



DEAKIN UNIVERSITY Information 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). o These days we frequently think first of web search, but there are many other cases: Searching your laptop Corporate knowledge bases Legal information retrieval

3

Unstructured (text) vs. structured (database) data in 1996

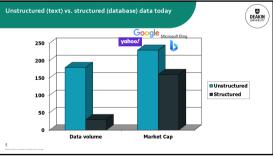
Data volume

Market Cap

250

200 150

Unstructured (text) vs. structured (database) data today DEAKIN UNIVERSITY



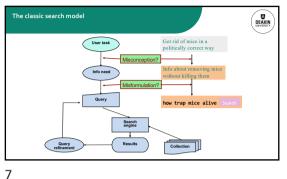
Basic assumptions of Information Retrieval DEAKIN UNIVERSITY • Collection: A set of documents o Assume it is a static collection for the moment • Goal: Retrieve documents with information that is relevant to the user's information need and helps the user complete a task

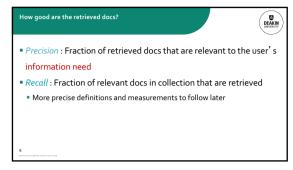
4 5 6

DEAKIN UNIVERSITY

■Unstructured

■ Structured

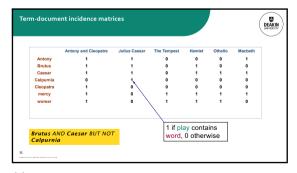




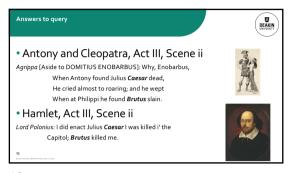


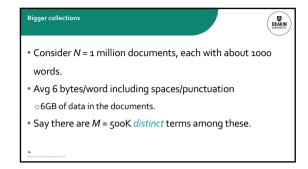
Which plays of Shakespeare contain the words Brutus AND Caesar but NOT Calpurnia?
 One could grep all of Shakespeare's plays for Brutus and Caesar, then strip out lines containing Calpurnia?

Why is that not the answer?
 Slow (for large corpora)
 NOT Calpurnia is non-trivial
 Other operations (e.g., find the word Ramans near countrymen) not feasible
 Ranked retrieval (best documents to return)
 Laterlectures



10 11 12





• 500K x 1M matrix has half-a-trillion o's and 1's.

• But it has no more than one billion 1's.

• matrix is extremely sparse.

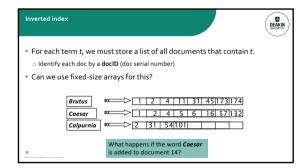
• What's a better representation?

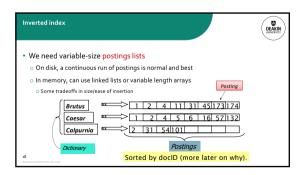
• We only record the 1 positions.

