

Assignment #2

Inverted Index : Inverted index is a mapping from content(words in this case) to their location in document. The purpose of this is to provide rapid full word searches at the cost of higher processing time at the time word is entered in the inverted index.

Process:

1. **Decompress:** To start with we are given a sum total for more than 3 million documents. The documents have been compressed using gzip so first we decompress this using “java.util.zip.GZIPInputStream” library. This allows us to uncompress the file in tsv(tab separated value) format. Note: I was unable to download the TREC format suggested as the download keeps on failing halfway through. Gzip is a lossless compression, so we can reconstruct original data.
2. **Program logic:**
 - In this code we will do a double pass of data. The first pass is to estimate the barrel size for lexicon partitioning while generating postings. In this pass we create a HashMap (frequencyGraph) and calculate frequency of each word as well as the number of documents in which the word occurs. This allows us to more precisely maintain allot a dynamic barrel size for each lexicon. We will be maintaining this table in the main memory. For all the terms this table seems to occupy around 1 GB of main memory.
 - Next is Term to TermID mapping: We are maintaining a HashMap(termIDS) which maps terms to term IDs. We are allotting term id sequentially. Each term has a unique term id. We are maintaining this data structure in main memory. Using Integer instead of string in the lexicon table allows for compression.
 - Next is MongoDB database which maintains document text mapped to document ID, this allows us to query data at last and highlight relevant(search query) data. We are doing this by forming a simple database URL and a collection in it called URLCollect. We are filling the database during the first pass of data.
 - Next is Lexicon data which is stored in a HashTable. The key for the hash table is the termID, which we get from HasMap TermIDS. While the value is an array of “Document frequency” i.e. number of documents the word has appeared in at least

once and the location in the inverted index file where the postings for that occurrence is stored.

3. The inverted Index Table:

- We can use RandomAccessFile to write to and read from any point the inverted index file.
- We also fill the lexicon HashMap in which barrel size is determined by the frequency details we gathered during the first pass.
- We iterate through each word in the collection of documents. We then split document ID, URL, title and data on the basis of tab(\t).
- We check if the word is in termsIDS HashMap. If not then, we assign it a new term ID sequentially. We add all the words to two HashMap,
 1. The wordposition has the term as key and an integer array which contains the position of the occurrences of that term in the document.
 2. The mapforpage HashMap which contains frequency of occurrence of each word.
- After iterating a document we have a Map containing each word that appears in it and another Map containing locations of terms each occurrence.

We store data in the form

DocID	frequency	position1	position2	position3	position4
-------	-----------	-----------	-----------	-----------	-----------

Difference postings are delimited by “ > ” (**greater than**), while (DocID,frequency) pairs and positions are delimited by a “ : ” (**a colon**). DocID and frequency is delimited by a comma and so are different positions.

Which would look like:

docID,frequency:locations1,location2,location3>docID2,frequency2:locations1,location2,location3;x

x denotes end for expression

- We start by iterating mapforpage Hashmap, for each term we find the relevant termID in the termIDS map. From this termID we find the location data in the

lexicon HashTable. We then seek this address using RandomAccessFile and read a number of bytes based on frequencyGraph. We process the bytes to string. If the string is empty we insert the DocID,frequency pair, getting frequency from mapforpage. This is followed by inserting position data by searching the key in wordposition Map. we add all the locations and delimit it with a **colon**.

If the String is not empty we split the data according to aforementioned delimiters and put it into a temporary HashMap foreachDocument with term as key and an integer array. The first element of the array is the frequency and the remaining the position. We also insert data from mapforpage and wordposition in the temporary hashmap. Now we sort the temporary HashMap according to the frequency in descending order(first element of value). This allows us to store documents with most frequent occurrence at the forefront. Now we insert back these elements.

- The HashTable *lexicon* and HashMap *termIDS* and *frequencyGraph* are stored in file “lexicon”, “termID” and “frequency” respectively.

