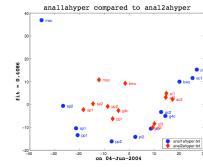


Data Mining and Information retrieval

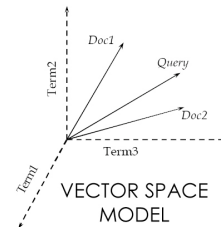


$$\begin{matrix} 1, 2, \dots, r, \\ r+1, r+2, \dots, n \end{matrix} \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & & \\ \vdots & & \ddots & \\ 0 & \dots & 0 & 1 \end{pmatrix} \begin{pmatrix} \square & \dots & * \\ \vdots & \ddots & \vdots \\ * & \dots & * \\ 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{pmatrix}$$

Pedro Contreras
pedro@cs.rhul.ac.uk

Department of Computer Science
Royal Holloway, University of London

27 February 2008



Overview, Lecture I

Data Mining

What's Data?

Record data, numerical data, data matrix,
document data, graph data, chemical data, etc.

What's Data Mining?

Why Data Mining?

Commercial viewpoint
Scientific viewpoint

Overview, Lecture I

Mining Data sets

Association Rules

Classification

Clustering

Forecasting

Challenges of Data Mining

References

Information Retrieval, Lecture II

What's Information Retrieval

Information Retrieval and Business Intelligence

Data preparation - Parsing

- Tokenisation
- Stop words removal
- Stemming
- Entity detection
- Part of speech

Data storage - Indexing

- Index construction

Exploitation - Querying

- Exploiting data repositories

Questions

Lecture II

Information Retrieval

27 February 2008

What's Information Retrieval

In the science of searching for information in heterogeneous data sources such as:

- Documents
- Images
- Audio
- Video

What I.R. systems do?

- Take documents in any format
 - Break into words
- Create an index
- Search it very quickly
- Update index when collection changes

Information Retrieval and B.I.

Business Intelligence transforms data into valuable information that can be accessed to help decision makers develop strategic, tactical and operational planning initiatives.

Business Intelligences deals with data integration from structured and unstructured data sources.

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

Enterprise Search comparison with Web search

Similarities

- Crawling
- Indexing and ranking
- User interface

Differences

- Diverse repositories
- Fewer documents
- Access control
- Diverse doc type
- Different user needs
- Less filtering

Data preparation and data source

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

Good Formatted file

```
docno="lists-046-11826843"
name="Ron Whitney"
email="RFW@math.ams.org"
sent="Sat, 20 Apr 1996 17:44:20 -0400 (EDT)"
inreplyto="199604201303.6479@uuea.wolfram.com"
* Mail-System-Version:
id="830036660.25042.RFW@MATH.AMS.ORG"
subject="Re: conference call Mon 22 April"
To: w3c-math-erb@w3.org
```

I was unable to attend last week, but am available this Monday. If the group plans to meet, I'll attend, but I have no specific items for the agenda now.

-Ron

Bad formatted file

```
name=""
email="mail137@163.com"
sent="Wed, 8 May 2002 14:39:51 -0400"
subject="øíÀÿÖΔπ, f μf Ö" æ°∞≤ΔΣ"
To: web-access@w3.org

ª•ï™ÖΣç'πμΩΩÖÄï" —æ±ªΨÜð‡μf...ÄªÄ¿™√Σ", «ÖΔπ,...
Äª≤ΔΣ≥...±æøÖμÖ"÷"—βμf"™œ
≤fl-ªΣμ"—μf»Äª%œ»ªΩ@çμfÖ"æΣμ'»¥√³∞—Ö"æÖΔπ,≥ª•
Σª"—μf»ÄÄïÄð°...œÖΣ'μ'»¥√³∞—ø
ªç
μf≤ΔΣΣç≤ªμΩÖ"...œ»ª°Σœ÷/ΣCE"√«Ö@πªª—©ø@μμf»ªç
Σfç—πÄÿμfΩ'f'μfÖ"æ°∞≤ΔΣÖΔπ,≥
ª•Σ°
---CE"√«μfΣCEÖ---
```

```
1°ç√Σ—ΣCEÖΣç
√Σ—Äπ@πÿ"ª•ï™Ö™œΣΣΣ@μf∞1ø...—ø
```

