

# Web Information Retrieval

*Eduardo Puglisi*

Web Information Retrieval = retrieve data/information from web resources

I.R. = retrieve information from a huge collection of unstructured docs (text) and analyze them  
An example is what Google does during a web search: within millions of sites/articles it shows the ones related to user's query

What we always have:

- 1) huge collection of docs
- 2) Information need: what the user is looking for

On Web there are some factors that make things harder

- links
- huge (very) and distributed sources
- Different kinds of files
- Web keep growing
- not homogeneous sources

#### COURSE OUTLINE

- collect a **Web corpus** (multiple documents)
- Pre-processing and organizing it
- Search for a document using a query
- Use web as platform to demand services

#### COLLECT A WEB CORPUS

Recursively explore all web pages iterating on links and collect informations: think of web as a graph with its necessities

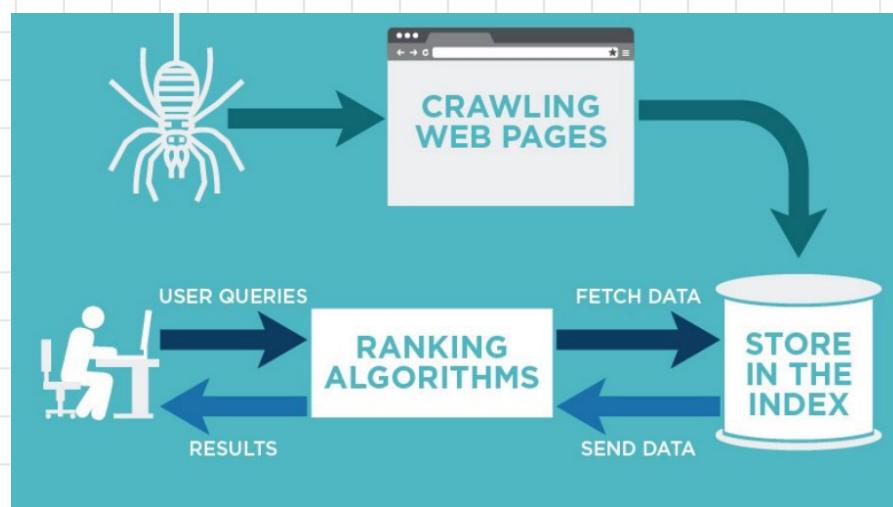
This takes problems such as huge space requirement, computing power and competitiveness (if we use more machines for the search). Need to be sure to visit "all" pages to avoid information missing (**Bias in data**)

#### ORGANIZE A WEB CORPUS

Main goals: efficiency and relevance at some time

- **Crawling**: visit web corpus graph's nodes
- **Indexing**: organize nodes → efficiency
- **Search Algorithms**: how to solve user's information need using query → relevance (only on runtime online)

Relevance must be compared with **authority** of the page: pages must give user real information

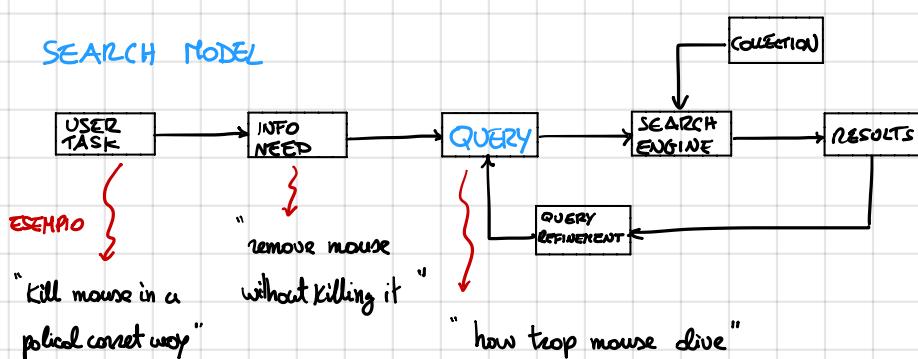


**META-SEARCH**: research using combined results from multiple search-engines

**COLLECTION**: set of documents, not always static

↳ goal: collect documents that are relevant to user's information need

### SEARCH MODEL



How to evaluate results:

• **Precision**: how many documents over the total are considered relevant to the user's information need

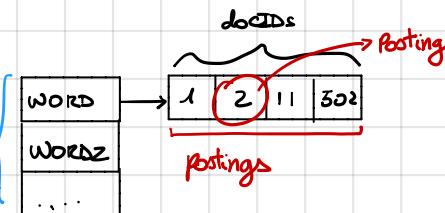
• **Recall**: how many documents, over the relevant ones, are shown (ex. First 10 links on Google)

More precision means less recall

Query based research use presence (or not) of keywords in documents: to do so we use **Inverted Index**: for every word we create a linked list with docIDs of documents that contain it.

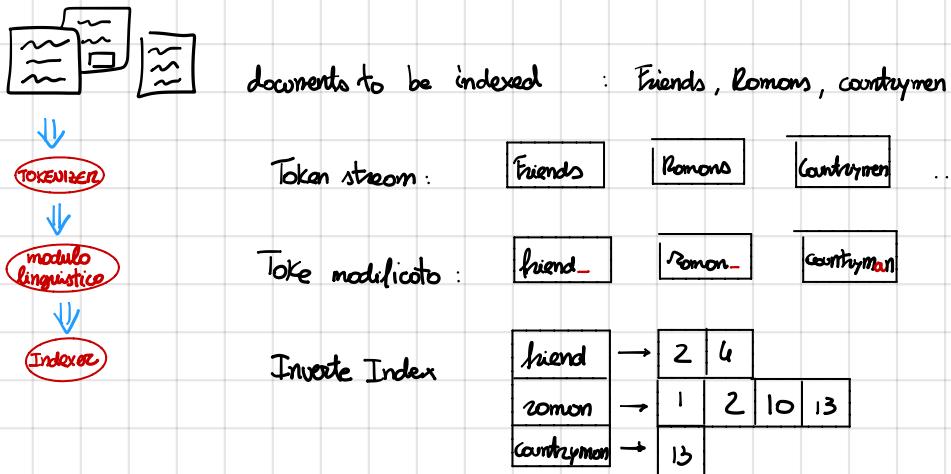
Inverted index save space and time:

- every document has a **docID** that identify it
- avoid using fixed length array      **hasmap / dictionary**
- docIDs are sorted

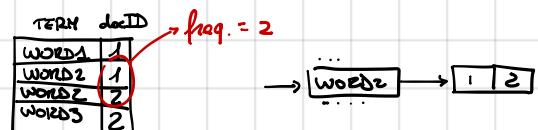


A "normal" array, actually a matrix, with docIDs as keyword and words as postings is way bigger and sparse

How to build on inverted index:



Note: after tokenization, we sort by word first then according to docIDs  
save doc frequency of every words and then sort postings.



Ex.

words1 AND words2 : locate words1 and words2 posting and intersect them . Sorted postings guarantee more efficiency.  
Complexity:  $O(x+y)$  with x and y as postings lengths

This is an example of **Boolean Query** which follow the boolean query model : AND , OR , NOT (with other tools)  
It's very precise: does match condition or not.

AND operand is very precise but recall is very low while OR has high recall and low precision making it hard to find a middle ground.

## QUERY OPTIMIZATION

Consider a query that is an AND of n-terms.

Brutus AND California AND Caesar → (Brutus AND California) AND Caesar

Intersect the two smallest postings list grants next intersection to be more efficient (less postings to check)

In cases like (mobbing OR crowd) AND (ignoble OR strike) AND... we use frequency to estimate which OR must be first: the one which terms's frequency sum is lower

If we want a query to match "Stanford university" or phrase <term: docs> we no longer suffices cause we need a **bi-term** combining the two words.

- friends romans

We could index pair of terms : Friends, Romans, Countrymen → "Romans countrymen" → too many dictionaries and false positive

We could save the position of term in docs :

- to : 2: [1, 17, 74] ; 4: [8, 16, 190, 429, 433] ; 7: [13...]
- be : 1: [17, 19]; 4: [17, 191, 291, 430, 434]; 5: [14, 19...]

So we can find docs in which we have "to be" with an extra iteration. Using position we can even do proximity searches.

But this use more space.

These methods can be used anyway in some cases and also combined, for example "name surname" can use bi-word method (ex search "Michael Jackson")

## QUESTIONS :

- which format ? pdf, word, excel ...
- which language ?
- which character set ?



**classification problems**: to identify the class of docs.

An index can contain docs with different format and language

## TOKENIZATION

**INPUT**: "Friends, Romans and Countrymen"

**OUTPUT**: • Friends • Romans • Countrymen

A Token is an instance of grouped together as a useful semantic unit for processing : depends on how we want to process the index (by words, bi-words etc)

After further processing they become index entries

Most simple strategy is splitting on white spaces , but "Finland's capital"? "Finland" and "s" ? "Finlands" ? "Finland's" ?

"San Francisco" is one or two tokens?

Similar problem with numbers : 3/20/91 and 20/3/91 is the same date but with different format → two different token or add "format" as meta-data to index entry.

Language: French → "L'ensemble" one or two token? → L? L'? le? → how to associate L'ensemble with un ensemble

Suppose we created the index: we now remove stop words like "the", "a" ... and create a separate list for them.

We need to normalize words: "U.S.A" is equal to "USA" → unique dictionary entry.

Most of time users doesn't use accent also and we have to know how they like to write queries.

Synonyms and homonyms? car = automobile: we can search for car too when query contains automobile or rewrite to form equivalence-class terms indexing automobile under "car - automobile".

Need to take care of spelling mistakes and so on...

Problems are a lot!

Lemmatization: reduce to base form . am, are, is → be

Stemming: reduce term to root before indexing. automobile, automatic, automotic → automot

↳ Algorithm: Porter's algorithm

- Mes → ss
- ies → i
- ational → ate
- tional → tion

- words weight sensitive
- ( $m > 1$ ) EMENT:
  - replacement → replac
  - cement → cement

} Typical rules (there are a lot of stemmers)

Back to "AND query" we can improve merge efficiency using jumps in postings list:

- [2] [4] [8] [61] [48] [64] [128] we have to take common docID so we check one-by-one and confront them
- [1] [2] [3] [8] [11] [17] [21] [31] when we reach 8 we match it and advance: we now have 11 and 11, but 11's successor is 31 which is smaller than 41 so we can skip to 31

Where do we put skips? Find a tradeoff.

## DATA STRUCTURES

We use dictionary for inverted index. But which data structures we use?

Most common choices are hash tables and trees

### Hashtable:

Pro: faster than tree finding postings list

Cons: - No prefix search **Tolerant retrieval**

- No easy way to find variants: judgement / judgment

- Need rehash if vocabulary keep growing

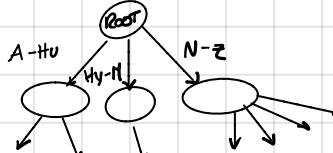
### Tree

More usual type is B-tree: self balancing tree which nodes have a number of child between  $[a, b]$  with a and b variable.

Pro: solve prefix problem

Cons: - Slower  $\mathcal{O}(\log M)$

- Rebalance is expensive



### WILD-CARD QUERIES

$mon^*$ : find all docs containing words beginning with "mon". Easy in B-tree

\* mon is way harder.  $\rightarrow$  use an alternative tree for terms backwards

$x^*y$  is an intersect of  $x^*$  AND  $y^*$   $\rightarrow$  very expensive

Now we look up the posting lists.

More efficient way: **permutation index**

Hello index under: hello\$, clo\$fh, clo\$he, clo\$hl, clo\$ll, \$hello (permutations)

- x lookup for  $x\$$        $\cdot x^* \rightarrow \$x^*$        $\cdot x^*y \rightarrow y\$x^*$

-  $x \rightarrow x\$$        $\cdot x^* \rightarrow x^*$

But quadruples size!!!

$x^*y^*z \rightarrow z\$x^* \rightarrow$  check every term if contains Y

### BIGRAM (K-gram) INDEXES

We maintain an index for bigrams (or K-grams)

\$m  $\rightarrow$  all words that start with "m"

mo  $\rightarrow$  all words that contain "mo"

$\rightarrow$  mon\*  $\rightarrow$  \$m AND mo AND on

we get "moon" too but it isn't a selection  $\rightarrow$  use post filters

### Spell Correction

- Correct indexed document }  $\rightarrow$  check isolated word or context in which it is used  
- Correct user query

A language specific way is to use **Soundex**: phonetic equivalents.

There are many ways to correct terms, for example we can search for terms with many K-grams in common with the query. (using k-gram index), we can calculate the "distance" of the query with terms in dictionary (Levenshtein distance).

## BSBI ALGORITHM

The steps we saw previously to create and sort an inverted index are done in memory (RAM) but this causes many problems:

- postings lists are incomplete until the end
- low space

Use disk (hard disk)? Too slow! We use it only to store intermediate results

We use then an **external** algorithm which uses very few **disk seeks** (operations that involve disk e.g. take back intermediate results)

E.g. To sort 100M postings we define 10 **blocks** of 10M postings. We can easily fit a block in memory.