15

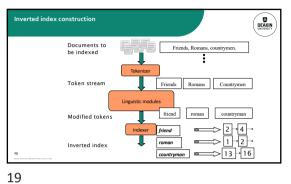
13 14

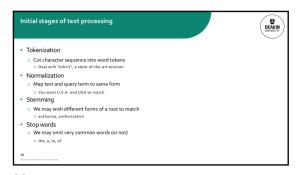


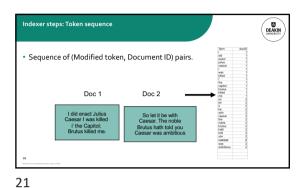


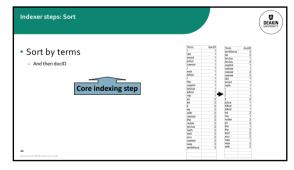


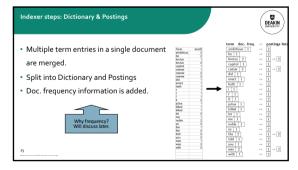
16 17 18

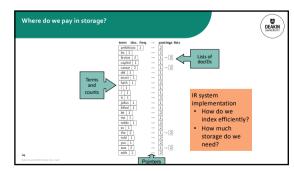




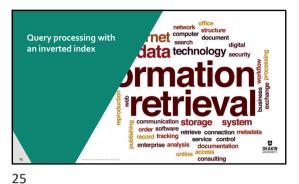


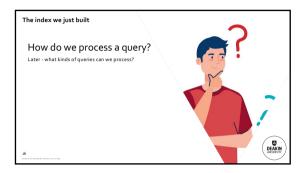


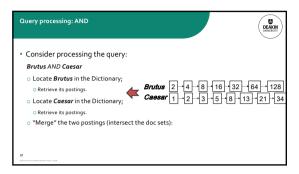




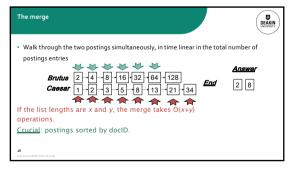
22 23 24

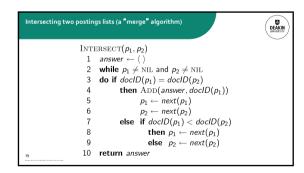


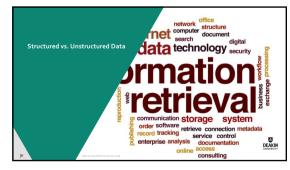




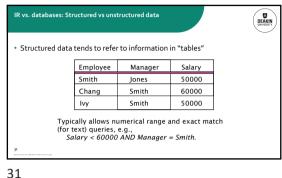
25 26 27

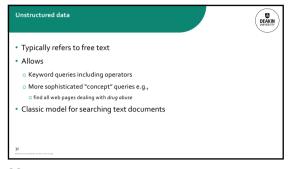






28 29 30





• In fact almost no data is "unstructured"

• E.g., this slide has distinctly identified zones such as the Title and Bullets

• ... to say nothing of linguistic structure

• Facilitates "semi-structured" search such as

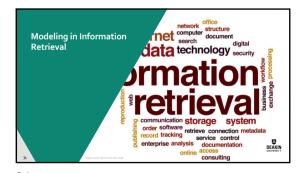
• Title contains data AND Bullets contain search

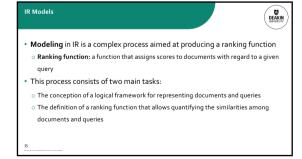
• Or even

• Title is about Object Oriented Programming AND Author something like strozup

• where * is the wild-card operator

32 33

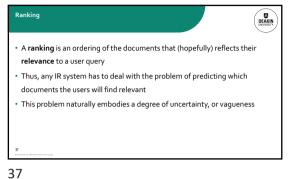


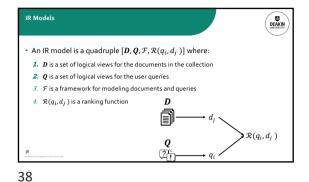


IR systems usually adopt index terms to index and retrieve documents
 Index term:
 In a restricted sense: it is a keyword that has some meaning on its own; usually plays the role of a noun
 In a more general form: it is any word that appears in a document
 Retrieval based on index terms can be implemented efficiently
 Also, index terms are simple to refer to in a query
 Simplicity is important because it reduces the effort of query formulation

35 36

6





A Taxonomy of IR Models

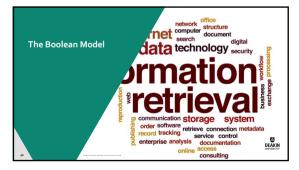
| Committee | Com

39

3/

In this lecture, we will discuss the following models:
 The Boolean Model
 The Vector Model
 Probabilistic Model

40



The Boolean retrieval model is being able to ask a query that is a Boolean expression:

Boolean Oueries are queries using AND, OR and NOT to join query terms

Views each document as are to fivords

Is precise document matches condition or not.

