

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Introduction to
Assignment Project Exam Help
Informa |
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
Lecture 5: Index on

Last lecture – in

Add WeChat edu_assist_pro

- Sort-based indexing
 - Naïve in-memory inversion
 - Blocked Sort-Based Indexing
 - Merge sort
- Single-Pass In
 - No global dictionary
 - Generate separate dictionary for each block
 - Don't sort postings
 - Accumulate postings in postings lists as they occur
- Distributed indexing using MapReduce
- Dynamic indexing: Multiple indices, logarithmic merge

Assignment Project Exam Help

Today

Add WeChat [edu_assist_pro](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Collection statistics in movie corpus (with RCV1)
 - How big will the dictionary and postings be?
- Dictionary compression
- Postings compression

Assignment Project Exam Help Why compression? (general)? Add WeChat edu_assist_pro

- Use less disk space
 - Saves a little money
- Keep more stuff in memory
 - Increases speed of data transfer to memory
- Increase speed of decompression
 - [read compressed data | decompress] faster than [read uncompressed data]
 - Premise: Decompression algorithms are fast
 - True of the decompression algorithms we use

Assignment Project Exam Help Why compression erated indexes? Add WeChat edu_assist_pro

- Dictionary
 - Make it small enough to keep in main memory
 - Make it so small that you can keep some postings lists in main memory
- Postings file(s)
 - Reduce disk space needed
 - Decrease time needed to read postings lists from disk
 - Large search engines keep a significant part of the postings in memory.
 - Compression lets you keep more in memory
- We will devise various IR-specific compression schemes

Assignment Project Exam Help Recall Reuters R

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

symbol	statistic	value
N	documents	800,000
L	avg	200
M	ter	400,000
	avg. # bytes per token (incl. spaces/punct.)	4.5
	avg. # bytes per term	7.5
	non-positional postings	100,000,000

Index

Assignment Project Exam Help
 (details IIR Table 5.1, p.

~~Add WeChat edu_assist_pro~~

size of	word types (terms)		non-positional postings		positional postings			
	dictionary		non-positional index		positional index			
	Size (K)	Δ			mul	Size (K)	Δ %	cumul %
Unfiltered	484			109,9		197,879		
No numbers	474	-2	-2	100,6		179,158	-9	-9
Case folding	392	-17	-19	96,969	-3	-12	179,158	0 -9
30 stopwords	391	-0	-19	83,390	-14	-24	121,858	-31 -38
150 stopwords	391	-0	-19	67,002	-30	-39	94,517	-47 -52
stemming	322	-17	-33	63,812	-4	-42	94,517	0 -52

Exercise: give intuitions for all the '0' entries. Why do some zero entries correspond to big deltas in other columns?

Assignment Project Exam Help

Lossless vs. lossy compression

Add WeChat edu_assist_pro

- Lossless compression: All information is preserved.
 - What we mostly do in IR.
- Lossy compression: Discard some information
- Several of the <https://eduassistpro.github.io/> viewed as lossy compression: case folding, words, stemming, number elimination
- Chap/Lecture 7: Prune postings entries that are unlikely to turn up in the top k list for any query.
 - Almost no loss quality for top k list.

Vocabulary vs. Collection size

Add WeChat edu_assist_pro

- How big is the term vocabulary?
 - That is, how many distinct words are there?
- Can we assume an upper bound?
 - Not really: A collection of length 20 can have words of length 20
- In practice, the vocabulary grows with the collection size
 - Especially with Unicode ☺

Vocabulary vs. size

Add WeChat edu_assist_pro

- Heaps' law: $M = kT^b$
- M is the size of the vocabulary, T is the number of tokens in the collection
- Typical values $k \approx 0.5$
- In a log-log plot of vocabulary size M vs. T , Heaps' law predicts a line with slope b :
 - It is the simplest possible relationship between the two in log-log space
 - An empirical finding ("empirical law")

Assignment Project Exam Help Heaps' Law 81 Add WeChat edu_assist_pro

For RCV1, the dashed line

$$\log_{10} M = 0.49 \log_{10} T + 1.64$$

is the best least squares fit.

