

CS4611

Tuesday 17th September 2013

Information Retrieval,

Prof. M.P. Schellekens (TA: Ang Gao)
Slides: P. Nayak and P. Raghavan.

lecture 1: Boolean Retrieval

- ◇ Information Retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers)
- ◇ Corpora: large collection of documents.

Ranking: Two ways = System \times Ask the user if it was helpful.

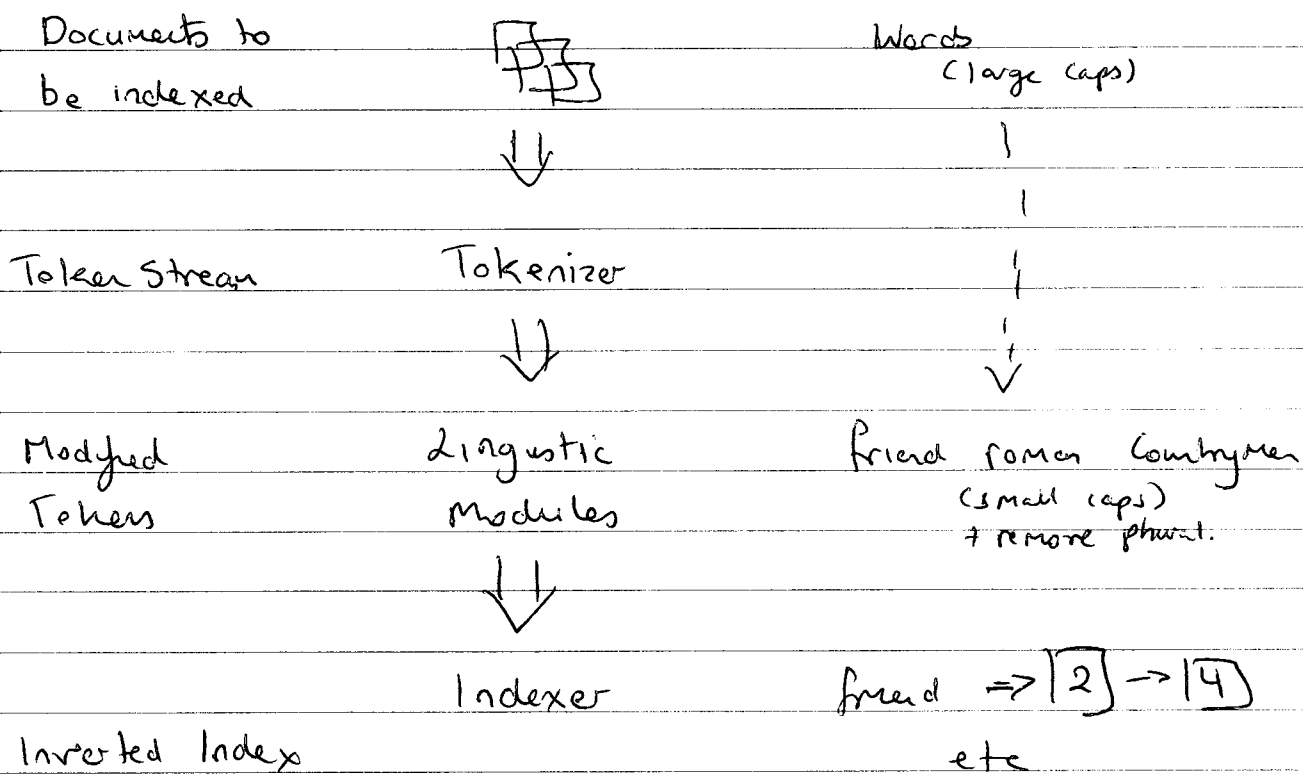
- ◇ Precision: fraction or proportion of documents that come back are useful to me.
- ◇ Recall: fraction of relevant docs in collection that are retrieved

Terms for now we can refer to them as words in a document. This is expanded as the words are processed and we will end up with terms.

Variable-Size Posting Lists

Dictionary \Rightarrow collection of all terms (processed words)

Sort by docID!



CS4611

Thursday 19th September 2013

Term frequency is handy for ranking and grouping results.

- o Book: Introduction to Information Retrieval (<http://informationretrieval.org/>)

lecture 2 : The term vocabulary and postings lists.

Postings list : Document ID's.

- ◇ A term is a (normalized) word type, which is an entry in our IR system dictionary.

Normalization

C.A.T. → CAT → cat

So when you type C.A.T. into ~~an~~ IR search engine 'cat' will be at the top instead of Caterpillar. We need to ensure that C.A.T. shows up at the top.

Synonyms & Homonyms.

- left (part of leave) } one word / two meanings
- left (opposite of right)

△ Porter Algorithm

- Used for stemming English

eg compress, compression \rightarrow compressed
 \rightarrow compress, compress \neq compress

sses \rightarrow ss

ies \rightarrow i

ational \rightarrow ate

tional \rightarrow tion

~~Rules~~ \Rightarrow

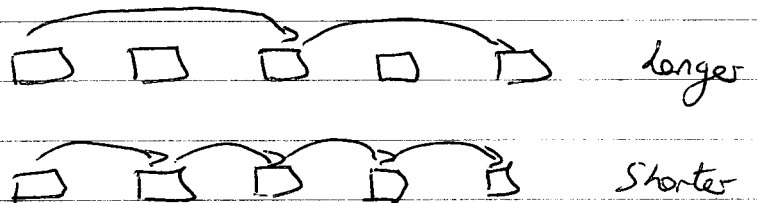
Other \Rightarrow Lovins Stemmer (approx 250 rules)

Rules are application & language specific

MIT.english

mit.german

Placing Skips = α use \sqrt{x} to ensure evenly spaced skip pointers



Phrase Queries & Positional Indexes

Search → Stanford University

Query → "I went to college at Stanford"

These are not a match.

Biwords

'friends romans' and 'romans countrymen'

each of the above biwords are not dictionary terms.

CS4611

Tuesday 24th September 2013

Dictionary (Dictionaries and tolerant Retrieval)

- o Hashtables
- o Trees

Hashtables - no prefix search \rightarrow tolerant retrieval is a disadvantage

Tree \rightarrow Binary Tree
 \rightarrow B-Tree

man* starts with man

*man ends with man

to find if man is in the middle then do an intersection of the two.

◇ Permuterm Index : permutation of a term

eg hello

hello\$, ello\$h, llo\$he, lobhel
o\$hell (where \$ is a symbol)

Bigram - little bit similar to bigrams

store two consecutive letters

example

April

\$a, ap, ri, il, l\$

Spelling Correction

- Correcting the doc(s) being indexed
- Correcting user queries to retrieve "right" answers

Will not catch for example from / form
but if it is part of a phrase "I came form school,"
clearly it should be "I came from school"
Google → Did you mean ...

o Lexicon

Alternatives

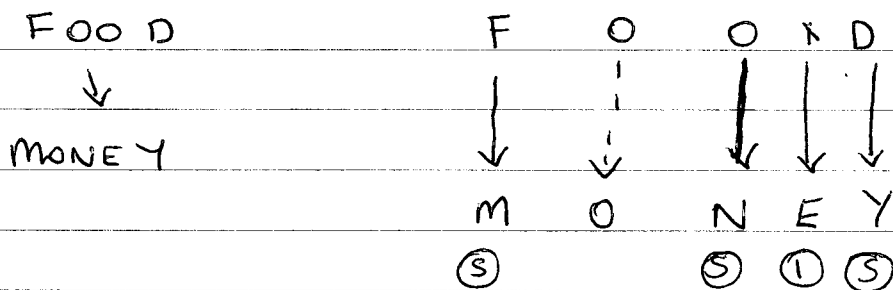
- ★ o Levenshtein Distance (Edit distance)
- o Weighted Edit Distance
- o n-gram overlap

