#### **Evaluation**

#### INFM 718X/LBSC 718X Session 6 Douglas W. Oard

## **Evaluation Criteria**

• Effectiveness

- System-only, human+system

- Efficiency
  - Retrieval time, indexing time, index size
- Usability

- Learnability, novice use, expert use

## **IR Effectiveness Evaluation**

#### User-centered strategy

- Given several users, and at least 2 retrieval systems
- Have each user try the same task on both systems
- Measure which system works the "best"
- System-centered strategy
  - Given documents, queries, and relevance judgments
  - Try several variations on the retrieval system
  - Measure which ranks more good docs near the top

# Good Measures of Effectiveness

- Capture some aspect of what the user wants
- Have predictive value for other situations

   Different queries, different document collection
- Easily replicated by other researchers
- Easily compared
  - Optimally, expressed as a single number

## **Comparing Alternative Approaches**

- Achieve a <u>meaningful</u> improvement – An application-specific judgment call
- Achieve reliable improvement in unseen cases
   Can be verified using statistical tests

## **Evolution of Evaluation**

- Evaluation by **inspection** of examples
- Evaluation by demonstration
- Evaluation by **improvised** demonstration
- Evaluation on data using a figure of merit
- Evaluation on test data
- Evaluation on common test data
- Evaluation on common, unseen test data

## **Automatic Evaluation Model**



These are the four things we need!

# **IR Test Collection Design**

- Representative document collection
   Size, sources, genre, topics, …
- "Random" sample of representative queries

   Built somehow from "formalized" topic statements
- Known binary relevance
  - For each topic-document pair (topic, not query!)
  - Assessed by humans, used only for evaluation
- Measure of effectiveness
  - Used to compare alternate systems

# Defining "Relevance"

- Relevance relates a <u>topic</u> and a document
   Duplicates are equally relevant by definition
  - Constant over time and across users
- **Pertinence** relates a <u>task</u> and a document Accounts for quality, complexity, language, ...
- Utility relates a <u>user</u> and a document
   Accounts for prior knowledge

Space of all documents



# Set-Based Effectiveness Measures

- Precision
  - How much of what was found is relevant?
    - Often of interest, particularly for interactive searching
- Recall
  - How much of what is relevant was found?
    - Particularly important for law, patents, and medicine
- Fallout
  - How much of what was irrelevant was rejected?
    - Useful when different size collections are compared

### **Effectiveness Measures**





# Balanced F Measure $(F_1)$

Harmonic mean of recall and precision

$$F_1 = \frac{1}{\frac{0.5}{P} + \frac{0.5}{R}}$$

# Variation in Automatic Measures

System

– What we seek to measure

- Topic
  - Sample topic space, compute expected value
- Topic+System
  - Pair by topic and compute statistical significance
- Collection
  - Repeat the experiment using several collections

