

Information Retrieval

→ Data :-
(structure-wise)
Structured
Semi - Structured
Unstructured

→ Data :-
(Nature wise)
image
text ✓
audio
video

Raw data : without interpretation only
text & number.

data
↓

information (data with some interpretation).
↓

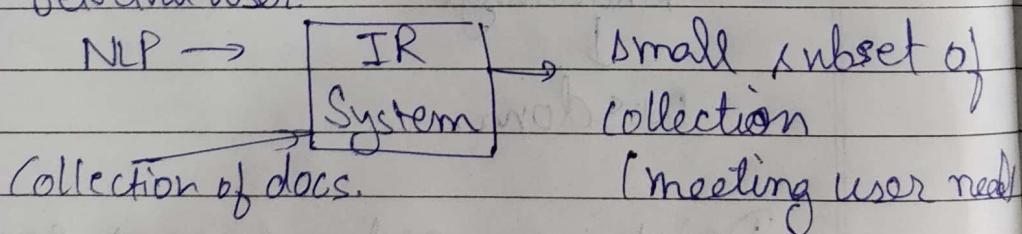
knowledge
↓

wisdom.

→ Database only support exact match.

<u>DB</u>	<u>IR/Web Search.</u>
1) Fixed format of Data	1) Free flowing data: text, audio, video, image.
2) Structured Data type	2) *Unstructured*
3) Knowledge of Schema.	3) No knowledge is required
4) Parameterized mode of querying.	4) Natural Language
5) Exact search.	5) Approximate search
6) Set based search.	6) rank-based retrieval.

Retrieval user:-



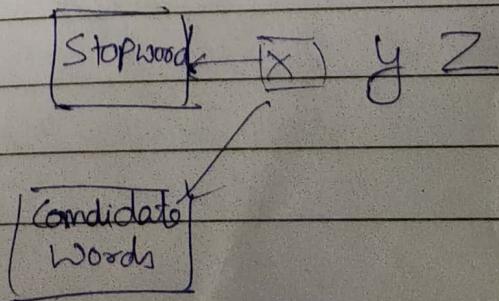
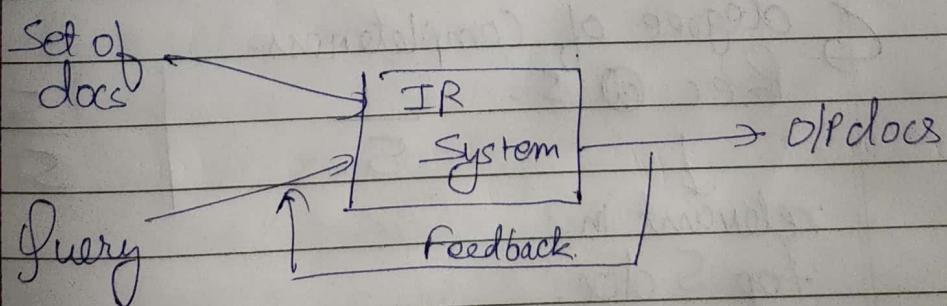
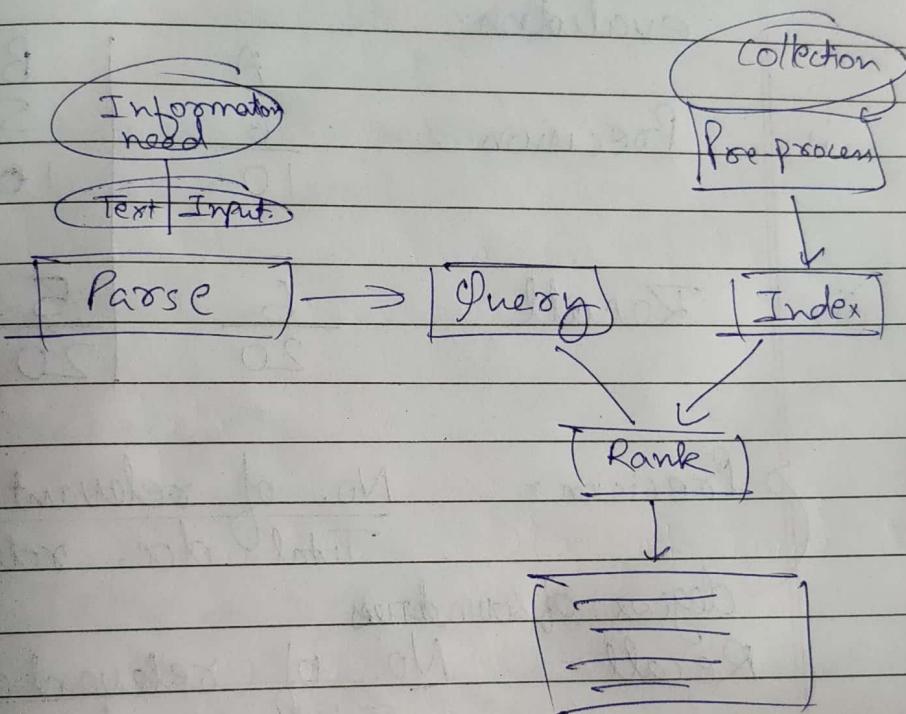
Google has accessed to several billions page.

20-30% publicly available
surface web

Small
subset of
collection - C
me...

Retrieval ← matching
ranking

Information Retrieval Process



search.

data:
video, image

*

is required

usage

file search
based
on al.

subset of

user needs

billions

Small
set of
on-C
meeting user need)

IR → Retrieval

Evaluation

2 basic technique for evaluation:

(i) Precision

A	B
6	5
10	10

(ii) Recall

A	B
6	5
20	20

Precision = $\frac{\text{No. of relevant doc.}}{\text{Total doc. returned}}$

degree of soundness

Recall = $\frac{\text{No. of relevant doc.}}{\text{Total rel. document}}$

(iii) degree of completeness.

Prec @ 5 =	A	B
	1	5

relevant in
top 5 doc.

Christopher D Manning

Prabhakar Raghavan.

PAGE NO.:
DATE: / /

IR applications :-

- i) Web search.
- ii) Desktop search
- iii) Email search
- iv) Enterprise search.
- v) Bibliographic.
- vi) Legal search
- vii) Medical search
- viii) Social search.

Filtering :-

Clustering

Categorization.

$$IRM = \langle D, Q, F, R(q_i, d) \rangle$$

Document ↗ System ↳ Ranking Function
Query

Measurable Properties.

