

Masovno otkrivanje dokumenata sličnih duplikata

Verzioniranje:

- Različite verzije jednog dokumenta
 - o Revizije, različiti formati ...

Zrcaljenje

- Objavljeno na više mjesta

Plagiranje

- Identična ili obrađena kopija

Štetni sadržaj

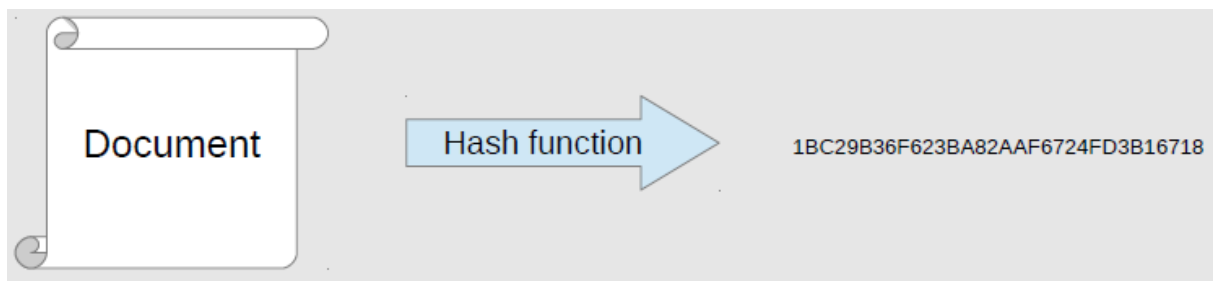
- Virusi, spam, ...

Skalabilno rješenje:

- Dokumenti i repozitoriji dokumenata su ogromni
 - o Internet
 - o Domena specifičnih tekstova
 - o Logovi
 - o ...

Checksuming

- Kriptografske hash funkcije
- MD5, SHA1, SHA2, ...



Hvata i najmanje uređivanje

- Dobro za detekciju identičnih kopija
- Nije dobro za detekciju sličnih duplikata

DETEKCIJA SLIČNIH DUPLIKATA

Dvije metode:

- Fingerprinting (identifikacija)
 - o Hashiranje dokumenta
- Rangiranje
 - o Tehnike pronalaženja informacija

FINGERPRINTING

Sličnost čuvajući raspršivanje

- X , skup ulaza
- d_x , funkcija udaljenosti nad X
 - o x_1, x_2 elements of X
- sličnost čuvajući hash funkciju
 - o $h: X \rightarrow Y$
 - o $|Y| < |X|$
- d_y , funkcija udaljenosti nad Y
- slični ulazi imaju slične hash vrijednosti

