

PRÉ-PROCESSAMENTO

Arquivo .db SQLite3
com texto em ascii

#include <sqlite3.h>

Importar texto p/ programa

Ler do arquivo "separadores.txt"
e usar os separadores p/ converter o
texto em vetor

Parecido com .split(...)

solução: strtok

Ler do arquivo "stopwords.txt"
e remover as stopwords.
Filtrá-las do vetor anterior.

#include <libstemmer.h>

Aplicar a função sb_stemmer_stem()
do SnowballStemmer em cada palavra
do vetor anterior.

Montar a matriz de frequência ou
índice de arquivo invertido.

Calcular o idf :

(<número de docs>)

PARALELISMO DE DADOS

BARREIRA →

"<word>": $\log_2 \left(\frac{\text{número de docs que}}{\text{<word> aparece}} \right)$



Calcular os vetores dos documentos:

"<docid>": $[tf(<word>, <docid>) * idf[<word>] \text{ for word in vocab}]$



Calcular as normas dos vetores anteriores

SALVAR
ESTRUTURA
RESULTANTE NUM
ARQ. BINÁRIO

CONSULTA

Consulta do usuário como argumento do programa:

cat consulta.txt | ./a.out
ou

./a.out "<consulta do usuario>"



Repetir as etapas de pré processamento.

PARALELISMO DE DADOS



Calcular similaridade:

Para cada docid:

PARALELISMO DE DADOS

$$\text{sim} = \frac{\langle \text{vetor do docid} \rangle * \langle \text{vetor da consulta} \rangle}{\langle \text{norma do docid} \rangle * \langle \text{norma da consulta} \rangle}$$



Retornar resultados em
ordem decrescente de
similaridade.