Web Data Compression and Search

Search, index construction and compression

Slides modified from Hinrich Schütze and Christina Lioma slides

Inverted Index

For each term *t*, we store a list of all documents that contain *t*.



Inverted index construction

1 Collect the documents to be indexed:

Friends Romans countrymen So

so

countryman

Friends Romans countrymen II.So. let it be with Caesar

2 Tokenize the text, turning each document into a list of tokens:

3 Do linguistic preprocessing, producing a list of normalized tokens, which are the indexing terms: [friend] roman

4 Index the documents that each term occurs in by creating an inverted index, consisting of a dictionary and postings.

Tokenizing and preprocessing



Data 1. I did enact julius caesar I was killed if the capitol brutus killed me Data 2. Iso let it be with caesar the noble brutus hath told you caesar was ambitious

Generate posting



Sort postings

term	docID		term	docID
i	1		ambitio	us 2
did	1		be	2
enact	1		brutus	1
julius	1		brutus	2
caesar	1		capitol	1
i	1		caesar	1
was	1		caesar	2
killed	1		caesar	2
i'	1		did	1
the	1		enact	1
capitol	1		hath	1
brutus	1		i	1
killed	1		i	1
me	1	\implies	i'	1
so	2		it	2
let	2		julius	1
it	2		killed	1
be	2		killed	1
with	2		let	2
caesar	2		me	1
the	2		noble	2
noble	2		so	2
brutus	2		the	1
hath	2		the	2
told	2		told	2
уоц	2		уоц	2
caesar	2		was	1
was	2		was	2
ambitio	us 2		with	2

Create postings lists, determine document frequency

term	docIC)			
ambitiou	is 2	2			
be	1	2	term doc. freg.	\rightarrow	postings lists
brutus	1	1	ambitious 1	_	2
brutus	2	2	he 1	_	2
capitol	1	L	be 1	_	E . A
caesar	1	l	brutus 2	-	$1 \rightarrow 2$
caesar	2	2	capitol 1	\rightarrow	
caesar	2	2	caesar 2	\rightarrow	$1 \rightarrow 2$
did	1	1	did 1	\rightarrow	1
enact	1	1	enact 1	\rightarrow	1
hath	1	L	hath 1	\rightarrow	2
i	1	1	i 1	\rightarrow	1
i	1	1	ř 1	\rightarrow	1
i'	1	\rightarrow	it 1	\rightarrow	2
it		2 1	iulius 1		1
julius	1	L	killed 1		1
killed	1	L	kined 1	-	1
killed	1	L	let 1	\rightarrow	
let	-	2	me 1	\rightarrow	1
me	1	1	noble 1	\rightarrow	2
noble	2	2	so 1	\rightarrow	2
so	2	2	the 2	\rightarrow	$1 \rightarrow 2$
the	1	1	told 1	\rightarrow	2
the	-	2	you 1	\rightarrow	2
told	-	2	was 2	\rightarrow	$1 \rightarrow 2$
you	2	2	with 1	_	
was	1	L	1111 1	-,	<u> </u>
was		2			

with

2

Split the result into dictionary and postings file



Simple conjunctive query (two terms)

- Consider the query: BRUTUS AND CALPURNIA
- To find all matching documents using inverted index:
 - 1 Locate BRUTUS in the dictionary
 - 2 Retrieve its postings list from the postings file
 - 3 Locate CALPURNIA in the dictionary
 - 4 Retrieve its postings list from the postings file
 - Intersect the two postings lists
 - 6 Return intersection to user

Intersecting two posting lists



- This is linear in the length of the postings lists.
- Note: This only works if postings lists are sorted.

Intersecting two posting lists

Inte	$\operatorname{CRSECT}(p_1, p_2)$
1	answer $\leftarrow \langle \rangle$
n an Anna	aliteran (2000 and 2000 and 20
	da li dodQ(mij = dodQ(mij
an a	dhen KDD(answer, diciQ(a))
. 5	and section
6	$p_2 \in aext(p_2)$
7	alse if doc/R(py) < doc/R(py)
	Cherry Art Coexc (A)
6	else $p_2 \leftarrow neuc(p_2)$
	Maluunin andraker

Typical query optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR
- Simple and effective optimization: Process in order of increasing frequency
- Start with the shortest postings list, then keep cutting further
- In this example, first CAESAR, then CALPURNIA, then BRUTUS

Recall basic intersection algorithm



- Linear in the length of the postings lists.
- Can we do better?

Skip pointers

- Skip pointers allow us to skip postings that will not figure in the search results.
- This makes intersecting postings lists more efficient.
- Some postings lists contain several million entries so efficiency can be an issue even if basic intersection is linear.
- Where do we put skip pointers?
- How do we make sure intersection results are correct?