### IIT CDIP v1.0 Collection

#### **Scanned**

* 3			
÷.	Philip Morris	U.S.A.	Inter-Office Correspondence
	Renefits Denartme	nt	Richmond, Virginia
	метеры ме		· · · · · · · · · · · · · · · · · · ·
	Te: Distributio		Date: May 30, 1997
	From: Lisa Halle		2 Start Contraction
	Subjects CIGNA V	Vell-Being News	etter - Fature Strategy
	During our last CIGNA J articles and discontinue s matter of discussion. I ha findings and preliminary I believe everyone's input whether you concur with	Action Plan meetin ending CIGNA W ave done some res recommendation it is valuable, and my recommendat	g, the insue of whether to stop previewing effil-Being gavelutent to our sampleyneen was a earch, and wanted to present you with my for FM's statingy regarding future newslotters, would appreciate hearing from each of you on loc.
	Background Informatio	in ,	
	CIGNA. Well-Being news using is to have one of the particular newsletter, and local articles can be repla- replace or modify a nation Since 1996, we have optor	sletters are sent or c analysis proview i then recommend used with another nal article, it costs of to either skip or	15 quarterity basis. The process we have been booth autional and local articles altend for a to either skip or send that issue. Offersive of similar length at no cost to PM. If we opt to PM 33,000 pre issue. send issues as follows:
	Date of Newsletter	Decision	Comments
	Spring 1996	Send	Deleted advertisement for CIONA Time- Life videos featuring Ex-ergeoen General Koop prior to sending.
<b>81/1</b> /144 3:	Spring 1996	Send : Norskoon 673 g Deskies	Datest advertisement for CIONA Time- Life video fasting Ex-surgeo General Koop prior to sending.
81/1.jees 3:	Spring 1996 <sup>1</sup> '00 OSECOUT <u>STRUCTURE</u> Date of Newsletter Pail 1996	Sand HEX-SHOOM 573 g Destrices Send	Dolest alvertisenen: for CIONA Time- Life videos featuring Ex-surgeon General Koop prior to sending.
01/1100 3	Spring 1996 1 '60' Ostron <u>en utternen</u> ove <u>Pase of Newsletter</u> Pall 1996	Send NORKON 672 g Zesteles Send	Comments Deleted advectionment: for CIONA Time- Life videos featuring Ex-express General Keep prior to sending. Comments A national article on locer structure which was demond to shore structure which was demond to sover that, what the ensurement of the structure of the structure which was demond to sover that, what the ensurement of the structure of the structure which was demond to sover that, what the ensurement of the structure of the structure of the structure of the video was demond to sover that, what the ensurement of the structure of the video was demond to sover that, what the ensurement of the structure of the video of the structure of the structure of the structure of the video of the structure of the structure of the structure of the video of the structure of the structure of the structure of the video of the structure of the structure of the structure of the video of the structure of the structure of the structure of the video of the structure of the structure of the structure of the structure of the video of the structure of the structure of the structure of the video of the structure of the structure of the structure of the video of the structure of the video of the structure of the structure of the str
01/11994 3	Spring 1996 1 - Gal Galerange <u>atternation</u> <b>Date of Newslatter</b> Patt 1996 Winner 1996	Send Poisson sra g Testiles Send	Deleted adversionment for CIONA Time- Life videos featuring Ex-express General Keep prior to sending. Command A national entities on baset statistic which was denaid on wear that what the wanning laid entits. Also, a baset statistic which was denaid on wears that what the wanning laid entits of the statistic of the wanning laid entits, which a baset which was denaid on wears that what the wanning laid entits, which a baset which was denaid on wears that what the wanning laid entits, which a baset which was denaid on wears that what the wanning laid entits, which a baset which was denaid on wears that what the wanning laid entits, which are the prophysical water was a baset concert, wear language and water. See the marketing entits water and water baset of the statistic statistics. Nearend a statistic water water and water and water and and a statistic statistics. Nearend a statistic water water and water and water and a statistic statistics.
01/11994 3	Spring 1996	Send Horstoon 673 g Datales Send Skip Send	Deleted adversionment: for CIONA Time- Life videos featuring Ex-express General Keep prior to sending. Command A antibud enfoid as lower studied and the sending of the sending video was denoid to seven that what the wanning label rates. Also, a these constant at the sendence of the sending video was denoid to seven that what the wanning label rates. Also, a these constant at the sendence of the sendence of the program, studentsford, program, studentsford, program, studentsford, as arentifier coaland as a sendence on the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the sendence of the

The proceed reviewing sticle and making a recommunicion to end or with a lates were, advancing on context. Typically, it is immediately clear if ecounting is objectionable. Other times, it may require disrussion with others and samagement. I would any breverse prince stall with COMA, perviewing of organization versions and full versions of cloud and mational articles, discussion if forecastry, subliding route of anyonic the starting organization. These works all full and the starting anyonic drive the starting organization. These provides the starting of anyonic drive the starting or principal from the discussion in starting anyonic drive period to most the starting any member there may discussion in starting mostings required the anyonic drive the starting.

