentru a construi sau a verifica existența unui astfel de separator. Această analiză este esențială pentru a înțelege limitele practice ale acestei probleme în contexte reale, cum ar fi verificarea software, învățarea automată sau procesarea limbajului natural.

##### 4.1 Decidabilitatea problemei

Primul aspect fundamental care trebuie discutat este dacă problema este decidabilă, adică dacă există un algoritm care determină în mod finit (într-un număr determinat de pași) dacă două cuvinte pot fi separate de un anumit tip de limbaj.

Pentru clasele mai simple de limbaje (de exemplu, limbajele regulate), problema separării este decidabilă și chiar eficientă din punct de vedere computațional. Există algoritmi bazați pe construcția DFA-urilor care pot verifica dacă două cuvinte sunt separate de un automat.

**Exemplu:** Se pot construi două DFA-uri: unul care acceptă doar primul cuvânt și altul doar pe al doilea. Dacă diferența dintre cele două DFA-uri generează un limbaj non-vid, atunci cuvintele pot fi separate.

Însă, pentru clase de limbaje mai puternice (ex. limbajele context-free, context-sensitive sau recursive), problema devine mult mai dificilă și adesea indecidabilă. De exemplu, dacă încercăm să separăm două cuvinte folosind o gramatică context-free, este posibil să nu existe un algoritm general care să poată stabili în toate cazurile separabilitatea, deoarece acest lucru ar implica rezolvarea unor probleme echivalente cu problema opririi (Halting Problem), care este indecidabilă.

##### 4.2 Complexitatea temporală și spațială

Chiar și în cazurile în care separabilitatea este decidabilă, complexitatea algoritmilor poate varia semnificativ:

- Pentru limbaje regulate, problema poate fi rezolvată în timp polinomial în raport cu lungimea cuvintelor implicate și dimensiunea alfabetului. Construcția automatului diferenței și verificarea acceptării pot fi realizate eficient.
- În cazul limbajelor context-free, separarea implică construirea de gramatici sau automate de tip pushdown (PDA), iar complexitatea tinde spre exponențială sau chiar nedeterministă în anumite cazuri.
- În domenii precum limbajele recursiv enumerabile, resursele necesare pentru verificare pot fi nelimitate, ceea ce face problema nerentabilă practic, chiar dacă în unele cazuri este teoretic rezolvabilă.

##### 4.3 Tehnici algoritmice de separare

Pentru a aborda această problemă în mod practic, au fost dezvoltate diverse metode algoritmice, fiecare adaptată la tipul de limbaj și complexitatea dorită:

###### 1. Automate finite deterministe (DFA)

Pentru limbajele regulate, se construiește un DFA pentru fiecare cuvânt, sau pentru o mulțime de exemple pozitive și negative. Algoritmi precum:

- Algoritmul Hopcroft pentru minimizarea stărilor
- Algoritmul Moore pentru determinarea echivalenței între stări
- Algoritmul de învățare RPNI (Regular Positive and Negative Inference)

sunt folosiți pentru a obține un automat care separă corect cuvintele date.

## 2. Automate pushdown (PDA)

În cazul în care limbajul este context-free, separarea necesită analiza stivei asociate PDA-ului. Verificarea dacă un PDA acceptă un cuvânt și respinge altul este mai costisitoare și implică algoritmi de parsing precum:

- Algoritmul CYK (Cocke–Younger–Kasami)
- Algoritmul Earley
- Construcția de gramatici echivalente

## 3. Machine learning & grammar induction

Metodele moderne de învățare automată utilizează tehnici inspirate din problema separării:

- Rețele neuronale recurente (RNN) pentru a învăța secvențe și a identifica diferențele semantice
- Algoritmi de tip decision tree care separă clase de cuvinte pe baza trăsăturilor lor
- Sisteme hibride care combină învățarea statistică cu reguli formale

## 4.4 Probleme deschise și limitări

Problema separării cuvintelor este încă subiect de cercetare activă, având numeroase direcții neexplorate sau incomplet rezolvate:

- Există o clasă minimă de limbaje în care orice pereche de cuvinte distincte este separabilă?
- Ce fel de complexitate apare în cazul separării în limbaje restricționate?
- Cum influențează alfabetul folosit dificultatea separării?