IDEA: for each block: ① accumulate postings ② sort in memory ③ write to disk  
then merge the blocks

This is called **Blocked Sort-Based Index (BSBI)**

This algorithm uses one more dictionary to map terms with termID (like docID) to reduce size

## SPINI ALGORITHM

BSBI use a dictionary to map terms to termID that grow dimensionally and this is a problem

**Single Pass in-memory index** solves it not using it

- Separated dictionary for each block ] → we create an inverted index for each block and then merge them
- Don't sort postings (they are already sorted anyway cause docIDs are incremental when received)

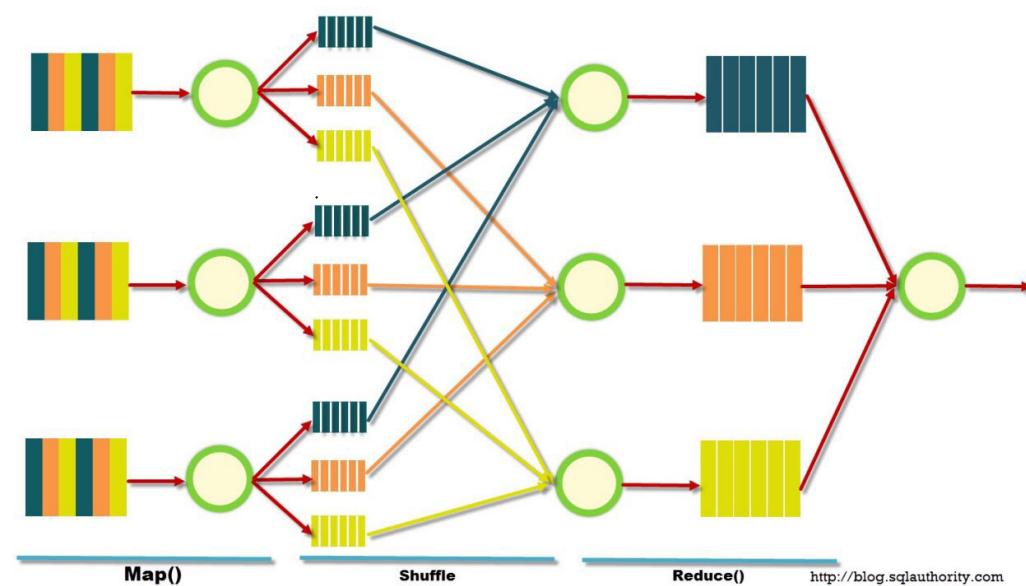
For web-scale indexing we (must) use distributed computer cluster. How?

Maintain a **master** machine considered "safe" which divide indexing into sets of **parallel** tasks assigned to other machines.

**MapReduce** is a framework that let task run in parallel

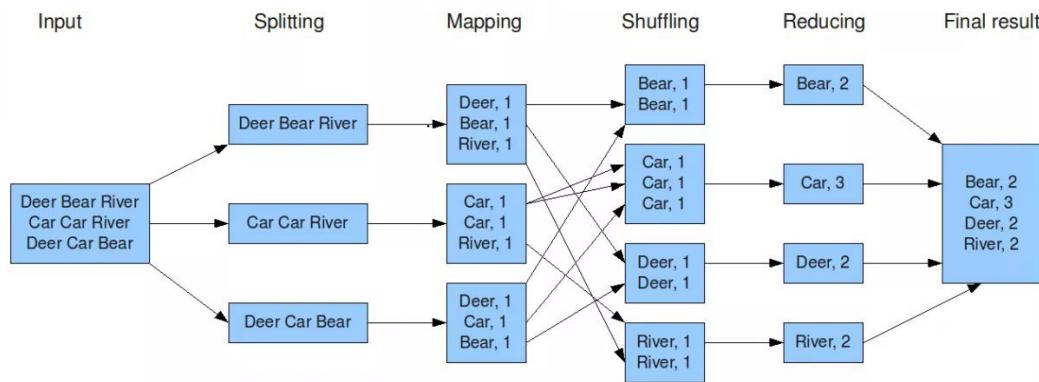
- Map: transforms an huge collection of documents (input) into a **stream** of **<key,value>** pairs, then it will group them by key (shuffle) (output)
- Reduce: takes in input Map's output and gives as result a set of distinct key and their corresponding combined value (any mathematical function can be used)

How MapReduce Works?



## Particular example

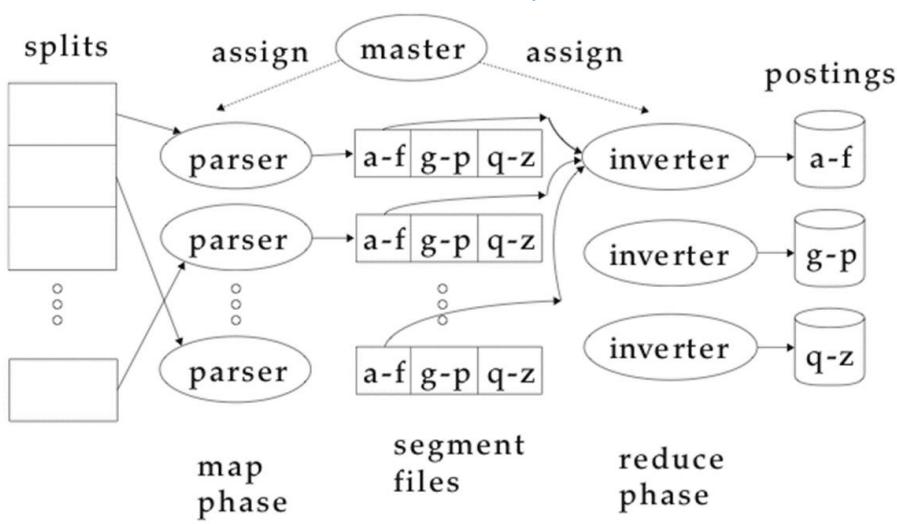
The overall MapReduce word count process



MapReduce defines two type of machine : **partitions** and **inverters** and break document collection into **splits** (same as block we saw in BSB1 / SP11)

- Partitions: read file and emits  $\langle \text{term}, \text{docID} \rangle$  pairs grouped in position of  $j$  terms
- Inverters: collect these pair (only one group), sort and writes to posting list

↓



We are assuming that collections are static which it's not true: posting must be **dynamically** modified.

The approach is simple :

- main big index on disk
  - new docs go to small auxiliary index in memory
  - search across both, merge results
  - periodically merge auxiliary to main
- Problems? Lot of merge and poor search performance during them

**Logarithmic merge**: amortizes the cost of merging index over time. It uses series of index each twice longer of the previous one. Smallest ( $z_0$ ) in memory and others ( $l_0, l_1, l_2, \dots$ ) on disk: when the last one is full we create on other one. When  $z_0$  is full, write it on  $l_0$  and empty  $z_0$ . When  $z_0$  is full again and  $l_0$  is too small, merge them in  $l_1$  and empty both. Then again: we full every index before  $l_1$  and then merge them together in  $l_1, \dots$  And so on.

It use a binary number to save which index is full: e.g. 1011 means  $l_2$  has space, others are full.

$O(\log T) = \# \text{ indexes}$  with  $T = \# \text{ postings}$  read so far  $\rightarrow$  query needs to merge  $O(\log T)$  indexes  $\rightarrow O(T \log T)$  because each of  $T$  postings is merged  $O(\log T)$  times, more efficient than  $O(n^2)$  of auxiliary index method

## COMPRESSION:

- Prefacing decompression algorithms are fast we want to compress index and posting lists to speed up the process e.g.  
Keep index in memory and save space on disk.
- **Lossy**: discard some info. (e.g. stop words) → what we "used" since now!
  - **Lossless**: all info is preserved.

**Heaps' law**: relationship between collection size and vocabulary size (empirical law, but works in general)  
 $\Rightarrow 44 \cdot N^{0.49} \sim \# \text{terms}$  ( $N = \# \text{tokens}$ )  $\Rightarrow$

**Zipf's law**: we have vocabulary size now → we want to know how which terms are more frequent or rare.  
 $\Rightarrow i^{\text{th}}$  most frequent term has frequency  $cf_i$  proportional to  $\frac{1}{i}$  ( $cf = \text{collection frequency}$ )  
NOTE: more rare a term in a query is, more relevant it is for the document  
 $cf_2 = \frac{1}{2} cf_1, cf_3 = \frac{1}{3} cf_1 \dots$

## DICTIONARY COMPRESSION

Dictionary as array of fixed-width entries need 11.2 MB for Reuters (1 year articles)

(20 bytes for term + 4 bytes for frequency + 4 bytes for pointer to posting list)  $\cdot 400,000 = 11.2 \text{ MB}$

Problems: average length is 8 → space wasted

cont handle very long terms.

- Dictionary as string: bytes for frequency and pointers are still the same, but with only 3 byte we have a pointer to the position of term in dictionary string: e.g. term1term2term3 ...  
+ 8 bytes (on average) for term in string:  $(4+4+3+8) \cdot 400,000 = 7.6 \text{ MB}$
- Dictionary as string with blocking: same as before but the dictionary-string is divided in blocks of k terms. We use an extra byte before every term to indicate the length: e.g. 4term9otherterm ...  
 $k=4$ : without blocks we use 3 bytes  $\times 4$  (terms) for pointers; now we use 3 bytes for 1 pointer + 4 bytes for 4 length indicators  $\rightarrow (4 \cdot 3) - (3+4) = 5$  bytes saved per block  $\rightarrow (400,000/k) \cdot 5 = 0.5 \text{ MB}$  total  
 $\rightarrow 7.6 \text{ MB becomes } 7.1 \text{ MB}$

NOTE: lookup of a term is slightly slower with blocking

It's possible to add **front coding**: e.g. 8 automata 8 automata 9 automatic → 8 automat \* a 1 ~~o e z o ic~~  
 $\rightarrow 5.9 \text{ MB}$

## POSTINGS COMPRESSION

Posting lists total size is about 10 times larger than total dictionary.

docIDs are stored in increasing order.

Gops: one idea is to store gops instead of docIDs e.g. COMPUTER → 283154, 283155, 283202  $\Rightarrow$  COMPUTER → 283154, 5, 43 ..

gops for frequent terms are much smaller than for rare terms so we have to make a division.

- Variable byte (VB) code:

docID	824	823	$\left[ \begin{array}{c} \\ \\ \end{array} \right]$	824 in binary is 1100111000 → we use 7 bits "blocks" + 1 bit
gop		5	$\left[ \begin{array}{c} \\ \\ \end{array} \right]$	$\xrightarrow{\text{gop}=5}$ 00000 <u>110</u> 10111000 0 and 1 are continuation bit
VB	00000110 10111000	10000101		$\xrightarrow{\text{gop}=43}$ 1100111000 = 824 0 means that the number continues in next block

- Other VB: we can use other types of alignment: 32 bit (words), 16 bits, 4 bits (nibbles) etc.

- Gamma codes (bitfield code)

- Unary code:  $n = n$  times 1 with one zero at the end → 3 = 1110

Gamma code uses **length** and **offset** of a gop G:

- offset = gop in binary without first bit e.g. 13 = 1101 → 101

- length = offset length e.g. 13 → 101 → 3

Then encode length in unary (1110)

Gamma code for 13 is concatenation of length in unary code and offset: 1110101

Length of offset is  $\log_2 G$  bits, length of length is  $\log_2 G + 1$  bits  $\rightarrow$  total =  $2 \cdot \log_2 G + 1$  bits, 2 times the optimal encoding length.

This doesn't depend on distribution of gaps and it's universal and parameter-free.

Anyway machines have boundaries: 8, 16, 32 bit and this kind of compression and manipulation can be slow, making VB code a better choice sometimes.

## RANKED RETRIEVAL

Boolean queries are too difficult to most of the users (not user-friendly).

They are not very optimized (too few or too many results: **famine or feast**).

With ranking we can give as result top documents, the most relevant ones.

To do so we assign a score to each doc-query pair in  $[0, 1]$  e.g. one-term query if doc doesn't contain it  $\rightarrow$  score = 0, more frequent the term is in a doc, the higher the score will be.

There are a lot of alternatives.

**Jaccard coefficient:** A and B are two sets  $\neq \emptyset \rightarrow \text{JACCARD}(A, B) = \frac{|A \cap B|}{|A \cup B|}$

e.g.  $Q = \underline{\text{idles of March}}$   $D = \underline{\text{Coerce died in March}} \rightarrow J(A, B) = \underline{1/6}$

Prob: doesn't take frequency in consideration.