Parsing, tokenisation

Tokenisation is the process of splitting a stream of words into units or “tokens”. Normally this process does not include 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

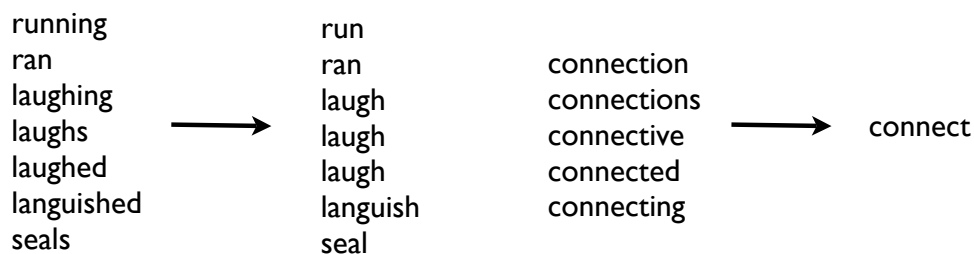
Parsing, stop words removal

Stop Word is the name given to a word that will be filtered and is not considered 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.

Parsing, Stemming

Stemming is the process to obtain a word's root (also called normalisation) through eliminating suffixes.



Benefits of Stemming

- Takes care of morphological variants
- Reduces index size

Known Limitations

- Impact on advanced syntax and exact match
- Accented characters are not supported
- Short words are not stemmed

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 Computational Linguistics (CoNLL). <http://cmts.uia.ac.be/signll/conll.html>

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	I think I 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 his .	pronoun
his	That is his book.	adjective
English	Can you speak English ?	noun
English	I am reading an English novel.	adjective

Data storage: indexing

After parsing

- Tokenisation
- Stop words removal
- Stemming
- Entity detection
- Part of speech

we can start the indexation process, which consist in storing the tokens in a DB, usually in a vector space fashion

Index design factors

- Merge factors
- Storage techniques
- Index size (compression)
- Lookup speed
- Maintenance
- Fault tolerance
- Scalability

Index construction, inverted index

T1: bab(y, ies, y's)
T2: child(ren's)
T3: guide
T4: health
T5: home
T6: infant
T7: safety
T8: toddler

D1: Infant & Toddler First Aid
D2: Babies & Children's Rooms (For Your Home)
D3: Child Safety at Home
D4: Your Baby's Health and Safety: From Infant to Toddler
D5: Baby Proofing Basics
D6: Your Guide to Easy Rust Proofing
D7: Beanie Babies Collector's Guide

$$A = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

Vector Space Model (VSM)

$$\mathbf{t}_i^T = [x_{i,1} \quad \dots \quad x_{i,n}] \quad \mathbf{d}_j = \begin{bmatrix} x_{1,j} \\ \vdots \\ x_{m,j} \end{bmatrix} \quad \mathbf{t}_i^T \rightarrow \begin{bmatrix} x_{1,1} & \dots & x_{1,n} \\ \vdots & \ddots & \vdots \\ x_{m,1} & \dots & x_{m,n} \end{bmatrix}$$

\mathbf{d}_j
↓

Binary	$\chi(f_{ij})$
Logarithmic	$\log(1 + f_{ij})$
Normal	$1/\sqrt{\sum_j f_{ij}^2}$
Inverse Document Frequency	$\log(n/\sum_j \chi(f_{ij}))$
etc.

Where:

f_{ij} : number of times term i in document j

Similarity

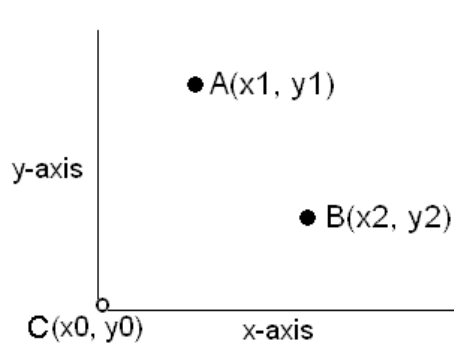
Similarity help to identifying the closeness between different vectors or in our case documents (or between a query and documents).