De asemenea, o problemă interesantă este identificarea lungimii minime a unui separator, sau a complexității minime a unui automat separator. Aceste întrebări au implicații directe în compresia datelor, criptografie, și în simularea limbajelor naturale.

## 4.5 Exemple concrete de analiză

Să considerăm două cuvinte:  $u = aab$  și  $v = abb$ .

Sunt acestea regulat separabile? Se poate construi un DFA care acceptă “aab” și respinge “abb”.

Un astfel de automat trebuie să distingă între succesiunea de simboluri aa urmată de b și ab urmat de b. Un DFA cu 4 stări este suficient pentru a realiza această diferențiere, ceea ce dovedește separabilitatea.

În schimb, considerăm cazul în care am avea două cuvinte de forma:  $u = a^n b^n$  și  $v = a^n b^{n+1}$ . Acestea nu sunt regulat separabile, deoarece limbajul lor nu este regulat - ar necesita memorarea unui număr arbitrar de apariții de a și b, ceea ce un DFA nu poate face. Pentru separare ar fi necesar un PDA.

## Concluzii parțiale

Complexitatea problemei separării cuvintelor reflectă gradul de sofisticare al clasei de limbaje în care operăm. Dacă în cazul limbajelor regulate putem vorbi despre o problemă computabilă și tractabilă, pe măsură ce ne deplasăm spre clase mai puternice, deciziile devin mai dificile, iar analiza implică metode avansate de calcul și teorie a complexității.

Separarea cuvintelor devine astfel o punte între teoria limbajelor formale și teoria complexității, deschizând direcții de cercetare în automatizare, logică, inteligență artificială și verificare formală.

## **5. Implicații în procesarea limbajului natural și recunoaștere vocală**

Problema separării cuvintelor are un impact major în domenii aplicate ale informaticii și lingvisticii computaționale, în special în procesarea limbajului natural (NLP) și în recunoașterea automată a vorbirii (ASR - Automatic Speech Recognition). În aceste domenii, abilitatea de a diferenția corect între unități lingvistice, adesea similare fonetic sau ortografic, este fundamentală pentru precizia și eficiența sistemelor automatizate.

### **5.1. Segmentarea și delimitarea cuvintelor în NLP**

Una dintre cele mai directe aplicații ale problemei de separare este în segmentarea textului, adică identificarea limitelor dintre cuvinte, propoziții sau fraze într-un șir continuu de caractere.

**Exemplu:** În limbi precum chineza sau japoneza, în care cuvintele nu sunt separate prin spații, segmentarea devine o problemă critică. Două secvențe de caractere pot fi interpretate diferit în funcție de context, iar sistemul trebuie să separe corect termenii pentru a evita ambiguitățile semantice.

**Aplicație:** Algoritmii de tokenizare, analiza morfologică și sintactică, recunoașterea entităților denumite sau traducerea automată depind direct de capacitatea unui sistem de a separa corect cuvintele în textul sursă.

### **5.2. Separarea fonetică în recunoașterea vorbirii**

În sistemele moderne de recunoaștere vocală, semnalele acustice sunt transformate în secvențe de foneme sau caractere. Însă uneori două cuvinte diferite pot avea reprezentări fonetice aproape identice sau chiar identice în contexte ambigue.

**Exemplu:** Cuvintele „pace” și „face” pot fi greu de separat fonetic într-un context ambiguu sau zgomotos. Problema devine și mai complexă în cazul cuvintelor omofone (de exemplu, „ceai” și „ce-i”).

**Soluție:** Se utilizează modele statistice de limbaj (n-gram, rețele neuronale de tip LSTM sau Transformer), care, pe baza contextului, stabilesc probabilitatea unei interpretări sau a alteia. Astfel, separarea cuvintelor este asistată de modele lingvistice și probabilistice care încorporează reguli implicite de separare.

### **5.3. Învățare automată și separabilitate**

În NLP, conceptele de „separabilitate” au fost extinse în metode de învățare automată:

- În clasificarea textelor, separarea documentelor în categorii presupune extragerea unor „cuvinte-cheie” care disting între clase. Astfel, problema devine una de identificare a cuvintelor „separator” între clasele de texte (spam vs. non-spam, recenzie pozitivă vs. negativă etc.).
- În vectorizarea semantică (Word2Vec, GloVe, BERT), spațiul vectorial în care cuvintele sunt proiectate trebuie să permită separarea semnificativă a vectorilor pentru a reflecta diferențele semantice și contextuale.