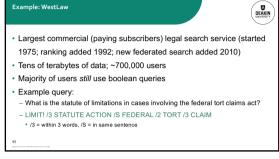
Perhaps the simplest model to build an IR system on

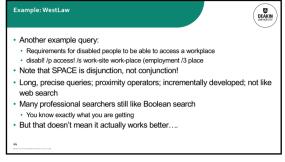
Primary commercial retrieval tool for 3 decades.

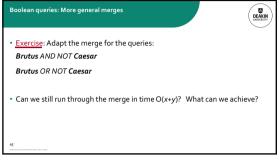
Many search systems you still use are Boolean:

Email, library catalog, Mac OS X Spotlight

41 42





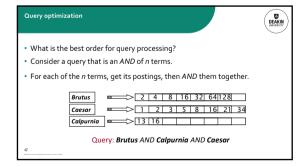


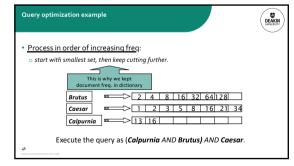
What about an arbitrary Boolean formula?
(Brutus OR Caesar) AND NOT (Antony OR Cleopatra)

Can we always merge in "linear" time?

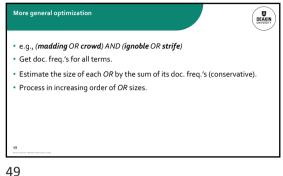
Linear in what?

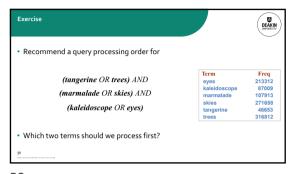
Can we do better?

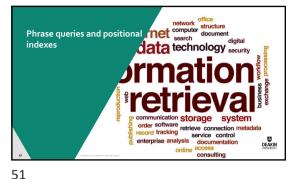




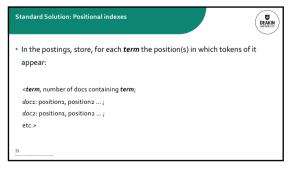
46 47 48







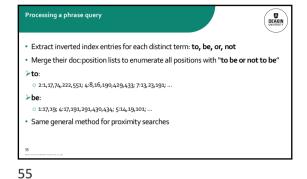
We want to be able to answer queries such as "stanford university" – as a phrase
 Thus the sentence "I went to university at Stanford" is not a match.
 The concept of phrase queries has proven easily understood by users; one of the few "advanced search" ideas that works
 Many more queries are implicit phrase queries
 For this, it no longer suffices to store only
 <term: docs> entries

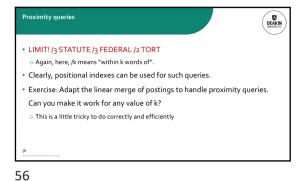


Positional index example

See: 993427;
Which of does 1,2,4,5
could contain "to be"?
to print to be"?

52 53 54





A positional index expands postings storage substantially
 Even though indices can be compressed
 Nevertheless, a positional index is now standardly used because of the power and usefulness of phrase and proximity queries ... whether used explicitly or implicitly in a ranking retrieval system.

Positional index size

Positional index size

Need an entry for each occurrence, not just once per document

Index size depends on average document size
Average web page has <1000 terms
SEC filings, books, even some epic poems ... easily 100,000 terms

Consider a term with frequency 0.1%

Document size
Postings
Positional postings
1000
1 100,000
1 1000

58

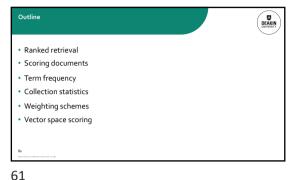
Rules of thumb

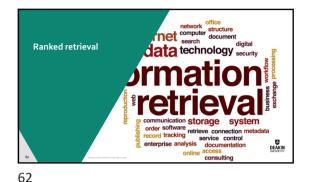
 A positional index is 2–4 as large as a non-positional index
 Positional index size 35–50% of volume of original text
 Caveat: all of this holds for "English-like" languages

