

Goldberg machine: is a mechanical machine that searched for a pattern of dots or letters across Catalog entries stored on a roll of micro film.

Information Retrieval: finding material of an unstructured nature that satisfies an information need from within large collections.

Data extraction: is a process that involves retrieval of data from various source in order to process it.

Information extraction: is the automated retrieval of specific information related to a selected topic bodies of text.

Data mining: is the process of analyzing of data to find patterns, discover trends and gain insight into how that data can be used.

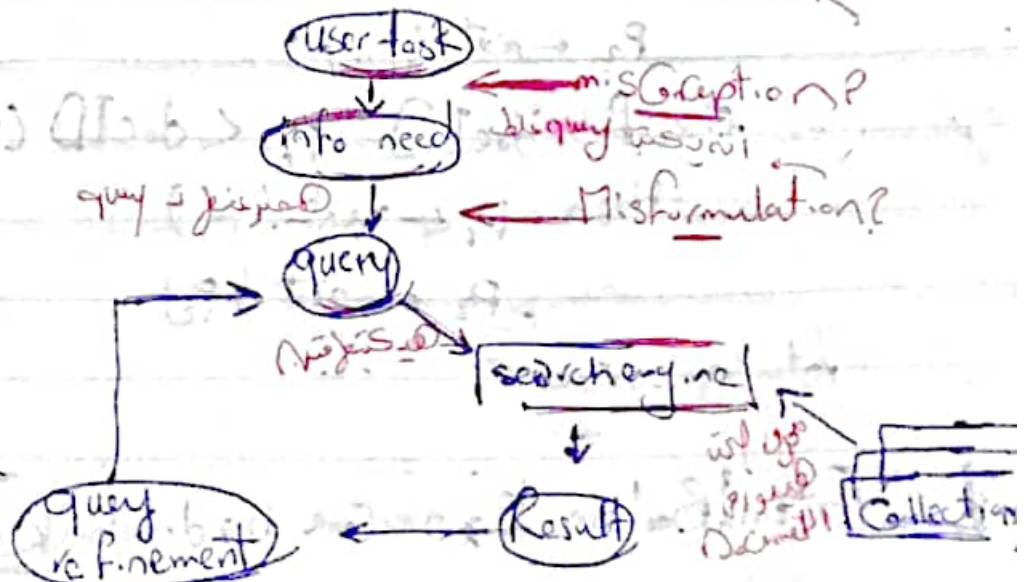
Web mining: is the process of using data mining techniques and algorithms.

web crawler: Called spider is a standalone bot that systematically scans the internet for indexing and searching for content, following internal links on web pages.

web scraper: is a process of extracting specific data.

Unlike web crawling, a web scraper searches for specific information on specific websites or pages.

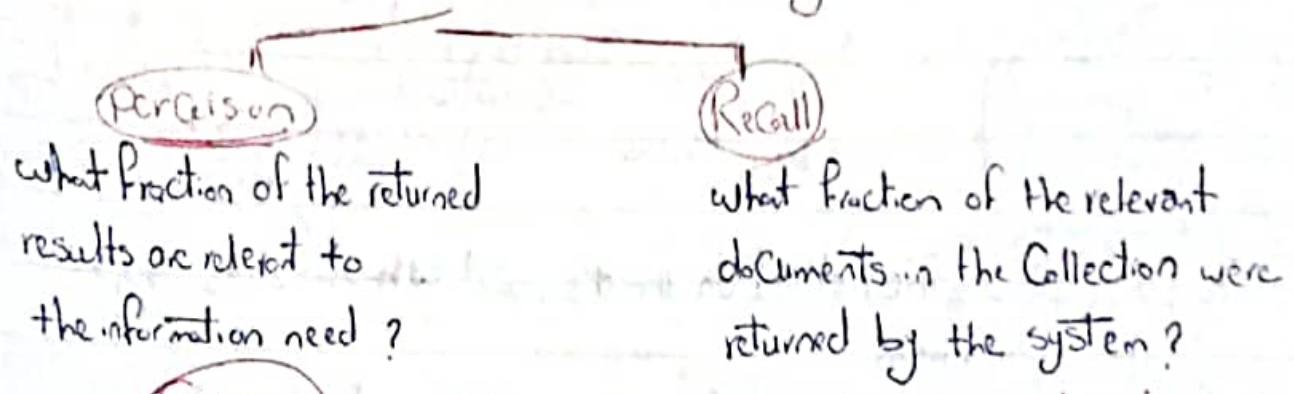
The classic search Model



lec 2

- information need: is the topic about which the user desires to know more
- query: is what the user conveys to the Computer in an attempt to communicate the information need
- Relevance: if it is one that the user perceives as containing information of value with respect to their personal information need.