### **5.4. Exemple concrete de utilizare**

- Asistenți vocali (ex. Siri, Google Assistant): Când utilizatorul spune „set alarm for six”, sistemul trebuie să distingă „six” de cuvinte precum „sex” sau „sick”.



- Traducere automată: Un cuvânt cu multiple sensuri trebuie interpretat corect, în funcție de context, prin separare semantică.

## 6. Generalizări ale problemei și cercetări actuale

Problema separării cuvintelor nu se limitează doar la perechi de cuvinte, ci poate fi extinsă și generalizată în mai multe direcții. Aceste generalizări au fost investigate intens în literatura de specialitate, deschizând noi căi de cercetare în teoria limbajelor formale, complexitatea algoritmică și inteligența artificială.

### 6.1 Separarea mulțimilor de cuvinte

O extensie naturală este problema separării a două mulțimi finite de cuvinte, denumite frecvent mulțimea pozitivă (cuvinte acceptate) și mulțimea negativă (cuvinte respinse). Scopul este de a construi un limbaj formal (regulat, context-free etc.) care să conțină toate cuvintele din mulțimea pozitivă și niciunul din cea negativă.

**Exemplu concret:** mulțimea pozitivă: {„abc”, „aab”, „aac”}; mulțimea negativă: {„abb”, „acc”, „aaa”}.

**Obiectiv:** determinarea unui limbaj regulat (sau a unui DFA) care separă aceste mulțimi.

Această problemă este centrală în domeniul gramaticilor inductive și al învățării automate a limbajelor formale, fiind adesea abordată prin algoritmi precum RPNI, EDSM sau prin tehnici de învățare bazate pe contraexemple.

### 6.2 Separabilitate în clase de limbaje

Există întrebări fundamentale în teoria limbajelor legate de capacitatea unei anumite clase de limbaje de a realiza separarea:

- Sunt toate perechile de cuvinte separabile în cadrul limbajelor regulate? **Răspuns:** Nu, dacă se impune ca separatorul să fie regulat, unele perechi pot necesita un limbaj mai expresiv.
- Ce limbaje pot separa toate perechile de cuvinte de o anumită lungime? Este o întrebare deschisă în legătură cu puterea expresivă a diferitelor clase formale.

### 6.3 Separabilitatea și logica formală

O direcție recentă de cercetare implică utilizarea logicii formale (ex. logica monadică a ordinului întâi) pentru a descrie proprietățile separabilității. Se urmărește găsirea unor formule logice care descriu un separator, ceea ce leagă teoria limbajelor de fundamentele matematicii și ale inteligenței artificiale.

**Exemplu:** În logica MSO (Monadic Second-Order Logic), se poate exprima o proprietate de forma „există o poziție în cuvânt unde simbolurile diferă”, ceea ce definește un separator formal.

### 6.4 Separabilitatea și securitatea informatică

Separabilitatea are aplicații și în **securitatea software**, în special în analiza fluxurilor de date și în detectarea potențialelor vulnerabilități. Două stări ale unui program pot fi considerate sigure dacă nu există „cuvinte” (șiruri de acțiuni) care să le confunde - cu alte cuvinte, trebuie să fie separabile.

**Aplicație practică:** În analiza taint, un flux de date periculos trebuie separat de un flux sigur. Dacă acest lucru nu este posibil, sistemul poate genera alerte.

### 6.5 Alte direcții de cercetare

- Separarea limbajelor infinite: Ce tip de automaton poate distinge între limbaje infinite? (ex: omega-automate pentru cuvinte infinite).

- Separabilitate aproximativă: Există situații în care două cuvinte sunt „aproape” identice. Aici intervine separabilitatea probabilistică sau fuzzy, utilă în recunoașterea vorbirii și în procesarea limbajului natural.
- Separarea în spații vectoriale semantice: Separarea cuvintelor prin învățare distribuită (embeddinguri vectoriale), în care cuvintele trebuie mapate în regiuni distincte ale unui spațiu de înaltă dimensiune.

## 7. Concluzii

Problema separării cuvintelor reprezintă un subiect fundamental în cadrul teoriei limbajelor formale, cu ramificații teoretice profunde și aplicații practice vaste în domenii precum inteligența artificială, procesarea limbajului natural, recunoașterea automată a vorbirii, securitatea informatică și analiza programelor.

