
Performance Evaluation of Information Retrieval Systems

Many slides in this section are adapted from Prof. Joydeep Ghosh (UT ECE) who in turn adapted them from Prof. Dik Lee (Univ. of Science and Tech, Hong Kong)

Why System Evaluation?

- There are many retrieval models/ algorithms/ systems, which one is the best?
- What is the best component for:
 - Ranking function (dot-product, cosine, ...)
 - Term selection (stopword removal, stemming...)
 - Term weighting (TF, TF-IDF,...)
- How far down the ranked list will a user need to look to find some/all relevant documents?

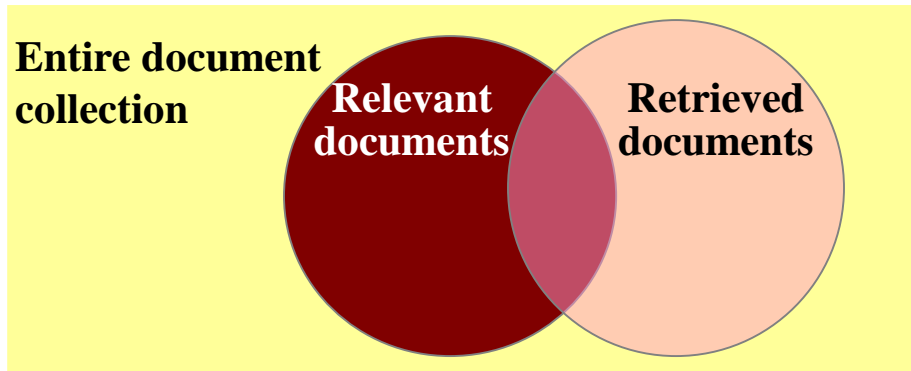
Difficulties in Evaluating IR Systems

- Effectiveness is related to the *relevancy* of retrieved items.
- Relevancy is not typically binary but continuous.
- Even if relevancy is binary, it can be a difficult judgment to make.
- Relevancy, from a human standpoint, is:
 - Subjective: Depends upon a specific user's judgment.
 - Situational: Relates to user's current needs.
 - Cognitive: Depends on human perception and behavior.
 - Dynamic: Changes over time.

Human Labeled Corpora (Gold Standard)

- Start with a corpus of documents.
- Collect a set of queries for this corpus.
- Have one or more human experts exhaustively label the relevant documents for each query.
- Typically assumes binary relevance judgments.
- Requires considerable human effort for large document/query corpora.

Precision and Recall



irrelevant	retrieved & irrelevant	Not retrieved & irrelevant
relevant	retrieved & relevant	not retrieved but relevant
	retrieved	not retrieved

$$\text{recall} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of relevant documents}}$$

$$\text{precision} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of documents retrieved}}$$

Precision and Recall

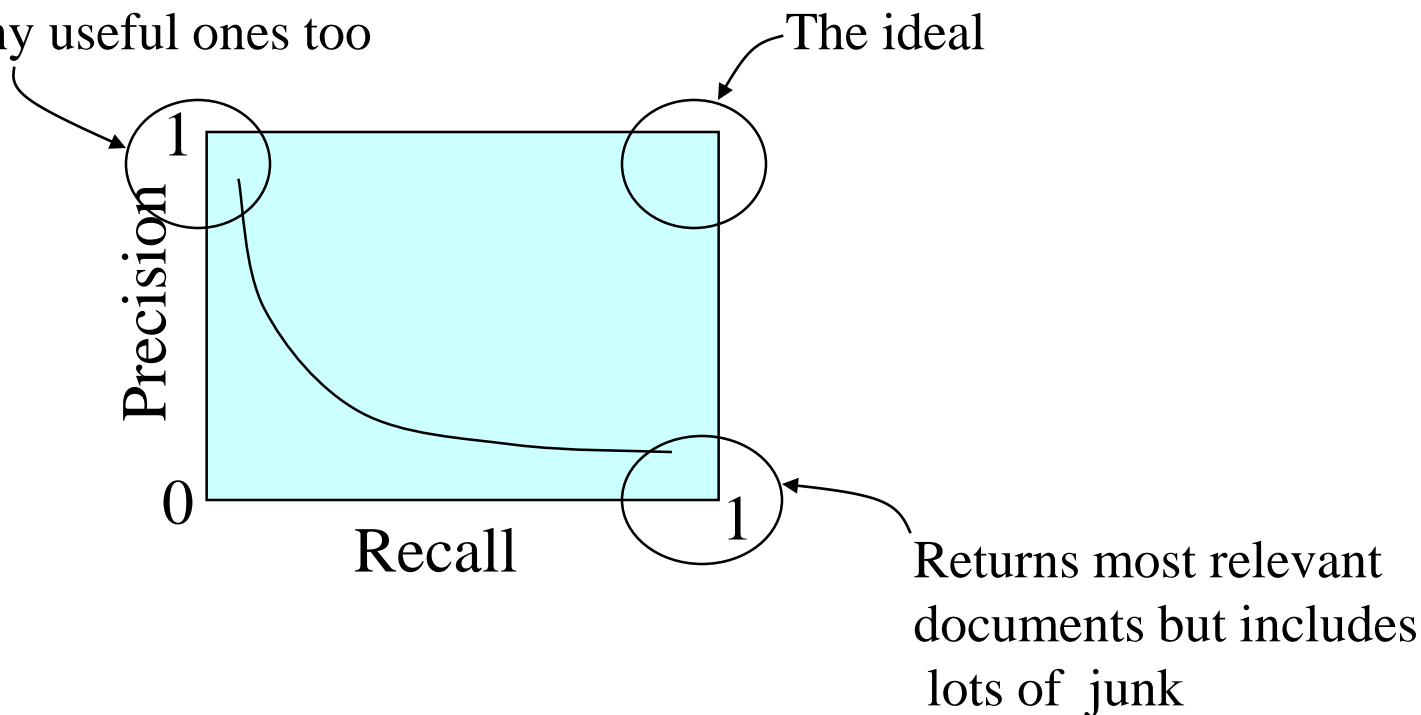
- Precision
 - The ability to retrieve top-ranked documents that are mostly relevant.
- Recall
 - The ability of the search to find *all* of the relevant items in the corpus.