□ Note : Must learn Levenshtein Distance algorithm

Edit Distance (Levenshtein)

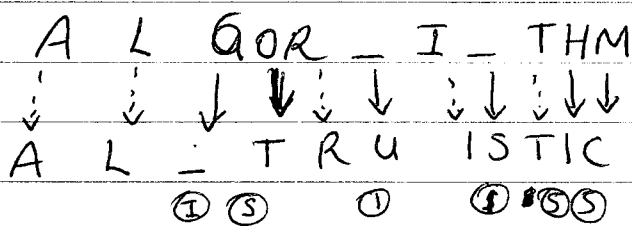
• Is the minimum number of letter insertions, letter deletions and letter substitutions to transform one word into another

Example 8- The distance between FOOD and MONEY



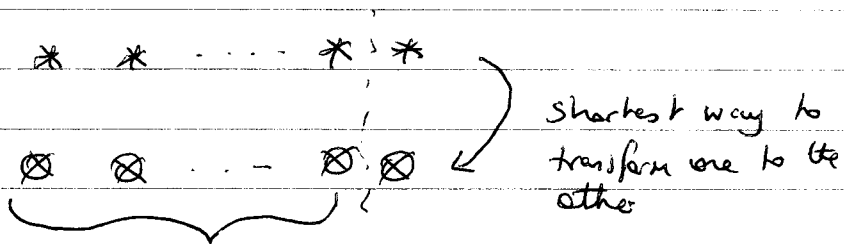
∴ The edit distance between MONEY and FOOD is 4
 $EditDistance(MONEY, FOOD) \leq 4$

Example 9-



6 operators - EditDistance = 6

Observation 3-



↳ Still shortest transformation
 is a subsequence. If you find the shortest way then it cannot

Recursive Definition of the Edit Distance between Two strings ϕ - $A [1 \dots m]$ and $B [1 \dots n]$

Denoted by $\text{Edit} (A [1 \dots m], B [1 \dots n])$

If neither string is empty (ϕ)

then there are 3 possibilities for the last column, in the shortest edit sequence we are adding up.

3 Possibilities

- o Insertion
- o Substitution
- o Deletion

Δ Insertion

The last entry in the top row is empty

$A [1] \dots A [m-1] \phi$
 $B [1] \dots B [n-1] B [n]$ \leftarrow Insertion

The edit distance = $\text{Edit} (A [1 \dots m], B [1 \dots n-1]) + 1$

Δ Deletion

The last entry in bottom row is empty.

$A [1] \dots A [m-1] A [m]$
 $B [1] \dots B [n-1] \phi$ \leftarrow Deletion.

The edit distance = $\text{Edit} ((A [1] \dots A [m-1]) (B [1] \dots B [n-1])) + 1$

Δ Substitution

Both rows have characters in the last column

If the characters are the same substitution is free (0 cost)
So the edit distance = edit \rightarrow

$$\rightarrow \text{edit}(A[1 \dots m-1], B[1 \dots n-1])$$

If characters are different

$$\text{edit}(A[1 \dots m-1], B[1 \dots n-1]) + 1$$

BASE-CASES

$$\text{EDIT}(A[1 \dots m], \emptyset) = m \text{ (deletions)}$$

$$\text{EDIT}(\emptyset, B[1 \dots n]) = n \text{ (insertions)}$$

Simplify Notation

Recursive subproblems are always prefixes of the strings A and B, so we can simply use the length of the substrings.

$$\text{EDIT}(i, j) \text{ (shorthand for edit}(A[1 \dots i], B[1 \dots j]))$$

$$\text{EDIT}(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \begin{cases} \text{EDIT}(i-1, j) + 1 \\ \text{EDIT}(i, j-1) + 1 \\ \text{EDIT}(i-1, j-1) + [A[i] \neq B[j]] \end{cases} \end{cases}$$

↳ is 1 if characters are different 0

SET UP A TABLE (Dynamic Programming) to compute the edit distance

	ϕ	f
ϕ	$+0$	$+$
c	$+$	

← Initialisation Row

→ cost to transform ϕ to ϕ is zero.

↑ Initialisation column

	ϕ	f
ϕ	$+0$	$+1$
c	$+1$	

→ End result.

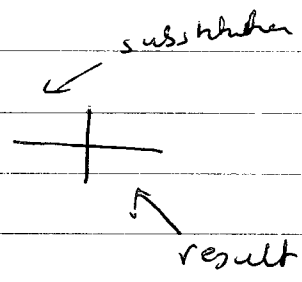
ϕ to F is insertion so cost is 1

c to ϕ is deletion so cost is 1

	ϕ	f		
ϕ	$+0$	$+1$	$+2$	$+3$...
c	$+1$			
e	$+2$			
\vdots				
$+3$				
\vdots				

	ϕ	P
ϕ	t_0	t_1
C	t_1	t_2

minimum
of three others



CS4611

Tuesday 8th October 2013

	A	L	T	R	O	I	S	T	I	C	
A	0/0	1/1	2/2	3/3	4/4	5/5	6/6	7/7	8/8	9/9	10/10
L	1/1	0/2	2/3	3/4	4/5	5/6	6/7	7/8	8/9	9/10	10/11
T	2/2	1/1	0/2	2/3	3/4	4/5	5/6	6/7	7/8	8/9	9/10
R	3/3	2/2	1/1	0/2	2/3	3/4	4/5	5/6	6/7	7/8	8/9
O	4/4	3/3	2/2	1/1	0/2	2/3	3/4	4/5	5/6	6/7	7/8
I	5/5	4/4	3/3	2/2	1/1	0/2	2/3	3/4	4/5	5/6	6/7
S	6/6	5/5	4/4	3/3	2/2	1/1	0/2	2/3	3/4	4/5	5/6
T	7/7	6/6	5/5	4/4	3/3	2/2	1/1	0/2	2/3	3/4	4/5
I	8/8	7/7	6/6	5/5	4/4	3/3	2/2	1/1	0/2	2/3	3/4
C	9/9	8/8	7/7	6/6	5/5	4/4	3/3	2/2	1/1	0/2	2/3

6 operations

d. Distance

0 1 2 3 4 5 6

1.000 - There are multiple paths that could have been taken.

□ = where both letters are the same, this is where we have changes.

Assignment 1 (DANGER → BALTIMORE)

- Compute the Edit Distance
- Find all shortest Edit Sequences that transform "DANGER" to "BALTIMORE"

Assignment Due @ 15th October 2013 10am

One Ophor - Jaccard Coefficient

$$X \cap Y / X \cup Y$$

Count amount from intersection and divide by count of union.

Soundex

Class of heuristics to expand a query into phonetic equivalents

Convert similar sounding letters to numbers.

Compression

- ◇ Heap's Law is linear in log-log space
It is the simplest possible relationship between collection size and vocab size in log-log space
Empirical law

CS4611

Thursday 10th October 2013

◇ Heaps Law: $M = kT^b$

$M =$ No. of terms

$T =$ No. of tokens

$$\log_{10} M = \log_{10}(kT^b)$$

$$\log_{10}(k) + \log_{10}(T^b)$$

* $\frac{b}{0.5} \log_{10}(T) + \frac{\log_{10} k}{c}$

↳ $\frac{b}{0.5} \log_{10}(T) + \frac{\log_{10} k}{c}$

$b + k$ are based on experiments

▲ Zipf's Law

Frequent Terms - short code

Infrequent Terms - long code

Few frequent terms and many rare terms.