Similarity is based in a metric that defines distance

Euclidean distances in the classical example

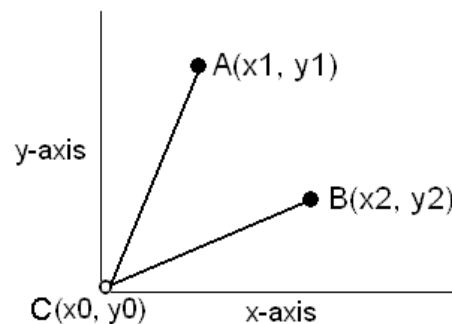
Similarity

Vector Space geometrically explained (2 dimensions)



Dot product

$$A \cdot B = x_1 \cdot x_2 + y_1 \cdot y_2$$

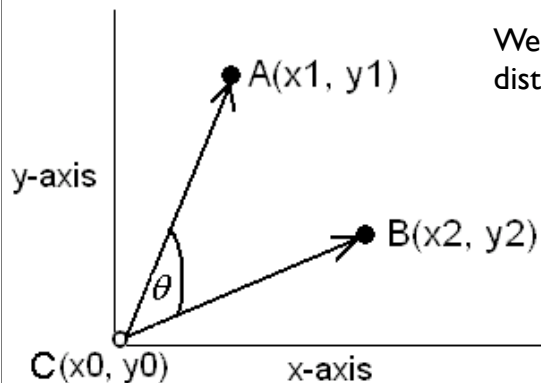


Pythagorean theorem

$$a^2 + b^2 = c^2 \quad c = \sqrt{a^2 + b^2}$$

Euclidean distance $a, b = d_{ab} = ((x_1 - x_0)^2 + (y_1 - y_0)^2)^{1/2} = (x_1^2 + y_1^2)^{1/2}$

Similarity, how far is A from B?



We normalise the dot product by the Euclidean distance. This ratio defines the cosine.

Dot product

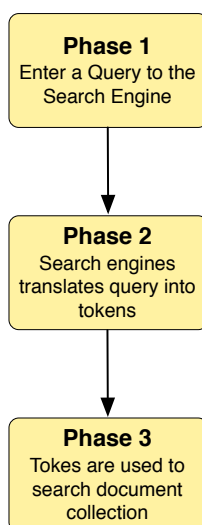
$$\text{Sim}(A, B) = \cosine \theta = \frac{A \bullet B}{|A||B|} = \frac{x1 \cdot x2 + y1 \cdot y2}{(x1^2 + y1^2)^{1/2} (x2^2 + y2^2)^{1/2}}$$

Distance

Similarity between a query and a document

$$\text{Sim}(Q, D_i) = \frac{\sum_j w_{Q,j} w_{i,j}}{\sqrt{\sum_j w_{Q,j}^2} \sqrt{\sum_i w_{i,j}^2}}$$

