#### Needles in Haystacks



Past, Present, and Future in Information Retrieval (and a little visualisation)

> Keith van Rijsbergen Zinal, January, 2012

"Overall, the impression must be of how comparatively little the nonnegligible amount of work done has told us about the real nature of retrieval systems"

KSJ, 1981.

- Ats
- Caveats
- Why we have survived.
- Where we were, where we are, where we are going.
- Challenging the status quo.

#### THEORY PRACTICE Has a lot to do with our methodology: Cranfield And Cross disciplinarity: CS/IS And Universality of methods And One search fits all

Why are we still here?

#### **Cranfield Paradigm**

- Document collection
- Relevance judgements in advance
- Run strategy A and B
- Evaluate A and B in terms of P & R
- Compare A with B statistically
- State whether  $A \sim B$ , A > B, B < A



Basic research in IR is concerned with the design of better IR systems.  $_{CH, 2012}$ 

#### What is Information Retrieval?

- General definition
  - Retrieval of unstructured data
  - Most often it is
    - Retrieval of text documents
      - Searching newspaper articles
      - Searching on the Web
  - Other types of retrieval
    - Image retrieval
    - Video retrieval
    - Music retrieval ....

#### **Definitions of Information Retrieval**

(Salton, 1968) – Information retrieval is a field concerned with the structure, analysis, organization, storage, searching, and retrieval of information.

(Needham, 1977).....the complexity arises from the Impossibility of describing the content of a document, Or the intent of request, precisely, or unambiguously

#### What is it, cont...

- Reference Retrieval?
- What is a document?
- What is a query?
- What is relevance?

#### The role of Information



**Pre-history** THEORY (1930) (see Buckland) Goldberg Bush (1945)(195 Bagley MIT Fairthorne 1945-52) RAE 1952` Mooers (1958) Luhn Maron, Kuhns, Swanson (1959+)



#### Time II

1978 1st SIGIR 1979 1<sup>st</sup> BCSIRSG **1980** 1<sup>st</sup> Joint BCS/ACM Symposium (SIGIR 3) **1981** KSJ book on IR Experiments 1982 Belkin et al ASK hypothesis 1983 - Okapi started 1985 RIAO-1 1986 CvR logic model 1990 Deerwester et al, LSI paper 1991 CoLIS 1 (in Tampere!) 1991 – Inquiry started 1992 Ingwersen's book 1992 TREC-1 1997 Kluwer series on IR 1998 Croft Ponte paper on language models 2012

# CIKM Quantum Interaction

Time III

#### **Experimental Methodology**

Cleverdor Lancaster Keen Saracevic Salton **Sparck Jones** Blair & Maron Harman/Voorhees

MENT THEORY Cranfield Medlars Cranfield/Smart **CWRU** Smart Ideal Test Collection **Stairs** TREC

EvaluationEvaluationEvent</t

Significance testing? Sampling distributions? Incomplete judgments? Multi-valued relevance?

Landmarks MENT THEORY Luhn's tf weighting Statistical weighting tf\* Architecture **Relevance Feedback** Stemming Poisson Model ->/BM25/DFR Various models (LUP -> LM)



In most cases of testimony the assertions of all witnesses do not agree. Some give evidence of one kind, some of the opposite, so that the evidence upon the same point is contradictory. In these, cases, the laws of induction and deduction are applied by the jury, to judge of the value of the testimony, and that which affords most probability, or that which most coincides with former knowledge is received.

Alfred Smee, 1851

Acknowledgment

The following material on 'Alias Smith and Jones' was taken from R. Jeffrey, Erkenntnis, 391-399, 1987. Also reprinted in his book 'Probability and the Art of Judgment' Alias Smith and Jones H: it is relevant, it is about..... E: Smith *says* it is F: Jones *says* it is

Their testimony can be clear: E, F are certain
Their testimony can be contradictory eg. E and -F
Their testimony may be uncertain

Smith and Jones contradict each other Cancellation?