**Bags of words:** each document is represented as a count vector with frequency for every term.

eg.	Anthony and Cleopatra	The tempter	Prob: no order (John is quicker than Mary = Mary is quicker than John)
Anthony	157	0	$tf_{t,d} = \text{term freq. of term } t \text{ in doc. } d$
Mary	2	3	How do we use it? A doc with $tf = 10$ is more relevant than a doc with $tf = 1$ BUT NOT 10 TIMES MORE RELEVANT (not proportional)

We add a **weight**:  $w_{t,d} = \begin{cases} 1 + \log tf_{t,d} & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$

Score for a doc-query pair:  $tf\text{-matching-score}(q, d) = \sum_{t \in q \cap d} (1 + \log tf_{t,d}) = \text{sum of weights of all terms which appear in both query and doc.}$

Rarity is extremely informative  $\rightarrow$  we want to use frequency of term **in the collection** too  $\rightarrow$  we want **high weights for rare terms** (higher than frequent terms)

$df_t = \# \text{ of docs in which term } t \text{ occurs}$

$idf_t = \log \frac{N}{df_t} = \text{measure of informativeness}$  ( $N = \# \text{ doc in collection}$ ), inverse docs frequency

$\Rightarrow \text{Score}(q, d) = \sum_{t \in q} tf_t \cdot idf_t$

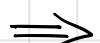
$idf$  affects ranking of docs for query with at least two terms e.g. "arachnacentric line":  $idf$  increases weight of "arach." and decreases weight of "line".

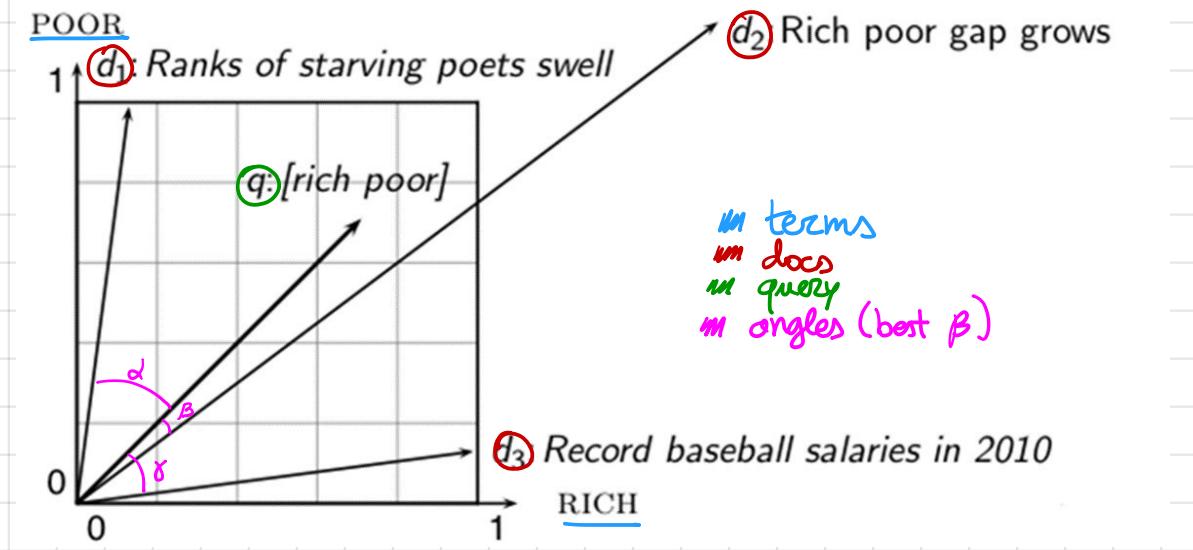
The best known weighting scheme is  $tf \cdot idf \rightarrow w_{t,d} = (1 + \log tf_{t,d}) \cdot \log \frac{N}{df_t}$ . It increases with number of occurrences within a doc and within rarity in collection.

We now modify the previous matrix using  $tf \cdot idf$  weight instead  $\rightarrow$  docs are now represented as real-valued vectors in high-dimensional space (terms axes).

To make it more complex we do the same for query and rank docs according to their proximity to the query. Proximity = Similarity.

To compare two (or more vectors) the best idea is to check the angle between them  $\rightarrow 0^\circ$  = same content





Convert angles to cosine. But how to compute it?

$$\cos(\vec{q}, \vec{d}) = \sin(\vec{q} \cdot \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| |\vec{d}|} = \frac{\sum_i q_i d_i}{\sqrt{\sum_i q_i^2} \sqrt{\sum_i d_i^2}}$$

two normalization

$$\begin{cases} q_i = t\text{-tf-idf of term } i \text{ in query} \\ d_i = t\text{-tf-idf of term } i \text{ in doc} \\ \|q\| \text{ and } \|d\| \text{ are } \vec{q} \text{ and } \vec{d} \text{ lengths} \end{cases}$$

Often we use different weightings for query and docs.

Most used combo is Inc. Itn (notation ddd.ooo) that stand for:

- does: logarithmic tf (l), no dtf weighting (n), cosine normalization (c)
  - query: logarithmic tf (l), idf (t), no normalization (n)

## Computing cosine scores

### COSINESCORE( $q$ )

1 float Scores[N] = 0

2 float Length[N]

3 for each query term  $t$

4 do calculate  $w_{t,a}$  and fetch postings list for  $t$

for each pair( $d$ ,  $tf_{t,d}$ ) in postings list

**do**  $Scores[d] \leftarrow w_{t,d} \times w_{t,a}$

## 7 Read the array *Length*

8 for each  $d$

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

10 **return** Top  $K$  components of Scores[]

Since we are not able to read user's mind, we can accept some error: taking K docs similar to top K docs.

How: find a set  $A$  of docs with  $k < A \ll n$ ,  $A$  does not necessarily contain top  $k$  docs but has many docs from top  $K$   $\rightarrow$  return top  $K$  in  $A$

How: find a set  $A$  of docs with  $k < A \ll N$ ,  $A$  does not necessarily contain top  $k$  docs but has many docs from top  $K$   $\rightarrow$  return top  $K$  in  $A$

Basic algorithm consider docs containing at least one query term → improve: consider only high-idf query terms and docs containing many query terms.

## CHAMPION LISTS

For every term  $t$ , compute  $r$  docs of highest weight  $\rightarrow$  champion list at index build time.

$r$  can be  $>$  than  $K$ !

At query time compute scores only for docs in champion list  $\rightarrow$  Pick the  $K$  top amongst these.

We want top-ranking docs to be both relevant and authoritative: relevance is modeled by cosine scores, authority is **query-independent!** (let's call it  $g(d)$  function)

## NET SCORE

$$\text{Net-score } (q, d) = g(d) + \cos(q, d)$$

It can use other linear combination  $\rightarrow$  now we seek the top  $K$  docs by net-score.

$$f(d, q) = \alpha g(d) + \beta \cos(q, d) \quad \text{with } \alpha + \beta = 1 \text{ in real world}$$

How i get top  $K$  by net-score? Order all postings by  $g(d)$ , in this way (usually) top  $K$  docs are in the first ones and so we can stop computation early

- Champion list +  $g(d)$ -ordering: each term has a champion list of  $r$  docs ordered by  $g(d) + tf \cdot idf_{t,d}$
- High and low list: each term has an high list (champion list) and a low one (others), when traversing postings only traverse high list. If we get more than  $K$  docs, select the best  $K$ , if not use the low list too.
- Impact-ordered postings: we compute only docs with  $wf_{t,d}$  high enough  $\rightarrow$  order every posting list by  $wf_{t,d}$   $\rightarrow$  all posting lists are now in different order. (a doc can be very relevant for a term and irrelevant for another)  $\rightarrow$  what we do?
  - ↳ ① Early termination  
When traversing postings stop after  $r$  docs or when  $wf_{t,d}$  is too low.  
Take the union of each sets of  $r$ -docs, one set for each query term  
Compute the score for docs in union
  - ↳ ② idf-ordered terms  
When considering postings of a query, look at them by decreasing idf.  
As we update score contribution from each query term, stop if doc scores relatively unchanged (not very clear)

## CLUSTER PRUNING

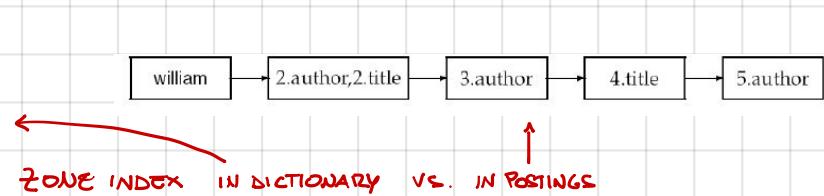
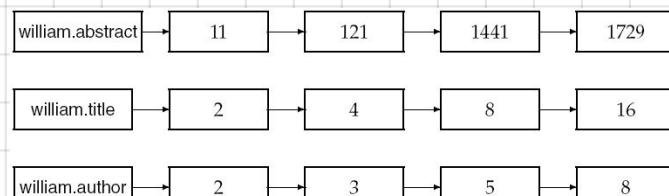
Preprocessing: pick  $\sqrt{N}$  docs random called **leaders**, for all other (**followers**) compute nearest leader. Every leader has about  $\sqrt{N}$  followers.

Query processing: given query  $Q$ , find its nearest leader  $L$  and seek  $k$  nearest docs in  $L$ 's followers.

## PARAMETRIC AND ZONE INDEXES

Some docs have specific term with special semantic called **metadata** (language, format, author...)

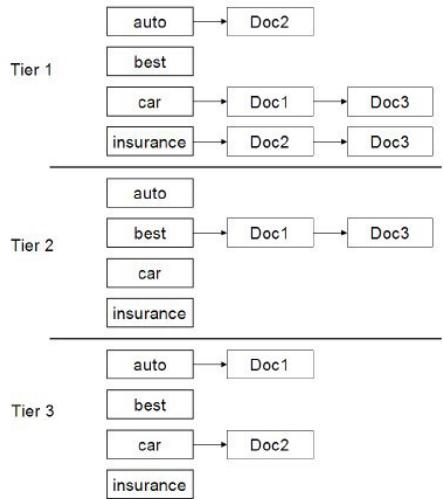
- Field: e.g. YEAR=1601  $\rightarrow$  postings for each field value  $\rightarrow$  doc must be from year 1601
- Zone: e.g. title  $\rightarrow$  build inverted index on zones as well to permit querying  $\rightarrow$  e.g. find docs with "merchant" in title



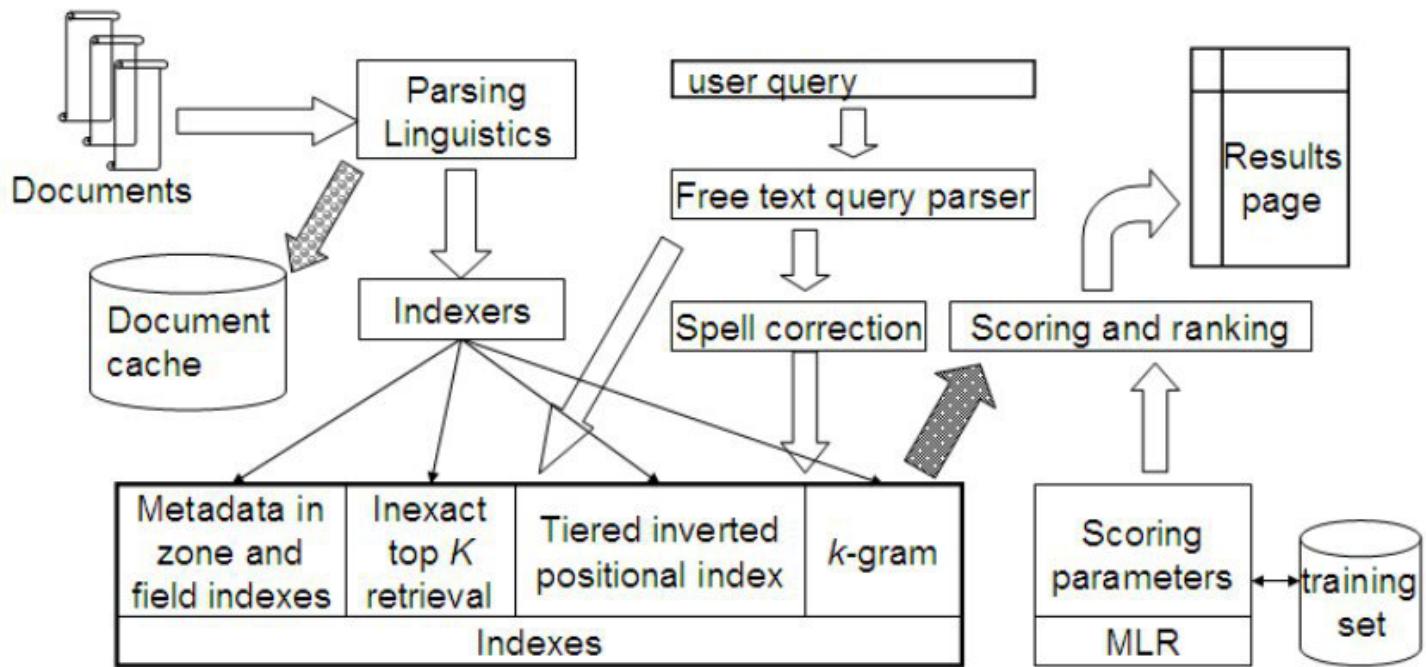
## TIERED INDEXES

Break postings up into a hierarchy of lists, from the most important to the least one. Can be done with  $g(d)$ .

Start finding  $K$  docs in top tier, if not enough go lower.



## PUTTING IT ALL TOGETHER



## USER HAPPINESS

How good are query results? User happiness is not a "factor" as speed of response or size of indexes.

