

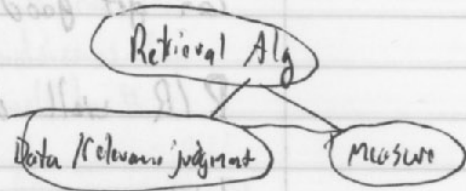
Information Retrieval

5/4/2005

IR models + IR evaluation

Evaluation in document retrieval

- dependent on model for relevance + task



TREC - Text REtrieval Conference

Relevance

- difficult to define
- a relevant document is one judged useful in the context of a query
- w/ real collections, never know full set of relevant documents
- Retrieval model incorporates notion of relevance
- individuals will disagree on individual documents, but will generally agree on avg. across data points

Test Collection

- set of documents
- set of queries
- set of relevance judgements (which docs relevant to each query)

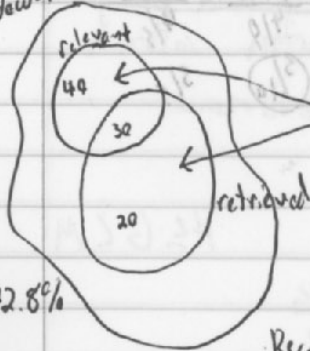
Precision and Recall

$|x|$ - size of x

Precision: proportion of a retrieved set that is relevant

$$\text{Precision} = \frac{|\text{relevant} \cap \text{retrieved}|}{|\text{retrieved}|} = P(\text{relevant} | \text{retrieved})$$

All documents: 12,000



2 types of errors

False-positives:

False negatives:

$$\text{Prec} = \frac{|\text{retrieved} \cap \text{relevant}|}{|\text{retrieved}|}$$

$$\text{Prec} = \frac{30}{50} = 60\%$$

Recall: proportion of all relevant documents in the collection included in the retrieved set

$$\text{Recall} = \frac{|\text{relevant} \cap \text{retrieved}|}{|\text{relevant}|} = P(\text{retrieved} | \text{relevant})$$

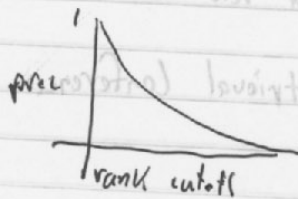
$$\text{Recall} = \frac{30}{70} = 42.8\%$$

can get good recall at expense of precision: return large (all) set
 can get good precision at expense of recall: return 1 relevant doc

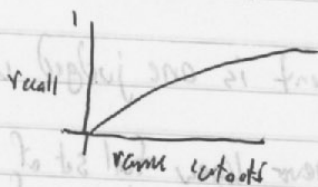
R/R well defined for sets

List

Rank	Relevance
1	R
2	N
3	N
4	R
⋮	R
⋮	⋮



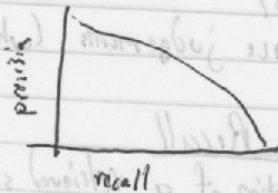
expect - precision to drop
 rank to rise



Average Precision

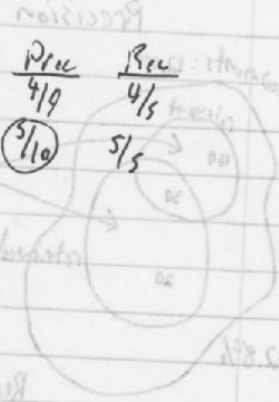
- average precision whenever recall increases (every time find a relevant document)
- highly correlated to what people actually feel
- highly weights early relevant documents in list.