Bayes' Theorem MENT THEORY you can derive: final odds = likelihood ratio \* prior odds  $\frac{P(H|E\neg F)}{P(\neg H|E\neg F)} = \frac{P(E\neg F|H)}{P(E\neg F|\neg H)} \times \frac{P(H)}{P(\neg H)}$ 

Cancellation of testimony: Likelihood ratio = 1

© CvR

### independence of assessors MENT THEORY P(EF|H) = P(E|H)P(F|H) $P(EF|\neg H) = P(E|\neg H)P(F|\neg H)$ P(EF) = P(E)P(F)?

Conflicting clear testimony of equally reliable independent assessors P(E|H) = P(F|H) = r $P(E|\neg H) = P(F|\neg H) = S$  $P(EF|H) = r^2, P(EF|\neg H) = s^2$ 

#### Cancellation

r = s means Smith's assessment nor Jones' changes prior odds. More interesting case:

## $P(E \neg F|H) = P(E \neg F|\neg H) \text{ iff}$ (r = s) v (r + s = 1)

Cancellation when :  $P(E|H) + P(E|\neg H) = 1$ 

# Reliability context dependent

Smith and Jones are more reliable when the document is relevant than when it is not.

$$P(E|H) = 0.9 \text{ but } P(\neg E|\neg H) = 0.8$$
$$\Rightarrow P(E|\neg H) = 0.2$$
$$\therefore P(E|H) + P(E|\neg H) = 1.1$$



### Maron's theory of indexing

....in the case where the query consists of single term, call it B, the probability that a given document will be judged relevant by a patron submitting B is simply the ratio of the number of patrons who submit B as their query and judge that document as relevant, to the number of patrons, who submit B as their search query





Probabilistic View



geometric view



#### The logical view

To evaluate a conditional, first hypothetically make the minimal revision of your stock of beliefs required to assume the antecedent. The evaluate the acceptability of the consequent on the bais of this revised body of beliefs.

#### LUP in Hilbert space

PRACTICE

Given any two propositions F and E; a measure of the uncertainty of  $E \rightarrow F$  relative to a given data set is determined by the minimal extent to which we have to add information to the data set, to establish the truth of  $E \rightarrow F$ .


# **Classical Probability** $P(E_0) = 0$ and $P(E_1) = 1$ $P(E_i \cup E_i) = P(E_i) + P(E_i)$ provided that $E_i \cap E_j = E_0$ **Quantum Probability** $\mu_{\varphi}(\Phi) = 0$ $\mu_{\omega}(\mathbf{H}) = 1$ For subspaces $L_i$ and $L_i$ , $\mu_{\omega}(L_i \oplus L_i) = \mu_{\omega}(L_i) + \mu_{\omega}(L_i)$ provided $L_i \cap L_i = \Phi$

Probability

Total Probability  $P(A) = P(A|B)P(B) + P(A|B^c)P(B^c)_{CH, 2}$  (interference)





# Cluster Hypothesis I

'Cluster-based retrieval has as its foundation a hypothesis, the *cluster hypothesis*, which states that closely associated documents tend to be relevant to the same requests' CvR, 1971



 $\{d(\mathbf{x}, \mathbf{y}) \leq \max\left[d(\mathbf{x}, \mathbf{z}), d(\mathbf{z}, \mathbf{y})\right]\}$ 

# But (Tversky says)

- Symmetry: 'Turks fight like tigers', not 'Tigers fight like Turks'
- 2. ∆ inequality: Jamaica is similar to Cuba, Cuba is similar to Russia, but Jamaica and Russia not similar
- Maybe what is needed is conditional probability  $P(A|B) \neq P(B|A)$  or, logic  $B \rightarrow A \neq A \rightarrow B$

#### Cluster hypothesis II



### Static Clustering

- 1. dependence on rank-ordering of dissimilarity
- 2. insensitive to small errors in DC
- 3. preservation of well marked clusters
- 4. stable under growth
- 5. labelling independence
- 6. invariance of ultrametric
- 7. subject to 3 minimises distortion