The effectiveness of an IR system



Grepping → the Unix Command → allow useful possibilities for wildcard pattern matching
 very effective process

Shortfalls of Grepping

- ↳ To process large document collections quickly
- ↳ To allow more flexible matching operations
- ↳ To allow ranked retrieval

Term document incidence matrix

Terms	documents
1	1 1 0 0
2	1 1 0 0
3	1 1 0 0
4	1 1 0 0
5	1 1 0 0

query

• 0 AND 0
 • 10111 AND 110001

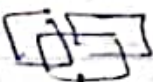
Boolean Retrieval Model: model for information retrieval in which we can pose any query which is in the form of a Boolean expression of terms

Lec 3

inverted index Construction

Processing Document

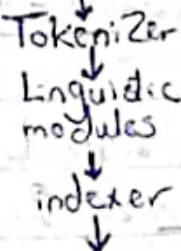
Document to be indexed



Token stream

Modified tokens

inverted index



Friends, Romans, Countrymen

Friends | Romans | Countrymen

Friend | roman | Countryman

Friend → [2] → [4]

roman → [13] → [2]

Countryman → [13] → [16]

tokenizer → Normalization → Stemming → Stopword

Note: Posting list are bounded by the number of Terms

Lec 4

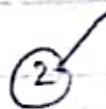
Ex: Brutus and Caesar

Locate Bruts → Retrieval position

Locate Caesar → Retrieval position



Large Two Position



Intersect (P_1, P_2)

answer ← < >

while $P_1 \neq NIL$ and $P_2 \neq NIL$

do if docID(P_1) = docID(P_2)

then ADD (answer, docID(P_1))

$P_1 \leftarrow \text{next}(P_1)$

$P_2 \leftarrow \text{next}(P_2)$

elseif docID(P_1) < docID(P_2)

