Understanding how search engines work, an Information Retrieval perspective

Pedro Contreras pedro@cs.rhul.ac.uk

Department of Computer Science Royal Holloway, University of London

26 February 2008



Enterprise Search

Enterprise search characteristics:

- Diversity of content sources and formats, and not necessarily HTTP based

- Secure access
- Combined structured and unstructured search
- Dedicated search (e.g. email search)
- Intranet search
- Ranking and categorisation problem
- Social forces behind the creation of Internet and Intranet content are quite different
- Deployment environments for these domains also differs

Department of Computer Science. Royal Holloway, University of London

3

Enterprise Search comparison with Web searchSimilaritiesDifferences• Crawling- Diverse repositories• Indexing and ranking- Rewer documents• User interface- Access control• Diverse doc type- Different user needs• Less filtering

Data preparation and data source

Different document types requiered different access methods. PDF, MS Word, HTML, Open formats, email, etc.

Good Formatted file

sent="Sat, 20 Apr 1996 17:44:20 -0400 (EDT)"

id="830036660.25042.RFW@MATH.AMS.ORG"

subject="Re: conference call Mon 22 April'

have no specificitems for the agenda now.

inreplyto="199604201303.6479@uvea.wolfram.com"

I was unable to attend last week, but am available this

Monday. If the group plans to meet, I'll attend, but I

docno="lists-046-11826843"

email="RFW@math.ams.org"

name="Ron Whitney"

* Mail-System-Version:

-Ron

To: w3c-math-erb@w3.org

Bad formatted file

name="" email="<u>mail137@163.com</u>" sent="Wed, 8 May 2002 14:39:51 -0400" subject="øΪÀŸÕΔπ,*f*'µfÕ⁺æ⁰∞≤Δ∑" To: <u>web-access@w3.org</u>

 $\begin{array}{l} \overset{a\bullet}{\to} i^{\mathsf{TM}} \widetilde{O} \Sigma \varphi^{*} \pi \mu \Omega \Omega \widetilde{O} \widetilde{A} \widetilde{I}^{*} - \mathfrak{B} \pm^{a} \widehat{V} \widehat{O} \pm \mu f \dots \widetilde{A}^{ou} \widetilde{A}^{*} \widetilde{c}^{**} \sqrt{\Sigma^{**}}, \quad \langle \widetilde{O} \Delta \pi, \dots \widetilde{A}^{ou} \leq \Delta \Sigma \geq \dots \pm \mathfrak{B} \Diamond \widetilde{\mu} \widetilde{O} \widetilde{I} \widetilde{D} \widetilde{I}^{*} = \mathcal{I}^{\mathsf{TM}} \widetilde{D} \widetilde{I} \widetilde{D} \widetilde{I} = \mathcal{I}^{\mathsf{TM}} \widetilde{I} = \mathcal{I}^{\mathsf{$

$$\begin{split} & \mu f \leq \Delta \sum \mathfrak{e} \leq \mathfrak{P} \Omega \tilde{\mathcal{O}}^{-} \dots \mathfrak{E}^{\circ} \mathfrak{E} \mathfrak{E}^{+} / \mathfrak{L}^{-} \mathbb{E}^{+} \sqrt{\mathfrak{O}} \mathbb{B} \pi^{\mathfrak{e}_{\mathfrak{q}}} - \mathbb{O} \wedge \mathbb{B}^{+} \mu f^{\mathfrak{s}})^{\mathfrak{d}_{\mathfrak{q}}} \\ & \mathfrak{L}^{+} f \leftarrow - \prod A \tilde{Y} \mu f \Omega^{-} f \mu f \tilde{\mathcal{O}}^{-} \mathfrak{B}^{\circ} \mathfrak{S}^{\circ} \mathfrak{S} \leq \Delta \sum \tilde{\mathcal{O}} \Delta \pi_{\mathfrak{s}} \geq \mathfrak{S}^{\circ} \\ & \rightarrow \mathfrak{L}^{\circ} \end{split}$$

5

---Œ"√«µƒ∑ॄœÒ---

1°¢√,∑−∑,ŒÒ£∫ √,∑−÷π©πӱ"^∙j™Õ™™œ∑Ω∑®μƒ∞ ι◊...−Ø

Department of Computer Science. Royal Holloway, University of London

Parsing, tokenisation Tokenisation is the process of splitting a stream of words into units or "tokens". Normally this process does not included the following symbols: - period (.) - comma (,) - semicolon (;) - quotation marks (") - colon (:) - brackets [] - braces { } - parentheses () - mathematical operators + - / * = < > - special characters | & ~ - the at sign @ - underscores and other rare characters e.g The "brown" fox jumps, quickly over the lazy dog*. The brown fox jumps quickly over the lazy dog Department of Computer Science. Royal Holloway, University of London 6

Parsing, stop words removal

Stop Word is the name given to a word that will be filtered and is not consider relevant by an information retrieval system. Some of the more frequently used stop words for English include: "a", "of", "the", "I", "it", "you", and "and". These are generally regarded as functional words that do not carry meaning for the system.

Stop List is the list or set of *stop words*, there are as many *stop lists* as there are languages. I.e. if a system processes text that includes English, French, German and Spanish it also should be a *stop list* for each of these languages.





Department of Computer Science. Royal Holloway, University of London

7

Parsing, entity detection The problem here is to find various structured data within unstructured documents, e.g. people's names project's names places amounts Algorithms for entity detection normally are either rule or statistical based. see: Special Interest Group on Natural Language Learning on the Association for

see: Special Interest Group on Natural Language Learning on the Association for Computational Linguistics (CoNLL). http://cnts.uia.ac.be/signll/conll.html

Department of Computer Science. Royal Holloway, University of London

Parsing, part of speech Part of speech is a categorisation process that take in consideration phrase function. Each part of speech explains not what the word is, but how the word is used. POS is very useful when dealing with

Natural Language Processing in IR.

Parts of speech: the verb, the noun, the pronoun, the adjective, the adverb, the preposition, the conjunction, and the article.

e.	g
----	---

can	l think l can do it.	verb
can	Don't open that can of beans.	noun
only	This is my only pen.	adjective
only	He was only joking.	adverb
his	That book is <mark>his.</mark>	pronoun
his	That is <mark>his</mark> book.	adjective
English	Can you speak English?	noun
English	I am reading an English novel.	adjective

Department of Computer Science. Royal Holloway, University of London

9





Department of Computer Science. Royal Holloway, University of London







System exploitation, querying



















Baire, or longest common prefix

Definition: a Baire space consists of countably infinite sequence with a metric defined in terms of the longest common prefix

Case of vectors x and y, with 1 attribute. Precision: digits 1, 2, ..., |K|

$$d_B(x_K, y_K) = \begin{cases} 1 & \text{if } x_1 \neq y_1 \\ \inf 2^{-n} & x_n = y_n & 1 \le n \le |K| \end{cases}$$

- each coordinate is normalised, so is a floating point value.
- then: we will define $d_B(x,y)$ based on sharing common prefix in all coordinates.

Baire, or longest common prefix

An example of Baire distance for two numbers (x and y) using a precision of 4