Basic idea

89 92 CAESAR



Skip lists: Larger example



Intersection with skip pointers



Where do we place skips?

- Tradeoff: number of items skipped vs. frequency skip can be taken
- More skips: Each skip pointer skips only a few items, but we can frequently use it.
- Fewer skips: Each skip pointer skips many items, but we can not use it very often.

Phrase queries

- We want to answer a query such as [stanford university] as a phrase.
- Thus The inventor Stanford Ovshinsky never went to university should not be a match.
- The concept of phrase query has proven easily understood by users.
- About 10% of web queries are phrase queries.
- Consequence for inverted index: it no longer suffices to store docIDs in postings lists.
- Two ways of extending the inverted index:
 - biword index
 - positional index

Positional indexes

- Postings lists in a nonpositional index: each posting is just a docID
- Postings lists in a positional index: each posting is a docID and a list of positions

Positional indexes: Example

```
Query: "to<sub>1</sub> be<sub>2</sub> or<sub>3</sub> not<sub>4</sub> to<sub>5</sub> be<sub>6</sub>"
то, 993427:
     < 1: <7, 18, 33, 72, 86, 231>;
      2: <1, 17, 74, 222, 255>;
      4: <8, 16, 190, 429, 433>;
      5: <363, 367>;
      7: <13, 23, 191>; . . . >
BE, 178239:
     < 1: <17, 25>;
      4: <17, 191, 291, 430, 434>;
      5: <14, 19, 101>; . . . > Document 4 is a match!
```

Inverted index



Dictionaries

- The dictionary is the data structure for storing the term vocabulary.
- Term vocabulary: the data
- Dictionary: the data structure for storing the term vocabulary

Dictionary as array of fixed-width entries

- For each term, we need to store a couple of items:
 - document frequency
 - pointer to postings list
 - • •
- Assume for the time being that we can store this information in a fixed-length entry.
- Assume that we store these entries in an array.

Dictionary as array of fixed-width entries

term	document	pointer to
	frequency	postings list
а	656,265	\longrightarrow
aachen	65	\longrightarrow
zulu	221	\longrightarrow

space needed: 20 bytes 4 bytes 4 bytes

How do we look up a query term q_i in this array at query time? That is: which data structure do we use to locate the entry (row) in the array where q_i is stored?

Data structures for looking up term

- Two main classes of data structures: hashes and trees
- Some IR systems use hashes, some use trees.
- Criteria for when to use hashes vs. trees:
 - Is there a fixed number of terms or will it keep growing?
 - What are the relative frequencies with which various keys will be accessed?
 - How many terms are we likely to have?

Hashes

- Each vocabulary term is hashed into an integer.
- Try to avoid collisions
- At query time, do the following: hash query term, resolve collisions, locate entry in fixed-width array
- Pros: Lookup in a hash is faster than lookup in a tree.
 - Lookup time is constant.
- Cons
 - no way to find minor variants (resume vs. résumé)
 - no prefix search (all terms starting with automat)
 - need to rehash everything periodically if vocabulary keeps growing

Trees

- Trees solve the prefix problem (find all terms starting with automat).
- Simplest tree: binary tree
- Search is slightly slower than in hashes: O(logM), where M is the size of the vocabulary.
- O(logM) only holds for balanced trees.
- Rebalancing binary trees is expensive.
- B-trees mitigate the rebalancing problem.
- B-tree definition: every internal node has a number of children in the interval [a, b] where a, b are appropriate positive integers, e.g., [2, 4].

Sort-based index construction

- As we build index, we parse docs one at a time.
- The final postings for any term are incomplete until the end.
- Can we keep all postings in memory and then do the sort inmemory at the end?
- No, not for large collections
- At 10–12 bytes per postings entry, we need a lot of space for large collections.
- But in-memory index construction does not scale for large collections.
- Thus: We need to store intermediate results on disk.

Same algorithm for disk?

- Can we use the same index construction algorithm for larger collections, but by using disk instead of memory?
- No: Sorting for example 100,000,000 records on disk is too slow – too many disk seeks.
- We need an external sorting algorithm.

"External" sorting algorithm (using few disk seeks)

- We must sort 100,000,000 non-positional postings.
 - Each posting has size 12 bytes (4+4+4: termID, docID, document frequency).
- Define a block to consist of 10,000,000 such postings
 - We can easily fit that many postings into memory.
 - We will have 10 such blocks.
- Basic idea of algorithm:
 - For each block: (i) accumulate postings, (ii) sort in memory, (iii) write to disk
 - Then merge the blocks into one long sorted order.