Determining Recall is Difficult

- Total number of relevant items is sometimes not available:
 - Sample across the database and perform relevance judgment on these items.
 - Apply different retrieval algorithms to the same database for the same query. The aggregate of relevant items is taken as the total relevant set.

Trade-off between Recall and Precision

Returns relevant documents but misses many useful ones too



F-Measure

- One measure of performance that takes into account both recall and precision.
- Harmonic mean of recall and precision:
- $1/F = 1/2(1/P + 1/R)$; $1/F = (P+R)/2PR$;

$$F = \frac{2PR}{P+R} = \frac{2}{\frac{1}{R} + \frac{1}{P}}$$

- Compared to arithmetic mean, both need to be high for harmonic mean to be high.

E Measure (parameterized F Measure)

- A variant of F measure that allows weighting emphasis on precision over recall:

$$E = \frac{(1 + \beta^2)PR}{\beta^2 P + R} = \frac{(1 + \beta^2)}{\frac{\beta^2}{R} + \frac{1}{P}}$$

- Value of β controls trade-off:
 - $\beta = 1$: Equally weight precision and recall (E=F).
 - $\beta > 1$: Weight recall more.
 - $\beta < 1$: Weight precision more.

Computing Recall/Precision Points

- For a given query, produce the ranked list of retrievals.
- Mark each document in the ranked list that is relevant according to the gold standard.
- Compute a recall/precision pair for each position in the ranked list that **contains a relevant document**.

Computing Recall/Precision Points: Example 1

n	doc #	relevant
1	588	x
2	589	x
3	576	
4	590	x
5	986	
6	592	x
7	984	
8	988	
9	578	
10	985	
11	103	
12	591	
13	772	x
14	990	

Let total # of relevant docs = 6
Check each new recall point:

$R=1/6=0.167$; $P=1/1=1$

$R=2/6=0.333$; $P=2/2=1$

$R=3/6=0.5$; $P=3/4=0.75$

$R=4/6=0.667$; $P=4/6=0.667$

$R=5/6=0.833$; $p=5/13=0.38$

Missing one
relevant document.
Never reach
100% recall

Computing Recall/Precision Points: Example 2

n	doc #	relevant
1	588	x
2	576	
3	589	x
4	342	
5	590	x
6	717	
7	984	
8	772	x
9	321	x
10	498	
11	113	
12	628	
13	772	
14	592	x

Let total # of relevant docs = 6
Check each new recall point:

