# IR System for E-commerce platform.

- search for products
- Whether relevant items are being recommended in the context of searched item.

- Relevant products being shown.

- are the searches getting materialized to money ( $ , ₹ )

- [ Repeat buyers ( in days/week/months )
  └→ Dwell time.

# Evaluation of IR systems.

→ How fast can we index?
 → # of documents /hour
 → 10K docs are added per day?

→ How fast can we search?
 = CPU speed      — query size
 = latency

→ Whether "related items" are featuring?

None of these are the correct metrics for judgement at the PRIMARY SPOT.

query → information need → IR system should be judged on the basis of this need.

50K sample queries.
5M products.

5M products

50K
sample
queries

relevance judgement

Binary
0/1

graduated
fair (0)
Gud (1)
Excellent (2)

Each judgement takes
2.5 secs. for a human.

✓ 10¹¹ secs ← 3000+ years.
✓ $10 per hour to a person. $300M

Cyril Cleverdon → Cranfield experiments

Make ~~some~~ generic assumptions about an IR systems.

A test collection:

↳ Three elements

(1) A benchmark collection of documents

(2) A benchmark suite of queries.

(3) An assessment whether a benchmark document is relevant / not relevant to a benchmark query. (human judges).

benchmark docs: products

benchmark queries: ?

Judgements whether a product is relevant to a query.

Relevance → as per the user
need/intent

&

NOT the
query itself

Information need — "My swimming
pool bottom is becoming
black & needs to be
cleaned."

Query → "pool cleaner".

crowd source the judgement task.

AMT → Amazon Mechanical
            Turk platform.

## Benchmark
## query

- query should be appropriate
  to the corpus.

- query $\equiv$ actual info
  need of the user.

- random query terms are
  not a good choice

- What is a good choice?
  "query logs" → Internet
                  capital
                  these day

# Binary assessments.

Precision : No. of <u>relevant</u> docs out of all the <u>retrieved</u> docs.

Recall . No. of <u>relevant</u> docs out of all the <u>relevant</u> docs.

| | Relevant | Not relevant |
|---|---|---|
| Retrieved | true positive (tp) | false (fp) positive |
| Not retrieved | false (fn) negative | true (tn) negative. |

→ P

← recall.

$$Precision = \frac{tp}{(tp + fp)}$$

$$Recall = \frac{tp}{(tp + fn)}$$

Why care about precision & recall?

When is precision important
— precision medicine "~~softe~~ safely critical system"

When is recall important
— covid testing.
— covid vaccination · (flu shot)

" CV screening"

## Rank based measures.

Precision @K
Mean Average Precision (MAP)
Mean Reciprocal Rank (MRR).

Graduated Relevance scores:

Normalized Discounted Cumulative Gain (NDCG).

Set a rank threshold = K.
% of relevant docs in top K.

Ignores all documents that are ranked below k.

R   NR   R   NR   R

■  ■  ■  ■  ■

1     3     5

$$\text{Avg}(\text{Pr@1}, \text{Pr@3}, \text{Pr@5})$$

$$\frac{1}{3}\left(1 + \frac{2}{3} + \frac{3}{5}\right) \simeq 0.76.$$

Estimate whether relevant items are at "<u>better</u>" rank positions.

<u>MAP.</u>

Mean of avg precisions across different queries.
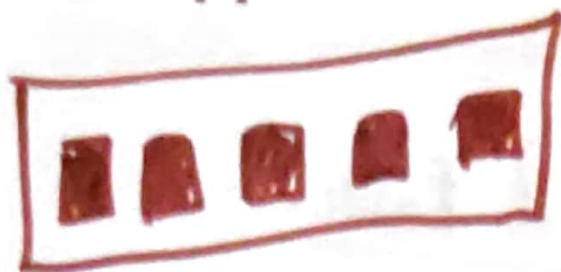
R NR R NR R
█ █ █ █ █

1 2 3 4 5

$$Pr@3 = \frac{2}{3}$$

$$Pr@5 = \frac{3}{5}$$

## Average Precision:

the rank positions of the relevant docs are at $K_1, K_2, K_3$

..., $K_r$.

