

IRS - Unit - 2 Imp.

1. BSB1 Algorithm.

With insufficient main memory, we need to use an external sorting algorithm, i.e., one that uses disk.

One solution is blocked-sort-based indexing algorithm or BSB1.

- BSB1 segments the collection into parts of equal size.
- BSB1 sorts term ID-doc ID pairs of each part in memory.
- BSB1 stores intermediate sorted results on disk and
- BSB1 merges all intermediate results into the final index.

Algorithm:

BSBINDEXCONSTRUCTION()

1. $n \leftarrow 0$
2. while (all documents have not been processed)
3. do $n \leftarrow n + 1$ \uparrow No space \uparrow No space
4. $block \leftarrow \text{PARSENEXTBLOCK}()$
5. BSB1-INVERT(block)
6. WRITEBLOCKTODISK(block, f_n)
7. MERGEBLOCKS($f_1, \dots, f_n; f_{\text{merged}}$)

2. Spelling-Correction Algorithm:

Ans: Indexing for spelling correction is a bit different than for document retrieval.

Most common spelling errors are:

1. Insertion: adding an extra letter. Ex: truly.
2. Deletion: missing a letter.
3. Transposition: Interchanging of a letter's position.
4. Substitution.

To overcome these errors & issues Spelling-Correction Algorithm is used.

Spelling Correction Algorithm:

1. Here, a Ternary Search tree stores dictionary.
2. For each stored word, we also keep frequency count, obtained from analysis of large corpora.
3. Queries entered in the search engine are parsed & individual terms are extracted with non word tokens ignored. Each word is then converted to lower case & checked to see if it is correctly spelled. These, correctly spelled words found in user queries are updated in the dictionary, by incrementing their frequency count.
4. Edit distance, k-gram indexes are the two types of basic Spelling correction.

Algorithm:

EDIT DISTANCE (s_1, s_2)

1. $\text{int } m[|s_1|, |s_2|] = 0$
2. for $i \leftarrow 1$ to $|s_1|$
3. do $m[i, 0] = i$
4. for $j \leftarrow 1$ to $|s_2|$
5. do $m[0, j] = j$
6. for $i \leftarrow 1$ to $|s_1|$
7. do for $j \leftarrow 1$ to $|s_2|$
8. do $m[i, j] = \min\{m[i-1, j-1] + 1 \text{ if } (s_1[i] = s_2[j]) \text{ then } 0 \text{ else } 1, m[i-1, j] + 1, m[i, j-1] + 1\}$.
- 9.
10. return $m[|s_1|, |s_2|]$.

Stemming Algorithm:

It is a process of linguistic normalisation, in which ~~various~~ variant forms of a word are reduced to a common form. for example, connection connections connective → connect connected connecting.

- It reduces a word to its base word using various approaches.
- It is simple to develop.
- There are two ^{types of} errors that occur. They are:

(i) Over-stemming:

It occurs when 2 words are stemmed from same root that are of different stems. It can also be regarded as false-positive.

(ii) Under-stemming:

It occurs when 2 words are stemmed from same root that are not of different stems. It can be interpreted as false-negatives.

→ There are various stemming algorithms - Porter's stemming algorithm, Lovins stemmer, Dawson stemmer, Krovetz stemmer etc.

APPLICATIONS:

Stemming is used in information retrieval systems like search engines.

It is used to determine domain vocabularies in domain analysis.