Recommendation

#### <u>OCR</u>

Philip Moxx's. U.S.A. x.dr~am~c.

cvrrespoaa.aa

- Benffrts Departmext Rieh>pwna, Yfe&ia Ta: Dishlbutfon Data aday 90,1997.
- From: Lisa Fislla

Sabj.csr CIGNA WeWedng Newsbttsr - Yntsre StratsU

During our last CIGNA Aatfoa Plan meadng, tlu iasuo of wLetSae to i0op per'Irw+ng

artieles aod discontinue mndia6 CIGNA Well-Being aawslener to om employees was a

msiter of disanision . I Imvm done somme reaearc>>, and wanted to pruedt you with my

Sadings and pcdiminary recwmmeadatioa for PM's atratezy leprding l4aas aewelattee\*.

I believe .vayone'a input is valuable, and would epproolate hoaring fmaa aaeh of you on

whetlne you concur with my reeommendatioa

#### <u>Metadata</u>

**Title:** *CIGNA WELL-BEING NEWSLETTER - FUTURE STRATEGY* 

**Organization Authors:** *PMUSA, PHILIP MORRIS USA* 

Person Authors: HALLE, L

**Document Date:** *19970530* 

**Document Type:** *MEMO, MEMORANDUM* 

Bates Number: 2078039376/9377

Page Count: 2

**Collection:** Philip Morris

...

#### "Complaint" and "Production Request"

...12. On January 1, 2002, Echinoderm announced record results for the prior year, primarily attributed to strong demand growth in overseas markets, particularly China, for its products. The announcement also touted the fact that Echinoderm was unique among U.S. tobacco companies in that it had seen no decline in domestic sales during the prior three years.

13. Unbeknownst to shareholders at the time of the January 1, 2002 announcement, defendants had failed to disclose the following facts which they knew at the time, or should have known:

a. The Company's success in overseas markets resulted in large part from bribes paid to foreign government officials to gain access to their respective markets;

b. The Company knew that this conduct was in violation of the Foreign Corrupt Practices Act and therefore was likely to result in enormous fines and penalties;

c. The Company intentionally misrepresented that its success in overseas markets was due to superior marketing.

d. Domestic demand for the Company's products was dependent on pervasive and ubiquitous advertising, including outdoor, transit, point of sale and counter top displays of the Company's products, in key markets. Such advertising violated the marketing and advertising restrictions to which the Company was subject as a party to the Attorneys General Master Settlement Agreement ("MSA").

e. The Company knew that it could be ordered at any time to cease and desist from advertising practices that were not in compliance with the MSA and that the inability to continue such practices would likely have a material impact on domestic demand for its products. ...

All documents which describe, refer to, report on, or mention any "in-store," "on-counter," "point of sale," or other retail marketing campaigns for cigarettes.

# An Ad Hoc "Production Request"

<ProductionRequest>

<RequestNumber>148</RequestNumber>

<RequestText>All documents concerning the Company's FMLA policies, practices and procedures.</RequestText>

<BooleanQuery>

<FinalQuery>(policy OR policies OR practice! or procedure! OR rule! OR
guideline! OR standard! OR handbook! OR manual!) w/50 (FMLA OR leave OR
"Family medical leave" OR absence)</FinalQuery>

<NegotiationHistory>

<ProposalByDefendant>(FMLA OR "federal medical leave act") AND (policies OR
practices OR procedures)</proposalByDefendant>

<RejoinderByPlaintiff>(FMLA OR "federal medical leave act") AND (leave w/10 polic!)</RejoinderByPlaintiff>

<Consensus1>(policy OR policies OR practice! or procedure! OR rule! OR
guideline! OR standard! OR handbook! OR manual!) AND (FMLA OR leave OR
"Family medical leave" OR absence)</Consensus1>

</NegotiationHistory>

</BooleanQuery>

<FinalB>40863</FinalB>

<RequestSource>2008-H-7</RequestSource>

### **Estimating Retrieval Effectiveness**



#### **Relevance** Assessment

• All volunteers

- Mostly from law schools

- Web-based assessment system
   Based on document images
- 500-1,000 documents per assessor
   Sampling rate varies with (minimum) depth