→ Response time of System.

how fast it process?

No. of doc./hr.

how fast it search?

High level view of IR process.

Back

End

Content

Management

Front

End

gi' ← 1
gi' ← 1
gi' ← 1

Case folding - Boing all the characters in same case.

Stemming - Provides the root of any word by chopping off prefix or suffix.

Education

Educating

educational

educative

educate

educat

Lemmatization

(ii)

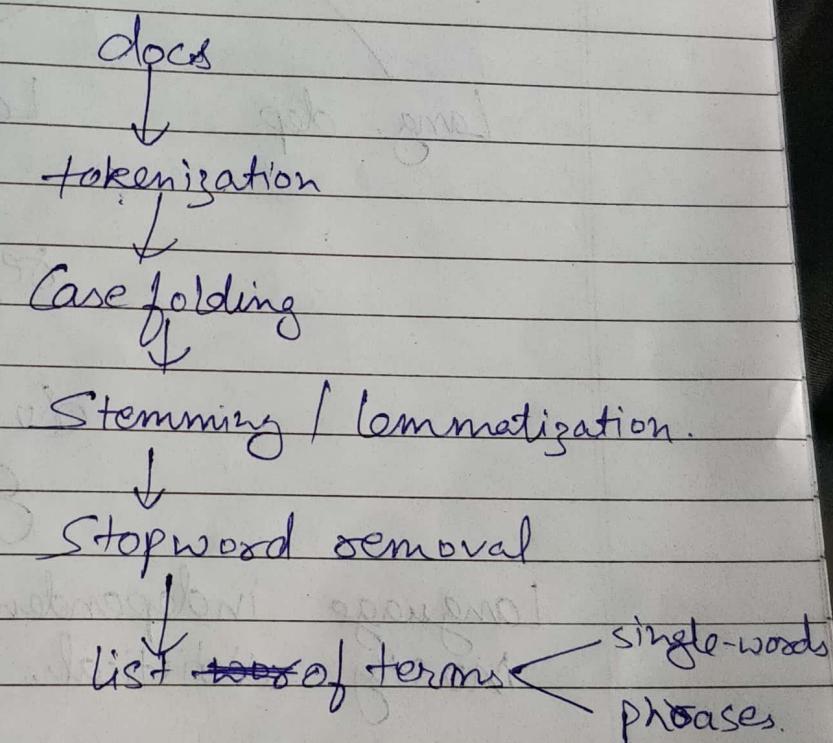
Getting the root from (lemm) of a word by applying rules of the lang. concerned.

Eg. Countrymen → Country man.

We either use Lemmatization or stemming not both.
Document Parsing.

Stopword Removal.

Phrases :- 'New York' → Newyork
Indexing.



Front End

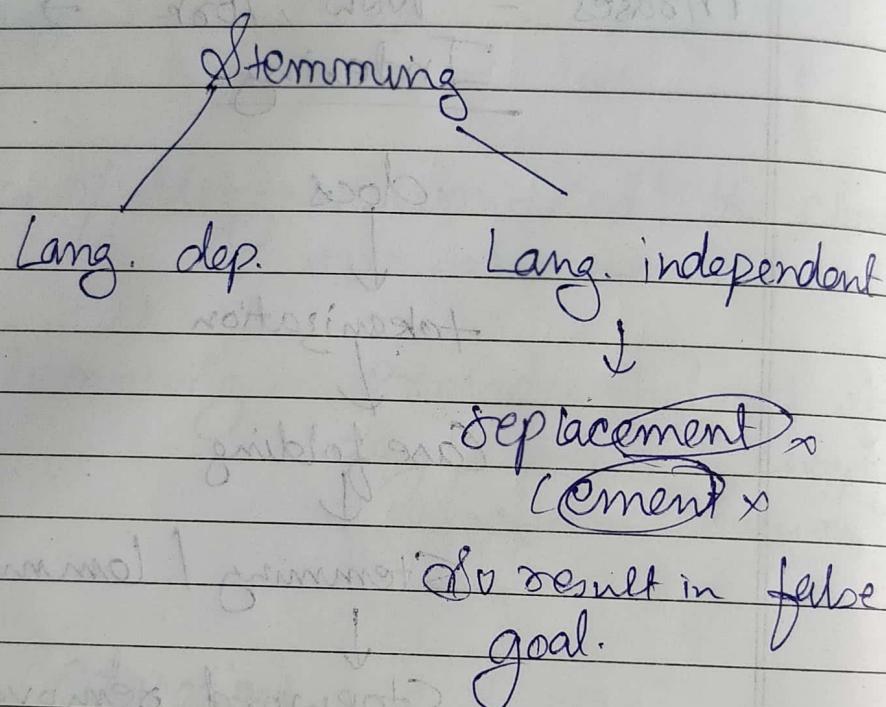
the character

root of
chopping
suffixes

(Lemma)
ules of
y man.

Inverted index construction.

index means list of words with document id.



Language independent technique are mostly statistical.

Download Porter's

saw ← Part of see
saw.

Lemmatization :- NLP tool.

$$10^6 \times 6 \times 1000$$

$$6 \times 10^9$$

$$100 \times 100 \\ 0.4$$

PAGE NO.: _____
DATE: / /

We need variable-size posting list.

On disk, a continuous run of postings is normal and best.

Boolean Retrieval.

Considering only presence (1) and absence (0) of terms are connected to each other by AND, OR, NOT

merge op. $O(m+n)$

And result $O(\min(m,n))$

OR result $O(m+n)$.

From documents we get term.
that's why invert indexing.

$$O(n^2)$$

$n = 1$ Gb dictionary is Mb.

Postings are usually larger $O(n)$.

- Normally kept on Disk.

Hegel law:- $\frac{1}{\text{dict. size}} \text{ const.}$

$$\beta \in [0.4, 0.6]$$

$$\sqrt{= K N}$$

dict. size
(length, fn).

↳ total no. of terms in collection.

Dictionary size < Posting size

When taking AND operation take
and min. of first two doc small document.

$$P = t_1 \& t_2 - \dots - \& t_n$$

Process in order of increasing freq.

Start with smallest one first.

Word doc-freq.

(t1) n₁ egg

(t2) n₂ bread

(t3) n₃ milk

Dict.