Thus, $M = 10^{1.64} T^{0.49}$ s

$10^{1.64} \approx 44$ and $b = 0.4$

Add WeChat edu_assist_pro

Good empirical fit for

Reuters RCV1 !

For first 1,000,020 tokens,
law predicts 38,323 terms;
actually, 38,365 terms

Assignment Project Exam Help Exercises

[Add WeChat edu_assist_pro](#)

- What is the effect of including spelling errors, vs. automatically correcting spelling errors on Heaps' law? [Assignment Project Exam Help](#)
- Compute the [this scenario:](https://eduassistpro.github.io/)
 - Looking at a collection of we find that there are 3000 different terms in the first 100 tokens and 30,000 different terms in the first 1,000,000 tokens.
 - Assume a search engine indexes a total of $20,000,000,000$ (2×10^{10}) pages, containing 200 tokens on average
 - What is the size of the vocabulary of the indexed collection as predicted by Heaps' law?

Zipf's law

Add WeChat edu_assist_pro

- Heaps' law gives the vocabulary size in collections.
- We also study the relative frequencies of terms.
- In natural language, there are very frequent terms and very infrequent terms.
- Zipf's law: The i th most frequent term has frequency proportional to $1/i$.
- $cf_i \propto 1/i = K/i$ where K is a normalizing constant
- cf_i is collection frequency: the number of occurrences of the term t_i in the collection.

Zipf consequence

Add WeChat [edu_assist_pro](#)

- If the most frequent term (*the*) occurs cf_1 times
 - then the second most frequent term (*of*) occurs $cf_1/2$ times
 - the third most frequent term (*and*) occurs $cf_1/3$ times ...
- Equivalent: $cf_i \propto i^{-\alpha}$ (with $\alpha \approx 1$)
 - $\log cf_i = \log K - \log i$
 - Linear relationship between $\log cf_i$ and $\log i$
- Another power law relationship

Zipf's law for Re

Assignment Project Exam Help
v1

Add WeChat [edu_assist_pro](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)

Assignment Project Exam Help Compression

[Add WeChat edu_assist_pro](#)

- Now, we will consider compressing the space for the dictionary and postings
 - Basic Boo
 - No study
 - We will consider com

[Assignment Project Exam Help](#)

<https://eduassistpro.github.io/>
s, etc.

[Add WeChat edu_assist_pro](#)

schemes

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

DICTIONARY COMPRESSION

Assignment Project Exam Help Why compress the dictionary?

Add WeChat edu_assist_pro

- Search begins with the dictionary
- We want to keep it in memory
- Memory footprint other applications
- Embedded/mobile devices very little memory
- Even if the dictionary isn't in memory, we want it to be small for a fast search startup time
- So, compressing the dictionary is important