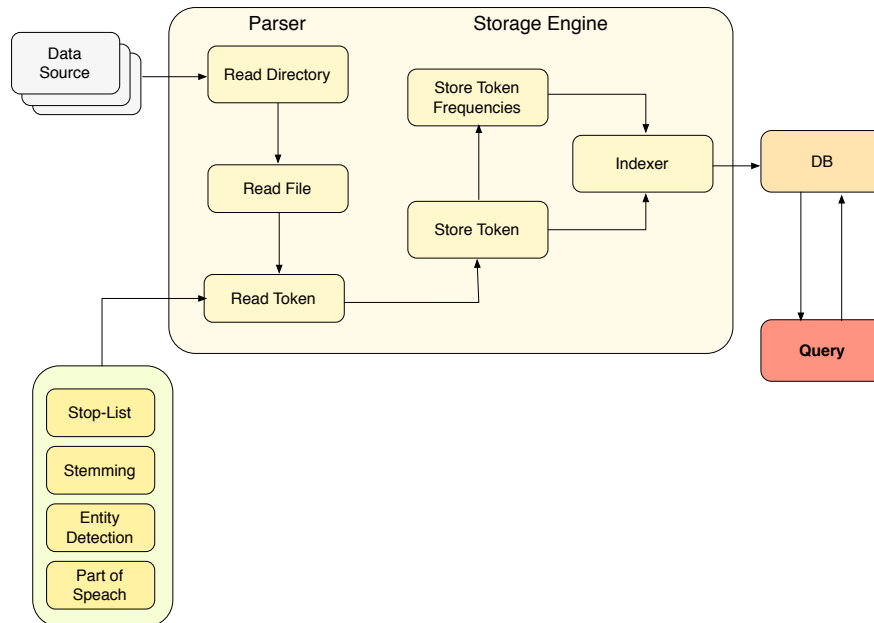
System exploitation, querying



Types of Queries

- **Boolean:** AND, OR, NOT
- **Natural Language Queries:** query is formulated as a question or a statement
- **Thesaurus Queries:** the user select the term from a previous term-set provided by the IR system
- **Fuzzy Queries:** threshold of relevance is expanded to include additional documents
- **Term Searches:** based in a few words or phrases provided by the user
- **Probabilistic Queries:** IR systems based in a computed probability to retrieve documents

An IR system overview



Sorting the query's result

To present a query's result to the user the output should be sorted in a meaningful way.

For example we can use the cosine similarity measure to rank result.

Also we can use clustering!

Ranking

For ranking we can use different mathematical functions based in word frequencies to determine the importance of a query's result.

Then we can sort this result in ascent order.

Clustering

Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that:

- Data points in one cluster are more similar to one another
- Data points in separate clusters are less similar to one another

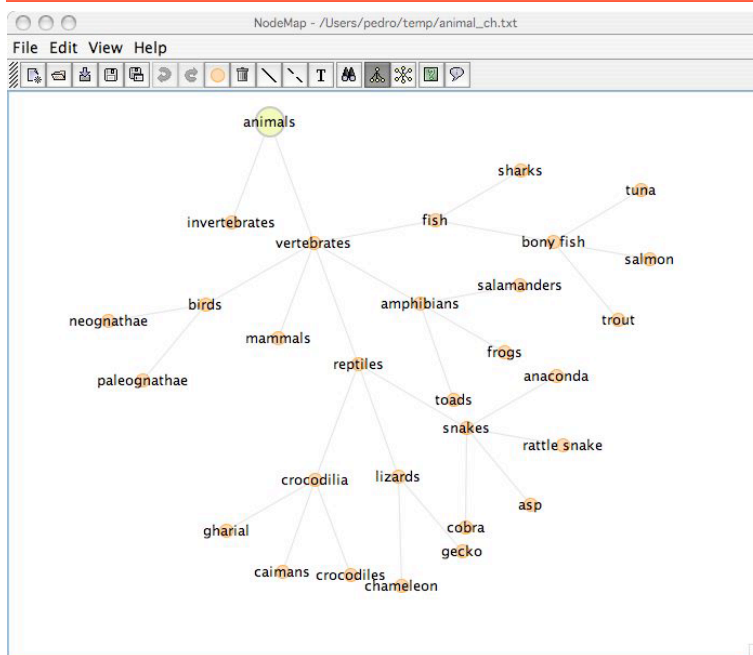
Unlike the classification problem here we do not know the labels or data categories, for this reason this is also called unsupervised learning.

Evaluation of I.R. systems

$$\text{Recall} = \frac{\text{No. of relevant docs retrieved}}{\text{Total relevant in the collection}}$$

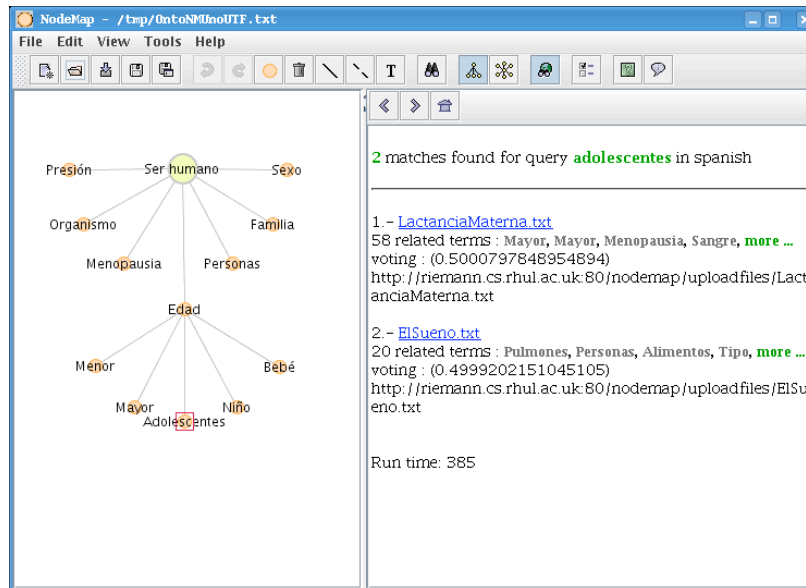
$$\text{Precision} = \frac{\text{No. of relevant docs retrieved}}{\text{Total retrieved in the collection}}$$