Need a way to quantify it!  $\Rightarrow$  RELEVANCE

Relevance measurement needs 3 elements:

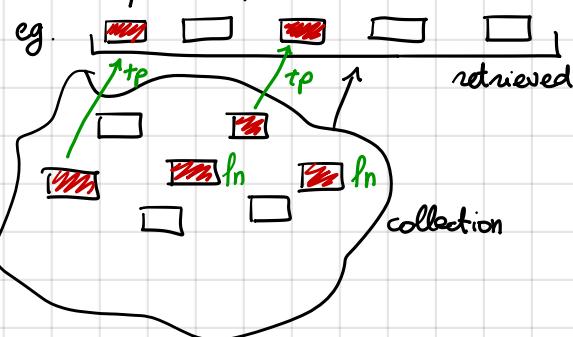
- A benchmark docs collection
- A benchmark suite of queries
- A (usually) binary assessment of either Relevant or Nonrelevant for each query and each doc.

Relevance is assessed relative to the information need (what the user wants for real in his mind) not the query.

A binary assessment is "composed" by Precision and Recall

Precision = fraction of retrieved docs that are relevant

Recall = fraction of relevant docs that are retrieved



	Relevant	Nonrelevant
Retrieved	tp <i>true positive (2)</i>	fp <i>false positive (2)</i>
Not Retrieved	fn <i>false positive (2)</i>	tn <i>true negative (0)</i>

$$\text{Precision} = P = \frac{tp}{(tp + fp)} = \frac{2}{5}$$

$$\text{Recall} = R = \frac{tp}{\frac{(tp + fn)}{\text{relevant items total}}} = \frac{2}{(2+2)} = \frac{1}{2}$$

$\Rightarrow$  Precision@K = precision at first K docs

Same for Recall@K

- Measure for a search engine

How fast does it index? (e.g. bytes/hour) How fast does it search? (e.g. latency as function of queries/second) How much does a query cost? (in dollars)

All these are measurable. But the key measure is user happiness we saw before. It includes speed of response, user-friendly UI and of course relevance (being free maybe more important)

The user change depending on the type of search engine and with it also change the measure used to evaluate its happiness.

e.g. E-commerce: user = buyer, success = buyer buys something, measure = time to purchase after query, fraction of "conversions" of searches to buyers

Back to relevance remember that given a query, a doc can match it very well but at same time can be not related to the user's information need!

About relevance, usually search engines make a tradeoff between recall and precision, lets call it F.

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1) PR}{\beta^2 P + R} \quad \text{with } \beta^2 = \frac{1-\alpha}{\alpha}, \alpha \in [0,1], \beta^2 \in [0, \infty] \rightarrow \text{harmonic mean: } \frac{1}{F} = \frac{1}{2} \left( \frac{1}{P} + \frac{1}{R} \right)$$

The most frequently used combination is  $\alpha=0.5$  and  $\beta=1$  called balanced F

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

Why we don't use accuracy  $= (Tp+Tn)/(Tp+Tp+Fn+Fn)$ ? cause Tn is very huge and a system tuned to maximize accuracy will give a lot of false positive docs (even only FP)

Why harmonic mean? We want to punish bad performance either on precision or recall, and when we have numbers that differ greatly, the harmonic mean is closer to their minimum than arithmetic or geometric mean

Example on next page

To explain, consider for example, what the average of 30mph and 40mph is? if you drive for 1 hour at each speed, the average speed over the 2 hours is indeed the arithmetic average, 35mph.

However if you drive for the same distance at each speed -- say 10 miles -- then the average speed over 20 miles is the harmonic mean of 30 and 40, about 34.3mph.

The reason is that for the average to be valid, you really need the values to be in the same scaled units. Miles per hour need to be compared over the same number of hours; to compare over the same number of miles you need to average hours per mile instead, which is exactly what the harmonic mean does.

Precision and recall both have true positives in the numerator, and different denominators. To average them it really only makes sense to average their reciprocals, thus the harmonic mean.

share improve this answer

answered Oct 14 '14 at 14:54



Sean Owen

58.5k • 18 • 125 • 159

Precision/recall / F score measure for unranked sets. For ranked lists just compute the set measure for each "prefix": top 1, top 2, top 3 ... results giving a precision-recall curve



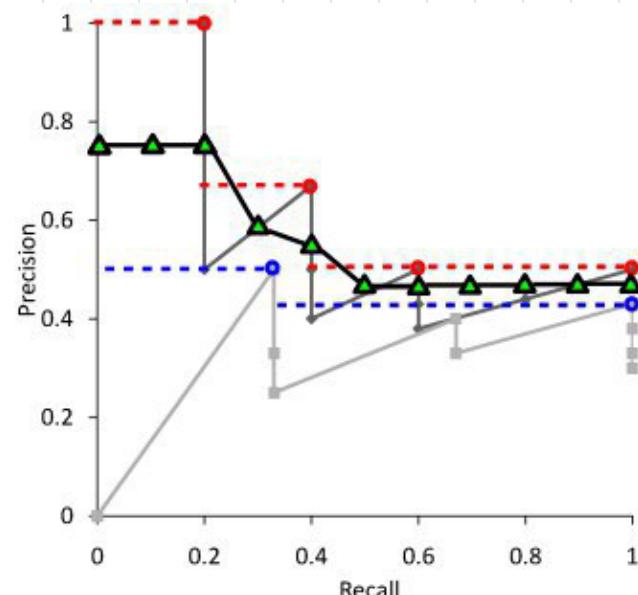
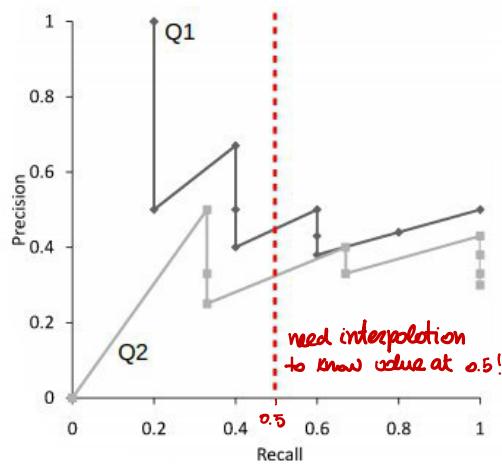
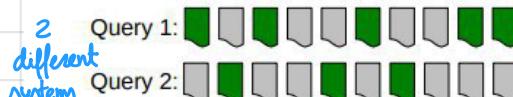
Ranking #1

	Recall	Precision
0.17	0.17	1.0
0.33	0.67	0.5

= the relevant documents

	Recall	Precision
0.5	0.67	0.75
0.67	0.83	0.8
0.83	0.83	0.83
0.83	0.83	0.71
1.0	0.6	0.63

How to compare systems' curves? Interpolation!



On average precision drops as recall increases and so we define interpolation to preserve monotonicity

M.A.P : mean average precision , its a single-number metric most frequently used in research paper

$$\text{avg. precision query 1} = (1 + 0.67 + 0.5 + 0.44 + 0.5) / 5 = 0.62 \quad \boxed{\text{M.A.P} = (0.62 + 0.44) / 2 = 0.53}$$

$$\text{avg P.q.2} = 0.44$$

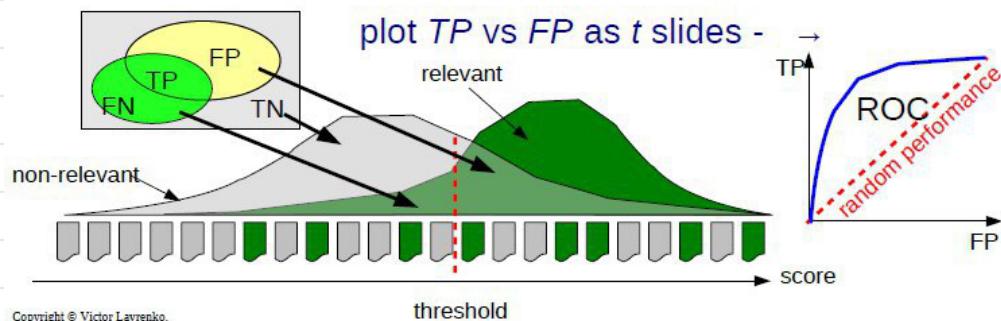
D.C.G. : discounted cumulative gain , its a measure of ranking quality.

Given  $n$  docs each has rating  $r_1, r_2, \dots, r_n$  (in ranked order)  $\rightarrow CG = r_1 + r_2 + \dots + r_n$

$$DCG = r_1 + \frac{r_2}{\log 2} + \frac{r_3}{\log 3} + \dots + \frac{r_n}{\log n}$$

## ERRORS CLASSIFICATION

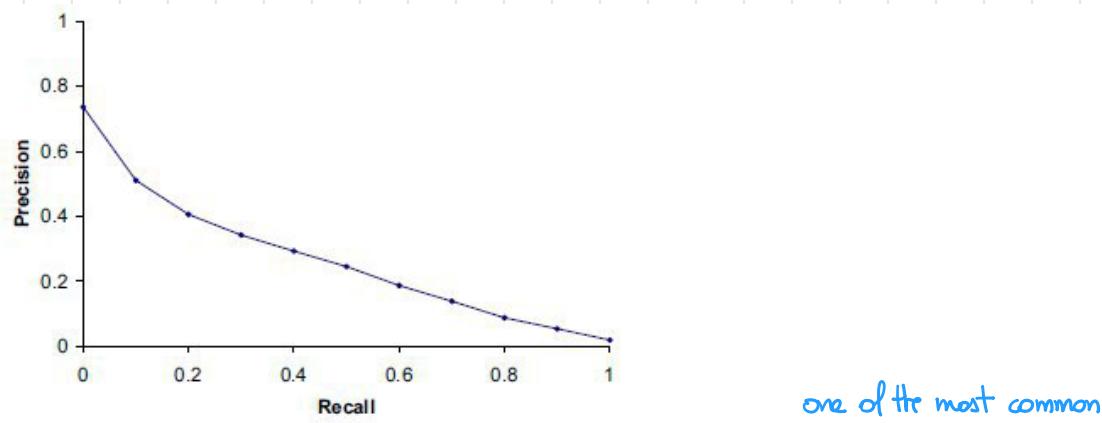
- False Pos. rate:  $\Pr(x > t | -) = FP / (FP+TN)$
- False Neg. rate =  $\Pr(x < t | +) = FN / (TP+FN)$
- True Pos. rate =  $\Pr(x > t | +) = 1 - \text{False Neg.}$
- Receiver Operating Characteristic (ROC):



Relevance assessment are usable only if consistent: how can we measure the consistency or agreement between who judge a doc relevant?  $\rightarrow$  [Kappa measure](#).

$K = \frac{P(A) - P(E)}{1 - P(E)}$  with  $P(A)$  = proportion of time judges agree and  $P(E)$  = what agreement would we get by chance.

change  $[\frac{2}{3}, 1]$  is considered acceptable.



- This curve is typical of performance levels for the [TREC](#) benchmark.
- 70% chance of getting the first document right (roughly)
- When we want to look at at least 50% of all relevant documents, then for each relevant document we find, we will have to look at about two nonrelevant documents.
- That's not very good.
- High-recall retrieval is an unsolved problem! [what to do?](#)  $\rightarrow$