$R=1/6=0.167$; $P=1/1=1$

$R=2/6=0.333$; $P=2/3=0.667$

$R=3/6=0.5$; $P=3/5=0.6$

$R=4/6=0.667$; $P=4/8=0.5$

$R=5/6=0.833$; $P=5/9=0.556$

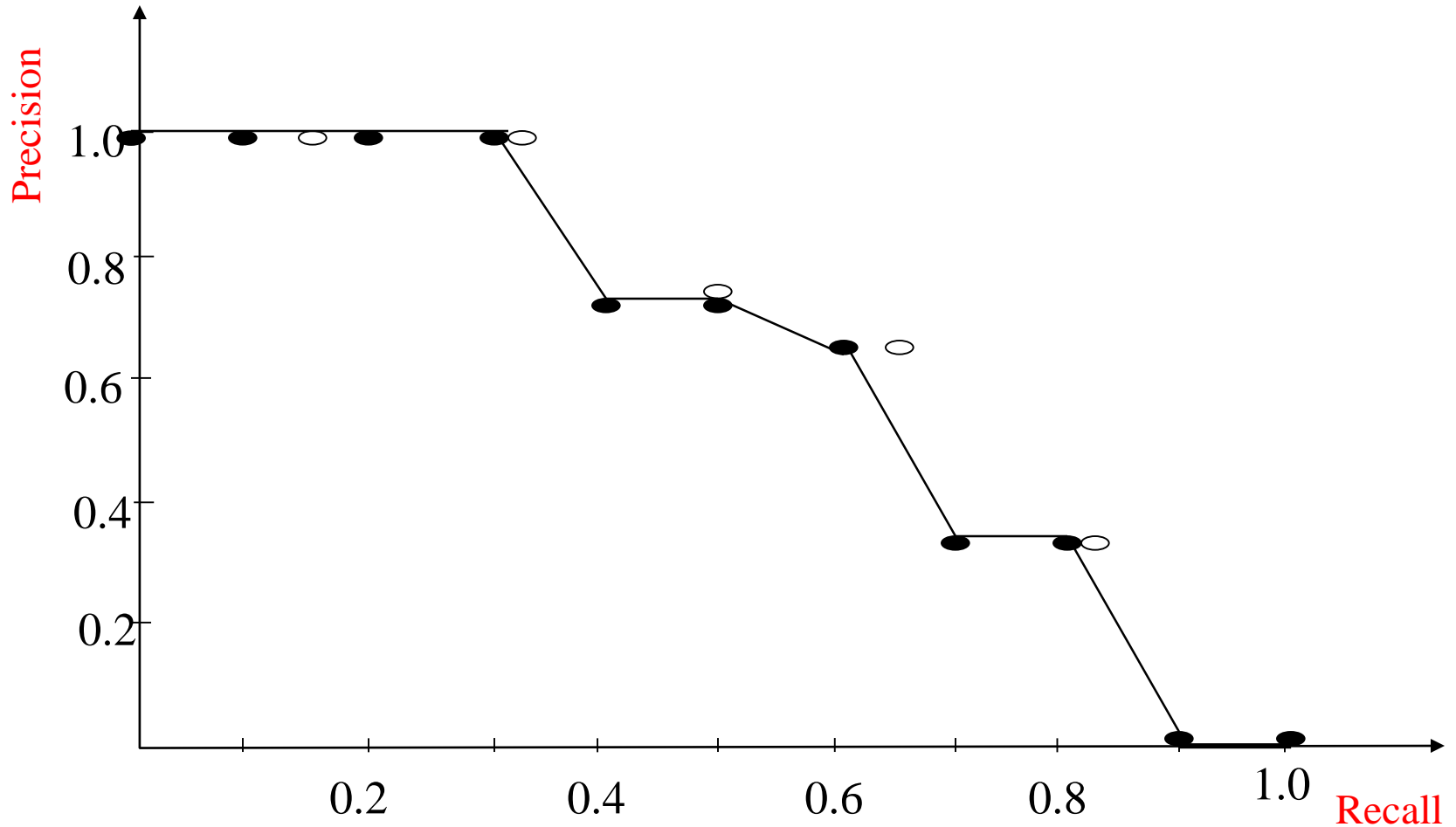
$R=6/6=1.0$; $p=6/14=0.429$

Interpolating a Recall/Precision Curve

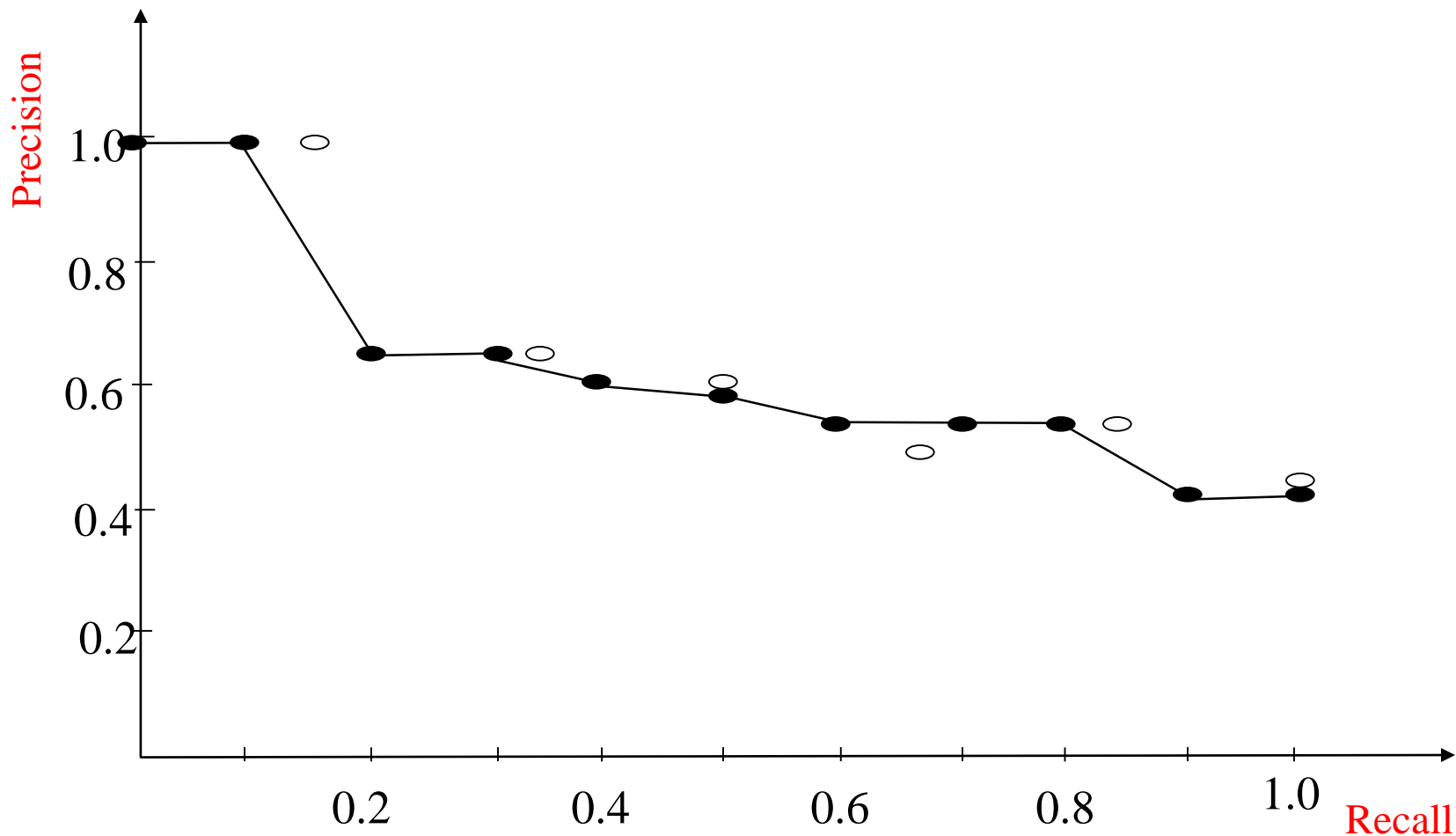
- Interpolate a precision value for each *standard recall level*:
 - $r_j \in \{0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$
 - $r_0 = 0.0, r_1 = 0.1, \dots, r_{10} = 1.0$
- The interpolated precision at the j -th standard recall level is the maximum known precision at any recall level between the j -th and $(j + 1)$ -th level:

$$P(r_j) = \max_{r_j \leq r \leq r_{j+1}} P(r)$$

Interpolating a Recall/Precision Curve: Example 1



Interpolating a Recall/Precision Curve: Example 2

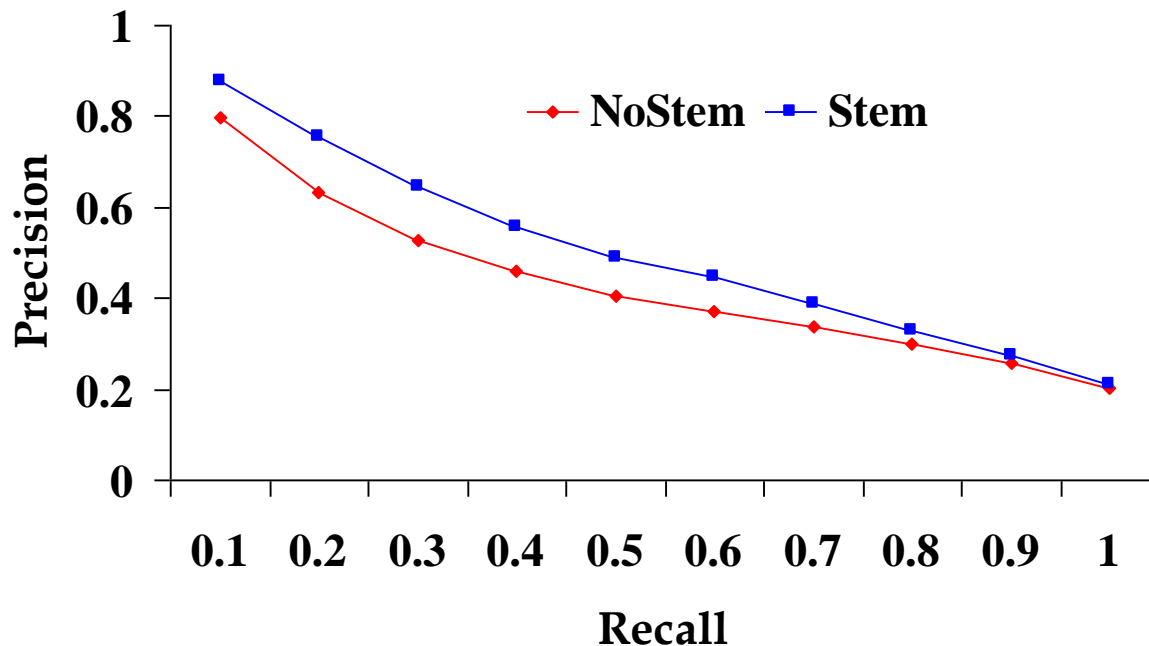


Average Recall/Precision Curve

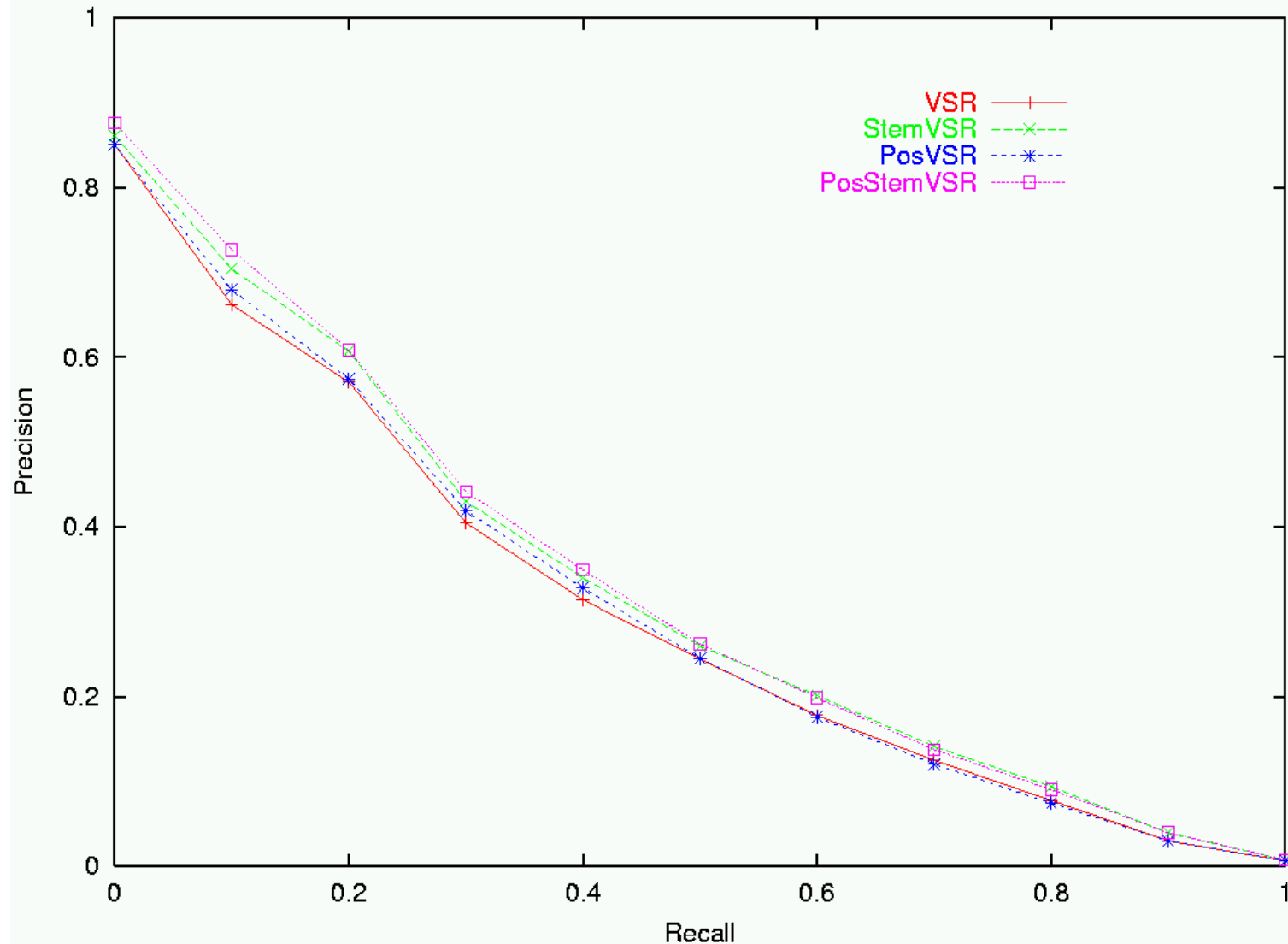
- Typically average performance over a large *set* of queries.
- Compute average precision at each standard recall level across all queries.
- Plot average precision/recall curves to evaluate overall system performance on a document/query corpus.

Compare Two or More Systems

- The curve closest to the upper right-hand corner of the graph indicates the best performance



Sample RP Curve for CF Corpus



R- Precision

- Precision at the R-th position in the ranking of results for a query that has R relevant documents.

n	doc #	relevant
1	588	x
2	589	x
3	576	