Demo I: text classification



Here we have an ontology and we would like to assign or classify documents on a given category

Demo 1: text classification



<http://thames.cs.rhul.ac.uk/wstalk/prototype.html>

Department of Computer Science, Royal Holloway, University of London

31

Demo 2: email search

The screenshot shows a web-based email search interface. On the left, a box contains search criteria:

- file: <keyword>
- from: <name>
- to: <email to>
- cc: <email cc>
- subject: <subject>
- <keyword>

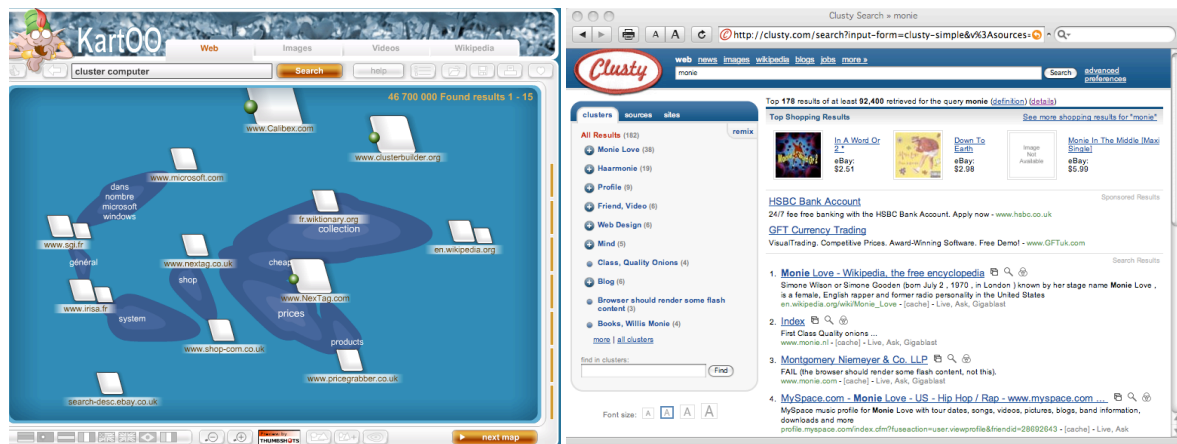
An arrow points from the 'from: <name>' field to a search bar that contains 'from: Bruce Smith'. To the right of the search bar is a 'Search' button. Above the search bar is a logo of a stylized orange tiger and the word 'Search' in orange. Below the search bar, a list of search results is displayed:

1. - <http://thames.cs.rhul.ac.uk/email/data/lists-046-11807292>
docno lists-046-11807292 nameFrom Bruce Smith emailFrom bruce@wolfham.com dateSent Sat, 20 Apr 1996 12:28:45 -0700
InReplyTo id 199604201228.6368@uuea.wolfham.com
subject Templates can work with precedences
2. - <http://thames.cs.rhul.ac.uk/email/data/lists-046-11812317>
docno lists-046-11812317 nameFrom Bruce Smith emailFrom bruce@wolfham.com dateSent Sat, 20 Apr 1996 12:36:32 -0700
InReplyTo id 199604201236.6403@uuea.wolfham.com
subject Re: Progress on Parsing
3. - <http://thames.cs.rhul.ac.uk/email/data/lists-046-11816340>
docno lists-046-11816340 nameFrom Bruce Smith emailFrom bruce@wolfham.com dateSent Sat, 20 Apr 1996 12:53:07 -0700
InReplyTo id 199604201253.6443@uuea.wolfham.com
subject Parsing of n-ary operators
4. - <http://thames.cs.rhul.ac.uk/email/data/lists-046-11820182>
docno lists-046-11820182 nameFrom Bruce Smith emailFrom bruce@wolfham.com dateSent Sat, 20 Apr 1996 13:01:35 -0700
InReplyTo id 199604201301.6465@uuea.wolfham.com
subject Re: scripted operators

Department of Computer Science, Royal Holloway, University of London

32

Search engines



Commercial

Google

YAHOO!

