

# Algoritmi de procesare a șirurilor de caractere. Trie și distanța Levenshtein

Grigore Iulia-Anita

Facultatea de Automatică și Calculatoare  
iulia\_anita.grigore@stud.acs.upb.ro

## 1 Descrierea problemei

Într-o lume extrem de digitalizată ca cea în care trăim acum, șirurile de caractere sunt peste tot: mesaje trimise online, căutări efectuate cu ajutorul diverselor browsere de internet și chiar în propriul organism (secvențele de ADN sunt uneori tratate ca o serie de simboluri cu semnificații bine cunoscute).

Din această cauză a apărut necesitatea implementării unor algoritmi din ce în ce mai eficienți de procesare a șirurilor de caractere.

Câteva necesități din viața reală sunt, spre exemplu, programele de spell checking implementate de majoritatea editoarelor de text, sau chiar sugestiile primite în momentul în care cauți ceva pe internet.

Un exemplu care nu este atât de la îndemână, dar care are o influență majoră în dezvoltarea științei, este alinierea secvențială, un termen întâlnit în bioinformatică. Aceasta constă în aranjarea unor secvențe de ADN și de ARN pentru a putea observa dacă au un predecesor comun (un șir de date comun) și pentru a depista mutațiile (de deleție sau de inserție a unor gene) apărute la nivelul acestora. [2]

## 2 Soluții alese

Ca soluții la problema prezentată anterior, am ales structura de date Trie și un algoritm de edit distance, distanța Levenshtein. Pe parcursul lucrării, voi compara, în linii mari, mai multe distanțe (precum Hamming distance), însă doar din punct de vedere al eficientizării sau în privința completărilor aduse soluției deja alese.

### 2.1 Trie

Trie este o structură de date implementată sub forma unui arbore, în care root-ul este gol, iar descendenții săi sunt reprezentați de un alfabet. Pentru a stoca un șir de caractere, este necesară parcurgerea arboreului în adâncime, pe copilul corespondent caracterului căutat.

Datorită acestei implementări, se folosește mai puțină memorie (în majoritatea cazurilor), nefiind necesar să salvăm fiecare șir individual, ci doar să adăugăm la șirurile existente caracterele lipsă. [3]

O utilizare des întâlnită a trie-ului este oferirea unor sugestii de căutare pe baza unor caractere scrise anterior. De exemplu, dacă vom căuta pe Google „tree“, vom primi sugestii precum „tree of life“ sau „treehouse“, toate acestea având un șir în comun.

## 2.2 Distanța Levenshtein

Distanța Levenshtein măsoară *diferența* dintre două șiruri de caractere. Această diferență este calculată conform unei formule matematice, iar operațiile care le calculează sunt ștergerea, adăugarea și înlocuirea unui caracter. [1]

Formula de calcul a distanței Levenshtein este:

$$\text{lev}_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{if } \min(i, j) = 0 \\ \min \begin{cases} \text{lev}_{a,b}(i-1, j) + 1 \\ \text{lev}_{a,b}(i, j-1) + 1 \\ \text{lev}_{a,b}(i-1, j-1) + 1 \end{cases} & \text{otherwise} \end{cases} \quad (1)$$

## 3 Criterii de evaluare a soluțiilor

Pentru evaluarea acestor soluții, îmi propun să creez 2 tipuri de teste. Pentru verificarea corectitudinii codului, voi folosi unit tests create cu date de intrare variate, care verifică atât cazuri generale, cât și edge case-uri, iar pentru a testa eficiența algoritmilor, voi scrie benchmark-uri.

Unit test-urile vor fi scrise astfel încât să verifice funcționalitatea Trie-ului, fără a ține cont de performanța algoritmilor. Cu toate acestea, toți algoritmii vor fi optimizați din punct de vedere al alocărilor dinamice, astfel încât memoria RAM să influențeze cât mai puțin acuratețea benchmark-urilor.

Pentru verificarea performanței operațiilor în Trie, voi construi benchmark-uri bazate pe cuvintele din dicționarul englezei americane. Dimensiunea dicționarul ce urmează a fi inserat într-o instanță de Trie este de aproximativ 4.5MB.

Fisierul în format CSV urmează să fie procesat cu un script de Python pentru a obține și compara:

- inserția a aproximativ 460 de mii de cuvinte de lungimi variate ordonate alfabetic, și
- inserția aceluiași set de cuvinte, dar ordonate aleatoriu.

Ulterior inserției voi măsura performanța căutării cuvintelor din același dicționar (intrări ce există în Trie) și a unor secvențe generate aleatoriu, astfel încât căutarea să nu întoarcă niciun rezultat.

Voi procesa dicționarul astfel încât să fie eliminate toate intrările ce contin alte caractere decât [a-zA-Z].

## References

1. Black, P.E.: "levenshtein distance", dictionary of algorithms and data structures. U.S. National Institute of Standards and Technology (2016)
2. DM, M.: Bioinformatics: Sequence and genome analysis (2nd ed.). Cold Spring Harbor Laboratory Press: Cold Spring Harbor, NY (2004)
3. Sedgewick, R., Wayne, K.: Digit-based sorting and data structures (???)