( Compute $Pr@K$ for each $K_1, K_2, K_3$

... $K_r$ )

↗ Averag

☐☐☐☐☐ = relevant docs for query 1.

| R | NR | R | NR | NR | R | AR | NR | R | R |
|---|----|---|----|----|---|----|----|---|---|
| ■ | ☐  | ■ | ☐  | ☐  | ■ | ☐  | ☐  | ■ | ■ |

Precision  1   0.5  0.66  0.5  0.4  .  .  .  .  .  .

Recall  0.2  0.2  0.4  0.4  .  .  .  .  .  .

☐☐☐ = relevant docs for query 2.

| NR | R | NR | NR | R | NR | R | NR | NR | NR |
|----|---|----|----|---|----|---|----|----|----|
| ☐  | ■ | ☐  | ☐  | ■ | ☐  | ■ | ☐  | ☐  | ☐  |

Precision  0  0.5  0.33  .  .  .  .  .  .

Recall  0  0.33  0.33  0.33  0.67  .  .  .  .  .  .

AP for q1 = 0.62

AP for q2 = 0.44

MAP = (0.62 + 0.44)/2 = 0.53.

# Mean Reciprocal Rank.

rank R of the first relevant
document in the list

$$RR = \frac{1}{R}.$$

MRR → RR across different
queries.

MRR@K.

↳ inspect upto rank K.

Graded relevance.

$\rightarrow$ fair $\Big\}$ —1

$\rightarrow$ good $\Big]$ — 2

$\rightarrow$ Ex — 3

$\rightarrow$ NR $\rightarrow$ 0.

$\longrightarrow$

Highly relevant docs
should be at the top.

$\hookrightarrow$ if a document
comes very low in
the rank list then it
is not so important.

# Discounted gain.

Discount the gain

$$\frac{1}{\log_2(\text{rank})}$$

↳ discount at rank

$$4 = \frac{1}{2}$$

$$8 = \frac{1}{3}$$

$$\vdots$$

$[0 \ldots k]$ , $k >= 2$.

relevance based ratings of 'n' docs

$r_1, r_2, \ldots, r_n$ (in ranked order)

$CG = r_1 + r_2 + r_3 \ldots + r_n$

$DCG = r_1 + r_2/\log 2 + r_3/\log 3 + r_4/\log 4 + \ldots + r_n/\log n$

$$DCG = rel_1 + \sum_{i=2}^{\to P} \frac{rel_i}{\log_2 i}$$

$\downarrow$ normalization

$\boxed{NDCG}$

Precision at a rank. (K)

$$Pr@3 : = \frac{1}{3}$$

| R | NR | NR | R | NR |
|---|----|----|---|----|
| ■ | ☐  | ☐  | ■ | ☐  |

1 2 3 4 5

## Avg. Precision

AP:

$Pr@K \rightarrow$ only those K's rank positions where we have a relevant doc

$Avg\left(Pr@1, Pr@4\right) = \frac{1}{2}\left(1 + \frac{2}{4}\right)$

$$\begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{bmatrix} \begin{array}{l} \rightarrow AP \\ \rightarrow AP \\ \\ \rightarrow \end{array}$$

MAP. (mean avg precision).

MRR. $\rightarrow$ RR = inverse of the rank of the first relevant doc $\rightarrow (L) = \frac{1}{R}$

E Mean of RR over multiple queries

# Discounted cumulative gain. (DCG).

$$R \rightarrow R_1 . \left(\frac{full}{score}\right)$$

$$\underset{\text{log discount.}}{\longrightarrow} R \rightarrow R_4 \qquad \left(\frac{score \leftarrow discounted)}{\log_2(rank)}\right)$$

$$DCG = rel_1 + \sum_{i=2}^{P} \frac{rel_i}{\log_2 i}$$

$$\underline{DCG@P} . (upto \; rank \; p).$$

Ideal DCG.

$$[0-3] \leftarrow \qquad \longrightarrow$$

IR.
$$\rightarrow \; \underset{\uparrow}{3}, \underset{\uparrow}{2}, \underset{\uparrow}{3}, \underline{0}, 0, 1, 2, 2, 3, 0.$$

Sort $\rightarrow$ 3, 3, 3, 2, 2, 1, 0, 0, 0.

Actual DCG: $\underline{3}, \underline{5}, \underline{6.89}, \underline{6.89} \cdots$

Ideal DCG: 3, 6, 7.89, 8.89 $\cdots$

NDCG : $\underset{1}{\underline{1}}, \underset{2}{\underline{0.83}}, \underset{3}{\underline{0.87}}, \underset{4}{\underline{0.76}} \cdots$

**Intrinsic eval technique.**

**Extrinsic eval technique.** ( task based evaluation ).

User judgements for that task at your disposal.