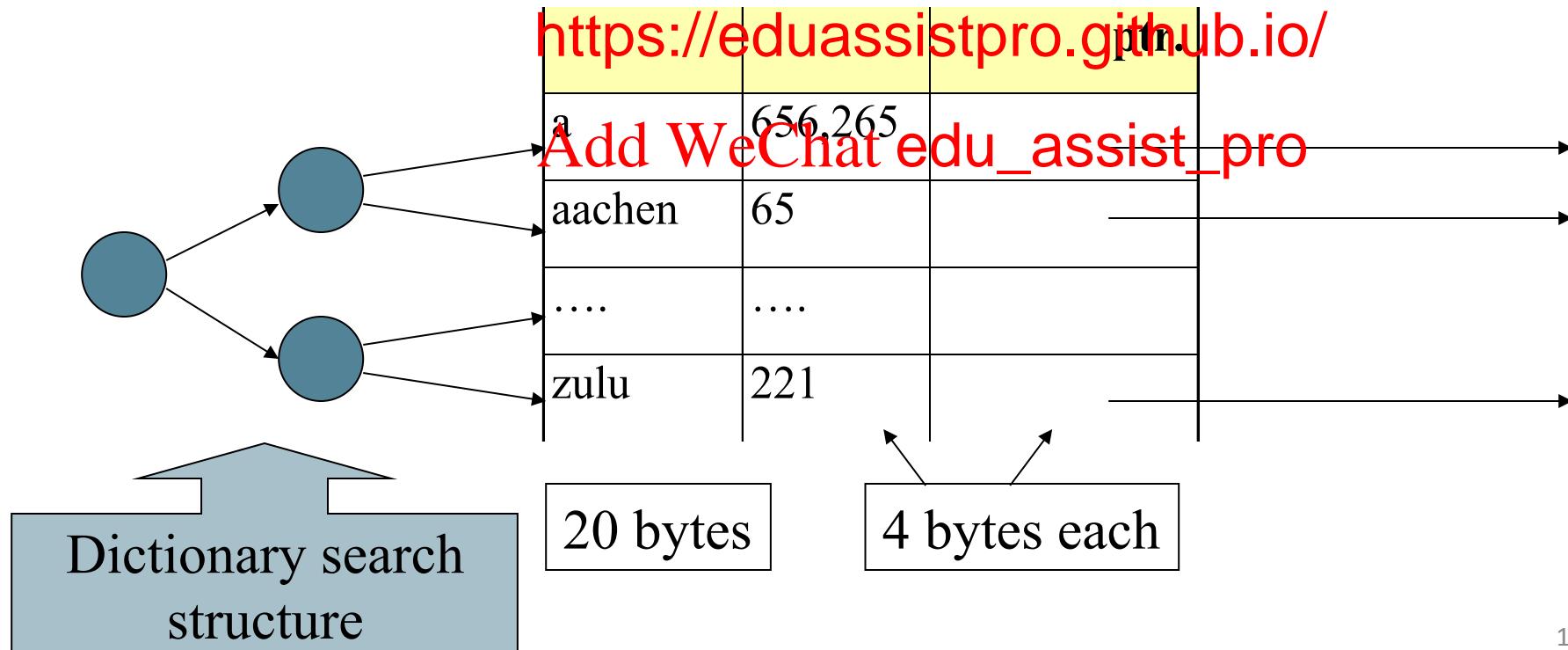
Assignment Project Exam Help

Dictionary storage

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- Array of fixed-width entries
 - ~400,000 terms; 28 bytes/term = 11.2 MB.

Assignment Project Exam Help



Assignment Project Exam Help

Fixed-width terms are wasteful

Add WeChat edu_assist_pro

- Most of the bytes in the **Term** column are wasted – we allot 20 bytes for 1 letter terms.
 - And we still can't handle *superfragilisticexpialidocious* or *hydrochlorofluo*
- Written English <https://eduassistpro.github.io/> characters/word.
 - Exercise: Why isn't this the best way for estimating the dictionary size?
- Ave. dictionary word in English: ~8 characters
 - How do we use ~8 characters per dictionary term?
- Short words dominate token counts but not type average.

Compre

t:

Assignment Project Exam Help
Dictionary-as-a-Str[Add WeChat edu_assist_pro](#)

- Store dictionary as ring of characters:
 - Pointer to next word shows end of current word
 - Hope to save up to 60% of dictionary space.

[Assignment Project Exam Help](#)

Freq.	Postings ptr.	Term ptr.
33		
29		
44		
126		

Total string length =
400K x 8B = 3.2MB

Pointers resolve 3.2M
positions: $\log_2 3.2M =$
22bits = 3bytes

Assignment Project Exam Help

Space for dictionary string

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- 4 bytes per term for Freq.
 - 4 bytes per term for pointer to Postings.
 - 3 bytes per te
 - Avg. 8 bytes p
 - 400K terms x 19 → 7.6 MB
- Now avg. 11 bytes/term, not 20.

Assignment Project Exam Help

Assignment Project Exam Help Blocking

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- Store pointers to every k th term string.

- Example below: $k=4$.

- Need to store term lengths (1 extra byte)

....7systyle9s <https://eduassistpro.github.io/> belyite8szczecin9szomo....

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

Freq.	Postings ptr.	Term ptr.
33		—
29		
44		
126		
7		—

Save 9 bytes
on 3
pointers.

Lose 4 bytes on
term lengths.