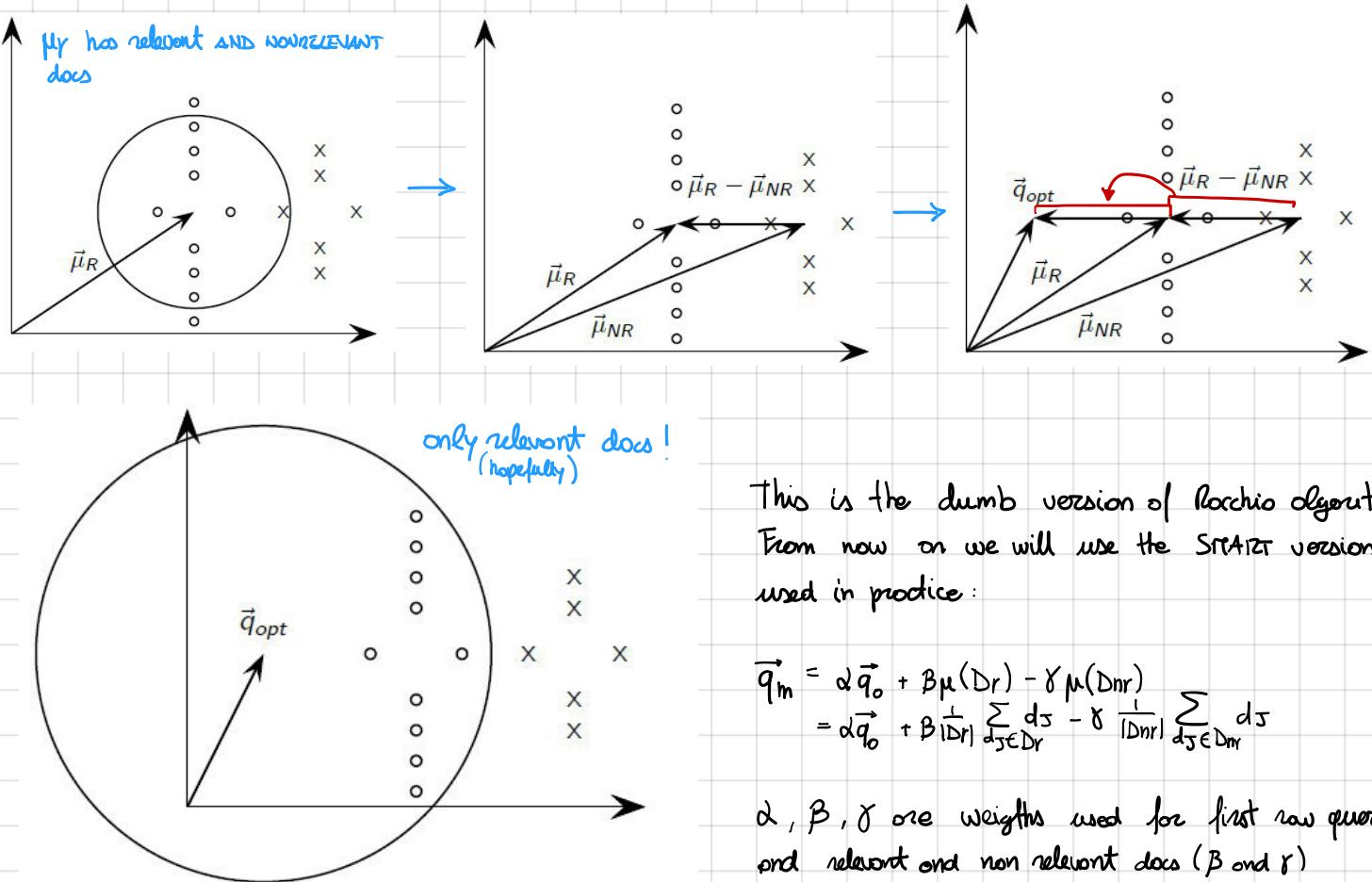
- **Interactive Relevance Feedback**: make a query, take the first results set and ask user to select relevant and nonrelevant docs. The user will see the result summaries (e.g. Google) and will click docs that he thinks are relevant → improve overall performance  
Best method: Rocchio feedback
- **Query expansion**: improve retrieved result by adding synonyms / related term to the query. Used for example on one-term query.

The main idea is to do a "local" on-demand analysis for an user query with relevance feedback and a "global" one (e.g. of collection) to produce **resources** and use it for query expansion.

PART 1: the user issues a (short, simple) query, get a first set of result and mark the docs as relevant or non relevant (assumption). Then search engine computes a new representation of information need, hopefully better, and runs a new query with new (better) results

**Centroid** is the center of mass of a set of points that summarize their characteristics  $\vec{\mu}(D) = \frac{1}{|D|} \sum_{d \in D} \vec{v}(d)$   
Remember we represent docs as vector/point in an highdimensional space.

Rocchio algorithm chooses the query  $\vec{q}_{opt}$  that maximizes  $\vec{q}_{opt} = \underset{q}{\operatorname{argmax}} [\operatorname{sim}(q, \mu(D_r)) - \operatorname{sim}(q, \mu(D_{nr}))]$   
 $\rightarrow \vec{q}_{opt}$  is the vector that separates relevant and nonrelevant docs maximally difference b/w q and  $D_r / D_{nr}$   
with additional assumptions:  $\vec{q}_{opt} = \mu(D_r) + [\mu(D_r) - \mu(D_{nr})]$



This is the dumb version of Rocchio algorithm  
From now on we will use the SMART version used in practice:

$$\begin{aligned}\vec{q}_m &= \alpha \vec{q}_0 + \beta \mu(D_r) - \gamma \mu(D_{nr}) \\ &= \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{d_j \in D_r} d_j - \gamma \frac{1}{|D_{nr}|} \sum_{d_j \in D_{nr}} d_j\end{aligned}$$

$\alpha, \beta, \gamma$  are weights used for first raw query ( $\alpha$ ) and relevant and non relevant docs ( $\beta$  and  $\gamma$ )

Positive feedback ( $\beta$ ) is more valuable than negative ( $\gamma$ ) and many systems only allow positive feedback ( $\gamma=0$ )

Relevance Feedback has 2 assumptions:

- A1: user knows the terms in collection for first query  
↳ violation: mismatch of user's vocabulary and collection one (e.g. cosmonaut / astronaut)
- A2: relevant docs contain similar terms.  
↳ Violation: docs with different terms but some query concept are not helped by the feedback

Of course "edited" query is better than raw one → a fair evaluation must be on docs not yet judged by user

Alternative: user revises and resubmits query (what a normal user does right now IRL)

Problems: expensive to process, users are reluctant to give feedback

⇒ **Pseudo-Relevance Feedback**: ① retrieve a ranked list of **hits** for the user's query ② assume top K docs are relevant ③ do relevance feedback.

Works very well on average but several iteration can cause **query drift**: no longer "close" to the query.

**PART 2**: in global query expansion, the query is modified based on some global query-independent resource.

A publication or DB that collect synonyms is called **thesaurus**: HOSPITAL → MEDICAL

- manual: maintained by editors
- automatic: generated automatically based on statistics (distribution of words)

Or can be used query-equivalence based on query log mining.

Manual th. is expensive and usually used in specialized search engines for science and engineering

Automatic th. has a fundamental notion: similarity between two words

① two words are similar if they co-occur with similar words e.g. "car" ≈ "motorcycle" cause they both occur with "road", "gas", "plate".

② two words are similar if they occur in a given grammatical relation with the same words.  
e.g. "apples" and "pears" are similar cause you can "harvest", "eat", "peel" both

① is more robust ② more accurate

Anyway query logs is the main source of query expansion: e.g. user searching "flower pics" and user searching for "flower clipart" both click on same link → "flower pics" and "flower clipart" are potential expansions of each other.

e.g. user issues "herbs" before than "herbal medicine" → "herbal medicine" is potential expansion of "herb"

## TEXT CLASSIFICATION / CATEGORIZATION

Classification means distinguish relevant doc from collection ( $\neq$  ranking)

STANDING QUERY: a query, used for classification, that is run periodically to retrieve new relevant doc according to an information need.

e.g. Google Alert  $\rightarrow Q = \text{"bike" AND "Yamaha"}$   $\rightarrow$  notify new docs relevant for Yamaha bikes.

Classification is used for spam filtering too: when you mark an email as spam, the system tries to create a query which describes it and use it to find new spam (+ training set for machine learning)

### RECAP

- BAYES THEOREM  $\rightarrow P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A|B) P(B)}{P(A)}$
- $P(A_1, A_2, A_3, A_4) = P(A_1|A_2, A_3, A_4) P(A_2|A_3, A_4) P(A_3|A_4) P(A_4)$

Given a representation of a doc, we want to identify its class (category)  $\rightarrow$  we need a classification function called **classifier**.

Classification methods:

- Manual: accurate when done by expert, difficult and expensive to scale
- Hand-coded rule-based: need to build very good rule (code) to categorize docs
- Supervised Learning: what we said for spam, data set can be built by amateurs (no need experts)

Keeping the last method in mind we define the **probabilistic relevance feedback**: if user has told us some relevant and some nonrelevant docs we can build a probabilistic classifier

e.g.  $P(\text{"dog"} | \text{Relevant})$  = probability that a relevant doc contains word "dog"  
 $\approx \frac{N_{\text{dog}}}{N_r} = \frac{\text{relevant docs with dog}}{\text{total relevant docs}} \rightarrow P(\text{"dog"} | \text{NR}) \approx \frac{N_{\text{non-dog}}}{N_{\text{nr}}}$

Bayes' rule for text classification:  $P(c, d) = \frac{P(d|c) \cdot P(c)}{P(d)}$

$P(c)$  = prob. that the doc I pick is class C

Assuming  $d = \langle x_1, x_2, x_3, \dots \rangle$  in position order then  $p(d|c) = P(x_1, x_2, x_3, \dots | c)$   $\rightarrow$  very complicated!

To "simplify" this, we assume that terms probability are independent:  $P(d|c) = \prod_i^n p(x_i | c)$   
 but this is still on huge number because we still keep track of position.

Then:  $\forall i \quad P(x_i = t | c) = P(x_j = t | c) \quad \forall i, j \neq i$

In this way "I go home" and "go I home" are computed in the same way because we don't keep track of position: maybe too much simplification but it's ok!

$$\Rightarrow \frac{P(d|c) \cdot P(c)}{P(d)} = \frac{\prod_i^n P(x_i | c) P(c)}{\prod_i^n P(x_i)}$$

independent variables and position

What we want to do something to maximize this probability:

$$g(d) = \arg \max_{c \in C} \frac{\prod_i^n P(x_i | c) P(c)}{\prod_i^n P(x_i)} = \arg \max_{c \in C} \prod_i^n P(x_i | c) P(c)$$

doesn't depend on classes

$$\forall c: \hat{P}(c) = \frac{N_c}{N}$$

$\forall t, \forall c: \hat{P}(t|c)$  = occurrences of term  $t$  over total number of terms in all docs of class  $c$ , easily done during dictionary creation

Problem: if  $P(x_i|c) = 0$  !!! (e.g. a misspell)

$$\hat{P}(t|c) = \frac{N_{t,c}}{\sum N_{t,c}} \Rightarrow \hat{P}(t|c) = \frac{N_{t,c} + 1}{\sum (N_{t,c} + 1)}$$

adding 1 doesn't make any problem even if a zero probability

### MULTIVARIATE BERNOULLI MODEL

$$P(d|c) = P(x_1, \dots, x_n | c)$$

but in this case every doc is a binary vector with length  $n = |V|$  (dictionary)

We can make some assumption as before:

$$P(d|c) = \prod_i P(x_i|c)$$

how many docs over the total have  $x_i$  in them (different from before!)

Naive Bayes (the edited version of Bayes' rule) is widely used for spam filtering and it's not so weak as someone could think: in fact it's very fast and robust and has low storage requirements

Note: evaluation must be done on test data that are independent of training data, sometimes it's possible to do a cross-validation (averaging results b/w multiple training and test sets)

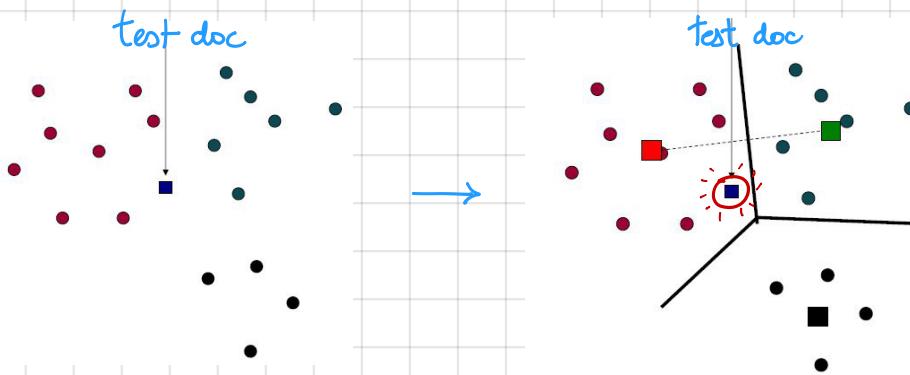
Evaluation is done according precision, recall, F<sub>1</sub> and classification accuracy ( $\frac{1}{n} = \frac{\text{docs correctly classified}}{\text{total docs}}$ )

Remember we are working with docs represented as vector in high dimensional space. How can we do classification?

Hypothesis 1: docs in some class form a contiguous region in space.

Hypothesis 2: docs from different classes don't overlap (much)

Learning a classifier means build surface to delineate classes in the space



How to find a good separator?