if $d_x(x_1, x_2) < \epsilon_x$, onda

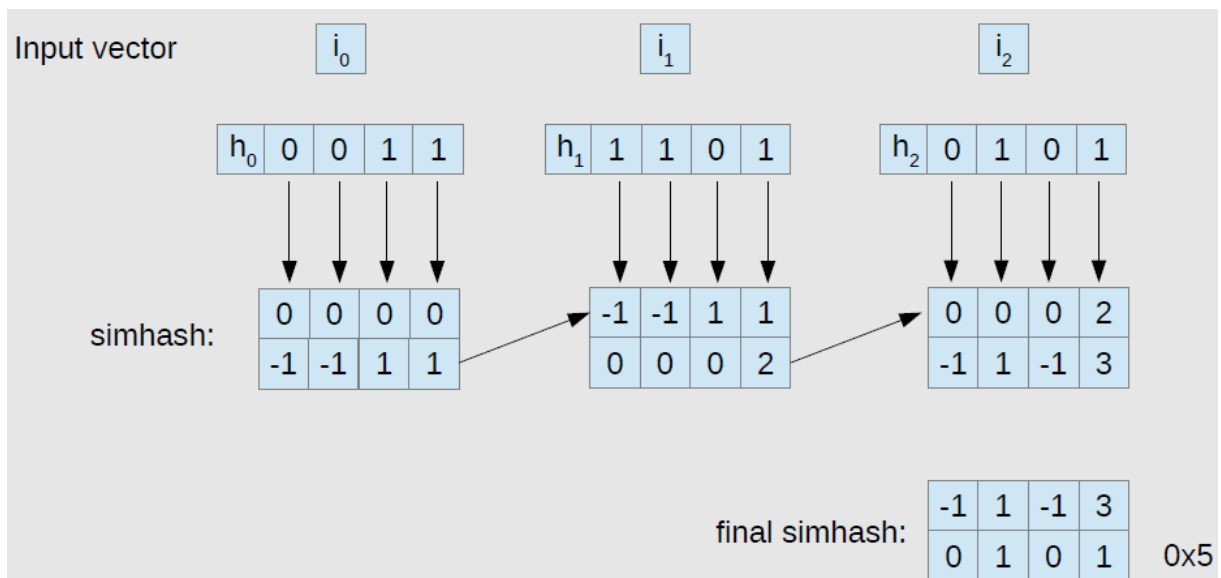
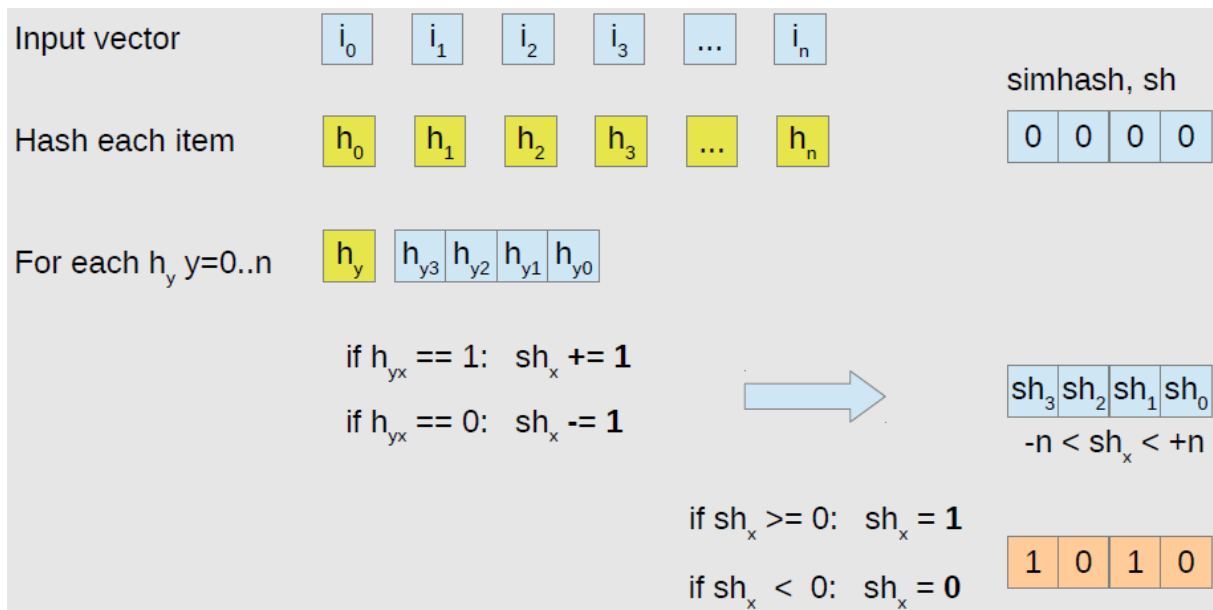
$$d_y(h(x_1), h(x_2)) < \epsilon_y$$

- primjer:
 - o $h(„text1“) = 0xaaaf$
 - o $h(„text2“) = 0xaaae$

Simhash algoritam

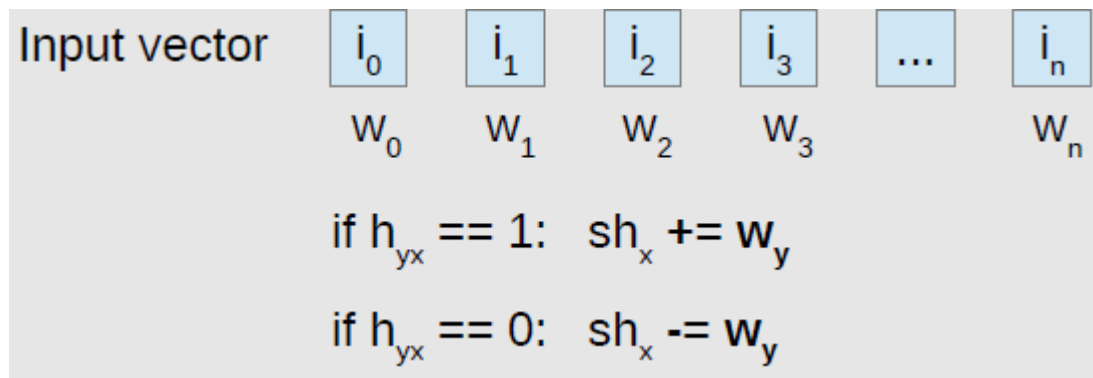
- Tehnika identifikacije
 - o Identifikacija sličnih duplikata razlikuje se u malom broj pozicija bit-ova (hamingova udaljenost)
- Smanjenje dimanzionalnosti
 - o Mape visoko-dimenzijskih vektora za male fingerprintove (fbits)
- Veličina fingerprinta (fsize)
 - o f je malen i proizvoljan