Net Assignment Project Exam Help

Add WeChat edu_assist_pro

- Example for block size $k = 4$
- Where we used 3 bytes/pointer without blocking
 - $3 \times 4 = 12$ bytes,

now we use 3 + <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro
Shaved another ~0.5MB. This the size of the
dictionary from 7.6 MB to 7.1 MB.

We can save more with larger k .

Why not go with larger k ?

Assignment Project Exam Help Exercise

[Add WeChat edu_assist_pro](#)

- Estimate the space usage (and savings compared to 7.6 MB) with blocking, for block sizes of $k = 4, 8$ and 16 .

[Assignment Project Exam Help](#)

<https://eduassistpro.github.io/>

[Add WeChat edu_assist_pro](#)

Assignment Project Exam Help
Dictionary search with king
~~Add WeChat edu_assist_pro~~

- Assuming each dictionary term equally likely in query (not really so in practice!), <https://eduassistpro.github.io/> number of comp = $(1+2+2+4+3+4)/8 \approx 2.6$

Exercise: what if the frequencies of query terms were non-uniform but known, how would you structure the dictionary search tree?

Assignment Project Exam Help

Dictionary search locking

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- Binary search down to 4-term block;
 - Then linear search through terms in block.
- Blocks of 4 (binary tree), avg. =
 $(1+2\cdot2+2\cdot3+2\cdot4+5)/8 = 3$ compares

Assignment Project Exam Help Exercise

[Add WeChat edu_assist_pro](#)

- Estimate the impact on search performance (and slowdown compared to $k=1$) with blocking, for block sizes of $k = 4, 8$ and 16 .

<https://eduassistpro.github.io/>

[Add WeChat edu_assist_pro](#)

Assignment Project Exam Help

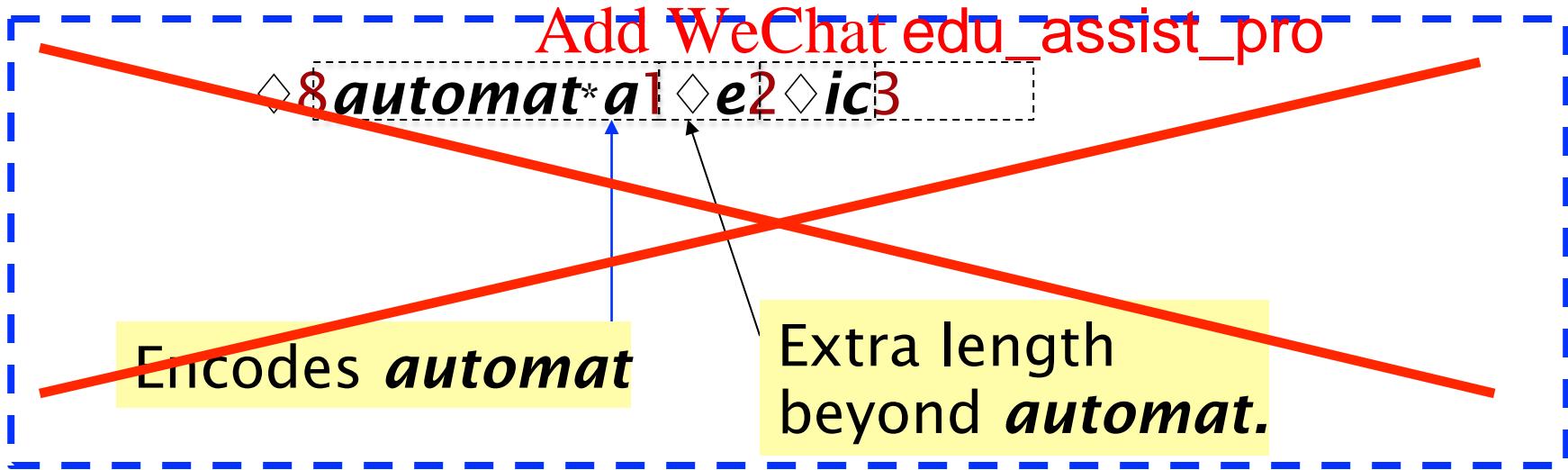
Front coding

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- Front-coding:

- Sorted words commonly have long common prefix – store differences only
- (for last $k-1$ i