4	590	x
5	986	
6	592	x
7	984	
8	988	
9	578	
10	985	
11	103	
12	591	
13	772	x
14	990	

R = # of relevant docs = 6

R-Precision = $4/6 = 0.67$

Mean Average Precision (MAP)

- **Average Precision:** Average of the precision values at the points at which each relevant document is retrieved.
 - Ex1: $(1 + 1 + 0.75 + 0.667 + 0.38 + 0)/6 = 0.633$
 - Ex2: $(1 + 0.667 + 0.6 + 0.5 + 0.556 + 0.429)/6 = 0.625$
- **Mean Average Precision:** Average of the average precision value for a set of queries.

Non-Binary Relevance

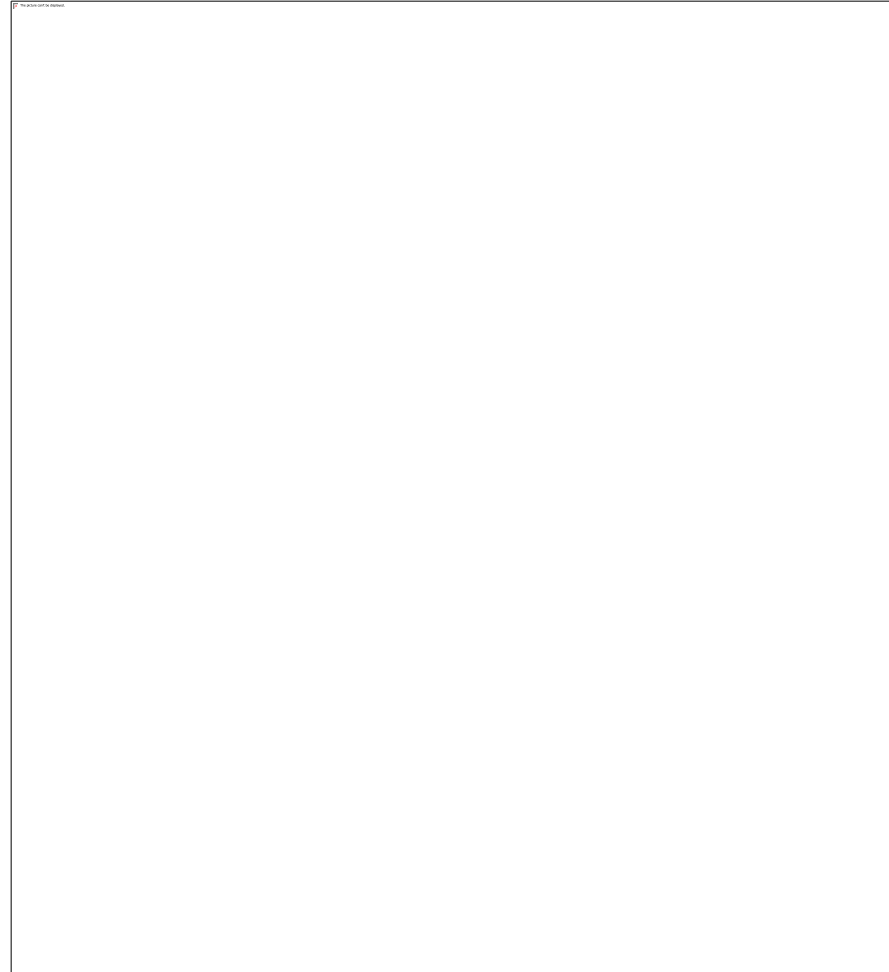
- Documents are rarely entirely relevant or non-relevant to a query
- Many sources of *graded relevance judgments*
 - Relevance judgments on a 5-point scale
 - Multiple judges
 - Click distribution and deviation from expected levels (but click-through != relevance judgments)