Merging two blocks



Blocked Sort-Based Indexing

BSBINDEXCONSTRUCTION()

- $1 \quad n \leftarrow 0$
- 2 while (all documents have not been processed)
- 3 do $n \leftarrow n+1$



Problem with sort-based algorithm

- Our assumption was: we can keep the dictionary in memory.
- We need the dictionary (which grows dynamically) in order to implement a term to termID mapping.
- Actually, we could work with term, docID postings instead of termID, docID postings . . .
- ... but then intermediate files become very large. (We would end up with a scalable, but very slow index construction method.)

Single-pass in-memory indexing

- Abbreviation: SPIMI
- Key idea 1: Generate separate dictionaries for each block no need to maintain term-termID mapping across blocks.
- Key idea 2: Don't sort. Accumulate postings in postings lists as they occur.
- With these two ideas we can generate a complete inverted index for each block.
- These separate indexes can then be merged into one big index.

SPIMI-Invert

6

SPIMI-INVERT(token_stream)

- 1 $output_file \leftarrow NEWFILE()$
- 2 dictionary ← NEwHASH()
- 3 while (free memory available)
- 4 **do** token ← next(token_stream)
- 5 **if** term(token) ∉ dictionary
 - then postings_list ← ADDTODICTIONARY(dictionary,term(token))
 - else postings_list ← GETPOSTINGSLIST(dictionary,term(token))

Reading the section of the section o

ട് പട്ടിന്റെ പ്രത്യായും പ്രത്യം നിന്നും പ്രത്യാം പട്ടിന്റെ പട്ടിയും പട്ടിന്റെ പട്ടിന്റെ പട്ടിന്റെ പട്ടിന്റെ പട്

- 16 ADTTEPEETINGSLIST(pastings-list,dack2(token))
- 11 seried_terms ← SORTTERMS[dictionary];
- 12 WRITEBLOCKTODISK(seried_terms, dictionary, autput_file)
- 13 reburn caspet_file

ⁿ la entre gable de la company de la company

Why compression in

Dictionary compression

- The dictionary is small compared to the postings file.
- But we want to keep it in memory.
- Also: competition with other applications, cell phones, onboard computers, fast startup time
- So compressing the dictionary is important.

Recall: Dictionary as array of fixed-width entries



Space needed: 20 bytes 4 bytes 4 bytes 4 bytes for Reuters: (20+4+4)*400,000 = 11.2 MB

Fixed-width entries are bad.

- Most of the bytes in the term column are wasted.
 - We allot 20 bytes for terms of length 1.
- We can't handle HYDROCHLOROFLUOROCARBONS and SUPERCALIFRAGILISTICEXPIALIDOCIOUS
- Average length of a term in English: 8 characters
- How can we use on average 8 characters per term?

Dictionary as a string



Space for dictionary as a string

- 4 bytes per term for frequency
- 4 bytes per term for pointer to postings list
- 8 bytes (on average) for term in string
- 3 bytes per pointer into string (need log₂ 8 · 400000 < 24 bits to resolve 8 · 400,000 positions)
- Space: 400,000 × (4 +4 +3 + 8) = 7.6MB (compared to 11.2 MB for fixed-width array)

Dictionary as a string with blocking



Space for dictionary as a string with blocking

- Example block size k = 4
- Where we used 4 × 3 bytes for term pointers without blocking . . .
- ...we now use 3 bytes for one pointer plus 4 bytes for indicating the length of each term.
- We save 12 (3 + 4) = 5 bytes per block.
- Total savings: 400,000/4 * 5 = 0.5 MB
- This reduces the size of the dictionary from 7.6 MB to 7.1 MB.

Lookup of a term with blocking: (slightly) slower



Front coding

One block in blocked compression $(k = 4) \dots$ **8** a u t o m a t a **8** a u t o m a t e **9** a u t o m a t i c **10** a u t o m a t i o n \downarrow ... further compressed with front coding. **8** a u t o m a t * a **1** ¢ e **2** ◊ i c **3** ◊ i o n

Dictionary compression for Reuters: Summary

data structure	size in MB
dictionary, fixed-width	11.2
dictionary, term pointers into string	7.6
~, with blocking, $k = 4$	7.1
~, with blocking & front coding	5.9

Postings compression

Key idea: Store gaps instead of docIDs

- Each postings list is ordered in increasing order of docID.
- Example postings list: COMPUTER: 283154, 283159, 283202, . . .
- It suffices to store gaps: 283159-283154=5, 283202-283154=43
- Example postings list using gaps : COMPUTER: 283154, 5, 43, ...
- Gaps for frequent terms are small.
- Thus: We can encode small gaps with fewer than 20 bits.