<https://eduassistpro.github.io/>
8automata8au



Begins to resemble general string compression. 29

Assignment Project Exam Help Front Encoding [Witten, Moffat, Bell]

Add WeChat edu_assist_pro

- Complete front encoding
 - (prefix-len, suffix-len, suffix)
- Partial 3-in-4 front encoding
 - No encoding <https://eduassistpro.github.io/> string in a block
 - Enables binary search

Add WeChat edu_assist_pro

Assume previous string is “auto”



String	Complete Front Encoding	Partial 3-in-4 Front Encoding
8, automata	4, 4, mata	, 8, automata
8, automate	7, 1, e	7, 1, e
9, automatic	7, 2, ic	7, 2, ic
10, automation	8, 2, on	8, , on

RCV1 dictionary c

Assignment Project Exam Help
sion summary

Add WeChat edu_assist_pro

Technique	Size in MB
Fixed width	11.2
Dictionary-as-String	7.6
Also, blocking $k = 4$	7.1
Also, Blocking + front coding	5.9

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

POSTINGS COMPRESSION

Assignment Project Exam Help Postings compression

[Add WeChat edu_assist_pro](#)

- The postings file is much larger than the dictionary, factor of at least 10.
- Key desideratum: store each posting compactly.
- A posting for <https://eduassistpro.github.io/>
- For Reuters (800,000 documents) would use 32 bits per docID when using 4 bytes.
- Alternatively, we can use $\log_2 800,000 \approx 20$ bits per docID.
- Our goal: use a lot less than 20 bits per docID.

Assignment Project Exam Help Postings: two co forces

[Add WeChat edu_assist_pro](#)

- A term like ***arachnocentric*** occurs in maybe one doc out of a million – we would like to store this posting using $\log_2 1M = 20$ bits.
- A term like ***the*** occurs in every doc, so 20 bits/posting is
 - Prefer 0/1 bitmap vector in t

Assignment Project Exam Help

Postings file entr

Add WeChat edu_assist_pro

- We store the list of docs containing a term in increasing order of docID.
 - *computer*: 33, 47, 154, 159, 202 ...
- Consequence <https://eduassistpro.github.io/>
 - 33, 14, 107, 5, 43 ...
- Hope: most gaps can be encoded with far fewer than 20 bits.

Assignment Project Exam Help
Three postings e
Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Variable length encoding

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- Aim:
 - For *arachnocentric*, we will use ~20 bits/gap entry.
 - For *the*, we will use ~1 bit/gap entry.
- If the average <https://eduassistpro.github.io/> is $\log_2 G$ bits/gap
- Key challenge: encode every gap) with about as few bits as needed for that integer.
- This requires a *variable length encoding*
- Variable length codes achieve this by using short codes for small numbers

Assignment Project Exam Help Variable Byte (V)

Add WeChat edu_assist_pro

- For a gap value G , we want to use close to the fewest bytes needed to hold $\log_2 G$ bits
- Begin with one byte to store G and dedicate 1 bit in it to be a [con](https://eduassistpro.github.io/)
- If $G \leq 127$, binary-encode it
- Else encode G 's lower-order 7 bits and then use 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 for the other bytes $c = 0$.

Example

Add WeChat edu_assist_pro

docIDs	824	829	215406
gaps		5	214577
VB code	00001101 1	00001101 00001100 10110001	00001101 00001100 10110001

Postings stored as the byte cion
00000110101110001000010100 110010110001

Key property: VB-encoded postings are uniquely prefix-decodable.

For a small gap (5), VB uses a whole byte.

Assignment Project Exam Help Other variable u s

[Add WeChat edu_assist_pro](#)

- Instead of bytes, we can use a different “unit of alignment”: 32 bits (words), 16 bits, 4 bits (nibbles).
- Variable byte wastes space if you have many small gaps in such cases.
- Variable byte
 - Used by many commercial formats
 - Good low-tech blend of variable-length coding and sensitivity to computer memory alignment matches (vs. bit-level codes, which we look at next).
- There is also recent work on word-aligned codes that pack a variable number of gaps into one word (e.g., simple9)

Assignment Project Exam Help

Simple9

- Encodes as many gaps as possible in one DWORD
- 4 bit selector + 28 bit data bits
 - Encodes 9 possible ways to “use” the data bits

Selector	#	Wasted bits
https://eduassistpro.github.io/		
0000	28	1
0001	14	2
0010	9	3
0011	7	4
0100	5	5
0101	4	7
0110	3	9
0111	2	14
1000	1	28

Unary code Assignment Project Exam Help

Add WeChat `edu_assist_pro`

- Represent n as n 1s with a final 0.