Not Relevant

→ Introduction to Information Retrieval:

Dictionary

Postings

<term, doc.freq> <doc.ids>

Single-word

Multiwords

Phrases.

→ Search: "Stanford University"

Not → I found Stanford Henry in the
Relevant univ. campus. ← Doc9

Stanford :- 1 3 5 (9)

University :- 3 (9) 11 25

Stanford University :-

Proximity Search

iit Varanasi

iit at Varanasi.

Store location of the word
in the document in our
posting list.

It will help in proximity
search.

fools rush in. brain glands

 2 4 7

fools 1 8 3,13

rush 2 9 4,14

in 3 10 5,15

 2,4,7,8,12

angels fear to tread:

angels X 12

fear X 13

to X 14

tread X 15

8 { 4 }

{ 4 }

positional index size

A positional index expands postings
storage substantially.

2 <1> 4 <8> 7 <513>
9

PAGE NO.: _____
DATE: / /

A Positional index is 2-4 times larger than a non-positional index.

Positional index size 35-50% of volume of original text.

For some popular phrases, we can use biword indexes.

Combination Combination

It requires 26% more space than having a positional index alone.

t_1 and t_2 and ... t_{n-1} and t_n match
 $n-1$ match = 0 match.

OR t_1 OR t_2 OR ... t_n .
 n matches = 1 match.

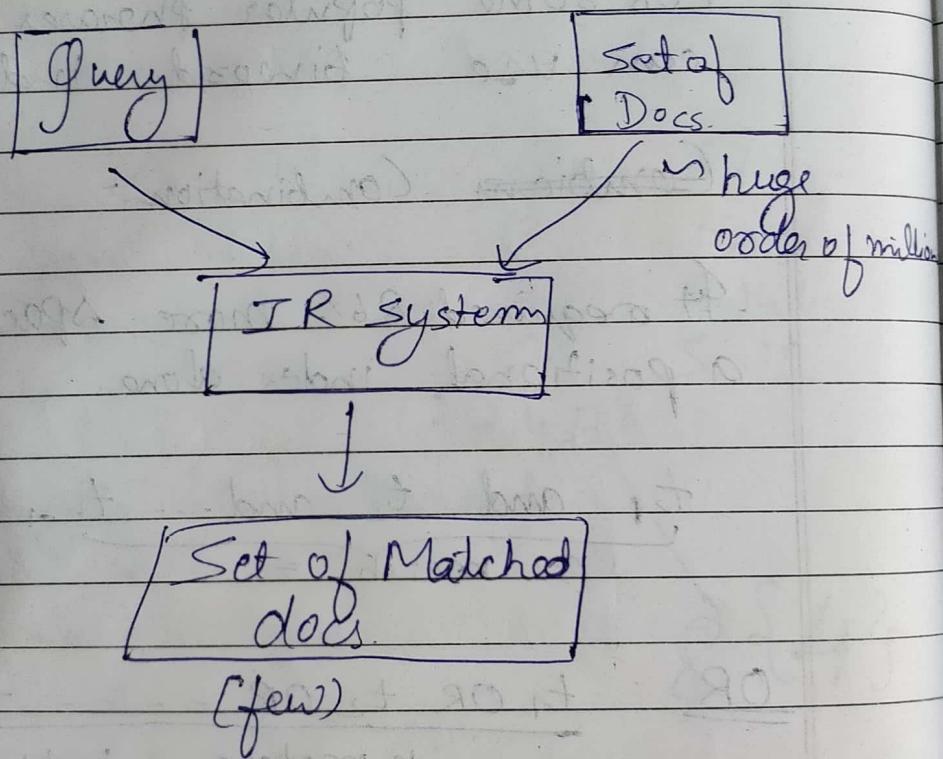
Skip pointers

$t_1 \rightarrow [3] \rightarrow [5] \rightarrow [100]$

$t_2 \rightarrow [1] \rightarrow [3]$

Retrieval Models

- Boolean → Set based retrieval / ranking not possible.
- Vector Space Models.
- Probabilistic



t_i 's connected by Boolean operators AND OR
= Boolean vector.

0.5 LPA

PAGE NO.: _____
DATE: / /

	t_1	t_2	\dots	t_n
D_1	1	0	0	\dots
D_2	0	1	0	0
D_3	1	0	0	0
D_n	0	0	0	1

Vector space model.

- Both docs & query are considered to be vectors with non-negative term weights. (importance of term in the doc.)
- Terms are dimensions in the vector space.

• Term-weight \propto term-frequency.

D_3 1000 words

\hookrightarrow 50 40 $t_i, 60$
 t_1 t_2 t_j

• Term-weight $\propto \frac{1}{\text{doc-length}}$.

• Term-weight \propto rarity of the term.

\propto inverse doc-frequency $\propto \frac{1}{\text{idf}}$

= # docs where a term is present.

tw \propto tf
 \propto idf
 \propto dl

$\phi_i \leftarrow \langle 0, w_{i1}, w_{i2}, \dots, 0 \rangle$

$w_{ij} \rightarrow$ weight of the term i in query j

$$D_1 = \langle w_{11}, w_{21}, \dots, w_{t1} \rangle$$

$$D_2 = \langle w_{12}, w_{22}, \dots, w_{t2} \rangle$$

$$D_j = \langle w_{1j}, w_{2j}, \dots, w_{tj} \rangle$$

Closeness check

D) Euclidean Distance.

Angle θ will represent the closeness b/w D & Q.

$$\text{Sim}(Q, D_i) = \frac{D_i}{\theta}$$

0° represent max. closeness

90° represent min. closeness.

Similarity

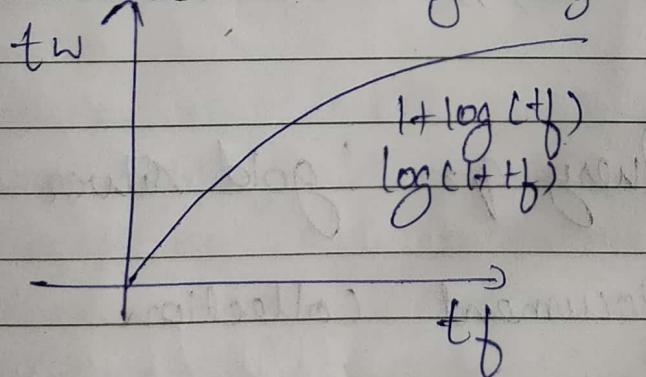
Similarity,

$$\text{Sim}(\vec{D}_j, \vec{D}) = \cos(\angle \vec{D}_j)$$

$$= \frac{\vec{D}_j \cdot \vec{D}}{|\vec{D}_j| \cdot |\vec{D}|}$$

$$\begin{aligned} & w_{1j} * w_{1g} + w_{2j} * w_{2g} + \dots w_{tj} \\ & \sqrt{w_{1j}^2 + w_{2j}^2 + \dots w_{tj}^2} \\ & \times \sqrt{w_{1g}^2 + w_{2g}^2 + \dots w_{tg}^2} \end{aligned}$$

We don't use raw term weights.
because it do symbolize real importance.