Računanje simhash-a za $f=4$



Ponderirana Simhash računanja

- Dodjela faktora težine na svaku novu značku



Izbor hash funkcije h

- Jedinstvena distribucija
- Brzina
- Kandidati: kriptografske hash funkcije
 - o MD5(128bit), SHA-1(256bit)
- Vlastite hash funkcije sa duljinom varijable

Ulazni vektor

- Fokus na sirovom tekstu dokumenata
- Konvertira dokument u oblikovni vektor
- Prikupljanje podataka
 - o Tokenizacija (shingling)
 - Unigrami, 2-grami, 3-grami
 - o Proizlazi
 - o Uklanjanje točke
 - o Detekcija fraza
 - o ...

“lorem ipsum dolor sit amet”

- “lorem”, “ipsum”, “dolor”, “sit”, “amet”
- “lorem ipsum”, “ipsum dolor”, “dolor sit”, “sit amet”
- “lor”, “ore”, “rem”, “em”, “m i”, ...

- Shingling
 - Hash k-grama
 - k-grami
 - Znakovi, riječi, rečenice
 - $k = ?$
 - malen k: različiti dokumenti se pojavljuju slično
 - velik k: slični dokumenti izgledaju različito
 - oblikovni vektor iz IR izlaza
 - Ponderirana sa IDF
 - inverzna učestalost dokumenta
 - može se promijeniti kada se promijene kolekcije
 - Informacije povezivanja
 - Struktura poveznica (slične stranice imaju zajedničke poveznice)
 - Sidreni tekst
 - Slični dokumenti bi trebali imati sličan sidreni tekst (poveznice)

BRZI UPITI

- F – kolekcija f-bit fingerprintova
- Q – upit
 - Jedinstveni ii skup fingerprintova
- Zadatak
 - Identificirati kad god se Q razlikuje od ijednog fingerprinta F u najviše k bitova
- Google numbers
- 8B 64-bitni fingerprintovi = 64Gb
- Online upiti
 - Q = jedan fingerprint
 - Ograničenja: nekoliko milisekundi
- Skupni upiti
 - Q = skup fingerprintova
 - e.g $|Q| = 1M$
 - Ograničenja: ~100sekundi
 - 1B upita po danu

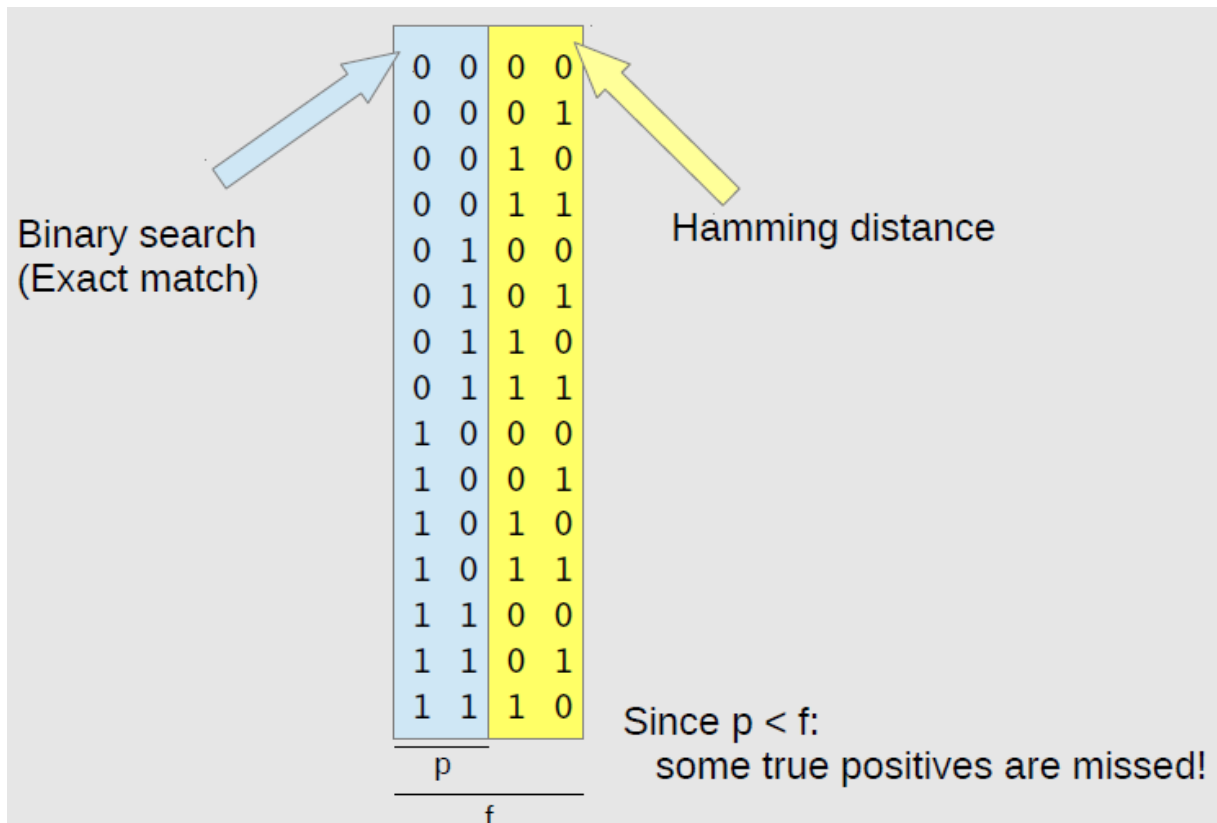
Prvi pristup:

- Napravi sortiranu tablu F-ova
- Napravi listu Q' sa svim fingerprintovima čija je hamingova udaljenost od Q najviše k

$$|Q| = \binom{64}{3} = 41664$$

Drugi pristup

- Napravi sortiranu tablicu F-ova
- Nađi skup fingerprintova (F') koji imaju jednako najviše značajnih dijelova (p bitova)
 - o Sortirana tablica – binarna pretraga $O(p)$.
- Provjeri hamingovu udaljenost za svaki fingerprint iz F'
- Ovaj pristup će locirati sve fingerprintove u F koji se razlikuju u najviše k bitova
 - o Ograničeno na najmanje značajnih $f-p$ bitova!



- Zamijenit žutu i plavu, + sort

Rješenje:

- Izgraditi dodatne tablice
 - o Svaka sa različitom permutacijom bitova
 - o Svaka tablica ima različit skup značajnih bitova
- Algoritam za brze (online) upite
 - o Napravi t sortiranih tablica od fingerprintova: T_1, T_2, \dots, T_t
 - o Svaka tablica T_i sadrži
 - P_i – broj značajnih bitova
 - Π_i – random permutacija
 - o Svaki fingerprint u T_i je permutiran sa permutacijom Π_i

Za dane Q i k

- Pročitaj svaku tablicu (u paraleli)
 - o Dohvati fingerprintove u T_i čiji se značajni P_i bitovi podudaraju sa značajnim P_i bitovima od $\Pi(Q)$
 - T_i'
 - $O(P_i)$ koraka (binarna pretraga)
 - o Za svaki fingerprint u T_i' , provjeri je li njegova hamingova udaljenost najviše k bitova od $\Pi(Q)$
- Primjer sa $t=20$, $f=64$, $k=3$, $|F| = 8B$
 - o Podjeli f u 6 blokova ($4 \times 11 + 2 \times 10$ bitova)
 - o Odaberi 3 od 6 blokova ($6 \text{povrh } 3 = 20$ putanja)
 - o Posloži te blokove kao značajne bitove
 - o $P = \text{zbroj tih bitova}$
 - 31,32 ili 33
 - o U prosjeku upit vraća $2^{34-31} = 8$ fingerprintova
 - o t i p parametri
 - $t \sim p$
 - vrijeme upita $\sim 1/p$
 - zahtjevi za pohranu $\sim p$
 - prostor/vrijeme
 - analitička solucija za t

Skupni upiti koristeći MapReduce i GFS

- F i Q su datoteke u GFS(sa replikacijom)
- $F \sim 64GB$, $Q \sim 8MB$
- F je spremljen u GFS komadima
- Broj mapera = broj F komada
- Map:
 - o Rješava hamingovu udaljenost za komad (64Mb) i emitira listu sličnih duplikata
- Reduce:
 - o Uklanja duplikate