

Come leggere la lista hackmeeting con TF-IDF

lucha@paranoici.org

24 giugno 2011

...un thread su lista hackmeeting

Tutti i tuoi movimenti su iphone?

DI COSA PARLA?

...un thread su lista hackmeeting

Tutti i tuoi movimenti su iphone?

DI COSA PARLA?

Ci viene in aiuto un algoritmo: TF-IDF.

- *TF* term frequency.
- *IDF* inverse document frequency.

come capiamo di cosa stiamo parlando

riconosciamo una *parola chiave* associando un punteggio di *rilevanza* ad ogni termine.

Ci viene in aiuto un algoritmo: TF-IDF.

- *TF* term frequency.
- *IDF* inverse document frequency.

come capiamo di cosa stiamo parlando

riconosciamo una *parola chiave* associando un punteggio di *rilevanza* ad ogni termine.

Il primo fattore del punteggio è la term frequency:

$$\text{tf}(\text{parola}) = \frac{\text{occorrenze parola nel testo}}{\text{numero di parole nel testo}}.$$

Da sola *non* è rilevante: congiunzioni, avverbi, etc.

Inverse Document Frequency

Consideriamo una collezione di documenti che costituiscono un insieme di riferimento *generico* (ovvero che non individuano un argomento specifico che vogliamo distinguere).

Il secondo fattore del punteggio è la inverse document frequency:

$$\text{idf}(\text{parola}) = \frac{\text{numero totale di documenti}}{1 + \text{numero di documenti nei quali compare parola}}.$$

Il valore di idf non dipende dal testo che vogliamo analizzare, ma solo dalla collezione di documenti che consideriamo come base neutra. Può quindi essere precalcolato.

tf-idf

$$\text{tf-idf} = \text{tf} * \text{idf} .$$

and the winner is

messico bada strano disgusto consensienti flat parlamento operatori elemosinante
rancido **sposarsi riscatto** location sud

machismo **vivano** disaccordo favori stronzi 10k **arma** schiava

normale interessava oppressione **lavoratrici**

unicamente figo **sincronizzare contatti spiegami**

ingravidarsi antifascista biologica **tendo rna** misogino prostituirsi dropbox

sicurezze gsm affidati database quindicenni zoccole parlamentari webmail
compro bieco cartello sceglie **condiviso terminale**

instabili guadagno crescono sciovinista volentieri **ruby**

gesto affidamento minetti codifica

appropriarsi silenzio linguistica **fascista** zona marrazzo **mese sic** sputo

brava moto ribrezzevole vessazioni elargitore bunga **aka** **preferiscono**

autobus **plausibile dogma** geografica **violentare** aperte ribrezzo maschile