HEAPS LAW & ZIPF'S LAW

Dictionary Compression

- Average length of a term is 8 characters (english)

o Entropy - enables one to compute the compressibility of data (kind of pre-scanning to see if it can be compressed)

▲ Huffman Code: ^{file} text replace each char. with binary number, end result is one list of binary numbers

Huffman's Code

Character A, ..., F

Num. of Char 100,000

Freq.	A	B	C	D	E	F
	45	13	12	16	9	5

Fixed 000 001 010 011 100 101

Variable 0 101 100 111 1101 1100

CASE A → 300,000 (100,000 x 3 bits)

CASE B → 45,000 (45,000 x 1 bit)

(ASIS
(Encoding size)

$$224,000 \text{ bits } (45^1 + 13^3 + 12^3 + 16^3 + 9^4 + 5^4 = 224)$$

Huffman's Code - gives the shortest code to the character that appears most frequently.

Two Codes

01011 010
└──────────┘
Prefix

one is a prefix of the other.

If you concatenate which do you decode? 01011010

Non-prefix code

A	B	C
1	00	001

⇓

1, 00, 00, 1

We want the Huffman Code to be leading, compare coding to entropy.

CS4611

Tuesday 15th October 2013

Entropy : A basic introduction

1) (See slides)

Entropy enables me to compute ... (see slides)

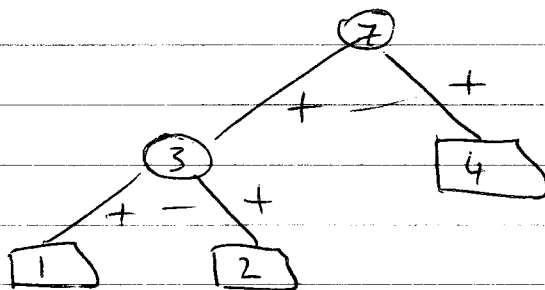
Huffman Encoding Example
• Fair Dice

$$p = 1/6$$

$$1/p = 1 / (1/6) = 6$$

$$\log(1/p) = \log(6) = 2.59$$

◆ HUFFMANN CODES

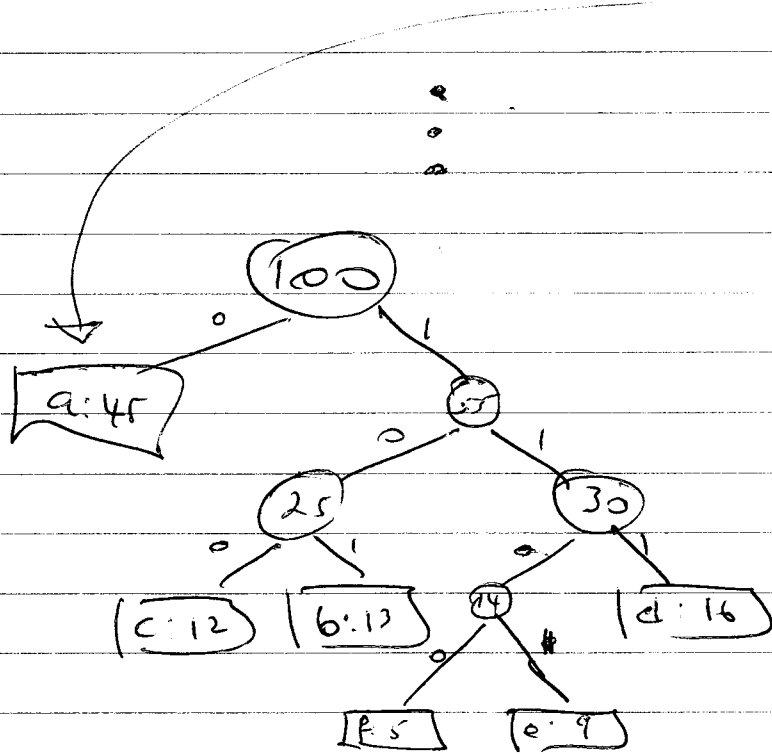
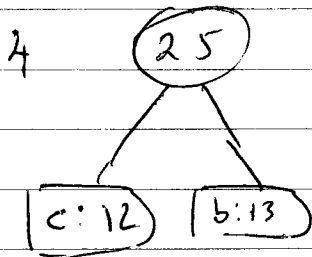
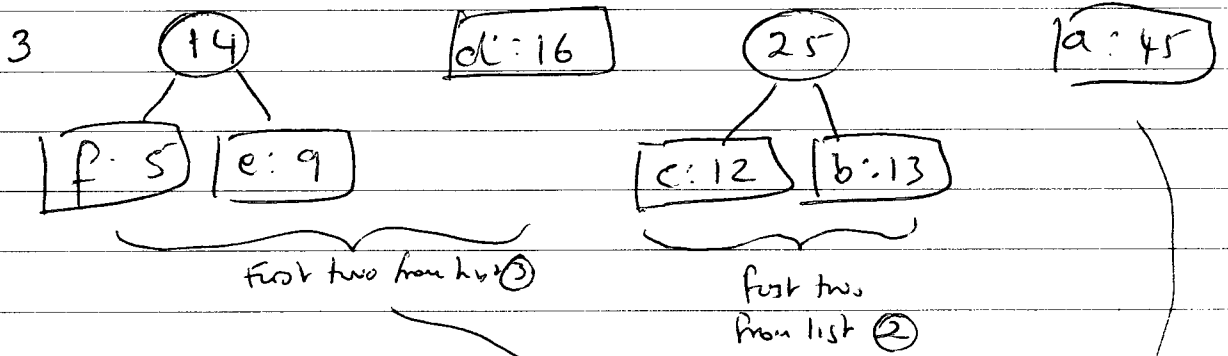
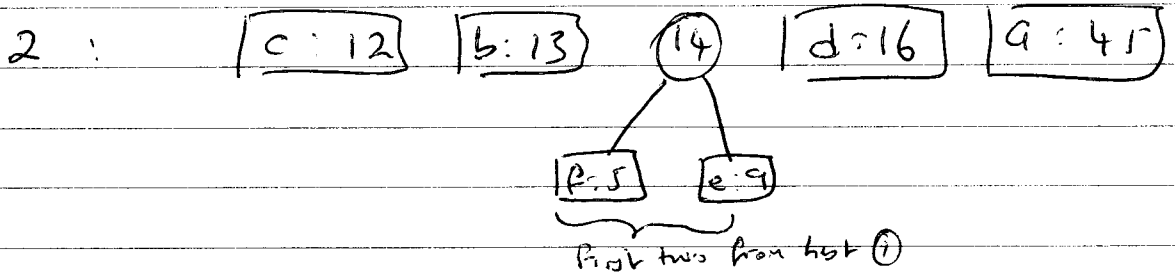


Once at the top of the tree, this = the final cost of the encoding.