The Vector Model

The Vector M

59 60





Ranked retrieval DEAKIN UNIVERSITY · So far, our gueries have all been Boolean o Documents either match or don't • Good for expert users with precise understanding of their needs and the collection Also good for applications: Applications can easily consume 1000s of results · Not good for the majority of users Most users incapable of writing Boolean queries (or they are, but they think it's too much work) o Most users don't want to wade through 1000s of results o This is particularly true of web search

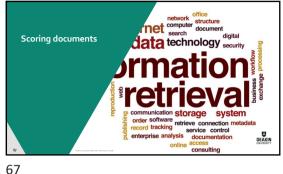
63

Problem with Boolean search: feast or famine DEAKIN UNIVERSITY • Boolean queries often result in either too few (=0) or too many (1000s) Query 1: "standard user dlink 650" → 200,000 hits • Query 2: "standard user dlink 650 no card found": o hits • It takes a lot of skill to come up with a query that produces a manageable number of hits. o AND gives too few; OR gives too many 64

Ranked retrieval models DEAKIN UMVERSITY • Rather than a set of documents satisfying a query expression, in ranked retrieval, the system returns an ordering over the (top) documents in the collection for a query • Free text queries: Rather than a query language of operators and expressions, the user's query is just one or more words in a human language • In principle, there are two separate choices here, but in practice, ranked retrieval has normally been associated with free text queries and vice versa

Feast or famine: not a problem in ranked retrieval DEAKIN UNIVERSITY • When a system produces a ranked result set, large result sets are not an o Indeed, the size of the result set is not an issue We just show the top k (≈ 10) results We don't overwhelm the user o Premise: the ranking algorithm works

65 66



Scoring as the basis of ranked retrieval DEAKIN UNIVERSITY • We wish to return in order the documents most likely to be useful to the • How can we rank-order the documents in the collection with respect to a • Assign a score - say in [o, 1] - to each document • This score measures how well document and query "match".

68

Query-document matching scores DEAKIN UNIVERSITY • We need a way of assigning a score to a query/document pair · Let's start with a one-term query • If the guery term does not occur in the document: score should be o • The more frequent the query term in the document, the higher the score · We will look at a number of alternatives for this

69

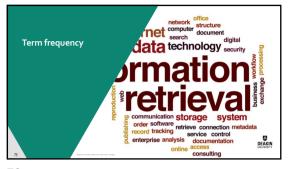
Take 1: Jaccard coefficient DEAKIN UNIVERSITY • $jaccard(A,B) = |A \cap B| / |A \cup B|$ • jaccard(A,A) = 1 • $jaccard(A,B) = o if A \cap B = o$ • A and B don't have to be the same size • Always assigns a number between o and 1

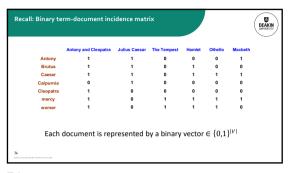
70

Jaccard coefficient: Scoring example DEAKIN UMVERSITY What is the guery-document match score that the Jaccard coefficient computes for each of the two documents below? Query: ides of march Document 1: caesar died in march • Document 2: the long march

Issues with Jaccard for scoring DEAKIN UNIVERSITY • It doesn't consider term frequency (how many times a term occurs in a document) · Rare terms in a collection are more informative than frequent terms. Jaccard doesn't consider this information • We need a more sophisticated way of normalizing for length • Later in this lecture, we'll use $|A \cap B|/\sqrt{|A \cup B|}$ • . . . instead of $|A \cap B|/|A \cup B|$ (Jaccard) for length normalization.