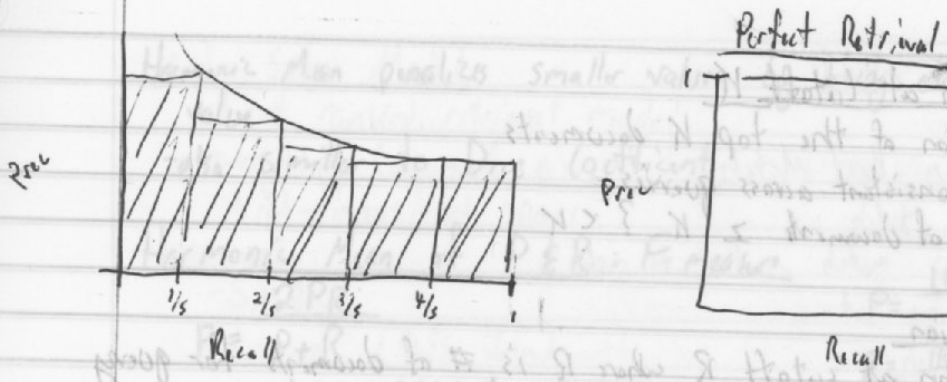
Recall / Precision Graphs



List

Rank	Relevance	Prec	Rec	Rank	Relevance	Prec	Rec
1	R	1	1/5	1	R	1/1	1/5
2	N	1/2	1/5	2	N	1/2	2/5
3	R	2/3	2/5	3	R	2/3	3/5
4	N	2/4	2/5	4	N	2/4	4/5
5	N	2/5	2/5	5	N	2/5	5/5
6	R	3/6	3/5	6	R	3/6	6/5
7	N	3/7	3/5	7	N	3/7	7/5
8	R	4/8	4/5	8	R	4/8	8/5

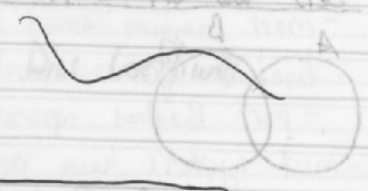




Area under P-R curve is approx. Avg Precision

List

Rank	Relevance	Prec	Recall	Interp Prec
1	R	1	1/5	1
2	N	1/2	1/5	
3	N	1/3	1/5	
4	R	2/4	2/5	2/3
5	R	3/5	3/5	2/3
6	R	4/6	4/5	2/3
7	N	4/7	4/5	
8	R	5/8	5/5	1/2
9	N	5/9	5/5	
10	N	5/10	5/5	



want to plot trend (Interpolated Recall-Precision)

- for a given recall, what is the highest precision anywhere beyond that
- Interpolated Precision always drops

H-LCM

$$H = \frac{1}{\frac{1}{p} + \frac{1}{r}} = \frac{pr}{p+r}$$

$$L = \min(p, r)$$

$$C = \max(p, r)$$

Precision at cutoff K

Precision of the top K documents

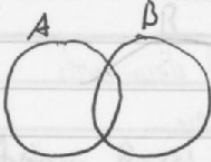
- not consistent across queries
- # of documents $\geq K$? $< K$

R-precision

- precision at cutoff R where R is # of documents for query
- value of 1 if and only if perfect retrieval

5/5/2005

Sets and Set Differences



Jaccard Coeff: $\frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$

Dice Coeff: $\frac{2|A \cap B|}{|A| + |B|}$

Given prec P
rec R

proposition: (arithmetic) average

$$M = \frac{P+R}{2}$$

Incorrect: bc for $M = 0.5$

1 $P = 0.5$ and $R = 0.5$ doing something interesting

2 $P \rightarrow 0$ + $R \rightarrow 1$ } trivial performance

3 $P \rightarrow 1$ + $R \rightarrow 0$

Geometric Mean: $G = \sqrt{P \cdot R}$

$$H \leq G \leq M$$

Harmonic Mean: $H = \frac{1}{\frac{1}{2}(\frac{1}{P} + \frac{1}{R})} = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P+R}$

Harmonic Mean penalizes smaller values by being closer to smaller value
- also similar to Dice coefficient

Harmonic Mean of P & R: F measure

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

$$P = \frac{|Ret \cap Rel|}{|Ret|}$$

$$R = \frac{|Ret \cap Rel|}{|Rel|}$$

$$= \frac{2 |Ret \cap Rel| \cdot |Ret \cap Rel|}{|Ret| \cdot |Rel|}$$

$$= \frac{|Ret \cap Rel|}{|Ret|} + \frac{|Ret \cap Rel|}{|Rel|}$$

$$= \frac{2 |Ret \cap Rel|}{|Rel| + |Ret|} \text{ which is Dice Coefficient}$$

can also parameterize to weight P or R

$$\frac{1}{\alpha \left(\frac{1}{P}\right) + (1-\alpha) \left(\frac{1}{R}\right)}$$

Vector Space Model

① Processing documents

② Remove stop words - "a", "and", "the", "is"

→ ③ Stemming: Reduce all words to their root stems

④ Keep track of statistics over root stems