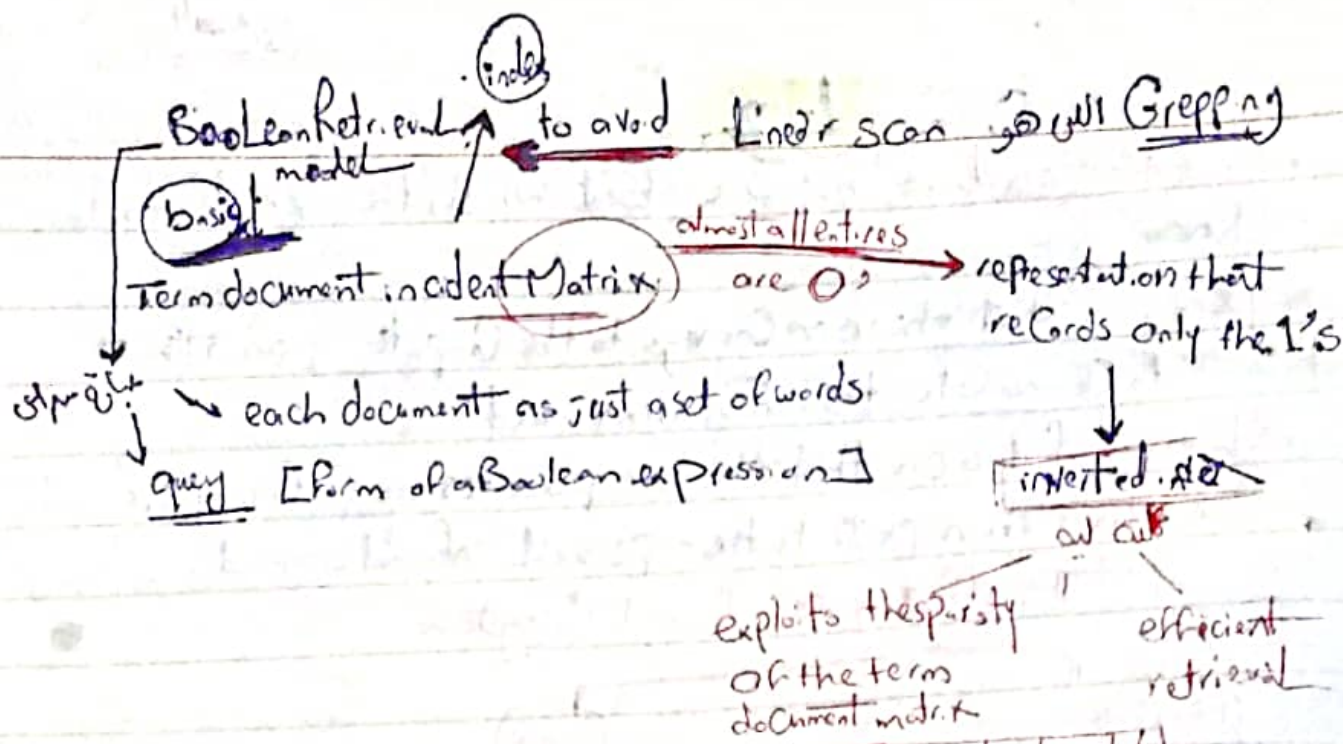
then $P_1 \leftarrow \text{next}(P_1)$

$P_2 \leftarrow \text{next}(P_2)$

return answer

Processing Query

Extend Boolean to overcome the drawback of the Boolean (doesn't term weight is)



Ex) Brutus →

1	2	3	4	5	6
---	---	---	---	---	---

Forward index: is the list of documents, and which words appear in them

inverted index: is the list of words, and the document in which they appear

westlaw → 13 within 3 words
 → 15 in the same sentence

Query optimization
 → best order for query processing (start with Smallest)

Lec 5

Phase Queries → want to be able to answer Stanford University
 (notably) I want to University at Stanford

first attempt

Biword indexes

each of biwords is a dictionary term

Problem Long phase queries
Stanford University palonito

Processed by breaking down

→ False Positives

• bigger dictionary

• not the standard solution

Combination Scheme

Solution 2

Positional indexes

= store for each term the positions

Positional intersect (P_1, P_2, k)

answer $\leftarrow \langle \rangle$

while $P_1 \neq \text{NIL}$, $P_2 \neq \text{NIL}$

do if $\text{docID}(P_1) = \text{docID}(P_2)$

then $L \leftarrow \langle \rangle$

$PP_1 \leftarrow \text{positions}(P_1)$

$PP_2 \leftarrow \text{positions}(P_2)$

while $PP_1 \neq \text{NIL}$

do while $PP_2 \neq \text{NIL}$

do if $|\text{Pos}(PP_1) - \text{Pos}(PP_2)| \leq k$

• then $\text{ADD}(L, \text{Pos}(PP_2))$

elseif $\text{Pos}(PP_2) > \text{Pos}(PP_1)$

then Break

• $PP_2 \leftarrow \text{next}(PP_2)$

while $L \neq \langle \rangle$ and $|\text{L}[0] - \text{Pos}(PP_1)| > k$

• delete $(L[0])$

foreach $P \in L$

do $\text{ADD}(\text{answer}, \langle \text{docID}(P), \text{Pos}(PP_1), \text{Pos} \rangle)$

• $PP_1 \leftarrow \text{next}(PP_1)$

$P_1 \leftarrow \text{next}(P_1)$

$P_2 \leftarrow \text{next}(P_2)$

elseif $\text{docID}(P_1) < \text{docID}(P_2)$

then $P_1 \leftarrow \text{next}(P_1)$

elseif $P_2 \leftarrow \text{next}(P_2)$

return answer

queries have all been Boolean

Document match or not

Problem

query 1: Stanford user drink 650 → 200,000 hits

query 2: Stanford user drink 650 no and found → 0 hit

Ranked retrieval → yes

system returns
an ordering over the top
documents in the collection
with respect to query

Free text query

* Fast or performance not a problem in ranked retrieval?

because produces a ranked result

* In the basic Ranking algorithm

document, query → Score

start with One term

document: Stanford
Score = 0

most frequent
higher the score

The Jaccard Coefficient

→ measure of overlap of two sets

$$Jaccard(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

issues

→ doesn't consider term frequency

are more information than frequent terms — rare terms

normalized =
$$\frac{|A \cap B|}{\sqrt{|A \cup B|}}$$

Bag of words model

vector representation → problem → doesn't consider the ordering of words

Ex: John is quicker than Mary = Mary is quicker than John

→ the same vector

able to distinguish these 2 documents

Positional index → fix

Term Frequency (TF)

the number of times that t occurs in d

Relevance doesn't increase with term TF

Use in Score

لا يزيد من مصداقية الوثيقة مع تكرارها
فقط اشارة

$$w_{t,d} = 1 + \log(t.f) \quad t.f = 70.1$$

$$\text{Score} = \sum (1 + \log_{10} t.f)$$

Lec 7

df_t is the document frequency of t : the number of documents that contain t

no effect on Ranking one term queries $\leftarrow idf = \log\left(\frac{N}{df_t}\right)$