Sort the characters by increasing freq.

f: 5	e: 9	c: 12	b: 13	d: 16	a: 45
------	------	-------	-------	-------	-------

In order
 12 → 13 → 14 → 16 → 45



How would you get num. of bits (AVG BIT LENGTH)

- 1 → Compute total cost 224,000 bits (100 + 55 + 25 + 30 + 10) / by number of character occurrences (100,000) (100) = 2.24

$$\Rightarrow \frac{224,000}{100,000} = \boxed{2.24 \text{ bits}}$$

OR (longer way)

$$2 \rightarrow \left(\frac{45}{100} \times 1 \right) + \left(\frac{12}{100} \times 3 \right) + \left(\frac{13}{100} \times 3 \right) + \left(\frac{5}{100} \times 4 \right) + \left(\frac{9}{100} \times 4 \right) + \left(\frac{10}{100} \times 3 \right) = 224 / 100 = 2.24$$

Compute the Entropy of the original file?

a	b	c	d	e	f
45	13	12	10	9	5

• Work in \log_2 (log base 2)

$$\rightarrow \frac{45}{100} \quad \frac{13}{100} \quad \frac{12}{100} \quad \frac{10}{100} \quad \frac{9}{100} \quad \frac{5}{100}$$

$$H(p_1 \dots p_n) = p_1 \log_2(1/p_1) + p_2 \log_2(1/p_2) + \dots + p_n \log_2(1/p_n)$$

The above is the Entropy Expression

In our case $n = 6$ as there are only 6 characters.

$$\log_2(k) = \ln(k) / \ln(2)$$

(where " \ln " is "log in base e")

This is if you do not have \log_2 on your calculator.

2.22

- Compute the Huffman encoding
- Compute the cost of the encoding

a) (prefix) codes for char

a: 0, b: 101, c: 100, d: 111, e: 1101, f: 1100

$$b) 45 \times 1 + 13 \times 3 + 12 \times 3 + 16 \times 3 + 9 \times 4 + 5 \times 4 = 224,000 \text{ bits}$$

$$c) 224,000 / 100,000 = 2.24 \text{ average encoding length}$$

e) Entropy

$$H\left(\frac{45}{100}, \frac{13}{100}, \frac{12}{100}, \frac{16}{100}, \frac{9}{100}, \frac{5}{100}\right) = \cdot$$

$$\begin{aligned} & 45/100 \log(45/100) - 13/100 \log(13/100) \\ & - 12/100 \log(12/100) - 16/100 \log(16/100) \\ & - 9/100 \log(9/100) - 5/100 \log(5/100) \end{aligned}$$

$$= 2.23 \quad \text{Answer.}$$

f) Conclusion

Excellent prediction of avg binary encoding length (some minor round off).

Other way

Instead of bytes we can use a different "unit of alignment": 32 bit (words), 16 bits.

△ Gamma Codes - More compression - bit level compression.

* We need unary code to be able to use gamma codes
Unary code for 3 ... (see slides)

- First encode it in binary
 - Then chop off the leading bit
 $13 \rightarrow 1101 \rightarrow 101 = \text{offset}$
 - Length is length of offset
 - $13 \Rightarrow 3$
 - Encode length in unary code: 1110
 - Gamma = 1110101
 - Concatenation of length and offset.
- $1110101 = \underbrace{111}_{\text{length}} \underbrace{0}_{\text{cont.}} \underbrace{101}_{\text{offset}}$

★ Exercise

Q:aps

1

2

3

4

9

13

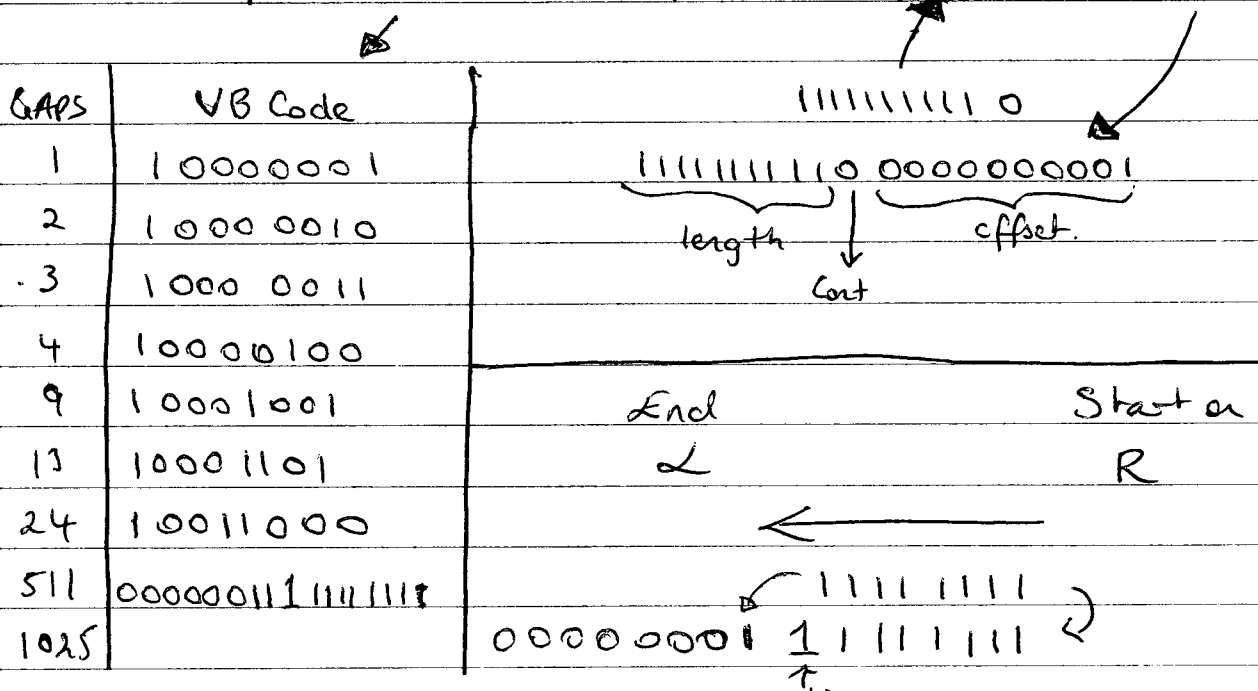
24

511

1 2 3 5

GAPS	BINARY		
1	0001	1	10000001
2	0010	2	10000010
3	0011	3	10000011
4	0100	4	10000100
9	1001	9	10001001
13	1110	13	10001101
24	1111 1	24	10011000
511		511	000000111111111
1025			

GAPS	BINARY	offset	length	⊗ Code
1	1	∅	0	0
2	10	∅	10	100
3	11	1	10	101
4	100	00	110	11000
9	1001	001	1110	1110001
13	1101	101	1110	1110101
24	11000	1000	11110	111101000
511	111111111	11111111	11111111	1111111101111111
1025	1000000001	0000000001		



$$130 = (2 \text{ to the power of } 7) \text{ plus } 2$$

Length of Gamma Code

- length of offset = $\lceil \log_2 a \rceil$ bits
- length of length = $\lceil \log_2 a \rceil + 1$ bits

RANKING

Scoring

Rank the documents by attaching a score to them.
Assign a score to each query-document (eg 0,1)
The more frequently the term occurs the higher the score.

- Jaccard Coefficient

$$\text{Jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$(A \neq \emptyset) (B \neq \emptyset)$$

Query

"ides of March"

Doc

"Caesar died in March"

$$\text{Jaccard} = 1/6$$

March = ①

7 words but March repeated so ⑥

How to rank search results.

- Term Frequency
- TF-idf Ranking
- Vector Space Model

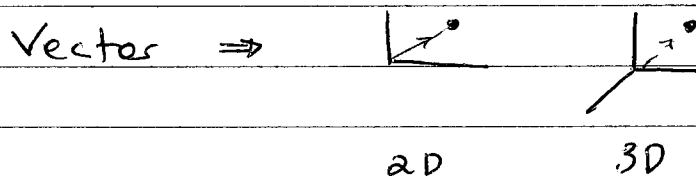
Boolean retrieval - All or nothing, a massive amount of results or no results.