Ask

alltheweb

HOTBOT

WebCrawler

Live Search

altavista

.....

Open source

Lucene

Terrier



Xapian

Estraier

and so on...

Department of Computer Science, Royal Holloway, University of London

33

The web needs a lot of servers



Department of Computer Science, Royal Holloway, University of London

34

Enterprise search



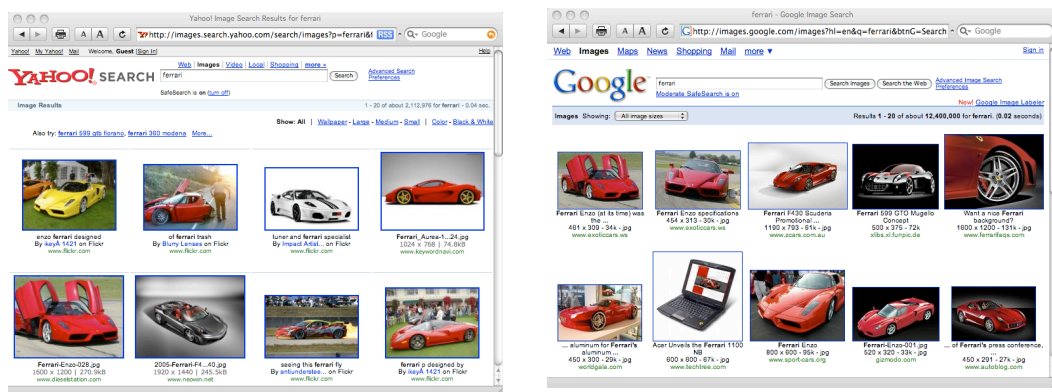
Not just docs. Multimedia I.R.

Content is not just text-based, there is more than text, i.e.

- Images
- Video
- Audio

But there again there isn't a perfect technique... let's see some examples

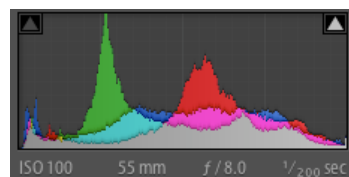
Text based image retrieval



e.g. Using information from the web site to annotate images

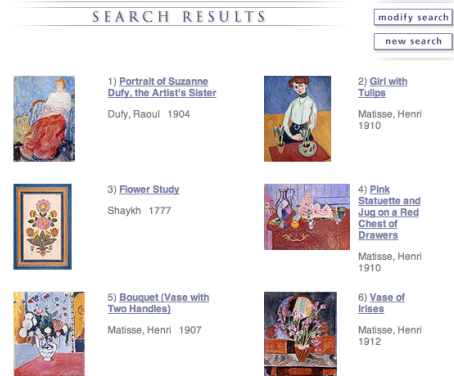
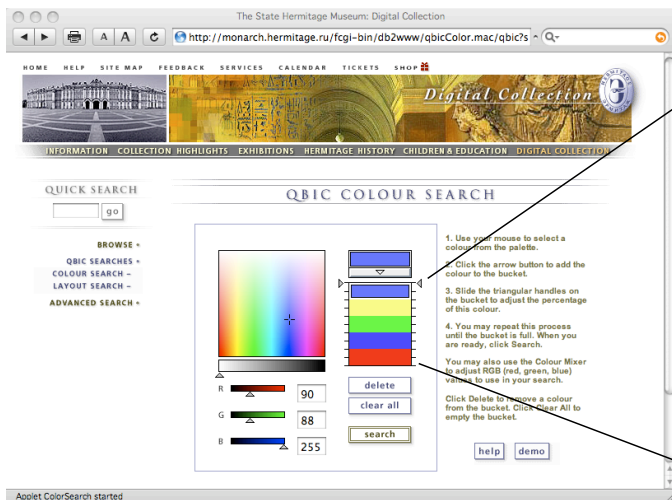
- Neighbouring text
- Same paragraph text
- Title
- Heading
- etc.