Similarly $\log(N)$

Other technique:

$$tf = 0.5 + 0.5 \times \frac{tf}{\text{max. tf.}}$$

weight = $t_f \times idf / \text{normalization}$

Q. Consider two Documents D_1, D_2 & query q .
 $D_1 = (0.5, 0.8, 0.3)$
 $D_2 = (0.9, 0.4, 0.2), q = (1.5, 1.0, 0)$.

$$\cos(D_1, q) = \frac{0.5 \times 1.5 + 0.8 \times 1}{\sqrt{1.5^2 + 1} \times \sqrt{0.5^2 + 0.8^2 + 0^2}}$$

$$\cos(D_2, q) = \frac{0.9 \times 1.5 + 0.4}{\sqrt{1.5^2 + 1} \times \sqrt{0.9^2 + 0.4^2 + 0.2^2}}$$

$$\cos(D_1, q) = \underline{0.87}$$

$$\cos(D_2, q) = 0.965 = 0.97$$

Q2. Query q_2 is 'gold silver truck'.

Document collection.

d_1 = 'shipment of gold damaged in a

d_2 = 'Delivery of silver arrived in a truck'

d_3 = 'Shipment of gold arrived in a truck'.

	$D_1 = 17)$	D_2	D_3
Shipment	$\frac{1}{2}$	0	$\frac{1}{2}$
of	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
gold	$\frac{1}{2}$	0	$\frac{1}{2}$
damaged	1	0	0
in	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
a	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
fire	1	0	0
delivery	0	1	0
silver	0	2	0
arrived	0	$\frac{1}{2}$	$\frac{1}{2}$
tough.	0	$\frac{1}{2}$	$\frac{1}{2}$

$$(1 + \ln(17)) \times \log \left(\frac{N}{2} \right)$$

No 3

Pros:-

(i) Impossibility of formulation of structure due to assumption of stochastic independence between terms. get rid off (all are orthogonal)

(ii) Terms are axes.

10.00

PAGE NO.:
DATE: / /

Evaluation

Relo

par

ap

IR
Retrieval

Pa

[Query] [Set of docs]

R

IR System

Retrieval

docs

SR

SN

(A) $\{d_1, d_2, d_3\}$

docs

met

req.

SR

SN

Relevance is decided by the user particularly the user who posed the query.

Precision = $\frac{\# \text{ relevant retrieved}}{\# \text{ retrieved}}$

Recall = $\frac{\# \text{ relevant retrieved}}{\# \text{ Total relevant}}$

3 System. Total relevant 50.

SR	9R	12R
SNR	11NR	18NR

PScore

5	9	12
10	20	30
S	9	12
50	50	50

No. of
docs

A Retrieved
docs
(Rel)

$$P = \frac{1}{1 + Rel}$$

$$R = \frac{1 + Rel}{1 + Rel}$$

No. of
docs

f_1 -measure.

\Leftarrow H.M of P & R.

$$f_1 = \frac{P + R}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P+R}$$

f_α -measure

$$f_\alpha = \frac{\alpha P + (1-\alpha)R}{\frac{1}{P} + \frac{1}{R}}$$

$0 \leq \alpha \leq 1$

$\alpha > 0.5$ means more importance
given to precision term.

$\alpha = 0.5$ balanced F -measure.

F_1

$$B^2 = \frac{1-\alpha}{\alpha} \sqrt{\text{ratio of importance}}$$

$f_1 \rightarrow$ means B is 1

$$F_2 = \frac{\frac{1}{R_1} + \frac{1}{R_2}}{R_1 + R_2} = \frac{0.2 \times \frac{2}{10} \times \frac{5}{50}}{\frac{2}{10} + \frac{5}{50}}$$

$$\frac{5}{1 + \frac{1}{5}} = \left[\frac{1}{6} \right]$$

F

1

F

2

$$F_1 = \frac{2 \times \frac{9}{10} \times \frac{9}{50}}{2 \times \frac{9}{10} + \frac{9}{50}} = \frac{9}{50}$$

obstacles

$$= \frac{9}{35} = 0.25$$

$$as we \quad F_1 = \frac{2 \times 0.4 \times 0.25}{0.64} = 0.3.$$

AM is influenced more by obstacles
Impacts **obstacles** **AM does not sense**
we take HM.

Sys	R	Ave
1 NR	0.1	0.150
2 R	1/2	1/50
3 NR	1/3	1/50
4 R	2/5	2/50
5 R	3/5	3/50
6 R	4/6	4/50
7 NR	4/12	4/50
8 NR	4/8	4/50
9 NR	4/9	4/50
10 R	5/10	5/50

$$P = \frac{1}{(Rel + Ret)} = \frac{1}{(Ret)}$$



$$P = \frac{1}{(Rel + Ret)}$$

$$R = \frac{Rel}{Rel + Ret} < 1$$

$$R = \frac{1}{Rel + Ret}$$

$$R = 1 \Rightarrow$$



other
other
other

$$R = \frac{\alpha}{\rho + \frac{1-\alpha}{R}}$$

$$\alpha = 0.5 \text{ (Balanced).}$$

$$\text{Average PRec} = \frac{1}{R_{\text{tot}}} \leq (1_2 + 2x_3 + \frac{3}{5}x_5 + \frac{4}{6}x_6)$$

$$AP_1 = \frac{1}{50}$$

Average PRec = $\frac{1}{R_{\text{tot}}} \sum_{i=1}^n P_{\text{rec}_i}$
(Total Rel.) (i = Ranks where Reloc is present)

$$AP_2 = \frac{1}{50} \left(1 + \frac{2}{5} + \frac{3}{6} + \frac{4}{7} + \frac{5}{9} \right)$$

AP is a top-heavy measure. High PRec gives high AP score.