- Good for applications
- Low Collector of documents

Query - relevance \rightarrow attach a number (eg 0,1)
score \uparrow

Take 1 - Jaccard Coefficient, can help us compute a score.

We don't measure how frequent a term occurs in a document using the Jaccard Coefficient - only that a term occurs in the document.



$$\Delta W_{e,d} = \begin{cases} 1 & \text{if } t_{e,d} > 0 \\ 0 & \text{otherwise} \end{cases} + \log_{10} t_{e,d}$$

$t_{e,d} \rightarrow W_{e,d} : 0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 1.3, 10 \rightarrow 2, 1000 \rightarrow 4$ etc.

Exercise

⊛ Compute the Jaccard matching score and the tf matching score for the following query-doc. pairs
(1. Jaccard coefficient, 2. term frequency rate)

1 [information on cars] D1: "all you've ever wanted to know about cars"

2 ["] D2: "information on trucks, information on planes, information on trains"

3 ~~[red cars and red trucks]~~ D3: "cars stop red cars more often"

$$\Delta \text{tf-matching-score}(q,d) = \sum t_{e,q} n_d (1 + \log t_{e,d})$$

Ans $J_1 = \frac{1}{11} = \frac{1}{11}$

$J_2 = \frac{2}{6} = \frac{1}{3}$

$J_3 = \frac{2}{8} = \frac{1}{4}$

TF - Matching - Score (Q, D)

TF - Matching - score (Q1, D1) = 1
TF - " - " (Q2, D2) = 2.954
" - " - " (Q3, D3) = 3

tf-idf

tf = term frequency

We want high weights for rare terms (Arachnoid^{example} tric)
For terms like (Good, Increase, line) we want positive weights

We are going to use document frequency for this.

idf-weight definition = $idf_t = \log_{10} \frac{N}{df_t}$

t = term

We measure the informativeness of a term to the entire collection.

example : $idf_t = \log_{10} \frac{1000,000}{df_t}$

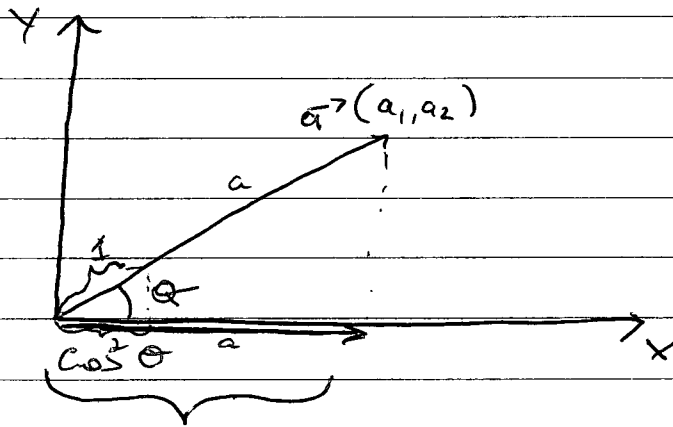
term	df _t	idf _t
animal	100	4

Word	collector freq.	document freq
insurance	10440	3997
try	10422	8760

We don't measure collector freq as the count could be the same but the document freq could be extremely different.

$$w_{t,d} = (1 + \log f_{t,d}) \cdot \log \frac{N}{df_t}$$

Product of its weight and its idf weight

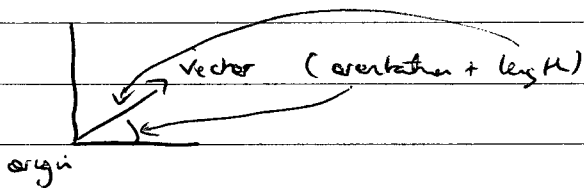


Vector Spaces

$\theta = \text{theta}$

1 x a
2 x a

$\|\vec{a}\| \cos \theta$
 > length of projection of \vec{a} on x-axis

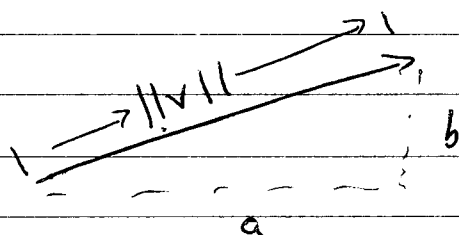


$$u \cdot v = u_1 v_1 + u_2 v_2 + u_3 v_3$$

$$u \cdot v = 1 \cdot 5 + 2 \cdot 6 + 3 \cdot 7 = 38$$

Length of a vector \Rightarrow

$$\|v\| = \sqrt{a^2 + b^2} = \sqrt{v \cdot v}$$



Distance between two vectors

$$\text{Dist}(P_1, P_2) = \sqrt{x \cdot x} = x$$

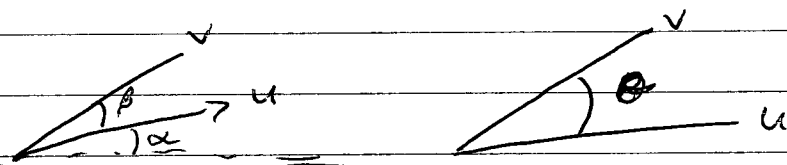
$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

↑ ↑
α α

Angle between two vectors

$$u = (\|u\| \cos(\alpha), \|u\| \sin(\alpha)) \text{ and}$$
$$v = (\|v\| \cos(\beta), \|v\| \sin(\beta))$$

$$u \cdot v = \|u\| \|v\| \cos(\beta - \alpha)$$



★

$$\cos(\theta) = \cos^{-1} \left(\frac{u \cdot v}{\|u\| \cdot \|v\|} \right)$$

↑ ↑
normal u + v

Required

①

Example:

Find the angle in degrees between the vectors
 $u = (-2, 2, 1)$ and $v = (2, 3, 6)$

Solution

$$\|u\| = \sqrt{4+4+1} = 3$$

$$\|v\| = \sqrt{4+9+36} = 7$$

$$3 \cdot 7 = 21$$

$$\begin{array}{ccc} (-2 \cdot -2) & (2 \cdot 2) & (1 \cdot 1) \\ (2 \cdot 2) & 3 \cdot 3 & (6 \cdot 6) \end{array}$$

$$\begin{array}{ccc} -2 & 2 & 1 \\ 2 & 3 & 6 \end{array}$$

$$\cos(\theta) = \frac{u \cdot v}{\|u\| \cdot \|v\|}$$

$$u \cdot v = -2 \cdot 2 + 2 \cdot 3 + 1 \cdot 6 = 8$$

$$\cos(\theta) = \frac{8}{21} = 67.6073^\circ$$

proximity = similarity

Distance is a bad idea as the distance could be very far apart, alternatively use the angle to measure the distance between two vectors.

Use angle instead of distance

example: $d + d'$ are quite similar but using the Euclidean distance, the distance is quite far away. The angle is close. in fact angle = 0.

△ The closer things are together the bigger the cosine is going to be.