Use centroid and classify docs according to nearest centroid (Rocchio)

Rocchio defines line or hyperplane as  $\sum_i w_i d_i = \theta$

$$\vec{w} = \bar{\mu}(c_1) - \bar{\mu}(c_2)$$

$$\theta = 0.5 \cdot (\|\bar{\mu}(c_1)\|^2 - \|\bar{\mu}(c_2)\|^2)$$

how?

click here → <https://youtu.be/eEEPAeweKqs?t=2516>

Rocchio is little used outside text classification and in general is worse than Naive Bayes but cheap to train and test docs.

$$\Rightarrow \vec{x}^\top \vec{w} = \frac{1}{2} \theta \Rightarrow \text{Rocchio is a linear classifier}$$

Why Naive Bayes is a linear classifier too? ( $n_d = \# \text{terms in doc } d$ )

We want to estimate  $\hat{P}(c|d) \propto \hat{P}(c) = \prod_{k=1}^{n_d} \hat{P}(t_k|c)$  where the probability to find  $t_k$  in class  $c$  is the same (for every some  $t_k$ ) regardless of position (only in the model, not IRL) odd ratio

Assume also  $\hat{P}(\bar{c}|d) \propto \hat{P}(\bar{c}) \prod_{k=1}^{n_d} \hat{P}(t_k|\bar{c})$

When we try to classify a doc we are doing the following:  $\frac{\hat{P}(c|d)}{\hat{P}(\bar{c}|d)} = \frac{\hat{P}(c) \prod_{k=1}^{n_d} \hat{P}(t_k|c)}{\hat{P}(\bar{c}) \prod_{k=1}^{n_d} \hat{P}(t_k|\bar{c})}$

We could decide which class is doc  $d$  checking odd ratio:  $c$  if  $\gg 1$ ,  $\bar{c}$  otherwise, but for numerical reason we use logarithm

$$\log \frac{\hat{P}(c|d)}{\hat{P}(\bar{c}|d)} = \log \frac{\hat{P}(c)}{\hat{P}(\bar{c})} + \sum_{i=1}^{n_d} \log \frac{\hat{P}(t_i|c)}{\hat{P}(t_i|\bar{c})}$$

$$= \log \frac{\hat{P}(c)}{\hat{P}(\bar{c})} + \sum_{i \in V} x_i \log \frac{\hat{P}(t_i|c)}{\hat{P}(t_i|\bar{c})}$$

cause we said probability is position independent

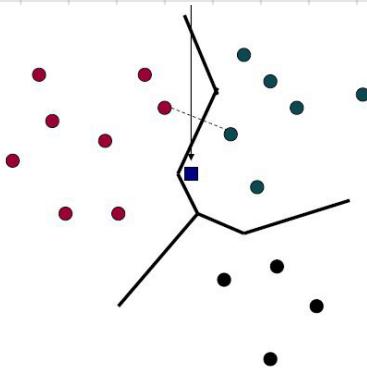
$$0 \leq \frac{\log \frac{\hat{P}(c)}{\hat{P}(\bar{c})}}{\theta} + \sum_{i \in V} x_i \log \frac{\hat{P}(t_i|c)}{\hat{P}(t_i|\bar{c})} \Rightarrow \sum x_i w_i \geq -\theta$$

e.g. "a rose is a rose"  $n_d = 5$   $|V| = 3$

$x_i = \# \text{occurrences of term } i \text{ in } d$  (e.g.  $x_{\text{rose}} = 2$ )

## KNN - K Nearest Neighbor Classification

To classify doc  $d$  we define its  $K$ -neighborhood as the  $K$  nearest neighbors and pick the majority class label in this neighborhood: for larger  $K$   $P(c|d) \approx \#(c)/K$



Voronoi diagram

As we can see this doesn't depend on a linear classifier!  
Usually  $K$  is odd to avoid ties; 3 and 5 are most common  
Naively finding nearest neighbors requires a linear search through  $|D|$  docs in collection that is the same as determining the  $K$  best retrievals using the test doc as a query to a database of training doc → use inverted index.

Pro: no training necessary, scales well, in most cases more accurate than NB and Rocchio (Bayes optimal)

Cave: expensive at test time, small changes can have ripple effect.

What we do with multiple-class docs? we have two types of answers:

• 1-of classifier: get only the most accurate class

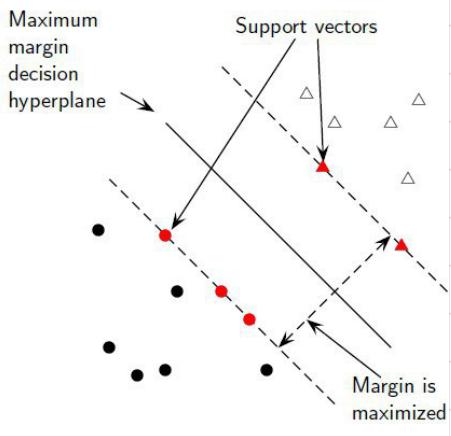
• only-of classifier: get a vector composed by scores given to each class (for a given doc) based on "accuracy"

Learning algorithms for vector space classification can be divided in two types: **simple** (N.B. Rocchio, KNN) and **iterative** (**Support Vector Machines**, Perception) that usually are the best choice. The simple ones often do only one linear pass.

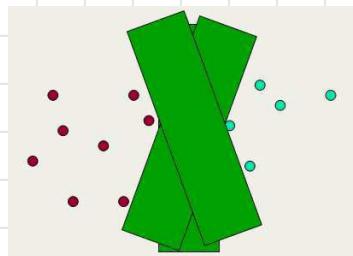
## SVM - Support Vector Machines

We still are in vector space (as previous algorithms) but we have now large margin classifier: we want to increase the division region btw classes and aim to find a separating hyperplane that is **maximally far** from any point in the training set.

In case of non-linearity we accept some points as noise (errors).



Why maximize the margin? Because points near the decision surface are less confident, uncertain classification decision and at the same time gives a classification safety margin respect to errors and random variation



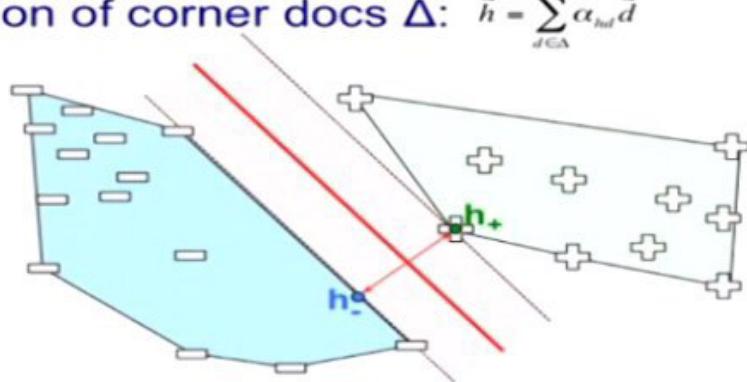
All points  $\vec{x}$  on the hyperplane satisfy  $\vec{w}^\top \vec{x} + b = 0$  ( $b$  is what we called  $\theta$ )

- Find boundary vector  $\mathbf{w}$  that satisfies conditions:

$$\left. \begin{array}{l} \sum_i w_j d_{j,i} \geq +1 \text{ all positives get scores } \geq +1 \\ \sum_i w_j d_{j,i} \leq -1 \text{ all negatives get scores } \leq -1 \\ \min \|\vec{w}\| \text{ maximum margin around } \mathbf{w} \end{array} \right\} \vec{w} = \sum_{j \in +} \alpha_j \vec{d}_j - \sum_{j \in -} \alpha_j \vec{d}_j \quad \left. \begin{array}{l} \text{"centroid" of positives} \\ \text{"centroid" of negatives} \\ +/- docs have "weights" \end{array} \right\} \begin{array}{l} \cdot \text{ similar to centroid / Rocchio} \\ \cdot \text{ some docs more important} \\ \cdot \text{ most don't matter: } \alpha_j = 0 \end{array}$$

- Convex hull: polytope around all positives / negatives

- any point  $\mathbf{h}$  is a linear combination of corner docs  $\Delta$ :  $\vec{h} = \sum_{d \in \Delta} \alpha_{hd} \vec{d}$
- find two nearest points  $\mathbf{h}_+$ ,  $\mathbf{h}_-$
- margin must be  $\|\mathbf{h}_+ - \mathbf{h}_-\|$
- boundary must be half-way
- perp. to:  $\vec{w} = \vec{h}_+ - \vec{h}_- = \sum_{d \in +} \alpha_d \vec{d} - \sum_{d \in -} \alpha_d \vec{d}$



Note: polytopes are the lines around positive/negative docs.

e.g. 2 classes  $y_i = +1$  and  $y_i = -1 \rightarrow f(\vec{x}) = \text{sign}(\vec{w}^\top \vec{x} + b)$

If a new doc is classified within the margin is possible to return "don't know" or use probability function to assign a class.

What if data is not linearly separable?

Allow to make **errors** and increase margin and reduce overfitting (accept the cost for every errors)

We introduce a **Slack variable**  $E_i$ , a non-zero value that takes in consideration these errors: the sum of  $E_i$  gives an upper bound on the number of errors.

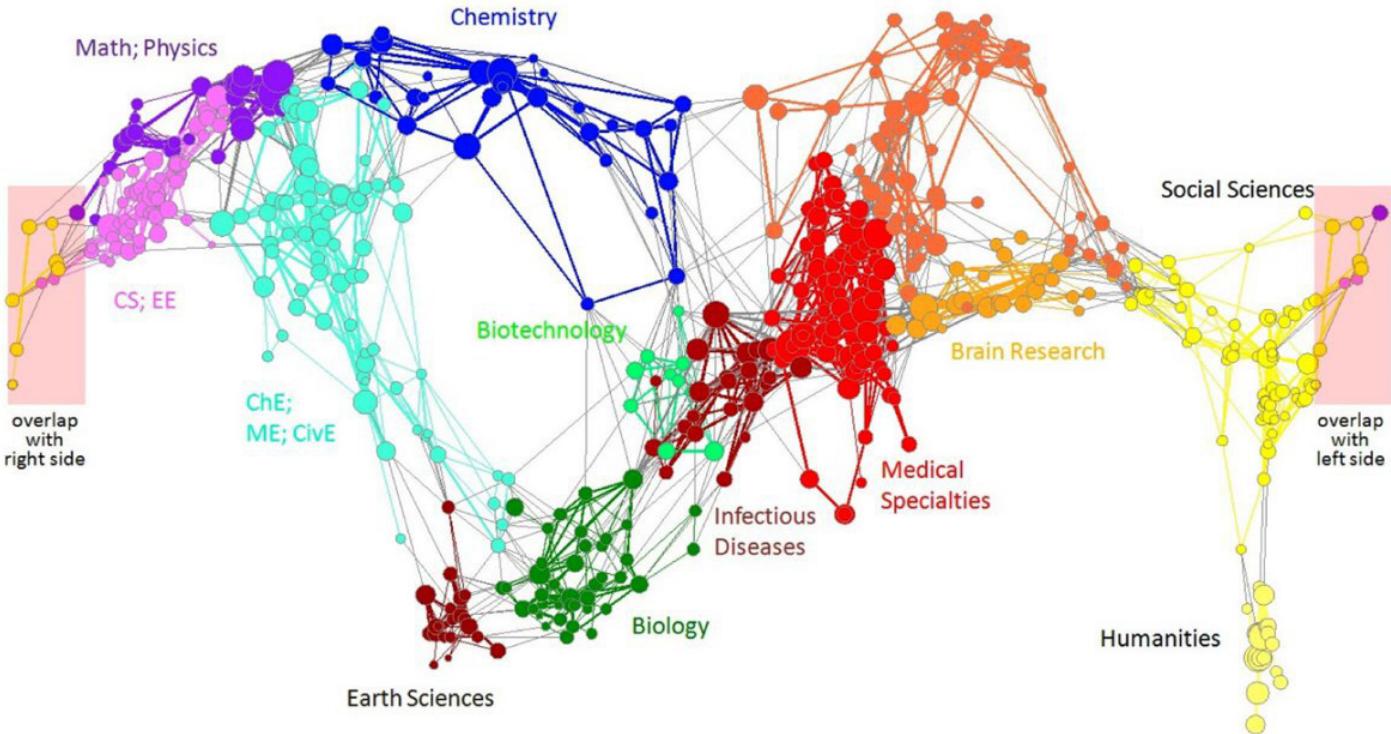
Ways to use SVM (for one-of classification)

- ① Train and run K classifier and then select the class with the highest confidence
- ② Build  $K(K-1)/2$  1vs1 classifier and choose the class selected by the most classifier: more classifier but less time of training due to data set is much smaller
- ③ Structured prediction

### CHOOSING CLASSIFIER

- How much training data is there currently available?
- No labeled training data? → Use hand written rules (even more complicated than boolean)

## LINK ANALYSIS



### Challenges:

- Web contains many sources, who trust? → Trustworthy page may point to each other
- What is the "best" answer to query e.g. "newspaper"? → Pages that actually know about newspapers might all be pointing to many "newspapers".

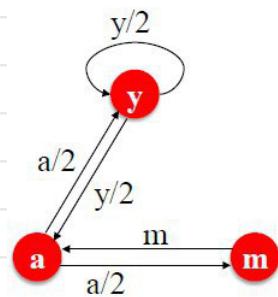
All webpages are not equally "important" and for this reason we will rank this graph's nodes by the link structure!

