

Detecció de similitud de documents amb hashing*

GRAU A CURS Q2-2024–2025

Departament de Ciències de la Computació
alg@cs.upc.edu

Resum

Aquest projecte té com a objectiu una validació experimental de l'efectivitat de diferents algorismes de hash per a la detecció de similitud de documents.

El projecte es farà en grups de 4 persones (excepcionalment 3, sota autorització expressa). Rebreu instruccions sobre la formació d'equips de projectes i com inscriure'ls mitjançant els canals habituals (Racó, canal #projecte de Slack).

El lliurament del projecte es farà mitjançant el Racó, tindreu temps fins a les 08:00 hores del dia 24 de març de 2025.

I. OBJECTIUS

L'objectiu d'aquest projecte és per una part analitzar l'efectivitat de la detecció de documents de text similars, en concret mitjançant algorismes de *Locality-Sensitive Hashing* (LSH), en funció de la precisió de la representació del text i de la selecció de les funcions de hash. Per això us proposem que considereu representacions dels documents basades en l'ús de k -shingles i de signatures *minhash* per tal d'avaluar l'efectivitat de la mesura de similitud computada. Una vegada determinats els paràmetres adients per a una col·lecció, volem un algorisme per determinar els documents més similars d'un corpus de documents donat.

Una descripció exacta de les definicions i els mètodes la podeu trobar al capítol 3 del llibre *Mining of Massive Datasets* [1]. En concret la representació d'un document per un conjunt de k -shingles a la secció 3.2 i la representació d'un document per una signatura minhash amb t components obtinguda utilitzant t funcions de hash a la secció 3.3.2. Una versió electrònica del capítol la podeu trobar a <http://infolab.stanford.edu/~ullman/mmds/ch3.pdf>.

Us demanem que implementeu algorismes per:

- Obtenir la similitud de Jaccard de dos documents representats per conjunts de k -shingles. Recordeu que el grau de similitud de Jaccard de dos conjunts A i B es defineix com

$$\text{sim}(A, B) = \frac{|A \cap B|}{|A \cup B|}.$$

*La versió més actualitzada d'aquest document, així com qualsevol material addicional relacionat es publicarà al Racó.

- Obtenir una aproximació del grau de similitud de Jaccard a través d'una representació via signatures minhash basades en t funcions de hash. La mesura de similitud de dues signatures a i b amb t components es defineix com:

$$\text{sim}(a, b) = \frac{1}{t} \sum_{i=1}^t \begin{cases} 1 & \text{if } a_i = b_i \\ 0 & \text{if } a_i \neq b_i \end{cases}.$$

- Implementar un algorisme de locality-sensitive hashing basat en signatures minhash.

L'objectiu és per una part veure experimentalment la dependència entre precisió i temps de còmput en funció dels valors dels paràmetres k i t per diferents col·leccions de documents i famílies de funcions de hash i per una altra, analitzar experimentalment la qualitat de les solucions computades per l'algorisme de LSH.

II. ENTREGA

El nivell de sofisticació i esforç dedicat al projecte és opcional i es tindrà en compte a l'hora d'avaluar-la. En la versió més senzilla (suficient per tenir una bona nota si està acompanyada d'un bon disseny d'experiments) implementeu programes en C++ per als problemes proposats seguint la versió més simple del capítol 3 de [1] en particular l'algorisme de LSH per signatures minhash descrit a la secció 3.4. Versions més sofisticades del projecte inclouran la implementació d'algorismes no bàsics o de més d'un algorisme de LSH o de variacions de l'algorisme bàsic.

Tingueu en compte que haureu de mesurar el temps dels algorismes. A més, potser haureu de fer un seguiment de diversos comptadors que permetin quantificar el tipus i la quantitat de treball que el programa fa i la qualitat de la mesura/solució computada.