Length Normalisation

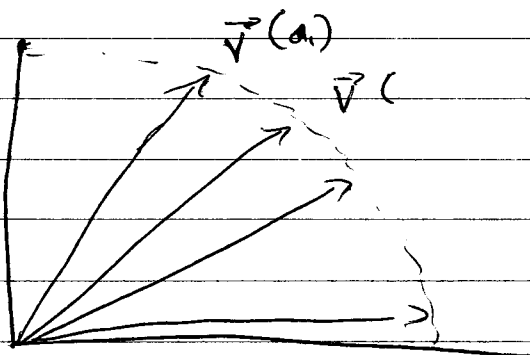
$$\|x\|_2 = \sqrt{\sum_i x_i^2}$$

This maps vectors onto the unit sphere

Query Vector Definition

$$\cos(\vec{q}, \vec{d}) = \text{sim}(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{\|\vec{q}\| \|\vec{d}\|} =$$

$$\frac{\sum_{i=1}^{|V|} q_i d_i}{(\text{see slides})}$$



TERMS	Q	D ₁	D ₂	D ₃	idf	Q idf	D ₁ idf	D ₂ idf	D ₃ idf
A	0	1	1	1	0	0	0	0	0
Arrived	0	0	1	1	0.176	0	0	0.176	0.176
Damaged	0	1	0	0	0.477	0	0.477	0	0
Delivery	0	0	1	0	0.477	0	0	0.477	0
Fire	0	1	0	0	0.477	0	0.477	0	0
Gold	1	1	0	1	0.176	0.176	0.176	0	0.176
in	0	1	1	1	0	0	0	0	0
of	0	1	1	1	0	0	0	0	0
Silver	1	0	2	0	0.477	0.477	0	0.954	0
Shipment	0	1	0	1	0.176	0	0.176	0	0.176
Truck	1	0	1	1	0.176	0.176	0	0.176	0.176

$$\text{Sim}(Q, D_1) = \frac{Q \cdot \vec{D}_1}{\|Q\| \cdot \|D_1\|}$$

$$\text{Sim}(Q, D_2) = \dots$$

$$\text{Sim}(Q, D_3) = \dots$$

Page Ranking

- D₁ - Shipment of Gold Damaged in a fire
- D₂ - Delivery of Silver arrived in a silver truck
- D₃ - Shipment of Gold arrived in a truck.

Q Gold Silver Truck.

→ VECTOR (query)

$$Q_{\#} = (0, 0, 0, 0, 0, 0.176, 0, 0, 0.477, 0, 0.1761)$$

Get this from table - idf column top → bottom.

Compute Similarity Between :

Q and D_1
 Q and D_2
 Q and D_3

} Largest Similarity best document
Match with Q .

$$\text{Sim}(Q, D_i) = \frac{\sum_{j=1}^3 w_{Q_j} w_{i,j}}{\sqrt{\sum_{j=1}^3 w_{Q_j}^2} \sqrt{\sum_{j=1}^3 w_{i,j}^2}}$$

$i = 1, 2, 3$

Vector (query) $\rightarrow \|Q\| = \sqrt{\quad}$

o Compute

(A) $\rightarrow \|Q\| = \sqrt{0^2 + 0^2 + 0^2 + 0^2 + 0^2 + (0.1761)^2 + 0^2 + 0^2 + (0.47)^2 + 0^2 + (0.1761)^2} =$

$\rightarrow \| \vec{D}_1 \|$
 $\| \vec{D}_2 \|$
 $\| \vec{D}_3 \|$

} Same as $\|Q\|$ - take column and \times^2
 then $\sqrt{\quad}$ the result of additions
 to get final result.

$$\begin{aligned} \rightarrow \|Q\| &= \sqrt{0.03101121 + 0.22762441 + 0.03101121} \\ &= \sqrt{0.2896} \\ &= \boxed{0.5382} \end{aligned}$$

(calculator $\Rightarrow \sqrt{\quad}$)

$$\rightarrow \| \vec{D}_1 \| = \sqrt{0.5173} = \boxed{0.7192}$$

$$\rightarrow \| \vec{D}_2 \| = \sqrt{1.2001} = \boxed{1.0955}$$

$$\rightarrow \| \vec{D}_3 \| = \sqrt{0.1240} = \boxed{0.3522}$$

(B) $\vec{Q} \cdot \vec{D}_1 =$ hold only one left, where ever you have a zero it disappears
 so $(0.1761) \cdot (0.1761) = \boxed{0.0310}$

$\vec{Q} \cdot \vec{D}_2 = (0.4771) \cdot (0.9542) + (0.1761) \cdot (0.1761) =$
 ~~$\vec{Q} \cdot \vec{D}_3 =$~~
 $(0.4771 \times 0.9542) + (0.1761)^2 = \boxed{0.4862}$

$\vec{Q} \cdot \vec{D}_3 = (0.1761 \times 0.1761) + (0.1761 \times 0.1761) = \boxed{0.062}$

(C) $\sin(\vec{Q}, \vec{D}_1) = \frac{\vec{Q} \cdot \vec{D}_1}{\|\vec{Q}\| \cdot \|\vec{D}_1\|} = \frac{0.0310}{0.5382 \times 0.7192} = \boxed{0.0801}$

$\left(\frac{B \text{ from above}}{A \cdot A} \right) \text{ How?}$
 $Q \cdot D_1 \text{ on previous page}$

$\sin(\vec{Q}, \vec{D}_2) = \frac{\vec{Q} \cdot \vec{D}_2}{\|\vec{Q}\| \cdot \|\vec{D}_2\|} = \frac{0.4862}{0.5382 \times 0.7192} = \boxed{0.824}$

$\sin(\vec{Q}, \vec{D}_3) \rightarrow \boxed{0.3271}$

From © which one will be presented?

\vec{D}_1 0.0801
 \vec{D}_2 0.8246
 \vec{D}_3 0.3271

} \vec{D}_2 will be presented because it has the largest similarity value

Next we are going to try and reduce the amount of computation. We are going to do this by filtering out certain documents and then doing the computations.

Efficient Ranking

- Solve K-nearest neighbour problem for a query vector.

$$\text{Sim}(\vec{Q}, \vec{D}_1) \quad \dots \quad \text{Sim}(\vec{Q}, \vec{D}_2)$$

↓ We are looking for the largest. ↓

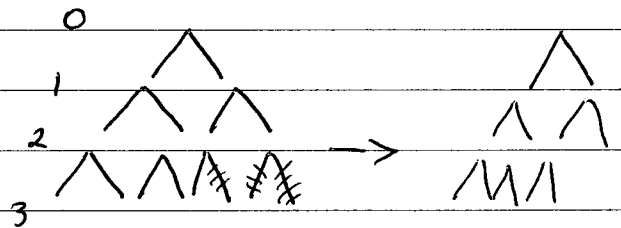
Heapify Exercise

[3, 1, 6, 2, 0, 9]

Heap: Complete Binary Tree, from L → R Order

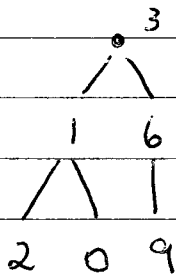
Create the heap produced for this list by the heapify algorithm.