- Unary code for 3 is 1110.

- ## ■ Unary code for

1111111111111 https://eduassistpro.github.io/11111111111110 .

- Unary code for ~~Add~~ is: WeChat edu_assist_pro

- This doesn't look promising, but....

Assignment Project Exam Help

Bit-Aligned Code

Add WeChat edu_assist_pro

- Breaks between encoders can occur after any bit position
- *Unary code*
 - Encode k by 0 at end making
 - 0 at end makes

Add WeChat edu_assist_pro

Number	Code
0	0
1	10
2	110
3	1110
4	11110
5	111110

Assignment Project Exam Help Unary and Binary

Add WeChat [edu_assist_pro](#)

- Unary is very efficient for small numbers such as 0 and 1, but quickly becomes very expensive
 - 1023 can be $2^{10} - 1$ bits, but requires 1024 bits in <https://eduassistpro.github.io/>
- Binary is more efficient for larger numbers, but it may be ambiguous

Assignment Project Exam Help

Add WeChat [edu_assist_pro](#)

Elias- γ Code

Assignment Project Exam Help

Add WeChat edu_assist_pro

- To encode a number

e

unary

Assignment Project Exam Help binary

- k_d is numb

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Elias- δ Code

Add WeChat [edu_assist_pro](#)

- Elias- γ code uses no more bits than unary, many fewer for $k > 2$
 - 1023 takes 19 bits instead of 1024 bits using unary
- In general, ta <https://eduassistpro.github.io/>
- To improve coding of large k_d , use Elias- δ code
 - Instead of encoding k_d in unary, encode $k_d + 1$ using Elias- γ
 - Takes approximately $2 \log_2 \log_2 k + \log_2 k$ bits

Elias- δ Code

- Split $(k_d + 1)$ into:

$$k_{dd} = \lfloor \log_2(k_d + 1) \rfloor$$

$$k_{dr} = (k_d + 1)^{\lfloor \log_2(k_d + 1) \rfloor}$$

- encode k_{dd} in binary

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Add WeChat **edu_assist_pro**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat **edu_assist_pro**

Assignment Project Exam Help Gamma code pr

Add WeChat edu_assist_pro

- G is encoded using $2 \lceil \log G \rceil + 1$ bits
 - Length of offset is $\lfloor \log G \rfloor$ bits
 - Length of length is $\lceil \log G \rceil + 1$ bits
- All gamma co <https://eduassistpro.github.io/> er of bits
- Almost within a factor of 2 $\frac{\text{number of bits}}{\log_2 G}$ Add WeChat edu_assist_pro

- Gamma code is uniquely prefix-decodable, like VB
- Gamma code can be used for any distribution
- Gamma code is parameter-free

Assignment Project Exam Help

Gamma seldom practice

Add WeChat [edu_assist_pro](#)

- Machines have word boundaries – 8, 16, 32, 64 bits
 - Operations that cross word boundaries are slower
- Compressing and manipulating at the granularity of bits can be slow
- Variable byte encoding is also thus potentially more efficient
- Regardless of efficiency, variable byte is conceptually simpler at little additional space cost

Assignment Project Exam Help Shannon Limit

Add WeChat edu_assist_pro

- Is it possible to derive codes that are optimal (under certain assumptions)?
- What is the optimal average code length for a code that encodes <https://eduassistpro.github.io/> independently?
Add WeChat edu_assist_pro
- Lower bounds on average code length: Shannon entropy
 - $H(X) = - \sum_{x=1}^n Pr[X=x] \log Pr[X=x]$
 - Asymptotically optimal codes (finite alphabets): arithmetic coding, Huffman codes

How to design an optimal code
for geometric distribution?