III. DADES

Aquest document és intencionadament vague. Per tant, a més d'analitzar i experimentar amb diferents funcions de hash i versions d'algorismes, haureu de documentar les fonts d'informació, les decisions preses (i el disseny d'experiments que els hi donen suport) i la selecció de conjunt de dades. Si feu servir un repositori públic, cal que proporcioneu l'adreça del repositori així com (si cal) les modificacions i/o simplificacions sobre els conjunts de dades triats. Per trobar informació addicional sobre funcions de hash aquests són bons punts de partida: <https://xxhash.com>, <https://github.com/jlumbroso/python-random-hash> (aquest darrer és una implementació en Python d'una família de funcions de hash senzilla i eficient, i és fàcil fer el *port* a C++).

Com a únic requisit respectat als experiments us demanem:

1. que un dels conjunts de documents sigui format per almenys 20 documents obtinguts mitjançant permutacions aleatòries d'un document de text “base” amb almenys 50 paraules diferents.
2. per un altre bloc d'experiments trebal·leu amb un conjunt de $D \geq 20$ “documents” virtuals: trebal·leu directament amb els D conjunts de k -shingles dels documents virtuals, cadascun d'aquests conjunts de k -shingles s'obté extraient una quantitat de k -shingles desitjats, de manera aleatòria, del conjunt de k -shingles “base” extret d'un document base amb almenys 100 paraules diferents, és a dir, tenim un cert grau de control sobre la similitud de dos documents: si dos “documents” tenen conjunts de k -shingles de n_i i n_j k -shingles extrets aleatòriament del conjunt base d' n k -shingles, llavors la seva similitud és, en mitjana, aproximadament

$$\approx \frac{p_i p_j}{p_i + p_j - p_i p_j},$$

on $p_i = n_i/n$ i $p_j = n_j/n$.

N.B. A l'hora d'extreure els k -shingles dels documents, i calcular-ne la similitud de documents considereu que es fa un preprocés filtrant i eliminant tots els anomenats *stopwords* (preposicions, verbs auxiliars, determinants, etc.) que són molt freqüents i formaran part de qualsevol text de llargària mitjana. Els *stopwords* matisen el sentit de les oracions, però no comporten significat, i distorsionen les mètriques de similitud. És molt senzill trobar llistes d'aquestes paraules a Internet, per exemple aquí: <https://github.com/stopwords-iso/stopwords-ca>. Quan diem que el document “base” que s'ha d'utilitzar en els dos blocs d'experiments obligatoris ha de tenir almenys 50 paraules diferents ens estem referint a paraules que **no** siguin *stopwords*.

IV. QUÈ CAL LLIURAR

Cal lliurar mitjançant el racó una carpeta comprimida (**.zip**) que contingui:

- Una documentació adequada dels algorismes i mètodes que heu implementat, les funcions de hash que heu fet servir, les proves que heu fet i la comparació dels resultats que heu obtingut. També és interessant que indiqueu altres idees que hagueu provat, encara que no hagin donat bons resultats, o d'altres que no heu explorat. El document en format **PDF** ha d'incloure les referències adients.
- Una carpeta amb totes les fonts necessàries per compilar i executar la vostra pràctica. S'han d'incloure les instruccions per a la compilació i execució, així com els conjunts de dades utilitzats o els programes per generar-los (si cal).

- Tingueu en compte que la documentació entregada ens ha de permetre valorar el nivell d'assoliment de la competència transversal que hem d'avaluar: **Capacitat d'autoaprenentatge**. En el context del projecte hi han uns quants aspectes rellevants relacionats amb aquesta competència: les funcions de hash, els algorismes per crear la representació de documents de text amb signatures minhash, els algorismes LSH, i el disseny i anàlisi dels experiments.
- La documentació ha de recollir i presentar la feina feta, les fonts que s'han consultat i els resultats de l'experimentació. Si no es compleix aquesta condició, la qualificació final del projecte reflectirà la qualitat de la presentació i no la del codi entregat.

REFERÈNCIES

- [1] Jure Leskovec, Anand Rajaraman and Jeffrey D. Ullman *Mining of Massive Datasets* Cambridge University Press