But it also takes into account recall. So it is reasonably balanced.
 Missing terms implies 0 PRec. In

other words, there will always be retrieved at rank infinity.

Interpolated Precision

→ To compare individual precision values at diff. ranks is difficult sometimes.

→ Better to compare prec. at diff. recall-values.

$$\text{Prec}(\delta) = \max_{\delta' \leq \delta} \text{Prec}(\delta')$$

$$0 \leq \delta, \delta \leq 1$$

$$\text{Pint}(0) = \text{highest p-value}$$

- Non increasing function.

Relevant Non-relevant

Rot.	TP	FP
Non-Rot.	FN	TN

$$\text{Prec.} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Fall-out =

$$FN +$$

~~accuracy
efficiency
difficulty~~

Interpolated Precision (δ) = max^{*} prec. (δ')

$$\delta' \geq \delta$$

Point δ' max prec. (δ')

$$\delta' \geq \delta$$

Interpolated
Precision δ' and
natural
Precision δ

If there are 2 queries Q_1, Q_2 = 10 relv.
 $\delta_1 = 2$ relv. we cannot compare them
~~as~~ by natural precision as they
 will have different recall values at
 that point but by interpolated precision
 we can do it.

11- Standard recall points.

$$= \{0, 0.1, 0.2, \dots, 1.0\}$$

(0)- Standard recall points.

$$= \{0, 0.01, \dots, 0.99, 1.0\}$$

AP@t = $\frac{1}{|Rel|} \sum_{i=1}^{|Rel|} \text{Rec}(S_i)$

$|Rel|$ no. of ranks where Roldoc is found.

Rel = set of relevant documents.

$$P_{\text{Rec}} = \frac{1}{|\text{Rel}|} \sum_{i=1}^{|\text{Rel}|} P(\text{rel}_i)$$

$$\text{AP}_{\text{Int}} = \frac{1}{|\mathcal{T}|} \sum_{t=1}^{|\mathcal{T}|} \text{P}_{\text{Int}}(t)$$

Mean Avg. Prec. (MAP)

$$= \frac{1}{|\mathcal{T}|} \sum_{t=1}^{|\mathcal{T}|} \text{AP}(t)$$

$\mathcal{T} \Rightarrow$ Set of queries/topics.

$$\frac{\text{MAP}_{\text{Int}}}{\text{MAP}_{\text{Rec}}} \quad |\mathcal{T}| = 50/100.$$

$P@K =$ Precision at rank K.

$P@5 =$ (r-precision) @ rank = 50

$$= P@50$$

$P@R =$ Recall at R.

↳

Mean average precision & gives priority to top heavy documents.

MAP cannot be used for web-search because we cannot get no. of relevant documents.

1.03

Cumulative Gain (Javelin & Kealihew)

To 13, 2002.

- 4 grades of rel 3 (Highly rel)
 2 (Moderately rel)
- 1 (Marginally rel)
- 0 (Non-relevant).

Ranks	Gain	C.G.	Ideal CG. (G.I.)
1	3	3/13	3/13
2	5	3/6	3/6
3	8	3/4	3/4
4	10	2/11	2/11
5	8	2/13	2/13
6	9	2/15	2/15
7	11	1/16	1/16
8	13	0/16	0/16
9	16	0/16	0/16
10	16	0/16	0/16

$$n \cdot C_G = C_G K_G$$

$$3/3 = 1$$

$$5/6$$

$$8/9$$

$$Vis & symbols$$

$$11/16$$

$$13/16$$

$$16/16$$

$$16/16$$

Discounted Cumulative Gain.

Logarithmic discount base

$$\text{Discounted gain } (\delta) = \frac{\text{Gain}(\delta) \text{ when}}{\log_2(\delta)}$$

D ₁	D ₂	D ₃	D ₄	D ₅	D ₆	D ₇	D ₈	D ₉	D ₁₀
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3

$$\text{DCG}_n(8) + \log \frac{\text{DCG}_n(2)}{\text{DCG}_n(1)}$$

NDCG

1

5/6

6.89

7.89

8.89

Evaluation

→ Set Based → Precision

→ Ranked → f-measure → P@K.

Binary relevance → AP
Reciprocal Rank → MRR

when

else

→ Gradient descent → CG

nCG

DCG

$\frac{n}{n \text{ DCG}}$

Mean Reciprocal Rank

$$MRR = \frac{1}{|S|} \sum_{t \in S} RRL_t$$

t is a topic in set S.

11/10

2/2

2/2

2/2

1/1

B. prof.
RBC
Stat AP.

Relevance Judgement

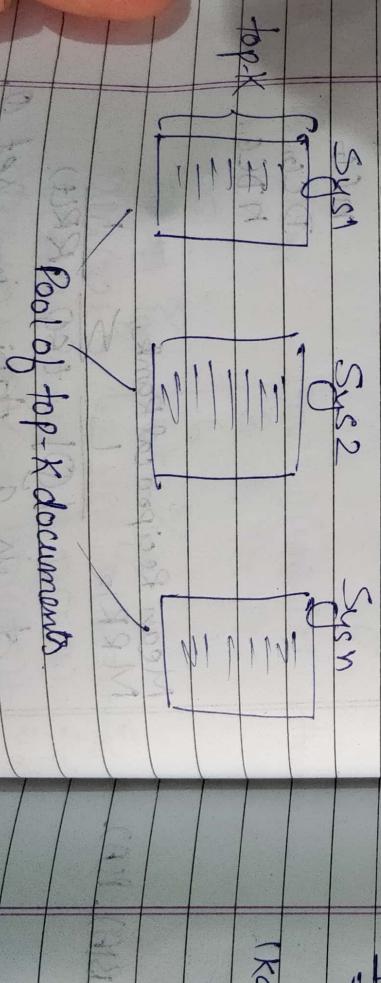
Test bed : 1) A set of documents (D)

2) A set of topics (Q)

3) Relevance Judged info
for each query (t) what are the docs
which are relevant.

So $\times 10,000$ manual labelling of infseal
queries.

Sampling : top-heavy sampling.



Assumptions of Orangified Evaluation

- 1) All the systems are arranging relevant docs in decreasing order of relevance (or preference of the system).
- 2)
 - a) All the systems can collectively identify all the relevant docs in the collection.
 - b) For a given query, within top-k ranks.
- 3) There is no relevant doc outside the pool created based on top-k pooling.