Cumulative Gain

- With graded relevance judgments, we can compute the *gain* at each rank.
- **Cumulative Gain** at rank n :

$$CG_n = \sum_{i=1}^n rel_i$$

(Where rel_i is the graded relevance of the document at position i)

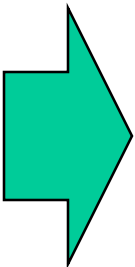


Discounting Based on Position

- Users care more about high-ranked documents, so we **discount** results by $1/\log_2(rank)$
- **Discounted Cumulative Gain:**

Normalized Discounted Cumulative Gain (NDCG)

- To compare DCGs, normalize values so that a *ideal ranking* would have a **Normalized DCG** of 1.0
- Ideal ranking:



n	doc #	rel (gain)	CG _n	log _n	IDCG _n
1	588	1.0	1.0	0.00	1.00
2	592	1.0	2.0	1.00	2.00
3	590	0.8	2.8	1.58	2.50
4	589	0.6	3.4	2.00	2.80
5	772	0.2	3.6	2.32	2.89
6	576	0.0	3.6	2.58	2.89
7	986	0.0	3.6	2.81	2.89
8	984	0.0	3.6	3.00	2.89
9	988	0.0	3.6	3.17	2.89
10	578	0.0	3.6	3.32	2.89
11	985	0.0	3.6	3.46	2.89
12	103	0.0	3.6	3.58	2.89
13	591	0.0	3.6	3.70	2.89
14	990	0.0	3.6	3.81	2.89

Normalized Discounted Cumulative Gain (NDCG)

- Normalize by DCG of the ideal ranking:

$$NDCG_n = \frac{DCG_n}{IDCG_n}$$

- $NDCG \leq 1$ at all ranks
- NDCG is comparable across different queries

n	doc #	rel			
		(gain)	DCG_n	$IDCG_n$	$NDCG_n$
1	588	1.0	1.00	1.00	1.00
2	589	0.6	1.60	2.00	0.80
3	576	0.0	1.60	2.50	0.64
4	590	0.8	2.00	2.80	0.71
5	986	0.0	2.00	2.89	0.69
6	592	1.0	2.39	2.89	0.83
7	984	0.0	2.39	2.89	0.83
8	988	0.0	2.39	2.89	0.83
9	578	0.0	2.39	2.89	0.83
10	985	0.0	2.39	2.89	0.83
11	103	0.0	2.39	2.89	0.83
12	591	0.0	2.39	2.89	0.83
13	772	0.2	2.44	2.89	0.84
14	990	0.0	2.44	2.89	0.84

Issues with Relevance

- ***Marginal Relevance***: Do later documents in the ranking add new information beyond what is already given in higher documents.
 - Choice of retrieved set should encourage **diversity** and **novelty**.
- ***Coverage Ratio***: The proportion of relevant items retrieved out of the total relevant documents ***known*** to a user prior to the search.
 - Relevant when the user wants to locate documents which they have seen before (e.g., the budget report for Year 2000).

Other Factors to Consider

- *User effort*: Work required from the user in formulating queries, conducting the search, and screening the output.
- *Response time*: Time interval between receipt of a user query and the presentation of system responses.
- *Form of presentation*: Influence of search output format on the user's ability to utilize the retrieved materials.
- *Collection coverage*: Extent to which any/all relevant items are included in the document corpus.

A/B Testing in a Deployed System

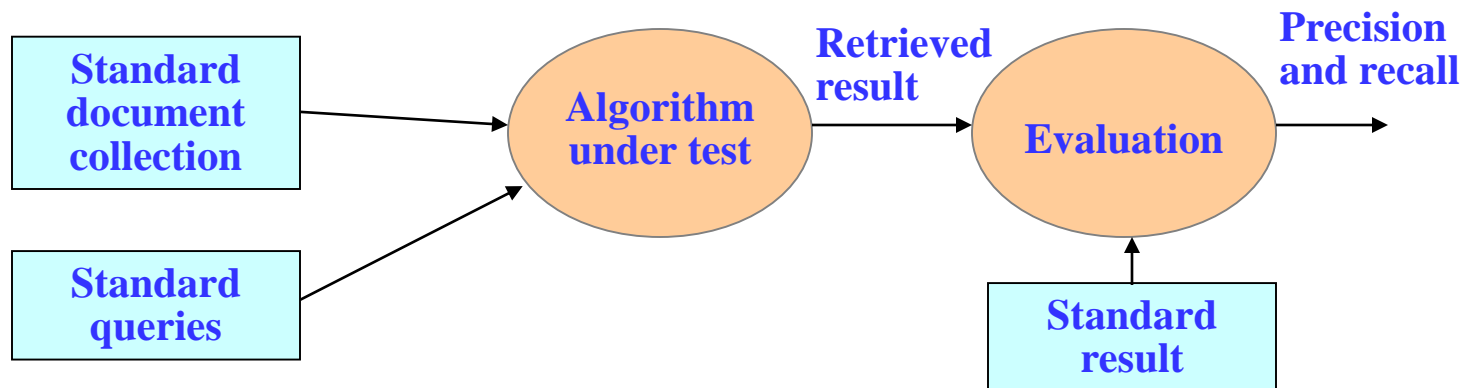
- Can exploit an existing user base to provide useful feedback.
- Randomly send a small fraction (1–10%) of incoming users to a variant of the system that includes a single change.
- Judge effectiveness by measuring change in ***clickthrough***: The percentage of users that click on the top result (or any result on the first page).