Prin natura sa, această problemă implică identificarea unui criteriu formal prin care două cuvinte, sau chiar două mulțimi de cuvinte, pot fi distinse în mod riguros. Această distincție nu este doar de ordin abstract, ci se reflectă în mod direct în construirea automatelor finite, în verificarea echivalenței limbajelor, în generarea de reguli gramaticale și în segmentarea textului.

Pe plan teoretic, separabilitatea cuvintelor stă la baza definițiilor de echivalență de stare în automatele finite, a algoritmilor de minimizare și a metodelor de inferență a limbajelor formale. În același timp, în plan practic, ea fundamentează numeroase aplicații din procesarea automată a datelor lingvistice, oferind soluții la probleme reale precum dezambiguizarea sensului, recunoașterea omofonelor sau delimitarea unităților lexicale.

Un aspect esențial evidențiat în această lucrare este faptul că separabilitatea nu este întotdeauna garantată, în funcție de clasa formală în care se caută separatorul. De exemplu, cuvinte separabile printr-un limbaj context-free pot să nu fie separabile printr-un limbaj regulat. Astfel, problema are și un caracter de analiză a puterii expresive a diferitelor clase de limbaje formale, constituind un instrument conceptual pentru înțelegerea profundă a limitelor teoretice ale automatelor și gramaticilor.

În contextul noilor direcții de cercetare, separarea cuvintelor a fost extinsă în diverse forme: separabilitate aproximativă, separabilitate probabilistică, separabilitate logică sau semantică (în spații vectoriale). Aceste extensii răspund provocărilor ridicate de complexitatea limbajelor naturale și de ambiguitățile inerente comunicării umane, evidențiind importanța continuă și evolutivă a acestei probleme.

Prin urmare, putem concluziona că problema separării cuvintelor, deși formulată aparent simplu, este una dintre cele mai profunde și fertile teme din știința calculatoarelor și lingvistica computațională. Ea contribuie nu doar la dezvoltarea teoriei limbajelor formale, ci și la soluționarea unor probleme practice esențiale pentru funcționarea eficientă a sistemelor informatice moderne.

## 8. Bibliografie

- Hopcroft, J.E., Motwani, R., Ullman, J.D. — Introduction to Automata Theory, Languages, and Computation, Pearson, 2006.
- Sipser, M. — Introduction to the Theory of Computation, Cengage Learning, 2012.
- Berstel, J., Perrin, D. — Theory of Codes, Academic Press, 1985.
- Grädel, E., Thomas, W., Wilke, T. (eds.) — Automata, Logics, and Infinite Games, Springer, 2002.
- De la Higuera, C. — Grammatical Inference: Learning Automata and Grammars, Cambridge University Press, 2010.

- Mohri, M. — Finite-State Transducers in Language and Speech Processing, Computational Linguistics, 2003.
- Jurafsky, D., Martin, J.H. — Speech and Language Processing, Pearson, 3rd Edition Draft, 2023.
- Klein, D., Manning, C.D. — Natural Language Grammar Induction with a Constituent-Context Model, ACL, 2002.
- Clark, A., Eyraud, R. — Polynomial Identification in the Limit of Subclasses of Context-Free Languages, Journal of Machine Learning Research, 2007.
- García, P., Ruiz, J.C. — Learning k-testable languages in the strict sense, Lecture Notes in Artificial Intelligence, 1999.
- Tan, Z., Smolka, S.A., et al. — Separability and Refinement Checking in Modal Transition Systems, CONCUR, 2013.
- Fijalkow, N., Zimmermann, M. — Separating Words with Automata, Proceedings of the 36th International Symposium on Theoretical Aspects of Computer Science (STACS), 2019.
- Yu, S. — Regular Languages, in Handbook of Formal Languages, Vol. 1, Springer, 1997.
- Eilenberg, S. — Automata, Languages, and Machines, Academic Press, 1974.
- ResearchGate, Academia.edu – articole academice diverse pe tema separabilității în limbaje formale.
- Arxiv.org – publicații recente privind automatele finite și învățarea automată a limbajelor.
- ChatGPT - Pentru aprofundări, explicații, căutare de surse și clarificări.