#### 2008 Est. Relevant Documents



Mean estRel = 82,403 (26 topics) • 5x 2007 mean estRel (16 904)

5x 2007 mean estRel (16,904)

Max estRel=658,339, Topic 131 (rejection of trade goods)

Min estRel=110 Topic 137 (intellectual property rights)

#### 2008 (cons.) Boolean Estimated Recall



Mean estR=0.33 (26 topics)

 Missed 67% of relevant documents (on average)

Max estR =0.99, Topic 127 (sanitation procedures)

Min estR=0.00, Topic 142 (contingent sales)

#### 2008 AestR@B: wat7fuse vs. Boolean



#### **Evaluation Design**



#### Interactive Task: Key Steps



#### **Interactive Task: Participation**

#### **2008**

4 Participating Teams (2 commercial, 2 academic)

□ 3 Topics (and 3 TAs)

Test Collection: MSA Tobacco Collection

#### 2009

11 Participating Teams (8 commercial, 3 academic)

7 Topics (and 7 TAs)

Test Collection: Enron Collection

#### 2010

□ 12 Participating Teams (6 commercial, 5 academic, 1 govt)

- □ 4 Topics (and 4 TAs)
- □ Test Collection: Enron Collection (new EDRM version)

UB	CI	H5	Pitt	AdHoc	Ν	n	а	r
R	R	R	R	R	5,727	46	46	38
R	R	R	R	N	24	5	5	4
R	R	R	N	R	11,965	98	98	78
R	R	R	N	N	995	9	9	9
R	R	N	R	R	131	5	5	3
R	R	N	R	N	0	0	0	0
R	R	N	N	R	1,547	13	13	2
R	R	N	N	N	220	5	5	2
R	N	R	R	R	1,901	15	15	11
R	N	R	R	N	46	5	5	2
R	N	R	N	R	17,082	145	145	111
R	N	R	N	N	10,291	84	84	61
R	N	N	R	R	176	5	5	1
R	N	N	R	N	19	5	5	2
R	N	N	N	R	7,679	62	61	23
R	N	N	N	N	9,531	77	77	17
Ν	R	R	R	R	8,068	65	65	49
Ν	R	R	R	N	101	5	5	2
Ν	R	R	N	R	73,280	541	540	393
Ν	R	R	N	N	28,409	235	235	146
Ν	R	N	R	R	1,185	10	10	4
Ν	R	N	R	N	37	5	4	3
Ν	R	N	N	R	23,688	193	193	84
Ν	R	N	N	N	20,078	171	164	57
Ν	N	R	R	R	5,321	43	43	33
Ν	N	R	R	N	371	5	5	2
Ν	N	R	N	R	151,787	800	795	552
Ν	N	R	N	N	293,439	1,100	1,095	621
Ν	N	N	R	R	2,253	18	18	6
Ν	Ν	N	R	N	456	5	5	2
Ν	Ν	N	Ν	R	526,099	1,100	1,087	234
Ν	Ν	N	Ν	N	5,708,286	1,625	1,579	111
TOTAL					6,910,192	6,500	6,421	2,663

### 2008 Interactive Topics

Торіс	Samples	Est Nrel Pre- adjudication	Est Nrel: Post- Adjudication	Relevance Density
102	4.500		562,402 ±73,000	8.1%
103	6,500	914,258 ±72,000	786,862 ±54,000	11.4%
104	2,500		45,614 ±25,000	0.7%

Topic 103

#### **Pre-Adjudication Results**



#### Topic 103

#### **Post-Adjudication Results**



#### Topic 103 Results on Good OCR



High OCR-accuracy documents only

Interactive Task - 2009

#### **TREC Enron Email Test Collection Version 1**

#### Enron Collection

- A collection of emails produced by Enron in response to requests from the Federal Energy Regulatory Commission (FERC)
- First year used in the Legal Track