Gap encoding

	encoding	postings list								
THE	docIDs		283042		283043		283044		283045	
	gaps			1		1		1		
ECONOPUENEZE.	ne: lBs		253947		262153		352159		DK\$292	
	用的现在			1 <i>97</i>		<u>1</u>		∂Q		
ABACONCERSIO	arc n≪De sicilities 26	.257505 2525055551251	500100							

Variable length encoding

• Aim:

- For ARACHNOCENTRIC and other rare terms, we will use about 20 bits per gap (= posting).
- For THE and other very frequent terms, we will use only a few bits per gap (= posting).
- In order to implement this, we need to devise some form of variable length encoding.
- Variable length encoding uses few bits for small gaps and many bits for large gaps.

Variable byte (VB) code

- Used by many commercial/research systems
- Good low-tech blend of variable-length coding and sensitivity to alignment matches (bit-level codes, see later).
- Dedicate 1 bit (high bit) to be a continuation bit c.
- If the gap G fits within 7 bits, binary-encode it in the 7 available bits and set c = 1.
- Else: encode lower-order 7 bits and then use one or more additional bytes to encode the higher order bits using the same algorithm.
- At the end set the continuation bit of the last byte to 1
 (c = 1) and of the other bytes to 0 (c = 0).

VB code examples

docIDs	824	829	215406
gaps		5	214577
VB code	00000110 10111000	10000101	00001101 00001100 10110001

VB code encoding algorithm

VBEncodeNumber(n)		VBENCODE(<i>numbers</i>)		
$\sim \frac{1}{\sqrt{2}} \sqrt{\frac{1}{2}} \frac{1}{\sqrt{2}} \frac{1}{$		nn nu syssk ⇔rig i – ° '		oyxestream
n∈ numbers		2 while true		2) far each i
WREDIGODS Y WIDER C.	١.,,,	warmen Baak of BESSIN About on a	g. mag.	1.98 march and brings
COM (Lysestream, Syles)	ζ_{i}	iî 22 < 1028	4	by test ream $\leftarrow \mathbb{E}\mathbb{X}^{*}$
	5	then Break	5	ntiurn bytestream
	6	$n \leftarrow n$ div 128		
	7	bytes[Length(bytes)] += 128		
	8	return <i>bytes</i>		

VB code decoding algorithm



Gamma codes for gap encoding

- You can get even more compression with another type of variable length encoding: bitlevel code.
- Gamma code is the best known of these.
- First, we need unary code to be able to introduce gamma code.
- Unary code
 - Represent n as n 1s with a final 0.
 - Unary code for 3 is 1110

 - Unary code for 70 is:

Gamma code

- Represent a gap G as a pair of length and offset.
- Offset is the gap in binary, with the leading bit chopped off.
- For example $13 \rightarrow 1101 \rightarrow 101 = offset$
- Length is the length of offset.
- For 13 (offset 101), the length is 3.
- Encode length in unary code: 1110.
- Gamma code of 13 is the concatenation of length and offset: 1110101.

Gamma code examples

number	unary code	length	offset	$\gamma \operatorname{code}$
0	0			
1	10	0		0
2	110	10	0	10,0
3	1110	10	1	10,1
4	11110	110	00	110,00
9	1111111110	1110	001	1110,001
13		1110	101	1110,101
24		11110	1000	11110,1000
511		111111110	11111111	111111110,11111111
1025		11111111110	000000001	11111111110,000000001

Properties of gamma code

- Gamma code is prefix-free
- The length of offset is $\lfloor \log_2 G \rfloor$ bits.
- The length of length is $\lfloor \log_2 G \rfloor + 1$ bits,
- So the length of the entire code is $2 \times \lfloor \log_2 G \rfloor + 1$ bits.
- Υ codes are always of odd length.
- Gamma codes are within a factor of 2 of the optimal encoding length log₂ G.

Gamma codes: Alignment

- Machines have word boundaries 8, 16, 32 bits
- Compressing and manipulating at granularity of bits can be slow.
- Variable byte encoding is aligned and thus potentially more efficient.
- Regardless of efficiency, variable byte is conceptually simpler at little additional space cost.

Compression of Reuters

data structure	size in MB
dictionary, fixed-width	11.2
dictionary, term pointers into string	7.6
\sim , with blocking, k = 4	7.1
~, with blocking & front coding	5.9
collection (text, xml markup etc)	3600.0
collection (text)	960.0
T/D incidence matrix	40,000.0
postings, uncompressed (32-bit words)	400.0
postings, uncompressed (20 bits)	250.0
postings, variable byte encoded	116.0
postings, gamma encoded	101.0