Experimental Setup for Benchmarking

- ***Analytical*** performance evaluation is difficult for document retrieval systems because many characteristics such as relevance, distribution of words, etc., are difficult to describe with mathematical precision.
- Performance is measured by ***benchmarking***. That is, the retrieval effectiveness of a system is evaluated on a *given set of documents, queries, and relevance judgments*.
- Performance data is valid only for the environment under which the system is evaluated.

Benchmarks

- A benchmark collection contains:
 - A set of standard documents and queries/topics.
 - A list of relevant documents for each query.
- Standard collections for traditional IR:
 - Smart collection: <ftp://ftp.cs.cornell.edu/pub/smart>
 - TREC: <http://trec.nist.gov/>



Benchmarking – The Problems

- Performance data is valid only for a particular benchmark.
- Building a benchmark corpus is a difficult task.
- Benchmark web corpora are just starting to be developed.
- Benchmark foreign-language corpora are just starting to be developed.

Early Test Collections

- Previous experiments were based on the SMART collection which is fairly small.
(<ftp://ftp.cs.cornell.edu/pub/smart>)

Collection Name	Number Of Documents	Number Of Queries	Raw Size (Mbytes)
CACM	3,204	64	1.5
CISI	1,460	112	1.3
CRAN	1,400	225	1.6
MED	1,033	30	1.1
TIME	425	83	1.5

- Different researchers used different test collections and evaluation techniques.

The TREC Benchmark

- TREC: **T**ext **RE**trieval **C**onference (<http://trec.nist.gov/>)
Originated from the TIPSTER program sponsored by Defense Advanced Research Projects Agency (DARPA).
- Became an annual conference in 1992, co-sponsored by the National Institute of Standards and Technology (NIST) and DARPA.
- Participants submit the P/R values for the final document and query corpus and present their results at the conference.

What we need for a benchmark

- A collection of documents
 - Documents must be representative of the documents we expect to see in reality.
- A collection of information needs
 - . . .which we will often incorrectly refer to as queries
 - Information needs must be representative of the information needs we expect to see in reality.
- Human relevance assessments
 - We need to hire/pay “judges” or assessors to do this.
 - Expensive, time-consuming
 - Judges must be representative of the users we expect to see in reality.

Standard relevance benchmark: Cranfield

- Pioneering: first testbed allowing precise quantitative measures of information retrieval effectiveness
- Late 1950s, UK
- 1398 abstracts of aerodynamics journal articles, a set of 225 queries, exhaustive relevance judgments of all query-document-pairs
- Too small, too untypical for serious IR evaluation today

Standard relevance benchmark: TREC

- TREC = Text Retrieval Conference (TREC)
- Organized by the U.S. National Institute of Standards and Technology (NIST)
- TREC is actually a set of several different relevance benchmarks.
- Best known: TREC Ad Hoc, used for first 8 TREC evaluations between 1992 and 1999
- 1.89 million documents, mainly newswire articles, 450 information needs
- **No exhaustive relevance judgments – too expensive**
- Rather, NIST assessors' relevance judgments are available only for the documents that were among **the top k returned** for some system which was entered in the TREC evaluation for which the information need was developed.

Standard relevance benchmarks: Others

- GOV2
 - Another TREC/NIST collection
 - 25 million web pages
 - Used to be largest collection that is easily available
 - But still 3 orders of magnitude smaller than what Google/Yahoo/MSN index
- NTCIR
 - East Asian language and cross-language information retrieval
- Cross Language Evaluation Forum (CLEF)
 - This evaluation series has concentrated on European languages and cross-language information retrieval.
- Many others

Validity of relevance assessments

- Relevance assessments are only usable if they are **consistent**.
- If they are not consistent, then there is no “truth” and experiments are not repeatable.
- How can we measure this consistency or agreement among judges?
- → Kappa measure

Kappa (K) measure

- Kappa is measure of how much judges agree or disagree.
- Designed for categorical judgments
- Corrects for chance agreement
- $P(A)$ = proportion of time judges agree (on the test set)
- $P(E)$ = what agreement would we get by chance

$$\kappa = \frac{P(A) - P(E)}{1 - P(E)}$$

- $\kappa = ?$ for (i) chance agreement (ii) total agreement

Kappa measure (2)

- Values of k in the interval $[2/3, 1.0]$ are seen as acceptable.
- With smaller values: need to redesign relevance assessment methodology used.