71 72





• Consider the number of occurrences of a term in a document:

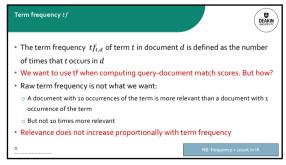
• Each document is a count vector in N': a column below

Autony and Cleopatra Not so to the sound vector in Structure in the sound vector in N': a column below

Autony and Cleopatra Not sound vector in Structure in the sound vector in Structure in the sound vector in Structure in the sound vector in t

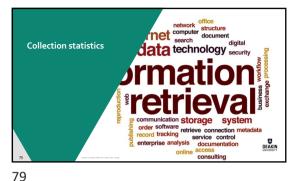
73 74 75

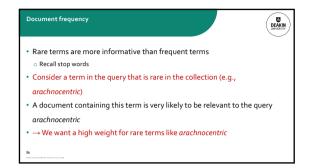
Vector representation doesn't consider the ordering of words in a document
John is quicker than Mary and Mary is quicker than John have the same vectors
This is called the bag of words model
In a sense, this is a step back: The positional index was able to distinguish these two documents
The IIR book considers "recovering" positional information
For now: bag of words model



Log-frequency weighting• The log frequency weight of term t in d is $w_{t,d} = \begin{cases} 1 + log_{10}(tf_{t,d}), & \text{if } tf_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}$ • Score for a document-query pair: sum over terms t in both q and d: $score(q,d) = \sum_{t \in Q \cap d} \left[1 + log_{10}(tf_{t,d}) \right]$ • The score is o if none of the query terms is present in the document

76 77 78





• Rare terms are more informative than frequent terms
• Recall stop words
• Consider a term in the query that is rare in the collection (e.g., arachnocentric)
• A document containing this term is very likely to be relevant to the query arachnocentric
• → We want a high weight for rare terms like arachnocentric

81

79

80

Frequent terms are less informative than rare terms

Consider a query term that is frequent in the collection (e.g., high, increase, line)

A document containing such a term is more likely to be relevant than a document that doesn't

But it's not a sure indicator of relevance.

For frequent terms, we want high positive weights for words like high, increase, and line

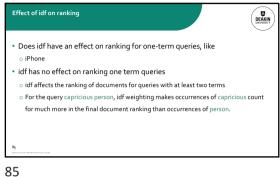
But lower weights than for rare terms

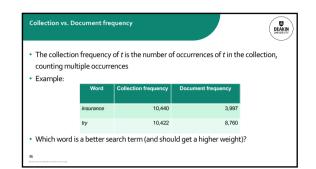
We will use document frequency (df) to capture this

82

 $df_t \text{ is the } \frac{\text{document } f_t \text{ equency of } t\text{: the number of documents that contain } t \\ o df_t \text{ is an inverse measure of the informativeness of } t \\ o df_t \leq N \\ \bullet \text{ We define the } idf \text{ (inverse document frequency) of } t \text{ by } \\ idf_t = log_{10} \binom{N}{df_t} \\ \bullet \text{ We use } log_{10} \binom{N}{df_t} \text{ instead of } \binom{N}{df_t} \text{ to "dampen" the effect of idf} \\ \hline \text{Will turn out the base of the log is immaterial.}$

83 84

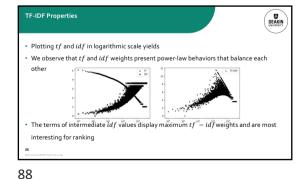




**Consider the tf, idf, and tf-idf weights for the Wall Street Journal reference collection
**To study their behavior, we would like to plot them together
**While idf is computed over all the collection, tf is computed on a per document basis. Thus, we need a representation of tf based on all the collection, which is provided by the term collection frequency
**This reasoning leads to the following tf and idf term weights: $w_t = 1 + log_{10} \sum_{j=1}^{N} tf_{i,j}, \quad ldf_t = log_{10} \binom{N}{df_t}$

87

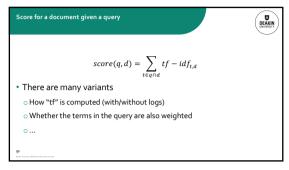
86

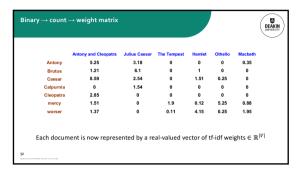




• The tf-idf weight of a term is the product of its tf weight and its idf weight $tf-idf_{t,d}=\left(1+log_{10}(tf_{t,d})\right)\times log_{10}\left({}^{N}\!/{}_{df_{t}}\right)$ • Best known weighting scheme in information retrieval
• Note: the "-" in tf-idf is a hyphen, not a minus sign!
• Alternative names: tf.idf, $tf\times idf$ • Increases with the number of occurrences within a document
• Increases with the rarity of the term in the collection