Importances of a node:

- Page Rank
- Hubs and Authorities (HITS)
- Topic-specified (Personalized) Page Rank
- Web Spam detection Algorithms

Idea: links as vote. More links more points. link to/from important pages count more. Recursive!

e.g. Page  $y$  with importance  $r_y$  has  $n$  out-links, each get  $r_y/n$  votes. Page  $y$  importance is the sum of in-links.  $\text{Rank } y = r_y = \sum \frac{r_i}{d_i} \quad (d_i = n_i)$



No uniq. solution! Gaussian elimination is not suitable for web-size graphs.  
New solution: Stochastic adjacency matrix.  $\Rightarrow$

"Flow" equations:

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

$\Gamma$	$J$	
$i$	$A_{ij}$	
	$\dots$	
	$A_{in}$	
		$\downarrow A$

if  $i \rightarrow j$   $A_{ij} = 1$  else  $A_{ij} = 0 \rightarrow \pi_j = \sum \frac{A_{ij} \cdot \pi_i}{d_i} = \sum \pi_{i,j} \cdot \pi_i$

$\pi^T$  (row-vector) =  $\pi^T M$

$\forall i \sum_{j=1}^n \pi_{i,j} = 1 \rightarrow$  stochastic property (row stochastic)

Note:  $r^T = r^T M$  means  $r^T$  is an eigenvector with eigenvalue 1

How to solve  $M \cdot r^T = r^T$ ? Power Iteration Method

To simplify formula we will start from  $\pi = M^T z = A^T z$  (not previous  $A$ ,  $A$  is column stochastic)

Initialize  $\pi^{(0)} = x$  and iterate  $\pi^{(t+1)} = A \cdot \pi^{(t)}$   $\rightarrow \pi^{(1)} = Ax$ ,  $\pi^{(2)} = A^2x$ , ...,  $\pi^{(t)} = A^t x$

$x$  is a vector  $[1/n, \dots, 1/n]$  (usually)

$(w_1, \lambda_1), \dots, (w_n, \lambda_n)$  are eigenvectors ( $w_i$ ) and corresponding eigenvalues ( $\lambda_i$ ) of  $A$

$\lambda_1$  in stochastic matrix is equal to 1 and all eigenvalues are real ( $\leq \lambda_1$ )

Since this is a system of eigenvectors, it's a base of  $\mathbb{R}^n$  so every other vector can be expressed as linear combination of these ones.

$$\Rightarrow x = \sum \alpha_i w_i$$

$$\hookrightarrow A^t x = A^t (\sum \alpha_i w_i) = \sum \alpha_i A^t w_i$$

$$A w_i = \lambda_i w_i, A^t w_i = A A w_i = A \lambda_i w_i = \lambda_i^t w_i, \dots$$

$$\Rightarrow A^t x = \sum \alpha_i \lambda_i^t w_i \quad \text{Note: } |\lambda_2| = \lambda_2 = 1 \quad |\lambda_i| < 1 \quad i > 1 \quad \text{ALWAYS (if no error occurs)}$$

$$= \alpha_1 w_1 + \sum_2^n \alpha_i \lambda_i^t w_i \quad (= \alpha_1 w_1, (n \rightarrow \infty, \sum_2^n \alpha_i \lambda_i^t w_i \rightarrow 0))$$

The main eigenvalue is 1!  $\rightarrow \pi = \text{main eigenvector} \rightarrow \alpha_1 w_1 = \alpha_1 \pi \quad \text{ALWAYS NOT NEGATIVE}$

So: to solve  $\pi^T - \pi^T M$  we do  $(\pi^{(t+1)})^T = (\pi^{(t)})^T M$

Term  $\sum_2^n \alpha_i \lambda_i^t w_i$  converges and how fast depends on how small  $|\lambda_i|$  are  $\rightarrow$  depends on  $\max_{i>1} |\lambda_i|$  exponentially

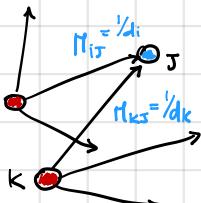
Let's assume  $\pi^{(0)} = x = \begin{pmatrix} 1/n \\ \vdots \\ 1/n \end{pmatrix}$ . This give all probability to be at time  $t=0$  in any given point.

$$(\pi^{(1)})^T = (\pi^{(0)})^T M \quad \Leftarrow \quad \pi_j^{(1)} = \sum_{i=1}^n \pi_{i,j}^{(0)}$$

Explain: The probability to be in node  $j$  at time  $t=1$  ( $\pi_j^{(1)}$ ) is equal to the probability to be in any incoming neighbours ( $\pi_i^{(0)}$ ) times the probability to go to  $j$  ( $\pi_{i,j}$ ) at time  $t=0$ .

So the probability distribution is given by  $(\pi^{(t+1)})^T = (\pi^{(t)})^T M$

This is called random walk and  $\pi$  is its stationary distribution.



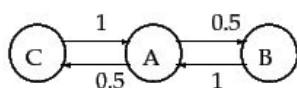
## MARkov CHAINS

A Markov chain is a discrete-time stochastic process: a process that occurs in a series of time-steps in each of which a random choice is made. A Markov chain consists of  $N$  states. Each web page will correspond to a state in the Markov chain we will formulate.

It's characterized by an  $N \times N$  transition probability matrix  $P$  each of whose entries is in the interval  $[0, 1]$ ; the entries in each row of  $P$  add up to 1. The chain can be in one of the  $N$  states at any given time-step; then, the entry  $P_{ij}$  tell us the probability that the state at next time-step is  $j$ , conditioned on the current state being  $i$ . Each entry  $P_{ij}$  is known as a transition probability and depends only on the current state  $i$ ; this is known as the Markov property:  $\forall i, j, P_{ij} \in [0, 1]$  and  $\forall i, \sum_j P_{ij} = 1$ .

The second property is the sum of the stochastic matrix.

$$\text{e.g. } \begin{matrix} & A & B & C \\ A & 0 & 0.5 & 0.5 \\ B & 1 & 0 & 0 \\ C & 1 & 0 & 0 \end{matrix}$$



In a Markov chain, the probability distribution of next state depends only on the current state and NOT on how we reached that state.

Starting from A we have 0.5 probability to proceed either to B or C.

From them we can only reach A, so probability is 1

We can view a random surfer on the web graph as a Markov chain, with one state for each page, and each transition probability representing the probability of moving from one page to another.

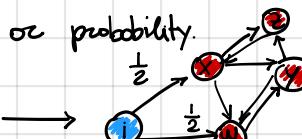
Back to page rank! We have  $Z_j^{(t+1)} = \sum_i \frac{r_i^{(t)}}{d_i}$  or  $Z^T = M Z^T$  but:

1) does this converge?

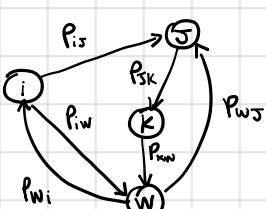
2) does it converge to what we want?

3) are results reasonable?

We saw previously that  $Z_j^{(t+1)}$  can be either zero or probability.  
 $(Z^{(t+1)})_i^T = (Z^{(t)})^T M$  is equal to  $Z_i^{(t+1)} = \sum_j \frac{r_j^{(t)}}{d_j}$   
 $M_{ii} = i^{\text{th}} \text{ row} = (\frac{1}{d_1}, 0, 0, \dots)$  e.g.  $(0, \frac{1}{2}, 0, 0, \frac{1}{2})$



The difference with a markov chain is that all non-zero entries (in the same row) are equal!



$$\forall i, \sum_j P_{ij} = 1$$

$X^{(t)}$  = node (state) in which M.c. is at step t  
 $\hookrightarrow Z_j^{(t)} = P(X^{(t)} = j)$

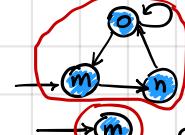
$$M = \begin{pmatrix} & & & & j \\ & & & & | \\ & & & & P_{j,j} \\ & & & & | \\ & & & & \end{pmatrix} \quad [(Z^{(t+1)})^T = (Z^{(t)})^T M]$$

$$\begin{array}{cccc} r_a & = & 1 & 0 & 1 & \dots \\ r_b & = & 0 & 1 & 0 & \dots \\ t & = & 0 & 1 & 2 & 3 \dots \end{array} \quad \left\{ \rightarrow a \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} = A \rightarrow \det(A - \lambda I) = 0 = P(\lambda) \right.$$

$$A - \lambda I = \begin{pmatrix} -\lambda & 1 \\ 1 & -\lambda \end{pmatrix} \rightarrow P(\lambda) = \lambda^2 - 1 = 0 \rightarrow \lambda = \pm 1$$

This is an example of periodic markov chain that doesn't converge!

Other problem: Spider traps



Similar problem is a DeadEnd

All pagerank score get "trapped" in this circle  
e.g. a link to a file without outgoing links

To fix dead ends we adjust matrix with a link to ALL other nodes with  $\frac{1}{n}$  probability

To fix spider traps we use **teleports**. at each time step, the random walker has two options

- with prob.  $\beta$  follow a link at random
- with prob  $1-\beta$  jump to some random page

$$\beta = [0.8, 0.9]$$

When only nodes is reachable from any other nodes, the Markov chain is called **irreducible**

$$\text{irreducible + periodic} = \text{ergodic} \rightarrow \begin{cases} \lim_{t \rightarrow \infty} \pi^{(t)} = \pi \\ \pi \text{ is unique} \end{cases}$$

A M.C. has the following property:  $P(X^{(t)}=j | X^{(t-1)}=a_{t-1}, \dots, X^{(0)}=a_0) = P(X^{(t)}=i | X^{(t-1)}=a_{t-1})$

with  $a_t$  = state at end of round  $t$

$$\text{Note: } P_{ij} = P(X^{(t)}=j | X^{(t-1)}=i)$$

We assume that each state  $i$  has at least one outgoing link and the matrix is stochastic (we enforce this property adding  $n$  links to dead ends)

After  $t$  step, M.C. is in each state with some probability.  $p^{(t)}$  corresponds to a probability distribution, an  $n$ -dimensional vector such that its  $j^{\text{th}}$  entry is  $p_j^{(t)} = P(X^{(t)}=j)$ : since at the end of each round M.C. has to be in some state we have  $\sum p_j^{(t)} = 1 \quad \forall t$

$$p_j^{(t)} = \sum_{i \in S} p_i^{(t-1)} P_{ij} \Leftrightarrow (p^{(t)})^T = (p^{(t-1)})^T \cdot P \Rightarrow (p^{(t)})^T = (p^{(0)})^T P^{(t)} \text{ iterating over different rounds.}$$

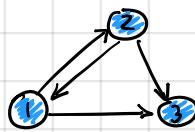
Steps for random walk:

① remove dead ends - a dead end corresponds to a row of zeros. We replace it with a row of  $\frac{1}{n}$  which corresponds to adding links to all nodes (even it self) [matrix  $\hat{P}$ ]

From now on we consider  $A = \hat{P} + \left(\frac{1}{n}\right)[1]^T$  where  $a$  is a column vector such that  $a_i = 1$  if  $i$  is a dead end,  $a_i = 0$  otherwise. ( $[1]^T$  is a row of 1's)

Now it is stochastic

e.g.



$$\hat{P} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

$$\frac{1}{n} \cdot \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}^T$$

② Teleportation - the random walk described by  $P$  is modified as follow: when at a generic node  $i$

- with probability  $\beta$  → follow one of  $i$ 's outgoing links
- with probability  $1-\beta$  → jump to any node at random (even  $i$ )

$$r_j = \sum_{i \in S} \beta \frac{r_i}{d_i} + (1-\beta) \frac{1}{n}$$

larger  $\beta$  is ( $< 1$ ) then more similar the graph is to the original one but without periodicity and irreducible

$$\rightarrow [\pi^{(t)}]^T = \beta \cdot [\pi^{(t-1)}]^T \cdot A + \frac{1-\beta}{n} \cdot [1]^T$$

$$\sum r_j^{(t-1)} = 1 \rightarrow = \beta \cdot [\pi^{(t-1)}]^T \cdot A + \frac{1-\beta}{n} \cdot [\pi^{(t-1)}]^T [1][1]^T = [\pi^{(t-1)}]^T (\beta A + \frac{1-\beta}{n} [1][1]^T) \quad (\rightarrow \pi^{(t)} = \pi^{(t-1)} M)$$

$\underbrace{1}_{\text{P stochastic}}$

Smaller  $\beta$  means faster convergence but lose original topology → find a compromise  $\alpha \sim [0.8, 0.8]$

## SOME PROBLEM WITH PAGERANK (this kind of pagerank)

- Measures generic popularity of a page and biased against topic-specific authorities  
↳ Topic Specific Page Rank
- Uses a single measure of importance when there can be other ones  
↳ Hubs-and-Authorities
- Susceptible to Link Spam: artificial in-link created in order to boost page rank  
↳ TrustRank

## TOPIC-SPECIFIC PAGERANK

Goal: evaluate web pages not just according to popularity, but by how close they are to a particular topic.  
Search engines use browser history, cache and cookies to profile the user and select the best topic for him.

$$\mathcal{Z}_S(t) = \beta \sum_{i \in S} \frac{r_i(t-1)}{d_i} + \frac{1-\beta}{n}, \text{ let assume we want to rank pages according to topic "sport".}$$

We update teleportation in the following way:

$$\begin{cases} \beta M_{ij} + \frac{1-\beta}{|S|} & \text{if } i \in S = \text{sport pages} \\ \beta M_{ij} + 0 & \text{otherwise} \end{cases} \rightarrow \text{this means that when (if) teleporting, jump only to "sport" pages}$$

$$\rightarrow \mathcal{Z}(t)^T = \mathcal{Z}(t-1)^T (\beta A + \frac{1-\beta}{|S|} [1] P^T) \quad P_i = \begin{cases} 1 & i \in S \\ 0 & \text{otherwise} \end{cases} \quad (\text{same as } \alpha \text{ in dead ends elimination})$$

$P^T$  is called personalization vector. How to compute it? Usually you don't have there one open source repository with million pages handmade categorized. (Dmoz)

$P$  can be composed by multiple topics with different percentage of "usage" by the user

e.g. Police =  $(0.3, \dots, 0 \dots, 0.2, \dots 0 \dots, 0.5, \dots 0 \dots)$  (not a personalization vector)

↑ Books      ↑ Movies      ↑ Music

$$\rightarrow \begin{bmatrix} P_{\text{music}} \\ P_{\text{movies}} \\ P_{\text{books}} \end{bmatrix} \rightarrow P = 0.3 P_{\text{books}} + 0.2 P_{\text{movies}} + 0.5 P_{\text{music}}$$

Let assume we have  $P_1$  and  $P_2$ , two personalization vector. Then we can compute two stationary distribution:

$$\mathcal{Z}_1^T = \mathcal{Z}_1^T (\beta A + \frac{1-\beta}{|S_1|} [1] P_1^T) = \mathcal{Z}_1^T (\beta A + (1-\beta)[1] P_1^T) \text{ with } \sum P_{1,i} = 1$$

$$\mathcal{Z}_2^T = \mathcal{Z}_2^T (\beta A + \frac{1-\beta}{|S_2|} [1] P_2^T) = \mathcal{Z}_2^T (\beta A + (1-\beta)[1] P_2^T) \text{ with } \sum P_{2,i} = 1$$

Note: Personalization vectors are stored somewhere

Assume, by user profile,  $P = \alpha_1 P_1 + \alpha_2 P_2$  and we have to compute  $\mathcal{Z}^T = \mathcal{Z}^T (\beta A + (1-\beta)[1] P^T)$ .