Assignment Project Exam Help Global Bernoulli

Add WeChat [edu_assist_pro](#)

- Assumption: term occurrence are Bernoulli events
- Notation:
 - n: # of documents, m: # of terms in vocabulary
 - N: total # of (<https://eduassistpro.github.io/>)
- Probability of a term t_i occurring in document d_i : $p = \frac{N}{nm}$
- Each term-document occurrence is an independent event
- Probability of a gap of length x is given by the geometric distribution $\Pr[X = x] = (1 - p)^{x-1} \cdot p$

Golomb Code

Add WeChat edu_assist_pro

It can also be deemed as a generalization of the unary code.

- Golomb Code (Golomb 1966): highly efficient way to design optimal Huffman-style code for geometric distribution
- Parameter b
- For given $x \geq 1$, find $q = \lfloor (x-1)/b \rfloor$ and remainder $r = (x-1) - q \cdot b$
- Assume $b = 2^k$
 - Encode q in unary, followed by r coded in binary
 - A bit complicated if $b \neq 2^k$. See wikipedia.
- First step: $(q+1)$ bits
- Second step: $\log(b)$ bits

Golomb Code & de

Assignment Project Exam Help
Add WeChat edu_assist_pro

- How to determine optimal b^* ?
- Select minimal b such that
 - Assignment Project Exam Help $b+1$
1
- Result due to <https://eduassistpro.github.io/> rms 1975:
generates an optimal prefix geometric distribution
- Small p approximation:
$$b^* \approx \ln 2 / p = 0.69 \cdot avg_val$$
- Rice code: only allow $b = 2^k$

Assignment Project Exam Help Local Bernoulli

[Add WeChat edu_assist_pro](#)

- If length of posting lists is known, then a Bernoulli model on each individual inverted list can be used
- Frequent words are coded with smaller b , infrequent words with l_a <https://eduassistpro.github.io/>
- Term frequency need to be μ use gamma-code)
- Local Bernoulli outperforms global Bernoulli model in practice (method of practice!)

Assignment Project Exam Help RCV1 compressi

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

Data structure	Size in MB
dictionary, fixed-width	11.2
dictionary, term postings into string	7.6
with blocking, k = 4	7.1
with blocking & front	5.9
collection (text, xml markup etc)	3,600.0
collection (text)	960.0
Term-doc 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, g-encoded	101.0

Assignment Project Exam Help

Google's Indexing

Add WeChat edu_assist_pro

- Index shards partition by doc, multiple replicates
- Disk-resident index
 - Use outer parts of the disk
 - Use different fields:
Rice_k (a specific encoding for positions) or gaps, and Gamma
- In-memory index
 - All positions; No docid
 - Keep track of document boundaries
 - Group-variant encoding
 - Fast to decode

Source: [Jeff Dean's WSDM 2009 Keynote](#)

Assignment Project Exam Help Other details

[Add WeChat edu_assist_pro](#)

- Gap = $\text{docid}_n - \text{docid}_{n-1} - 1$
- Freq = freq – 1
- Pos_Gap = pos

<https://eduassistpro.github.io/>

- C.f., Jiangong Zhang, Xiaohui Torsen, Suel: Performance of Compressed Inverted List Caching in Search Engines. WWW 2008.

Assignment Project Exam Help

Index compression

Add WeChat edu_assist_pro

- We can now create an index for highly efficient Boolean retrieval that is very space efficient
- Only 4% of the total size of the collection
- Only 10-15% of the total size of the collection
- However, we've ignored position information
- Hence, space savings are less for indexes used in practice
 - But techniques substantially the same.

Assignment Project Exam Help

Resources for to cture

Add WeChat edu_assist_pro

- *IIR* 5
- *MG* 3.3, 3.4.
- F. Scholer, H. Compression [FastIO](https://eduassistpro.github.io/FastIO/) I. 2002.
Evaluation. *Proc. ACM-SIGI*
- Variable byte codes
- V. N. Anh and A. Moffat. 2005. Inverted Index Compression Using Word-Aligned Binary Codes. *Information Retrieval* 8: 151–166.
■ Word aligned codes