X

informs another task ($T_2$)

unevaluated task with no user judgement $T_1$ output

user judgement is there

Based on this addl. info → does $T_2$'s performance improve?

## Improve the recall of an IR System?

- Relevance feedback.
- Query expansion.

term.

q: [aircraft].

= [plane]
= [airplane].

Can we leverage this info. — to improve
the overall performance?
┌ RF ← local method.
│        some feedback from the
│        user to change the
└ QE      query.

↳ Global method.
  ↳ Dictionary/
    Thesaurus.
    (to expand
    your query)

## Relevance feedback.

$u \xrightarrow{\text{fires}} q_0$ (initial query).

SE → returns some documents
(search engine) related to $q_0$
(relevant to) $q_0$.

feedback
$\Big[$ $u$ → marks some A the retrieved
docs as R OR NR.

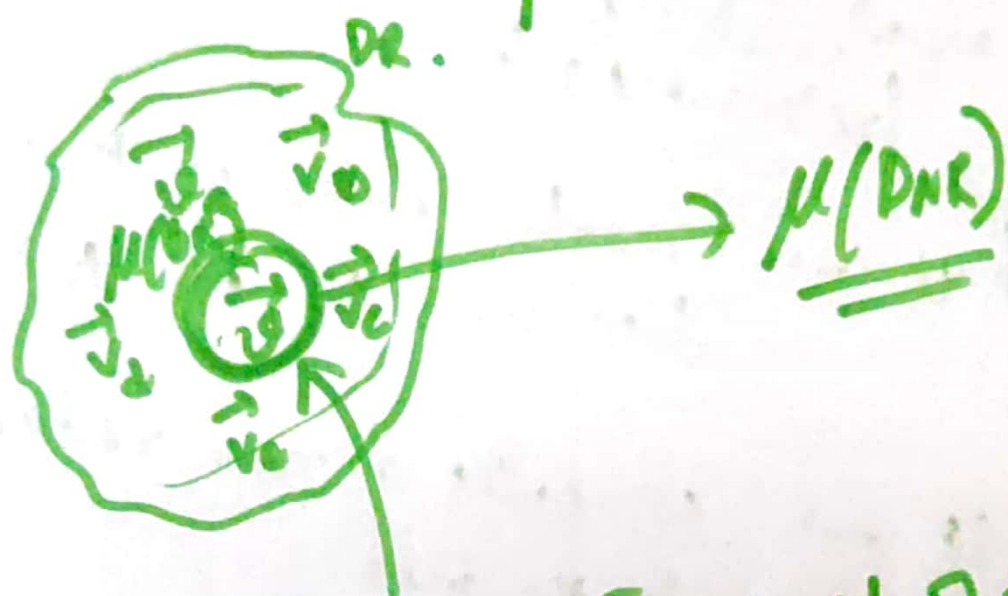SE → computes a new representation
of the info.

A $q_1$ is run on SE

Expectation: recall should improve.

---

## Rocchio's Technique.

- Initial query $q_0$
- $u$ ← relevant documents ( $D_R$ )
  non-relevant documents ( $D_{NR}$ )

→ having a vector ($\vec{V}_{opt}$) that maximally
separates $D_R$ from $D_{NR}$ in your vector space.

$$\vec{v}_{opt} = \boxed{\operatorname{argmax}_{\vec{v}}} \left[ \sin(\vec{v}, D_R) - \sin(\vec{v}, D_{NR}) \right].$$

cosine sim.



$\mu(D_{NR})$

$\vec{\mu}(D_R)$ [Centroid of the vectors in $D_R$].

$$= \frac{1}{|D_R|} \sum_{d_j \in D_R} d_j$$

Compute: $\vec{\mu}(D_{NR})$. $\rightarrow \vec{v}_{opt} = \underline{\vec{\mu}(D_R) - \vec{\mu}(D_{NR})}$

Shift $q_0$ with $\vec{v}_{opt}$ to get the optimal query.

$$\vec{q}_{opt} = \boxed{\vec{\mu}(D_R)} + \left[ \vec{\mu}(D_R) - \vec{\mu}(D_{NR}) \right]$$