Calculating the kappa statistic

		Judge 2 Relevance		
		Yes	No	Total
Judge 1 Relevance	Yes	300	20	320
	No	10	70	80
	Total	310	90	400

*Observed proportion of
the times the judges agreed*

$$P(A) = (300 + 70)/400 = 370/400 = 0.925 \text{ [probability of agreement by judges]}$$

Pooled marginals

$$P(\text{nonrelevant}) = (80 + 90)/(400 + 400) = 170/800 = 0.2125 \text{ [average for judges]}$$

$$P(\text{relevant}) = (320 + 310)/(400 + 400) = 630/800 = 0.7878 \text{ [average for judges]}$$

Probability that the two judges agreed by chance

$$P(E) = P(\text{nonrelevant})^2 + P(\text{relevant})^2 = 0.2125^2 + 0.7878^2 = 0.665$$

$$\text{Kappa statistic } \kappa = (P(A) - P(E))/(1 - P(E)) =$$

$$(0.925 - 0.665)/(1 - 0.665) = 0.776 \text{ (still in acceptable range)}$$

Interjudge agreement at TREC

Information need	number of docs judged	disagreements
51	211	6
62	400	157
67	400	68
95	400	110
127	400	106

Impact of interjudge disagreement

- Judges disagree a lot. Does that mean that the results of information retrieval experiments are meaningless?
- No.
- Large impact on absolute performance numbers
- Virtually no impact on ranking of systems
- Suppose we want to know if algorithm A is better than algorithm B
- An information retrieval experiment will give us a reliable answer to this question . . .
- . . . even if there is a lot of disagreement between judges.

Evaluation at large search engines

- Recall is difficult to measure on the web
- Search engines often use precision at top k , e.g., $k = 10 \dots$
- \dots or use measures that reward you more for getting rank 1 right than for getting rank 10 right.
- Search engines also use non-relevance-based measures.
 - Example 1: clickthrough on first result
 - Not very reliable if you look at a single clickthrough (you may realize after clicking that the summary was misleading and the document is nonrelevant) \dots
 - \dots but pretty reliable in the aggregate.
 - Example 2: Ongoing studies of user behavior in the lab
 - Example 3: A/B testing

Critique of pure relevance

- We've defined relevance for an isolated query-document pair.
- Alternative definition: marginal relevance
- The **marginal relevance** of a document at position k in the result list is the additional information it contributes over and above the information that was contained in documents $d_1 \dots d_{k-1}$.
- Exercise
 - Why is marginal relevance a more realistic measure of user happiness?
 - Give an example where a non-marginal measure like precision or recall is a misleading measure of user happiness, but marginal relevance is a good measure.
 - In a practical application, what is the difficulty of using marginal measures instead of non-marginal measures?

How do we present results to the user?

- Most often: as a list – aka “10 blue links”
- How should each document in the list be described?
- This description is crucial.
- The user often can identify good hits (= relevant hits) based on the description.
- No need to “click” on all documents sequentially

Doc description in result list

- Most commonly: doc title, url, some metadata . . .
- . . . and a summary
- **[but others exist: on mouseover display page!]**
- How do we “compute” the summary?

Summaries

- Two basic kinds: (i) static (ii) dynamic
- A **static summary** of a document is always the same, regardless of the query that was issued by the user.
- **Dynamic summaries** are **query-dependent**. They attempt to explain why the document was retrieved for the query at hand.

Static summaries

- In typical systems, the static summary is a subset of the document.
- Simplest heuristic: the first 50 or so words of the document
- More sophisticated: extract from each document a set of “key” sentences
 - Simple NLP heuristics to score each sentence
 - Summary is made up of top-scoring sentences.
 - Machine learning approach
- Most sophisticated: complex NLP to synthesize/generate a summary
- For most IR applications: not quite ready for prime time yet

Dynamic summaries

- Present one or more “windows” or **snippets** within the document that contain several of the query terms.
- Prefer snippets in which query terms occurred as a ***phrase***
- Prefer snippets in which query terms occurred jointly in **a small window**
- The summary that is computed this way gives the entire content of the window – all terms, not just the query terms.

A dynamic summary

*Query: “new guinea economic development” Snippets (in bold) that were extracted from a document: . . . **In recent years, Papua New Guinea has faced severe economic difficulties and economic growth has slowed, partly as a result of weak governance and civil war, and partly as a result of external factors such as the Bougainville civil war which led to the closure in 1989 of the Panguna mine (at that time the most important foreign exchange earner and contributor to Government finances), the Asian financial crisis, a decline in the prices of gold and copper, and a fall in the production of oil. PNG’s economic development record over the past few years is evidence that governance issues underly many of the country’s problems. Good governance, which may be defined as the transparent and accountable management of human, natural, economic and financial resources for the purposes of equitable and sustainable development, flows from proper public sector management, efficient fiscal and accounting mechanisms,***

Google example for dynamic summaries

Generating dynamic summaries

- Where do we get these other terms in the snippet from?
- We cannot construct a dynamic summary from the positional inverted index – at least not efficiently.
- We need to cache documents.
- The positional index tells us: query term occurs at position 4378 in the document.
- Byte offset or word offset?
- Note that the cached copy can be outdated
- Don't cache very long documents – just cache a short prefix

Dynamic summaries

- Real estate on the search result page is limited ! Snippets must be short . . .
- . . . but snippets must be long enough to be meaningful.
- Snippets should communicate whether and how the document answers the query.
- Ideally: linguistically well-formed snippets
- Ideally: the snippet should answer the query, so we don't have to look at the document.
- Dynamic summaries are a big part of user happiness because . . .
 - . . .we can quickly scan them to find the relevant document we then click on.
 - . . . in many cases, we don't have to click at all and save ⁵⁵time.