A heap is a complete binary tree with leaves removed in Right to Left Order



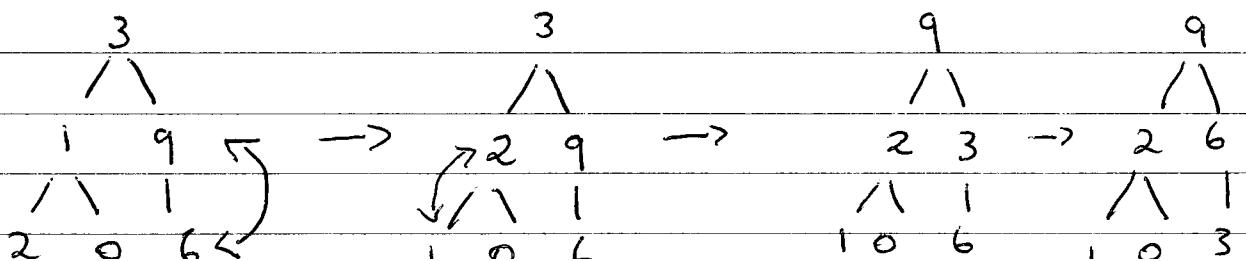
Label of Parent ⇒ labels of children (in particular: largest label is at root)

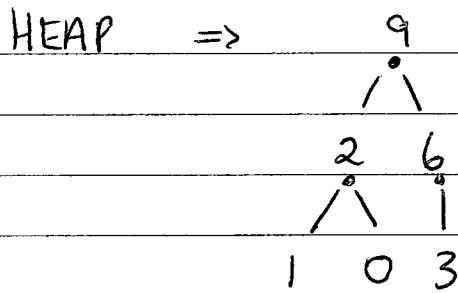
A) Create nice tree labeled with list element.



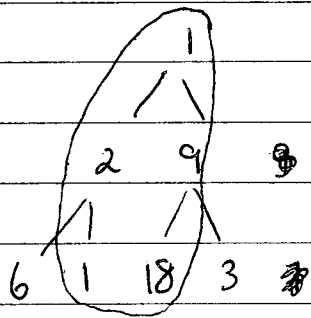
* Last Time / Near Heap

B) For loop :





Use Binary Tree - Heap for selecting top K



Generic Approach - Find set A of candidates
 $K < |A| \ll N$

o Static Quality Scores

- Top ranking documents to be both relevant and authoritative.

Relevance - Cosine Scores

Authority - query-independent property of a document
 eg Wikipedia / Newspapers / paper citations.

Net score

$$\text{net-score}(q, d) = g(d) + \text{cosine}(q, d)$$

Top K by Net Score

1/ Order all posts by $g(d)$

2/

CS4611

Tuesday 12th November 2013

Assignment 2: Cosine Similarity

Term Freq. (don't use \log)
Find closest match. $\frac{\text{term freq.}}$

(See Notes)

Measures for a search engine

- Measurable
 - User Happiness - Can we measure this?
 - Free
 - Who is the user?
- Measure
- Click-through rate
 - Rate of return to this search
 - Time to purchase, fraction of "conversions" of searches to buyers

User Happiness \rightarrow feedback is relevant.

How to you measure relevance? \rightarrow
(proximity - cosine similarity)

Relevance to what?

- Information need.
- Query q
- Some results returned may match the query but is not what we are looking for.

User happiness can only be measured by relevance to the information required.

$$\Delta \text{ precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant} | \text{retrieved})$$

$$\Delta \text{ recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = R(\text{retrieved} | \text{relevant})$$

Accuracy

This is the fraction of decisions (relevant | nonrelevant) that are correct.

$$\text{accuracy} = \frac{TP + TN}{(TP + FP + FN + TN)}$$

\Rightarrow roughly = 1, close to 100%

\Rightarrow Not right way to go about it.

FP = False Positive	- Small - Retrieves nothing - zero
FN = False Negative	- Small
TP = True Positive	- Small
TN = True Negative	- Large

} relevance

$$\frac{1}{2} P + \frac{1}{2} R = \text{Arithmetic Mean.}$$

(Harmonic mean)

Returns everything - answer is there for sure - 50% accuracy
this is still very bad.

$$\frac{1}{2} P + \frac{1}{2} R = \frac{1}{2} \frac{20}{1,000,120} + \frac{1}{2} \frac{20}{20} \approx \frac{1}{2}$$

F: Harmonic Mean

We can view the HM as a kind of smooth minimum.

$$F = \frac{1}{\frac{1}{2P} + \frac{1}{2R}}$$

$$F_1 = \frac{2PR}{P+R}$$

FORMULA

$$= \frac{1}{\frac{1}{2P} + \frac{1}{2R}}$$

$$= \frac{1}{\frac{2P+2R}{2P \cdot 2R}} = \frac{4PR}{2P+2R}$$

What we need for a benchmark

- Collection of documents
- Collection of information needs (i.e. queries)
- Human relevance assessments.

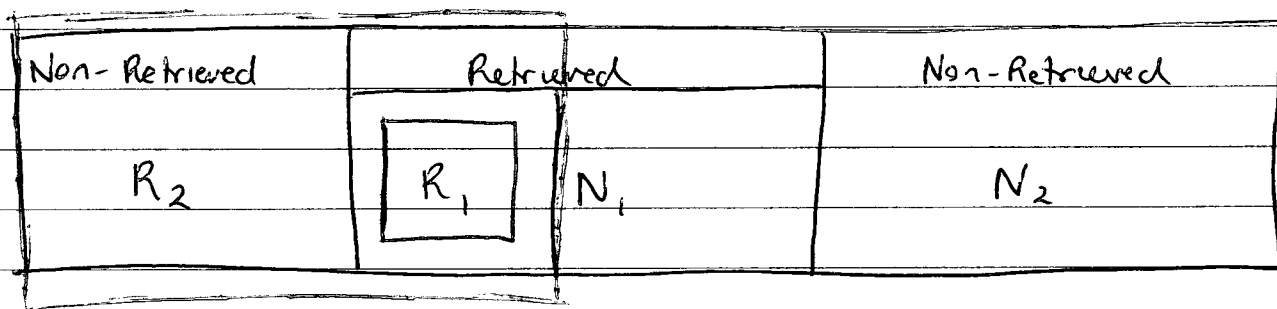
Benchmarks

- Cranfield (quite old, too small & too untypical for today)
 - TREC (Text Retrieval Conference) - Too expensive
 - CQV2
 - NTCIR
 - CLEF
- } others

$$\text{Precision (P)} = \frac{\# \text{ relevant items retrieved}}{\# \text{ retrieved items}}$$

$$\text{Recall (R)} = \frac{\# \text{ relevant items retrieved}}{\# \text{ relevant items}}$$

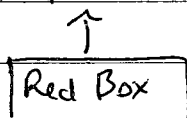
Enumerator
Same



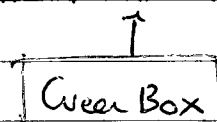
Not Relevant = N

Relevant = R

$$P = \frac{\# R_1}{(\# R_1 + \# N_1)}$$



$$R = \frac{\# R_1}{(\# R_1 + \# R_2)}$$



△ Aim : To represent P and R by via single number

* Decided : Use something like $\text{Min}(P, R)$

$$\begin{array}{l} \text{System 1 : } P = \frac{1}{8} \quad R = \frac{5}{9} \\ \text{System 2 : } P = \frac{1}{8} \quad R = \frac{7}{9} \end{array}$$

S2 choose because
it has slightly higher
recall

$$M(P, R) \text{ for } S1 = M(P, R) \text{ for } S2 = \text{Same } \left(\frac{1}{8}\right)$$

Exercise

What happens with our measure F (Harmonic Mean),
say $\alpha = \frac{1}{2}$, $F = \frac{1}{\frac{1}{2}P + \frac{1}{2}R} = \frac{1}{\frac{1}{2}P + \frac{1}{2}R}$

$$\frac{1}{\frac{1}{2}P + \frac{1}{2}R} = \frac{2PR}{P+R} = F$$

$$\begin{array}{l} F(\text{System 1}) = 0.204 = 0.2 \\ F(\text{u 2}) = 0.215 = 0.22 \quad \checkmark \text{ Better} \end{array}$$