Topic creators are the best assessors.

Inter-assessor agreement

$$\text{Kappa} = \frac{P(\text{Agreement}) - P(\text{Agreement by chance})}{1 - P(\text{Agreement by chance})}$$

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

$$\frac{325}{400} - \frac{\left(\frac{275}{400} \times \frac{300}{400} + \frac{125}{400} \times \frac{100}{400} \right)}{1 - \frac{825}{1600}} = \frac{1300 - 825}{1600 - 825} = \frac{475}{775}$$

Retrieval Model

\rightarrow Boolean
 \rightarrow VSM
 \rightarrow Probabilistic
 \rightarrow Information retrieval

\oplus

Probability

Event

Random Variable \rightarrow It's a function that maps an event to a real number.

$$0 \leq P(X = \dots) \leq 1$$

Joint Probability: $A, B \rightarrow P(A \cap B) = P(AB)$

$$P(AB)$$

$$A, B, C, \dots = P(A \cap B \cap C \dots) = P(ABC \dots)$$

Conditional Probability:

$$P(A|B) := \frac{P(A \cap B)}{P(B)}$$

Occurred
of A

has occurred

$$P(A \cap B) = P(A|B) P(B) = P(B|A) P(A)$$

$$P(A_1, A_2, \dots, A_n) = P(A_1|A_2, \dots, A_n) \times P(A_2|A_3, \dots, A_n) \times \dots \times P(A_n).$$

Bayes' Rule:

A_1, A_2, \dots, A_n are a set of mutually exclusive exhaustive set of events if B is another event then

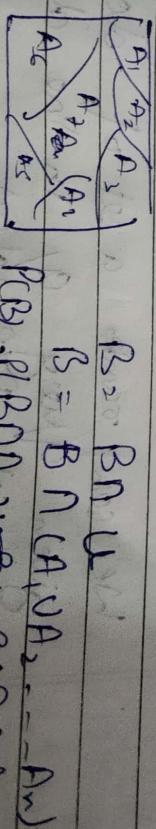
$$P(A_i|B) = \frac{P(A_i \cap B)}{P(B)}$$

$$\begin{aligned} (A \cap B) &= P(B|A_i) \times P(A_i) \\ &\stackrel{i}{=} P(B \cap A_i) \end{aligned}$$

$$ABC \dots \stackrel{i}{=} P(B|A_i) \times P(A_i)$$

$$B = B_1 \cup B_2 \cup \dots \cup B_n$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B|A) P(A) + P(B|A') P(A')}$$



$$P(A) = P(A_1) + P(A_2)$$

$$P(B) = P(B \cap A_1) + P(B \cap A_2)$$

$$P(B|A) = \frac{P(B \cap A_1)}{P(A_1)} + \frac{P(B \cap A_2)}{P(A_2)}$$

odd of an event.

$$O(A) = \frac{P(A)}{P(\bar{A})}$$

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

Probability Ranking Principle :- (PRP)

For any given query (q) we can consider a set of selected ($R=1$) and a set of non-selected docs ($R=0$) in the collection.

Each document will belong to ($R_{\geq 1}$)

With some prob. $P(C, R=1 | d, q)$ And to

$R=0$ with prob. $P(C, R=0 | d, q)$

$$P(C, R=1 | d, q) + P(C, R=0 | d, q) = 1.$$

We can rank a set of docs for a given query q in the decreasing order of $P(C, R=1 | d, q)$.

~~BEST~~ Binary Independence model
(BJM)

A doc in our ordered set of terms
 $\vec{x} = (x_1, x_2, \dots, x_m)$ where each x_i is
0 or 1.

Similarly $\vec{q} = (x_1, x_2, \dots, x_n)$.

$$\begin{aligned} P(R=1 | \vec{x}, \vec{q}) &= \frac{P(R=1 | \vec{x}, \vec{q})}{P(\vec{x}, \vec{q})} = \\ &\rightarrow P(\vec{x} | R=1, \vec{q}) \times P(R=1 | \vec{q}) \\ &= P(\vec{x} | \vec{q}) \cdot P(\vec{q}). \end{aligned}$$

$$P(R=0 | \vec{x}, \vec{q}) = \frac{P(R=0 | \vec{x}, \vec{q})}{P(\vec{x}, \vec{q})}$$

$$P(R=0 | \vec{x}, \vec{q}) = \frac{P(R=0 | \vec{x}, \vec{q})}{P(R=0 | \vec{x}, \vec{q})}$$

$$P(R=1 | \vec{x}, \vec{q}) = \frac{P(R=1 | \vec{x}, \vec{q})}{P(R=0 | \vec{x}, \vec{q})}$$

$$P(R=1 | \vec{q}) = P(R=1 | \vec{q}) \times \frac{P(\vec{x} | R=1, \vec{q})}{P(\vec{x} | R=0, \vec{q})}$$

[x_i 's are occurring independently]

PAGE NO.: _____
DATE: / /

$P(X_{i+1} = R=1, q_i)$
Probability term x_i in Relevant class
is present

$$O(R=1|q_i) \times \prod_{x_{i+1}} P(X_i|R=1) \prod_{x_{i+1}}$$

$$\begin{cases} R=1 & x_{i+1} \\ R=0 & u_i \\ p_i & x_{i+1} \\ u_i & x_{i+1} \\ 1-p_i & x_{i+1} \\ 1-u_i & x_{i+1} \end{cases}$$

$$O(R=1|q_i) = O(R=1|q_i^*) =$$

$$\prod_{x_{i+1}} p_i \times \prod_{x_{i+1}} (1-p_i)$$

$$\prod_{x_{i+1}} u_i \times \prod_{x_{i+1}} (1-u_i)$$

When $q_i = 0$: $p_i = u_i$ (Assumption)

