

Motor de căutare semantic

Anca-Nicoleta Ciubotaru

Rezumat

Cantitatea datelor nestructurate a avut de-a lungul timpului o creștere exponențială, iar regăsirea informațiilor de interes s-ar realiza cu dificultate fără o metodă eficientă de indexare, ordonare și clasificare.

Tema proiectul constă în crearea unui motor de căutare ce ia în considerare caracterul semantic al cuvintelor, facilitând procesul de regăsire și clasificare a informațiilor pe baza sensului interogării efectuate de utilizator. Diferența față de implementările motoarelor de căutare clasice, care redau informațiile independent de sensul termenilor din interogare, constă în analiza contextului în care aceștia apar. Unul dintre avantajele acestei abordări este faptul că în urma unei interogări date de utilizator, informațiile sugerate de motorul de căutare vor aparține unui întreg domeniu de referință, documentele cu caracter similar fiind considerate relevante pentru utilizator.

Motorul de căutare semantic realizează următoarele sarcini: vizitarea paginilor de pe World Wide Web, salvarea conținutului, acolo unde este permis, indexarea și generarea sensurilor cuvintelor pe baza contextului în care acestea apar. În urma introducerii unei interogări, utilizatorul va obține rezultate constând în pagini relevante din punct de vedere a sensului, dar și reprezentarea vectorială a cuvintelor reprezentative.

Arhitectura este modulară, fiecare modul având un scop clar definit în cadrul ansamblului din care face parte. Una dintre componentele de bază este un modulul ce realizează căutarea, selecția și indexarea documentelor în vederea minimizării timpului de regăsire. Pentru reducerea timpului de prelucrare, etapele sunt implementate prin intermediul modelului de programare *Map Reduce*, ce realizează eficient procesarea paralelă a unor cantități mari de date. Modulul ce realizează procesarea informațiilor prin generarea sensurilor cuvintelor raportate la contextul în care apar utilizează modelul arhitectural *Skip Gram*, ce determină reprezentările vectoriale ale cuvintelor prin intermediul unei rețele neuronale. Astfel, fiecare cuvânt are asociat un set de vectori n -dimensionali pe baza cărora se vor putea identifica sinonimele și sensurile cuvintelor, luând ca unitate de măsură distanța cosinus dintre aceștia. Modulul ce realizează legătura dintre utilizator și nucleul de procesare îi permite acestuia să obțină informații de interes. Astfel, acesta va vizualiza câmpul vectorial aferent rezultatelor obținute implementat prin intermediul algoritmului *t-distributed Stochastic Neighbor Embedding* (t-SNE) ce are ca scop reducerea numărului de dimensiuni ale vectorilor n -dimensionali și generarea grafică în două sau trei dimensiuni.

Modul de căutare a informațiilor pe World Wide Web are un caracter desprins din limbajul natural și deși cuvintele cheie sunt importante în vederea obținerii informațiilor de interes, contextul în care acestea apar este necesar a fi luat în considerare în procesul de construire a unui motor de căutare care să răspundă cerințelor date de utilizator.

Îndrumător: Conf. dr. ing. Mihai Horia Zaharia