89 90

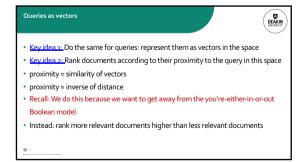




Vector space scoring

91 92 93

So we have a |V|-dimensional vector space
Terms are axes of the space
Documents are points or vectors in this space
Very high-dimensional: tens of millions of dimensions when you apply this to a web search engine
These are very sparse vectors - most entries are zero



Formalizing vector space proximity

• First cut: distance between two points
o (= distance between the end points of the two vectors)
• Euclidean distance?
• Euclidean distance is a bad idea . . .
• . . . because Euclidean distance is large for vectors of different lengths

94 95 96

• The Euclidean distance between q and d_2 is large even though the distribution of terms in the query q and the distribution of terms in the document d_2 are very similar.

Thought experiment: take a document d and append it to itself. Call this document d'

"Semantically" d and d' have the same content

The Euclidean distance between the two documents can be quite large

The angle between the two documents is o, corresponding to maximal similarity

Key idea: Rank documents according to angle with query

The following two notions are equivalent

Rank documents in decreasing order of the angle between query and document

Rank documents in increasing order of cosine(query,document)

Cosine is a monotonically decreasing function for the interval [0°, 180°]

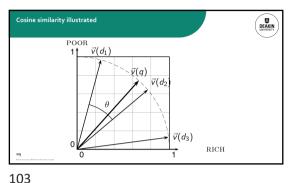
But how — and why — should we be computing cosines?

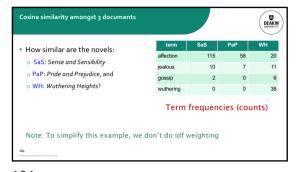
97 98

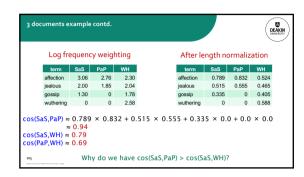
100

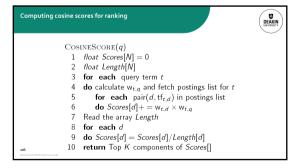
cosine(query, document) $cos(\vec{q}, \vec{d}) = \frac{\vec{q} \bullet \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\vec{q}}{|\vec{q}|} \bullet \frac{\vec{d}}{|\vec{d}|} = \frac{\sum_{i=1}^{|\vec{r}|} q_i d_i}{\sqrt{\sum_{i=1}^{|\vec{r}|} q_i^2} \sqrt{\sum_{i=1}^{|\vec{r}|} d_i^2}}$ $q_i \text{ is the tf-idf weight of term } i \text{ in the query } d_i \text{ is the tf-idf weight of term } i \text{ in the document}}$ $cos(q, d) \text{ is the cosine similarity of } q \text{ and } d \dots \text{ or, } equivalently, the cosine of the angle between } q \text{ and } d.$

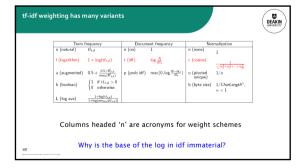
• For length-normalized vectors $\cos(\vec{q},\vec{d}) = \vec{q} \cdot \vec{d} = \sum_{l=1}^{|V|} q_l d_l$ o for q,d length-normalized.











Weighting may differ in queries vs documents

Many search engines allow for different weightings for queries vs. documents

SMART Notation: denotes the combination in use in an engine, with the notation ddd.qqq, using the acronyms from the previous table

A very standard weighting scheme is: Inc.ltc A bad idea?

Document: logarithmic tf (I as first character), no idf and cosine normalization

Query: logarithmic tf (I in leftmost column), idf (t in second column), no normalization ...

106 107 108