#### Spanning tree?



Hilbert-Schmidt: (A,B) = trace(A'B)

#### Pseudo Relevance Feedback

#### How does it differ from Relevance Feedback?

Also, think of Context



















 $\bigcirc$ 













Dependence e.g Unified Probabilistic Model Co-relevance Stochastic Processes Brouwerian Logics Error Analysis

**Buried Treasure** 

MENT

C.T Yu Maron/Cooper/SER Ivie Mandelbrot/Herdan Hillman Hughes/Cover/Duda

THEORY

#### Hypotheses/Principles

THEORY

Items may be associated without apparent meaning but exploiting their association may help retrieval

P & R trade-off – ABNO/OBNA Exhaustivity/Specificity Cluster Hypothesis Association Hypothesis Probability Ranking Principle Logical Uncertainty Principle ASK Polyrepresentation

#### Laws of Retrieval?

Inverse relationship of Precision/Recall
Perfect retrieval is impossible
....

CvR, 1979

THEORY

#### What cannot be done!

perpetual motion machine no free lunch fooling all of the people all of the time Godel's theorem

Failure Analysis Examine why B << A, instead of, why of why A >>B.

### Postulates of Impotence

(according to Swanson, 1988)

324

- An information need cannot be expressed independent of context
- It is impossible to instruct a machine to translate a request into adequate search terms
- A document's relevance depends on other seen documents
- It is never possible to verify whether all relevant documents have been found
- Machines cannot recognise meaning -> can't beat human indexing etc

# • Word-occurrence statistics can neither represent meaning nor substitute for it

.more postulates

THEORY

- The ability of an IR system to support an iterative process cannot be evaluated in terms of single-iteration human relevance judgment
- Thus, consistently effective fully automatic indexing and retrieval is not possible

# Future Theory I THEORY Relevance: R + R = NR, or NR + NR = RContext: R + Context = NR, NR + Context = RContextual Cluster Hypothesis

#### ССН



**Future Theory II** MENT THEORY P(X|Y) - sampling space  $P_{\rm V}({\rm X})$  - conditionalisation  $P(Y \rightarrow X)$  - probable inference

#### Areas of Research

- •How does the brain do it
- •How do we see to retrieve?-
- •How do we map IR onto Quantum Computation? (QM)
- •How do we reduce dimensionality in dynamic fashion? (Statistics)
- •What is a good logic for IR?
- •What is a good theory of uncertainty?
- •How do we model context?
- •How do we formally capture interaction?
- •How do we capture implicit/tacit/information?
- •Is there a theory of information for IR?

on? (Statistics) (mathematical logic) (frequency/geometry) (HCI)

(neuroscience)

(computer vision)


## Readings I

- Search Engines: Information Retrieval in Practice, W. Bruce Croft, et al, Addison Wesley, 2010.
- Introduction to Information Retrieval,
  C.D. Manning et al, Cambridge University Press, 2008.

" THEOBY

- 3. *Finding out About*,R.K. Belew, Cambridge University Press, 2000.
- 4. *Readings in Information Retrieval*,K. Sparck Jones and P.Willett, Morgan Kauffmann, 1997.
- 5. The Turn, P. Ingwersen and K. Jarvelin, Springer, 2005.



- 6. Information Retrieval: Implementing and Evaluating Search Engines S. Buttcher, C.L.A Clarke, and G.V. Cormack, MIT, 2010
- 7. Information Retrieval, C.J. van Rijsbergen, Butterworths, 1979.
- 8, *Modern Information Retrieval*, R. Baeza-Yates and B. Ribero-Neto Addison Wesley, 2010

http://www.cse.ucsd.edu/~rik/foa/

http://www-csli.stanford.edu/~hinrich/information-retrievalbook.html

http://www.search-engines-book.com/

http://searchuserinterfaces.com/book/

http://en.wikipedia.org/wiki/Information\_retrieval

http://people.ischool.berkeley.edu/~hearst/irbook/