~~$\bullet R=1$~~

$$O(R=1|q_i) \propto \prod_{x_{i+1}} p_i \times \prod_{x_{i+1}} (1-p_i)$$

$$= x_{i+1}, q_{i+1}^{-1}$$

$$\prod_{x_{i+1}} u_i \times \prod_{x_{i+1}} (1-u_i)$$

$$(O(R=1|q_i)) = x_{i+1}, q_{i+1}^{-1}$$

$$R=1 \quad x_{i+1}, q_{i+1}^{-1}$$

$$= O(R \cdot 1/q_f)$$

$$\frac{\prod_{i=1}^n R_i}{x_{i,1} q_{i,1}} \times \frac{\prod_{i=1}^n (1-u_i)}{x_{i,2} q_{i,2}}$$

$$\frac{\prod_{i=1}^n (1-p_i)}{x_{i,1} q_{i,1}} \times \frac{\prod_{i=1}^n u_i}{x_{i,2} p_{i,2}}$$

$$\frac{\prod_{i=1}^n (1-p_i)}{q_{i,1}} \times \frac{\prod_{i=1}^n (1-u_i)}{q_{i,2}}$$

constant independent of document.

$$O(R \cdot 1/d_f q_f) = \text{const.} \times \prod_{i=1}^n \frac{R_i}{1-p_i} \times \frac{1-u_i}{u_i}$$

$$\log(O(R \cdot 1/d_f q_f)) = \sum_{(i,1) \in q_f} \log\left(\frac{R_i}{1-p_i}\right) + \log\left(\frac{1-u_i}{u_i}\right)$$

for a query t/q_f .

	R_2	$R_2 \cdot 0$
x_i present	$\sum_{d_f \in S} d_f$	$d_f \rightarrow \text{no. of documents in which term } i \text{ occurs}$
$x_i = 1$	$\leq -S$	$N - d_f \rightarrow \text{no. of documents in which term } i \text{ does not occur}$
x_i absent	$S - d_f$	$N - d_f \rightarrow \text{Total no. of documents}$
$x_i = 0$	S	$N - S$

To

$$\hat{p}_i = \frac{s}{N}$$

$$u_i = \frac{d_{fi} - s}{N - s}$$

$$C_i = \frac{\hat{p}_i}{1 - \hat{p}_i} \times \frac{1 - u_i}{u_i}$$

$$= \frac{s}{N-s} \times \frac{(N-s-d_{fi}+s)}{(d_{fi}-s)/N-s}$$

$$= \frac{s}{N-s} \times \frac{N-s-d_{fi}+s}{d_{fi}-s}$$

$$s \leq N$$

$$s \leq d_{fi}$$

{ Statistical
assumption }

$$\cdot \frac{s}{N-s} \times \frac{N-d_{fi}}{d_{fi}}$$

\hat{p}_i represent
frequency term.

d_{fi} represent
frequency term.

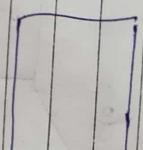
s represent
frequency term.

24

2

To get an idea of Σ .

init q_f



$V \approx 10 - 100$

$VR : VR_i$



$V \approx 10 - 100$

$VNR = V - VR$

V

$$\hat{p}_i = \frac{VR_i + k_2}{VR + k_1}$$

for Smoothing.

Probabilistic Reliance Feedback.

$$p_i^{(k+1)} = \frac{VR_i + \alpha p_i^{(k)}}{VR + \alpha}$$

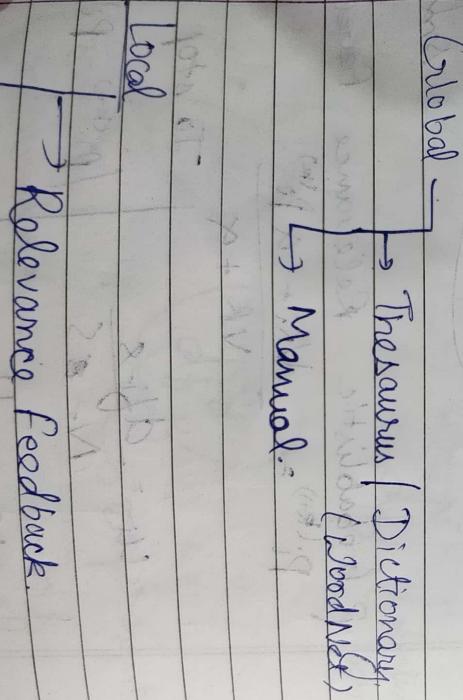
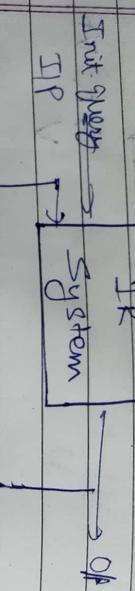
To stop when

$$c_i = \frac{d_{i-1}}{N-1} \quad |p_i^{(k+1)} - p_i^{(k)}| < \text{threshold}$$

$$\propto \frac{d_i}{N}$$

Retrieval Models

Set of docs.



(VSM)

$$\vec{q}_{opt} = \underset{\vec{q}}{\operatorname{argmax}} \left[\text{Sim}(\vec{q}, D_{\text{tr}}) - \frac{\text{Sim}(\vec{q}, D_{\text{te}})}{|D_{\text{te}}|} \right]$$

Modify \vec{q} and find out q_{opt} .

Reactive Relevance Feedback.

$$Q_{\text{mod}} = \vec{q}_{\text{mod}} + \beta \times \frac{\sum_{d_r \in D_r} d_r}{|D_r|}$$

$$= \vec{q}_x + \frac{\sum_{d_r \in D_r}}{|D_r|}$$

α, β, γ tuning param. for estimation
use "n" best in selection docs.

ideally

$\alpha = 1$
 $\beta = 0.6 - 0.85$ should be
 $\gamma = 0 - 0.15$ taken from whole set.

Pseudo-Relevance Feedback (Blind)
Query expansion.

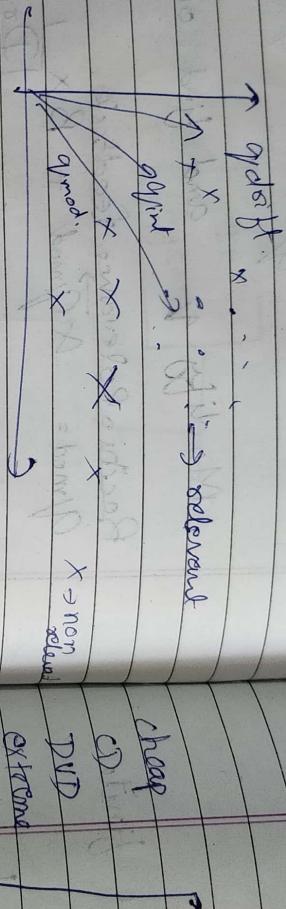


Initial 20 docs (say)

(Rank 1 - Rank 2) - d₁₈

The rest are done.

