

IR System Evaluation

Why do we perform evaluation?

Example for problems of evaluation

$$\Delta = (\pm \mid _ _ _ _)$$

$$\Delta = (+_ +_ + \mid _ + _ -)$$

Above notation succinctly represents how a system ranks documents and what user likes / dislikes

Relevance

+ user would like to retrieve

- user would not like to retrieve

vertical lines separate different ranks generated by the retrieval system (by means of RSVs)

Evaluation of systems by evaluation of tasks

2* 2 Table

	rel.	nonrel.
retr.	a	b
not retr.	c	d

$$\text{Recall} = a / (a+c) = R$$

$$\text{Precision} = a / (a+b) = P$$

$$\text{Fallout} = b / (b+d) = F$$

$$\text{Generality} = (a+c) / (a+ b + c +d) = G$$

$$\begin{aligned} \text{GR}/[\text{GR}+(1-\text{G})\text{F}] &= [(a+c)/(a+ b + c +d)]*[a/(a+c)] / \{ [(a+c)/ (a+ b + c +d)]*[a/(a+c)] +[(b+d)/(a+ b + c +d)]*[b/(b+d)] \} \\ &= a / (a+b) \\ &= P \end{aligned}$$

Ranked output (according to retrieval status values)

1	2	3			P
...
...
...

\longrightarrow
 user Rank 1,..., P

Assumption: User retrieves full ranks

R_v recall after having retrieved v ranks

P_v precision after having retrieved v ranks

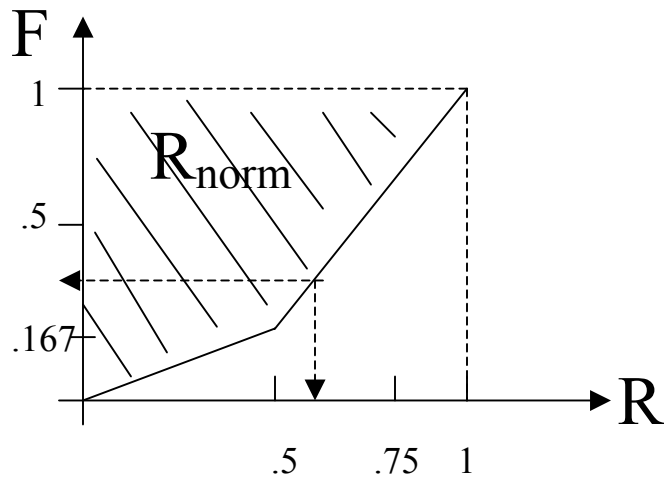
F_v fallout after having retrieved v ranks

Commonly used measure in this case:
recall_precision (R-P) graph

First consider R-F graph

Example:

($\begin{smallmatrix} + & + \\ _ & \end{smallmatrix}$ | $\begin{smallmatrix} _ & + \\ _ & \end{smallmatrix}$ | $\begin{smallmatrix} _ & + & _ \\ _ & \end{smallmatrix}$)



R_{norm} Interpretation:

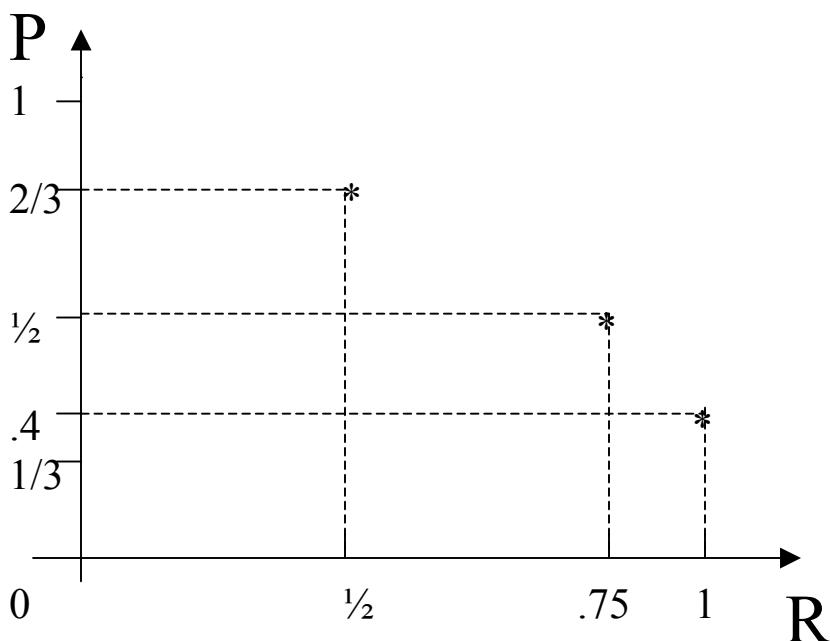
How much better is the system compared to worst case (Fallout is 1 for all search depths)

$R_{\text{norm}} = \text{Rocchio 1965}$

Interpretation of points on R-F graph:

The pair (expected recall after retrieving n documents, expected fallout after retrieving n documents) lies on the Recall_Fallout Graph (n is not necessarily an integer).

Next, we consider R-P graph



How to interpolate?

Straight line interpolation is not meaningful

Interpolation:

Map every point of the Recall_Fallout graph on the Recall_Precision graph using the following formula:

$$(R, F) \rightarrow (R, GR/[GR+(1-G)F])$$

Properties of R-P-graph:

- Every pair (expected Recall after retrieving n documents, expected precision after retrieving n documents) is on the graph.
- Multiplication (creating certain number of copies of every document) of collection gives the same graph
- If considered as graph of a representative sample of very large document collection,

the graph can be interpreted as:

- i) expected precision for a given Recall
- ii) $P(\text{rel/retr.})$ for a given Recall

Multilevel preference (relevance)

$$R_{\text{norm}} = \frac{1}{2} [1 + (I^+ - I^-) / I_{\text{max}}^+]$$

I^+ number of pairs where a better document precedes a worse document.

I^- number of pairs that are inverted (worse document precedes a better one)

I_{max}^+ maximum number of correct pairs

Example:

r relevant

m medium relevant

n non-relevant

r	r		r		m
m			m m		n n n
n			n n		

3 rel.

4 medrel.

6 nonrel.

$$I^+ = 2*8+5+4+6 = 31$$

$$I^- = 2+2+3=7$$

$$I_{\max}^+ = 30+24=54$$

$$R_{\text{norm}} = 1/2[1+(31-7)/54]$$

$$=1/2[(54+31-7)/54]$$

$$=78/108 =.72$$

Special case of rel and nonrel (i.e.
two-level relevance)

N collection size
n number of relevant

$$I_{\max}^+ = n(N-n)$$

Example: ($+_-+$ | $-_+_-$ | $-_-+$)

$$I^+ = 2*5+3= 13$$

$$I^- = 1+3= 4$$

$$I_{\max}^+ =n(N-n) = 4*6 =24$$

$$R_{\text{norm}} =1/2[1+(13-4)/24] = \frac{1}{2} [(24+13-4)/24] = 33/48 = 0.69$$

Theorem: For binary (two-level) relevance the two definitions of R_{norm} coincide.

Practical Problems:

n number of relevant documents may not be known

→ no R – P graph
no R_{norm}

other evaluation options:

- Just use precision, since it is known
- If more relevant documents are retrieved by a system for the same collection → Recall is higher for that system
- Other measures like expected search length, denoted, esl_k , can be used. It indicates the number of nonrelevant documents that can be expected to be retrieved in order to retrieve k relevant documents.

- “disgust rule” kraft et al.
(apprx.1982) how many relevant
document can we expect after
having retrieved a certain number, k ,
of nonrelevant documents.