

CHAPTER #3 DICTIONARIES & TOLERANT RETRIEVAL

★ Data Structures for Dictionary =>

↳ Hashtables =>

↳ Every vocab is hashed to an integer.

↳ Lookup is fast $O(1)$

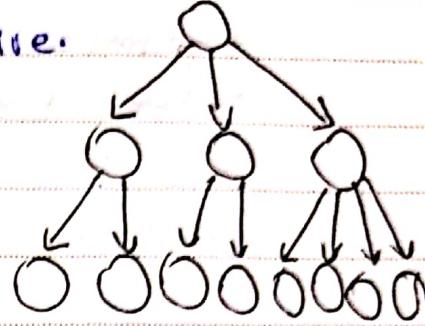
↳ Problems =>

↳ No easy way to find variants.

↳ No prefix search.

↳ may need to rehash over growing vocabulary.

- ↳ Trees => ~~to make for a better presentation of the material~~
- ↳ Simplest => Binary Trees
- ↳ Rebalancing binary trees is expensive.
- ↳ Usual Approach => B-Trees
- ↳ mitigate the rebalancing problem
- ↳ Pros =>
- ↳ solve the prefix retrieval problem



- ↳ Cons =>

↳ Slower $O(\log M)$ {M is the height}

★ Wildcard Queries => There are two situations for wildcard queries:

- ↳ user is uncertain of the spelling
- ↳ user is aware of multiple variants & seeking all.
- ↳ Trailing wildcards =>
 - ↳ Queries are of the form "mon*"
 - ↳ find all documents containing words starting with mon.
 - ↳ easy with binary & b-trees.

- ↳ Leading wildcards =>

↳ Queries are of the form "*mon"

↳ find all documents containing words ending with mon.

↳ harder since a reverse B-tree is required additionally.

- ↳ center wildcards =>

↳ Queries of the form "co*tion"

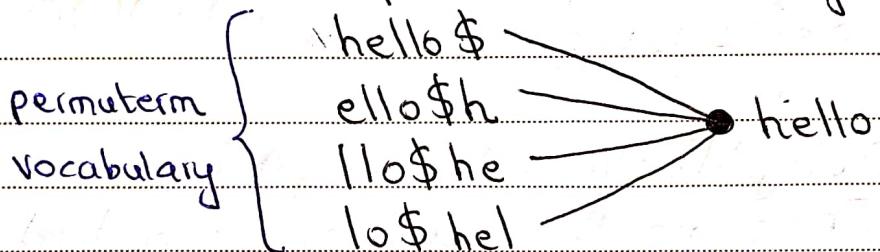
↳ first lookup "co*", then *tion" in a B-tree and then intersect the two sets.

↳ It is extremely expensive.

- ↳ Techniques to handle general wildcards => The strategy is to express the wildcard query q_{W} as a Boolean query Q on a specially constructed index such that the answer to Q is a superset of the set of vocabs matching q_{W} .

Then, we check each term in Q against q_{W} , discarding mismatches, then resort to inverted indexes.

↳ Permuterm Indexes ⇒ Firstly, we introduce a new symbol \$ which marks the end of the term.
 e.g: hello → hello\$. Next, construct a permuterm index in which various rotations of each term all link to the original vocabulary term. We call the rotations as permuterm vocabulary.



↳ How are queries looked up?

Consider the query "m*n", firstly rotate the query to turn this into a trailing wildcard query e.g: "n*m\$". Next, we lookup this string into a permuterm index. Finally, we can lookup in B-Trees for finding the original vocabulary like "man" and "moron". Finally, we can use inverted index to retrieve the results.

Wildcard query → rotate → permuterm index → B-Tree → Inverted

↳ Multiple wildcards ⇒ Index

For a query such as "fit*master". Firstly, we enumerate the terms involving "er\$fi*". Not all dictionary terms will have "mo" in the middle, so we filter such terms by "exhaustive enumeration", terms like fishmanger would survive but fillbuster will not. Then inverted index can be used to find the results.

↳ Problems ⇒

↳ Quadruples Lexicon size due to permutative rotations.

efficiency ++

- ↳ Bigram (k-gram) indexes => In the k-gram indexes, the dictionary contains all k-grams that occur in any term in the vocab. Each posting list points from a k-gram to the terms containing that k-gram.
e.g: "cruelst month" will have the bigrams → \$c, cr, ru, ue, el, le, es, st, t\$, \$m, mo, on, nt, th, h\$"
- \boxed{sm} → mace → madden → ...
- \boxed{mo} → among → amortize → ...

↳ How are queries looked up? \rightarrow always keep + attend
Consider the query "man*", the query can now be run as "\$m AND mo AND man". But we'd end up enumerating "moon" as well. So a post-filtering step is introduced so string compare the results with the original query. Survivors are checked in inverted indexes.

* Spelling Corrections =>

↳ Principle uses

↳ correcting documents being indexed

↳ Will / won't catch typos & correcting user queries

↳ Isolated Word Correction => Given a lexicon and a character sequence Q, return the words in lexicon closest to Q. We don't care about context and words originally in. Correct spellings won't be fixed.