In pseudo relevance feedback
 query drift ~ If you've modified
 query goes to non-relevant doc
 side.



RF increases recall. \Rightarrow MAPP.
 also it can increase prec.

Relevance Feedback

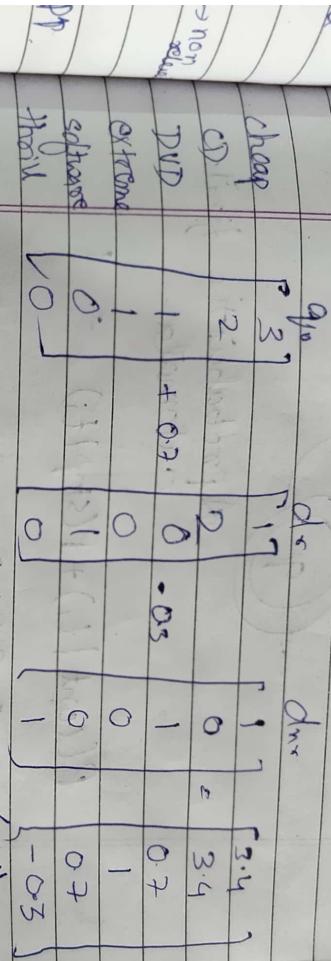
gives cheap DVDs, cheap CDs, extremely
 cheap CDs

$d_1 = \text{CDs, cheap software CDs}$

$d_2 = \text{cheap thrill DVDs, (NR)}$

$$\vec{q}_{mod} = A \vec{q}_{mod} + B \sum_{D_\alpha \in D_\alpha} \vec{d}_\alpha -$$

$$Y \cdot \frac{1}{(D_{n\alpha})} \sum_{D_{n\alpha} \in D_{n\alpha}} \vec{d}_{n\alpha}$$



for parameter α - λ_{mod} $\vec{d}_{n\alpha}$

2) Language Models:

$S = \text{Set of Symbols} = \{a, b, \dots\}$

P = Set of rules.

$S^* = \text{Set of strings}$

R)

$$P(t_1, t_2, t_3) = P(t_1) P(\text{continue} | t_1) P(t_2 | t_1) P(\text{stop} | t_2) P(t_3 | t_2)$$

Q1

A probabilistic finite automata.

$$P(\text{cont} | t_i) + P(\text{stop} | t_i) = 1$$

$$\sum P(t_i) = 1$$

We can rank the document by the probability by which a document generates a query.

Types of Lang Models.

Each document is a model of a universe for now.

Ingram: Consider a particular set of words.

$$P_{\text{in}}(t_1, t_2, t_3, t_4) = P(t_1) P(t_2 | t_1) P(t_3 | t_2) P(t_4 | t_3)$$

Bigram:

$$P_{\text{in}}(t_1, t_2, t_3, t_4) = P(t_1) P(t_2 | t_1) P(t_3 | t_2) P(t_4 | t_3)$$

multi
let
let
let

Multinomial

Probability PAGE NO.:
of document DATE: / /
language language
tokens tokens
very permutation

$$P(d) = \frac{t_1! t_2! t_3!}{(t_1+t_2+t_3)!} \times P(t_1)^{t_1} \times P(t_2)^{t_2} \times P(t_3)^{t_3}$$

find

Query likelihood Model.

$$P(q/d) = P(q/d) \times P(d)/P(q)$$

$$P(q/d) = \frac{K_q}{t_q} \prod_{t \in q} P(t/d)^{t_q}$$

$$\text{Where } K_q = \frac{q!}{(t_1+1)! \times (t_2+1)! \times \dots \times (t_q+1)!}$$

$$\text{by the } P(q/d) = \frac{P(q/d) \times P(d) \times P(q/d)}{P(q)}$$

a universe
for now.

Estimating Probabilities :-

Maximum likelihood estimate :-

$$\hat{P}(q/d) = \prod_{t \in q} \hat{P}_{\text{ML}}(t/d) = \prod_{t \in q} \frac{f_{t,d}}{\sum_t f_{t,d}}$$

Smoothing if $d=0$ then

$$\hat{P}(t/d) \leq c/d$$

$$\hat{P}(t|ld) = \lambda P_{\text{true}}(t|M_d) + (1-\lambda) P_{\text{model}}(t|M_d)$$

Where $0 < \lambda \leq 1$ and M_d is the model built from entire document collection.

Bayesian Smoothing

$$\hat{P}(t|ld) = \frac{f(t|ld) + \alpha P(t|M_c)}{ld + \alpha}$$

- Small lambda \Rightarrow large alpha means more smoothing.

- Low smoothing good for short periods, high smoothing for large periods

Example

$d_1 = \text{Profit}$ reports a profit but revenue is down
 $d_2 = \text{Losses}$ numerous quarters lost but revenue decreases further

11. -

PAGE NO.:
DATE: / /

probable

11
12

10010
00010

but will

10010
00010

01101

11011
11111
11111
11111
11111

12

110
110

X 10333
deposit
a

110
110

profit
but

110
110

lungs

deserve
is

110
110

series,
ies.

clown

110
110

quarrel
harrows

110
110

quarrel
loss
decreases
further

110
110

8

use is due
but

110
110

8

$$P(d_1 | \text{dev}) = \\ [\lambda P(t_1 | d_1) + (1-\lambda) P(t_1 | M_1)] \\ \times [\lambda P(t_2 | d_1) + (1-\lambda) P(t_2 | M_2)] \\ \times P(d_1)$$

$$= \frac{1}{8} \times \frac{1}{8} + \frac{1}{2} \times \frac{2}{16}] \times [[\frac{1}{8} \times \frac{1}{16}] \times \frac{1}{2}] \\ = \frac{1}{8} \times \frac{3}{32} = \frac{3}{256}$$

$$P(d_1 | \text{dev}) = \\ [[\left(\frac{1}{8} + \frac{2}{16} \right) \times \frac{1}{2}] \times [[\frac{1}{8} + \frac{1}{16}] \times \frac{1}{2}]] \\ = \frac{1}{8} \times \frac{1}{32} = \frac{1}{256}$$

So ranking $d_1 > d_2$.

(a) $\lambda = 0.5$

(b) $\lambda = 0.7$