5. For querying data:

- We read HashTable *lexicon* and HashMap *termIDS* and *frequencyGraph* from file “lexicon”, “termID” and “frequency” respectively.
- For querying data we are provided with a search query, we lookup this query term in the termIDS map and find corresponding term ID. We then look at the termID in the lexicon map and find the address of posting information in the inverted index file. We read the data (how much according to frequencyGraph), split it according to the delimiters set and get the first 10 results. Since these postings were sorted during time of insertion we don't have to sort it. Just take the top 10 results, query DocID in the MongoDB database to find the data stored and use location data to highlight the query where we found it.

Search term “there” being encolde by ←** **----->

mavenproject1 - Apache NetBeans IDE 12.1

File Edit View Navigate Source Refactor Run Debug Profile Team Tools Window Help

Search (Ctrl+F)

299.2/688.0MB

Output - Run (Web)

Enter a string:-----
there

CSF Analysis Share this page: Was this page helpful? Also known as: Spinal Fluid Analysis Formal name: Cerebrospinal Fluid Analysis Related tests: Glucose Total Protein
-----** there **----- other reasons to do a lumbar puncture? Why do I need a spinal tap? Why can't my blood or urine be tested? What other tests may be done in addition?
special needle is inserted through the skin between two vertebrae and into your spinal canal An opening or initial pressure reading of the CSF is obtained The healthcare provider
most patients it is a moderately uncomfortable procedure The most common sensation is a feeling of pressure when the needle is introduced Let your healthcare provider know
traumatic tap which just means that a small amount of blood may leak into one or more of the samples collected While this is not ideal it may happen a certain percentage
is often the best sample to use for conditions affecting your central nervous system because your CSF surrounds your brain and spinal cord Changes in the elements of your
to detect the source of the infection that led to meningitis or encephalitis Blood glucose total protein to compare with the concentration of CSF glucose and protein
Time for query execution: 3ms

In World War 2 the three great Allied leaders against Germany were President Roosevelt Winston Churchill and Joseph Stalin of the Soviet Union Three finer war leaders
at the start of the war But none of the other Axis leaders (certainly not Mussolini!) were of the caliber of the three Allied leaders Hitler was outnumbered three to one
The only thing we have to fear is fear itself President Roosevelt Roosevelt was one of America's greatest presidents To me he has always been the perfect example
Although he often had to compromise with his rich and powerful friends he was a true champion of the forgotten man President Roosevelt was conscious of the power of
and that Germany was wrong in its forced acquisition of surrounding nations Even though he knew which side we must support he could make no overt maneuvers because this
Nazi Germany's hands in 1939 1940 and 1941 he was unable to attempt a rescue because American public support was not there United States Prepares For War All in all the
capacity to produce 10 000 planes per year thereafter This advice did not set well with some advisors who wanted a more balanced approach But Roosevelt was trying to lead
It also magnified the power of the military chiefs which was a good thing as the nation prepared for war 3 In September of 1940 Roosevelt made a destroyer deal with Great
fire one does not haggle over the price to put it out the hose is readily loaned and the price is figured later America was slowly getting ready for war! U S Enters
could not help but feel a great relief (actually he was ecstatic) that America was now in the war He later wrote No American will think it wrong of me if I proclaim that
after seventeen months of lonely fighting We had won the war England would live; Our history would not come to an end Hitler's fate was sealed
doomed as the nations who made up the commonwealth began to fly the coop after the war and finally Britain was the only significant power left Although Britain ended up
end because of Roosevelt's hesitancy in joining Churchill in moving to limit Soviet Union expansionism This close relationship between the two leaders of course spilled
as the real Commander-in-Chief of the U S armed forces This is an authority that all presidents have a constitutional right to exercise but few do Roosevelt was part
the initiative and a master plan had been developed so there was no need for President Roosevelt to be so heavily involved Also as the war proceeded the military command

Run (Web) 50% 360:3 INS

Type here to search

6:48 PM 10/21/2020