The following strategies can be used to implement this:

↳ Levenshtein Distance =>

Given two strings S_1 and S_2 , the Levenshtein distance is the minimum number of operations to convert one to the other. The operations are:
Insert, Replace, Delete.

e.g. \Rightarrow dof → dog cost = 1

cat → act cost = 2

↳ ↳ ↳ EDIT DISTANCE (s_1, s_2):

- 1 int $m[i][j] = \{0\}$
- 2 for $i \leftarrow 1$ to $|s_1|$ } s_2 is empty so increment
 $m[i, 0] = i$
- 3 for $j \leftarrow 1$ to $|s_2|$ } s_1 is empty so increment
 $m[0, j] = j$
- 4 for $i \leftarrow 1$ to $|s_1|$: diagonal + 0 if same
- 5 for $j \leftarrow 1$ to $|s_2|$: diagonal + 1 if different
- 6 $m[i, j] = \min \{m[i-1, j-1] + s_1[i] = s_2[j] ? 0 : 1,$
- 7 $m[i-1, j] + 1, \} \text{ above} + 1$
- 8 $m[i, j-1] + 1\} \} \text{ left} + 1$
- 9
- 10
- 11 return $m[|s_1|, |s_2|]$ } last index is the solution

ex:

	\emptyset	f	a	s	t	
\emptyset	0	1	2	3	4	\emptyset
f	1	1 + 1 → 2	2 + 1 → 3	3 + 1 → 4		c
a	2	2 + 1 → 1	1 + 1 → 2	2 + 1 → 3		a
s	3	3 + 1 → 2	2 + 1 → 3	3 + 1 → 2		t
t	4	4 + 1 → 3	3 + 1 → 2	2 + 1 → 3		s
cost = 3						

↳ How to use?

↳ The simplest idea is to find all strings with minimum edit distance and then return them but this is extremely expensive.

One efficient heuristic is to restrict the search dictionary terms beginning with the same letter assuming the first letter is always correct.

A sophisticated variant is to use permuterm index while omitting \$. We can also refine the rotation scheme by omitting 'l' characters before each rotation so to have best matches. Then, perform B-Tree traversal.

↳ ↳ ↳ \hookrightarrow k-gram overlaps \Rightarrow We can improve the use of edit distance with k-gram indexes. The method is to use k-gram indexes to find the relevant k-grams then use edit distances over them to find the most relevant. Consider bi-grams of lord;

lo \rightarrow alone \rightarrow lord \rightarrow sloth

or \rightarrow border \rightarrow lord \rightarrow morbid

rd \rightarrow ardent \rightarrow border \rightarrow card

then merging can be performed to retrieve relevants.

Linear merging is still troublesome since for queries like "bord" will retrieve "boardroom", so a measure called "Jaccard Coefficient" is adapted. This is a measure of similarity b/w two terms. It is the measure of overlap $|A \cap B| / |A \cup B|$.

To compute the $\frac{\text{co-eff}}{\text{k-gram}}$, we need to find the k-grams of query and terms. Consider, q₁=bord and t=boardroom, the jaccard co-efficient can be calculated by dividing the number of k-gram hits by the sum of total bigrams - posting hits which

$$\text{is} \Rightarrow \{bo, or, rd\} \rightarrow 2 \\ \{bo, oa, ar, rd, di, ro, oo, om\} \rightarrow 8+3-2 = 0.222.$$

Finally, we can check if the jaccard co-efficient exceeds a certain threshold, we add the term to the answer.

\hookrightarrow Context Sensitive Spelling Correction \Rightarrow

Isolated words will fail to solve typographical errors when all terms are spelled correctly such as: "flew form Heathrow". We can try all permutations for each term but that would be computationally expensive. The other heuristics are to find hit-based spelling collections and then among

the retrieved results we can use the alternative with most hitting docs or by query log analysis to find relevant user queries to this query.

↳ Phonetic Correction =>

The main idea is to generate a phonetic hash so that similar sounding terms are always together.

↳ Soundex Algorithm =>

This algorithm is used to generate the phonetic hashes.

↳ Scheme =>

↳ Turn every term into a 4-character reduced form, then build an inverted index from these reduced forms to original words called the Soundex Index.

↳ Do the same with query terms.

↳ Run the queries over the soundex index.

↳ Encoding => The first character is the first letter of the term, rest are digits from 0-9.

↳ How to convert?

↳ Retain the first letter of the term.

↳ change all occurrences of A,E,I,O,U,H,Y,W to 0.

↳ change following

B,F,P,V → 1

C,G,J,K,Q,S,X,Z → 2

D,T → 3

L → 4

M,N → 5

R → 6

↳ Remove all pairs of consecutive digits.

Dated: _____

↳ Remove all zeros.

↳ Pad trailing zeros and return the first four characters.

↳ Query Processing

↳ A query like "Hermann" is processed like

Hermann → HOrmDna → H065055 → H6555

return H655 [first four]

↳ Observations

↳ Vowels are viewed as interchangeable in transcribing names.

↳ Consonants with similar sounds are put in equivalence classes.

↳ While it works to some extent, it is accurate to some extent but not completely.
e.g.: Rafi, Rafee have different soundex codes.

↳ Limitations

↳ It can handle only English words.

↳ It gives false +ves and false -ves.

↳ It does not consider semantic meaning.