KAPPA Measure

$$K = \frac{P(A) - P(E)}{1 - P(E)}$$

$P(A)$ = proportion of time judges agree
 $P(E)$ = what agreement would we get by chance

Example:

$$\frac{(300 + 70)}{400} = \frac{370}{400} = 0.925$$

$$P(\text{non-rel}) = (80 + 90) / (400 + 400) = \frac{170}{800} = 0.2125$$

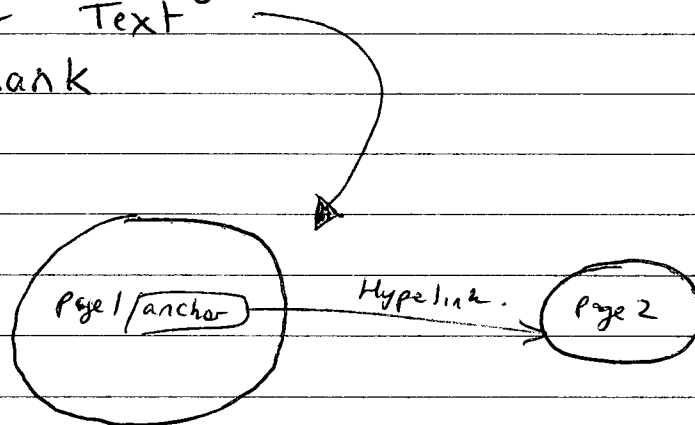
$$P(\text{relevant}) = (320 + 310) / (400 + 400) = \frac{630}{800} = 0.7875$$

Summary

- △ Static Summary (subset of document, eg 1st 50 words)
- △ Dynamic Summary (snippets of info, that contain a no. of query elements - also occur together "the ball", hit "the ball")

★ Evaluating webpages - links that go to them.
Google Page Ranking.

- △ Citation Analysis
- △ Anchor Text
- △ Page Rank
- △ Hubs



Anchor Text is after a better description than the text itself.

Matrix Multiplication

$$\begin{bmatrix} 1 & 3 & 2 \\ 0 & 4 & 0 \end{bmatrix} \times \begin{bmatrix} 1 \\ 6 \\ 2 \end{bmatrix} = \begin{bmatrix} 23 \\ 24 \end{bmatrix}$$

\downarrow
 2×3
 \downarrow
 rows cols

\downarrow
 3×1

\rightarrow row 1 \times col 1
 $1 \times 1 = 1$
 $3 \times 6 = 18$ add up = 23
 $2 \times 2 = 4$

\rightarrow row 2 \times col 1
 $0 \times 1 = 0$
 $4 \times 6 = 24 \rightarrow 24$
 $0 \times 2 = 0$

Page Ranking

(Compute important of webpage - links there are)

Consider $N \times N$ matrices store probability of moving from page i to page j

$$\begin{matrix} N & \begin{bmatrix} X & P_{12} & P_{13} \\ \cdot & \cdot & \cdot \\ P_{N1} & & P_{NN} \end{bmatrix} \\ & N \end{matrix}$$

P = Probability that user will jump from Page 1 to Page N
 X = Does not jump - stays on page.

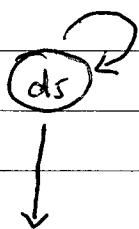
Slides : Example web Graph

High Probability Page is d_3 as all other can reach it through links, d_5 is isolated.

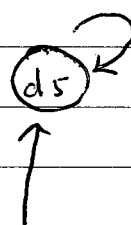
d_0 = no link so 0
 d_2 = link to itself so 1

	d_0	d_1	d_2
d_0	0	.	.
d_1	.	.	.
d_2	.	.	1

* The web graph must corr. with ergodic Markov chain
 The web graph must not contain any dead ends.



No dead end



Dead end.

If you reach a dead-end the probability to jump to another page is equal. $(1/N)$

A non-dead-end - probability 10% will jump to a random webpage. $(0.1/N)$

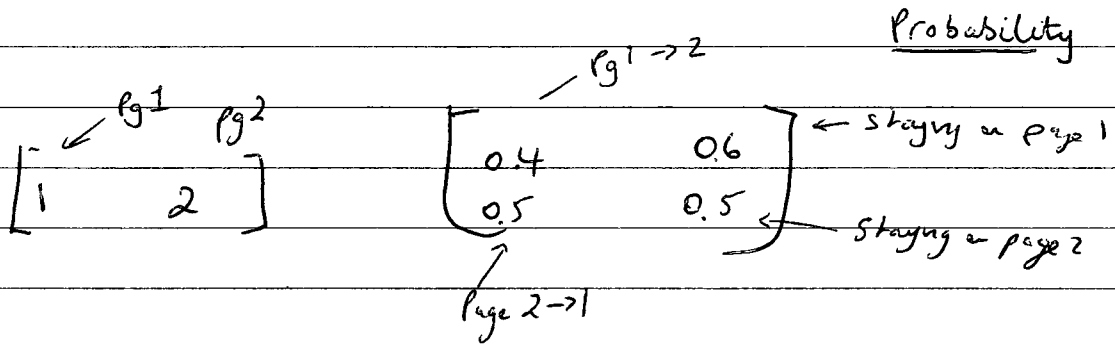
90% keep going $(0.9/N)$ ie following links.

Δ 10% is a parameter, the Teleportation Rate (TR)
 * "Jumping" from a dead end is independent from the TR

* Irreducibility }
 * Aperiodicity } cheer slides - Intro. to Inf. Ret. (115)

Formalisation of "visit"

How Random!



As X_m grows they will always stabilise - go to a steady state, ie when it reaches 1,000 or 1 million we will stop here.

Pg 170
Slides

$$P = \begin{bmatrix} 0.1 & 0.9 \\ 0.3 & 0.7 \end{bmatrix} \leftarrow \text{stay on Pg 2}$$

only give the bottom values - 0.3 0.7



$$E_0 \quad 0 \quad 1$$

⏟

Random choice will always give the same outcome

0.25 0.75
chances of landing
on page 1 + 2

0.75 high probability of moving to Pg 2 and staying on Pg 2.



$$\vec{\pi} = (x_1, x_2)$$

$$\vec{\pi} = \vec{\pi} P$$

$$(x_1, x_2) = (x_1, x_2) P$$

$$P = \begin{bmatrix} 0.1 & 0.9 \\ 0.3 & 0.7 \end{bmatrix}$$

$$(x_1, x_2) = \begin{matrix} 0.1 x_1 + 0.3 x_2 \\ 0.9 x_1 + 0.7 x_2 \end{matrix}$$

$$x_1 = 0.1 x_1 + 0.3 x_2$$

$$x_2 = 0.9 x_1 + 0.7 x_2$$

- Solve system of simultaneous equations
- Should give the same answer (0.25, 0.75)
 - There are two ways to do this.

Starting Vector $(x_1, x_2) = (0, 1)$

Step 1: Solve power method

Step 2: Solve $\vec{\pi} = \vec{\pi} P$ ($\vec{\pi} = (x_1, x_2)$)

Solution on pg 185 of slides

$$\vec{\pi} P = (0.7 x_1 + 0.2 x_2, 0.3 x_1 + 0.8 x_2)$$

$$x_1 = 0.7 x_1 + 0.2 x_2$$

$$x_2 = 0.3 x_1 + 0.8 x_2$$

SOLVE

$$(x_1, x_2) = 0.4, 0.6$$