#### Size of Collection (post-deduplication)

- 569,034 messages
- 847,791 documents
- Over 3.8 million pages

#### Distribution Format

- Extracted Text (in EDRM XML interchange format)
- Native .msg files

#### 2009 Results (pre-adjudication)



#### 2009 Results (post-adjudication)





#### 2009 Results (pre- to post-adj)



# EDRM Enron V2 Dataset

Email from ~150 Enron executives

1.3M records captured by FERC

Processed to several formats by ZL/EDRM

- EDRM XML (text+native) ~100GB
- PST ~100GB

Deduped, reformatted by U. Waterloo

- 455,449 messages + 230,143 attachments = 685,592 docs
- Text (1.2 GB compressed; 5.5GB uncompressed)
- Mapping from PST docs to EDRM document identifiers

Used for both Learning and Interactive tasks

• Participants submitted EDRM document identifiers

#### Topic 301 (2010)

#### Document Request

- All documents or communications that describe, discuss, refer to, report on, or relate to onshore or offshore oil and gas drilling or extraction activities, whether past, present or future, actual, anticipated, possible or potential, including, but not limited to, all business and other plans relating thereto, all anticipated revenues therefrom, and all risk calculations or risk management analyses in connection therewith.
- Topic Authority
  - Mira Edelman (Hughes Hubbard)

### 2010 Post-Adj Relevance Results



## 2010 Post-Adj Privilege Results



#### 2009 Change in F1



#### 2010 Change in F1



## **User Studies**

- Goal is to account for interface issues
  - By studying the interface component
  - By studying the complete system
- Formative evaluation
  - Provide a basis for system development
- Summative evaluation
  - Designed to assess performance

# Blair and Maron (1985)

- A classic study of retrieval effectiveness
  - Earlier studies used unrealistically small collections
- Studied an archive of documents for a lawsuit
  - 40,000 documents, ~350,000 pages of text
  - 40 different queries
  - Used IBM's STAIRS full-text system
- Approach:
  - Lawyers wanted at least 75% of all relevant documents
  - Precision and recall evaluated only after the lawyers were satisfied with the results

## Blair and Maron's Results

- Mean precision: 79%
- Mean recall: 20% (!!)
- Why recall was low?
  - Users can't anticipate terms used in relevant documents

"accident" might be referred to as "event", "incident", "situation", "problem," ...

- Differing technical terminology
- Slang, misspellings
- Other findings:
  - Searches by both lawyers had similar performance
  - Lawyer's recall was not much different from paralegal's

#### Additional Effects in User Studies • Learning

- Vary topic presentation order
- Fatigue

- Vary system presentation order

• Topic+User (Expertise)

Ask about prior knowledge of each topic

## Batch vs. User Evaluations

- Do batch (black box) and user evaluations give the same results? If not, why?
- Two different tasks:
  - Instance recall (6 topics)

What countries import Cuban sugar? What tropical storms, hurricanes, and typhoons have caused property damage or loss of life?

#### – Question answering (8 topics)

Which painting did Edvard Munch complete first, "Vampire" or "Puberty"? Is Denmark larger or smaller in population than Norway?

## Results

- Compared of two systems:
  - a baseline system
  - an improved system that was provably better in batch evaluations
- Results:

	Instance Rec	call	Question Answering		
	Batch MAP	User recall	Batch MAP	User accuracy	
Baseline	0.2753	0.3230	0.2696	66%	
Improved	0.3239	0.3728	0.3544	60%	
Change	+18%	+15%	+32%	-6%	
p-value (paired t-test)	0.24	0.27	0.06	0.41	

# **Qualitative User Studies**

- Observe user behavior
  - Instrumented software, eye trackers, etc.
  - Face and keyboard cameras
  - Think-aloud protocols
  - Interviews and focus groups
- Organize the data
  - For example, group it into overlapping categories
- Look for patterns and themes
- Develop a "grounded theory"