to dampen the effect of idf

- Collection Frequency of t is the number of occurrence of t in the Collection, Counting multiple occurrence

$$w_{t,d} = (1 + \log t.f) \cdot \log\left(\frac{N}{df_t}\right)$$

increase — # occurrences within a document
rarity of term in the Collection

$$\text{Score}(q, d) = \sum t.f \cdot idf$$

Queries as vector \rightarrow represent query as vector in the space

Euclidean distance

bad idea $\sqrt{x^2 + y^2}$

Because Euclidean distance is large for vectors

of different lengths

Use angle

angle between 2 doc is 0 corresponding to maximal similarity

Rank decrease

kind of inverse of angle

Cosine Rank increase

decreasing function for interval $[0, 180]$

q_t is the $t.f \cdot idf$ weight of term t in the query

$$\cos(q, d) = \frac{q \cdot d}{\|q\| \|d\|} = \frac{\sum q_t \cdot d_t}{\sqrt{\sum q_t^2} \sqrt{\sum d_t^2}}$$

Cosine Similarity

Lec 8

باب ستراتیژی
اسماء و محمد و یونس

• Term Frequency ($1 + \log tf$)

• Document Frequency ($\log \frac{N}{df}$)

• Normalization

$$\begin{array}{ccccccc} \textcircled{1} & \textcircled{2} & \textcircled{3} & \textcircled{4} & \textcircled{5} & \textcircled{6} & \textcircled{7} \\ 1 + \log(tf) & df & \log\left(\frac{N}{df}\right) & \textcircled{5} & \textcircled{6} & (w) & \text{normalize} = \frac{w}{\sum w} \\ & & & & & \frac{1}{\sum w} & \end{array}$$

• Computing Cosine Scores

Cosine SCORE (q)

Float Scores [N] = 0

Float Length [N]

for each query term t

do calculate $w_{t,q}$ and Fetch postings list for t

for each pair (d, $tf_{t,d}$) in Postings List

do $Scores[d] += w_{t,d} \times w_{t,q}$

Read the array Length

for each d

do $Scores[d] = Scores[d] / Length[d]$

return Top k Components of Scores[]

Lec 9

- Evaluation of a result (Measuring relevance)

if we have ① benchmark document Collection

② benchmark set of queries

③ assessor judgments of whether documents are relevant to queries

Then we can use Precision, Recall, F-measure

- Evaluation of ranked results \rightarrow the system can return any number of results
 \rightarrow By taking various numbers of the Top returned documents (level of recall), the evaluator can produce a precision-recall curve.

Measure For a search engine Q How fast does it index

② How fast does it search

③ UI

④ Free

⑤ Expressiveness of query Language

How do you tell: Users are happy?

↳ ① Search returns product's relevant to users

b(2) search results get clicked alot

b) users buy after using the search engine

↳ 4 Repeat visitors / buyers

precision: fraction of retrieved docs that are relevant

Recall: Fraction of relevant docs that are retrieved

	Relevant	Non-relevant
Retrieved	TP True Positive	FP
Not Retrieved	FN	TN

$$\text{Precision} = \frac{TP}{TP + FP}$$

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

$$F_{\text{measure}} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

Map is Average precision across multiple queries / rankings

relevant document

Average precision

Recall (17)

precision

1

2

3

90

4

0

0

2

1

0

Nera
Dici

92

100

$$\frac{0}{6} = \text{average Precision}$$
[illegible]

Map [mean average precision]

$$\left(\frac{P_{rank_1} + P_{rank_2}}{2} \right) \cdot 2$$

- if relevant document never gets retrieved, we assume the precision

Corresponding to that relevant doc to be Zero

- Map is macro-averaging: each query Counts equally

- lec 10

Basic crawler operation

- Begin with known seed URLs

- Fetch and Parse them

 - Extract URLs they point to

 - place the extracted URLs on a queue

- Fetch each URL on the queue and repeat

Complications

Malicious Pages

- Spam pages

- Spider traps

Spider trap

non-malicious pages pose challenges

- bandwidth to remote servers vary

- web masters stipulations

- site mirrors and duplicate pages

it is a set of web pages that may intentionally or unintentionally be used to cause a web crawler or search bot to make an infinite number of requests or cause a poorly constructed crawler to crash.

What any crawler must do

Be Robust

Be Polite

Respect implicit

and explicit politeness

Considerations

no specification

الرجوع إلى صفحة
الرجوع إلى صفحة
الرجوع إلى صفحة