Colour histogram based image retrieval



Take an image, convert to histogram and used to query image index

Colour based image retrieval



<http://www.hermitagemuseum.org/>

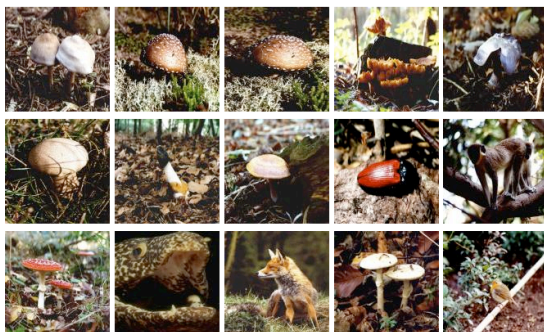
Features based image retrieval

Query Image



Weights: Perceptual Grouping = 0.1, Color = 0.3, Texture = 0.6, L, A, B channels.

Retrieved Images

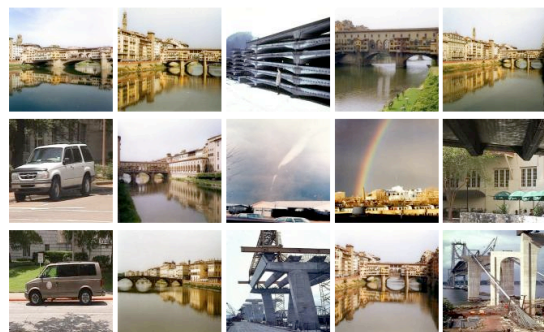


Query Image



Weights: Perceptual Grouping = 0.35, Color = 0.2, Texture = 0.45, L, A, B channels.

Retrieved Images



<http://amazon.ece.utexas.edu/~qasim/cires.htm>

Questions



References

- [1]** Michael Berry and Murray Browne, Understanding Search Engines. Mathematical Modeling and Text Retrieval. SIAM. 2nd Ed. 2005.
- [2]** Michael Berry, Zlatko Drma and Elizabeth Jessup. Matrices, Vector Spaces, and Information Retrieval. SIAM review. Vol. 41, No 2, pp 335-362. 1999.
- [3]** Ian Witten, Alistair Moffat and Timothy Bell. Managing Gigabytes. Compressing and Indexing Documents and Images. 2nd. Ed. M. Kaufman. 1999.
- [4]** Christopher Bishop. Pattern Recognition and Machine Learning. Springer, 2007.
- [5]** Fionn Murtagh, Geoff Downs and Pedro Contreras. Hierarchical Clustering of Massive, High Dimensional Data Sets by Exploiting Ultrametric Embedding. SIAM Journal on Scientific Computing, in press.
- [6]** Ella Bingham and Heikki Mannila. Random projection in dimensionality reduction: Applications to image and text data. Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-01). 2001.
- [7]** Survey of Text Mining, Clustering, Classification, and Retrieval. Michael Berry, editor. Springer. 2004
- [8]** Sholom Weiss, Nitin Indurkha, Tong Zhang, and Fred Damerau. Text Mining, Predictive Methods for Analysing unstructured Information. Springer. 2005.
- [9]** David Grossman and Ophir Frieder. Information Retrieval, Algorithms and Heuristics. 2nd. Ed. Springer. 2005.
- [10]** Ronen Feldman and James Sanger. The Text Mining Handbook, Advances Approaches in Analyzing Unstructured Data. Cambridge press. 2007.
- [11]** Clustering and Information Retrieval. Weili Wu, Hui Xiong, and Shashi Shekhar Ed. Kluwer. 2004.
- [12]** Boris Mirkin. Clustering for Data Mining, A Data Recovery Approach. Chapman and Hall. 2005.
- [13]** Keith van Rijsbergen. The Geometry of Information Retrieval. Cambridge press. 2004.