Can we compute it without applying power method again?

$P_1$  and  $P_2$  of course are two different MC.

$$\begin{cases} \mathcal{Z}_1^T = \mathcal{Z}_1^T P_1 \\ \mathcal{Z}_2^T = \mathcal{Z}_2^T P_2 \end{cases} \rightarrow \alpha_1 \mathcal{Z}_1^T + \alpha_2 \mathcal{Z}_2^T = \alpha_1 \mathcal{Z}_1^T P_1 + \alpha_2 \mathcal{Z}_2^T P_2 = \alpha_1 \mathcal{Z}_1^T \beta A + \alpha_2 \mathcal{Z}_2^T \beta A + (1-\beta) \alpha_1 P_1^T + (1-\beta) \alpha_2 P_2^T = (\alpha_1 \mathcal{Z}_1^T + \alpha_2 \mathcal{Z}_2^T) \beta A + (1-\beta)(\alpha_1 P_1^T + \alpha_2 P_2^T) \quad \text{yes } \mathcal{Z}^T \quad P^T \quad (\mathcal{Z}^T \cdot [1] = 1)$$

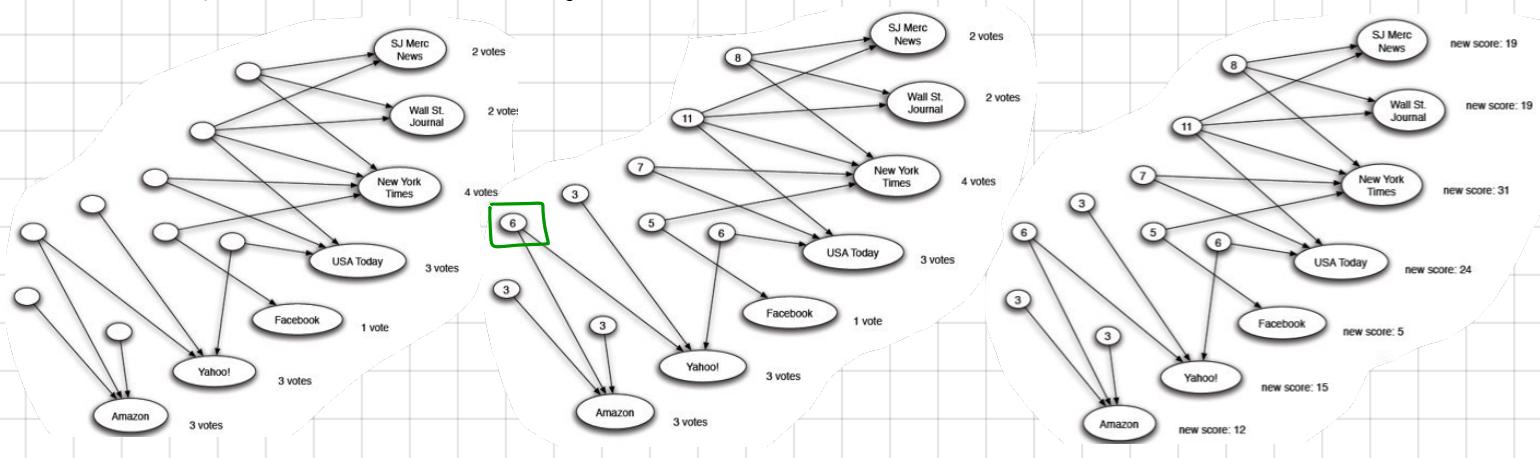
So we can use the linear combination on the fly. It's also true that this operation is not done for every user: users are classified in groups.

## HUBS AND AUTHORITIES

**HTS:** hypertext-induced Topic Selection, a measure of importance similar to PageRank. How it works: let's say we want to find good newspapers → don't just find newspapers but also "experts" who link in a coordinated way to good newspapers. Votes are always links, but this time also outgoing links count! "I'm a good page if trusted pages link to me, but also if I link to trusted pages"

Each page has 2 scores:

- as **hub** = total sum of votes of authorities pointed to
- as **authority** = total sum of votes coming from experts.



① Each page starts with hub score 1 → authorities collect their votes e.g. 3 hubs point to Amazon → vote 3 as authority

② Hubs collect authorities score e.g. green hub score now is  $3(\text{Amazon}) + 3(\text{Yahoo!}) = 6$

③ Authorities again collect hub scores.

Note: this is an example. In real graph is not bipartite and each page has both the scores

- good hub links to many good authorities

- good authority is linked from many good hubs

$$h(x) = \sum_{y: x \rightarrow y} a(y)$$

$$= \sum_y A_{xy} a(y) = A_{*x} \cdot a$$

$$a(x) = \sum_{y: y \rightarrow x} h(y) = \sum_y A_{yx} \cdot h(y) = A_x^T \cdot h$$

} two vectors  $h$  and  $a$

$$A = \begin{pmatrix} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{pmatrix}$$

$$\begin{cases} h^{new} = A a^{old} \\ a^{new} = A^T h^{old} \end{cases}$$

and so on ...

$$h(0) = [1] = 1$$

for  $t \geq 1$ :

$$\begin{aligned} a(t) &= A^T h(t-1) \\ h(t) - A a(t) &\rightarrow h(t-1) - A a(t-1) \end{aligned}$$

$$\begin{aligned} &\Rightarrow \begin{cases} a(t) = A^T A a(t-1) \\ h(t) = A A^T h(t-1) \end{cases} \\ &[c(t) = \Pi \cdot c(t-1)] \text{ power method} \end{aligned}$$

Entries in  $h$  and  $a$  tend to grow, so:

$$\begin{aligned} a(t) &= \frac{a(t)}{\|a(t)\|_2} \\ h(t) &= \frac{h(t)}{\|h(t)\|_2} \end{aligned}$$

## RECOMMENDATIONS

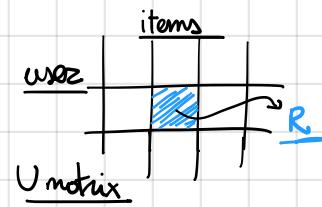
When a user search for something on the web usually he get recommendations from different sites. We come from scarcity (low amount of content in the past) to abundance and so we need to filter results.

Recommendations are very important also for "items" that are rare and without recommendations would never be found.

Different types:

- Editorial and hand curated e.g. list of favorites
- Simple aggregates e.g. top10, most popular, recent..
- Tailored to individual user e.g. Amazon, Netflix - that's what we will talk about

$X$ : set of users }  $\rightarrow$  Utility function:  $v = X \times S \rightarrow R$   
 $S$ : set of items       $R$ : set of ratings, totally ordered, 0-5 stars, real numbers  $[0, 1]$



Problems:

- 1) How to collect ratings
- 2) How to extrapolate unknown ratings from the known ones
- 3) Evaluate success of recommendations

1 Explicit way  $\rightarrow$  ask users to rate items (not so good)

Implicit way  $\rightarrow$  learn from past: purchase = high rate  $\rightarrow$  What about low ratings?

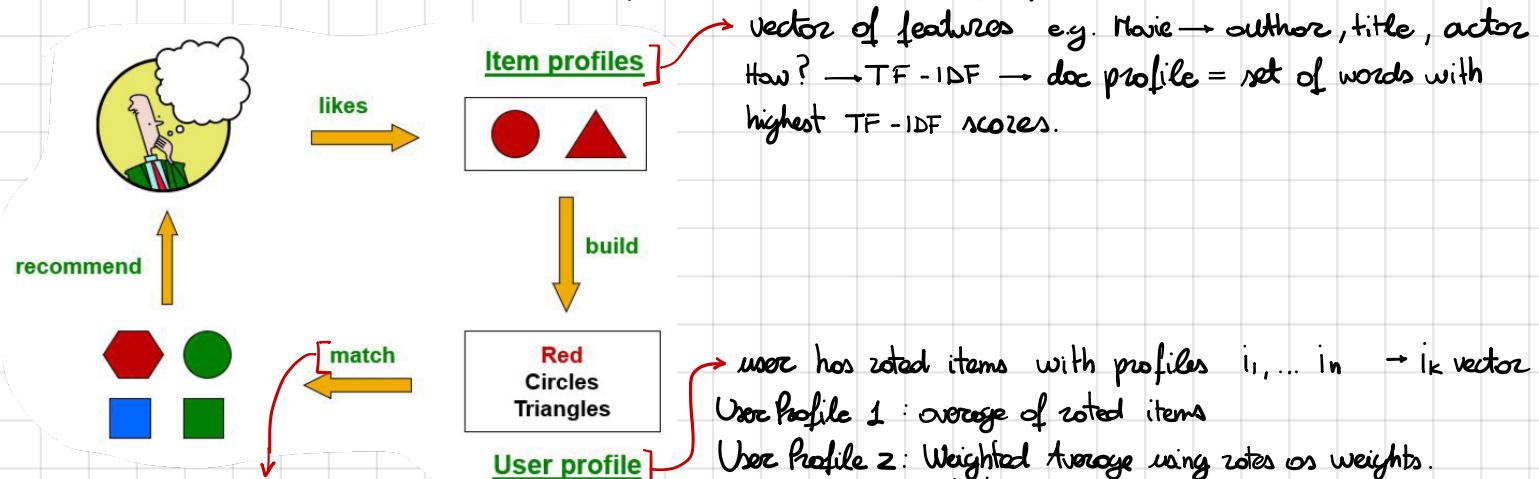
2 U matrix is sparse  $\rightarrow$  new items have no ratings, new users have no history

3 approaches:

- content based  $\xrightarrow{\text{used today!}}$
- collaborative
- latent factor based

### CONTENT-BASED

Recommend items to customer  $\propto$  similar to previous items rated highly by  $x$



Given user profile  $x$  and item profile  $i$ , estimate  $u(x, i) = \cos(x, i) = \frac{x \cdot i}{\|x\| \cdot \|i\|} \rightarrow$  cosine similarity

**PRO:** No cold start, able to recommend to users with unique tastes, able to recommend new and unpopular items, able to provide explanations (e.g. list of common features).

**CONS:** hard to find appropriate features, no recommendations for new users, too specific (overspecialization), unable to exploit quality judgments of other users.

## COLLABORATIVE FILTERING

Consider customer  $x$

Find set  $N$  of other users whose ratings are similar to  $x$ 's ratings.

Estimate  $x$ 's ratings based on ratings of users in  $N$

How to find "similar" users?

$\mathbf{r}_x$  = vector of user  $x$ 's ratings

- Jaccard similarity  $\rightarrow$  No, cause it ignores the value of ratings (it only see if both users rated the same item)

- Cosine similarity  $\rightarrow$  No, cause treats missing ratings as "bad ratings" (sparse matrix)

- Pearson Correlation coefficient:  $S_{xy} = \text{items rated by both users } x \text{ and } y$

$$\text{sim}(x, y) = \frac{\sum_{i \in S_{xy}} (r_{xi} - \bar{r}_x)(r_{yi} - \bar{r}_y)}{\sqrt{\sum_{i \in S_{xy}} (r_{xi} - \bar{r}_x)^2} \sqrt{\sum_{i \in S_{xy}} (r_{yi} - \bar{r}_y)^2}}$$

$\bar{r}_i = \text{avg rating of } r_i$

e.g.

ITEMS

	HP1	HP2	HP3	TW	SW1	SW2	SW3
Users	A 4 5	B OK 5	C 4	D not similar 2	E 1 4	F 5	G 5
A	4	OK	5	1	1	5	5
B	5	5	4	2	4	5	5
C				2	4		5
D				3			

As we can see  $\text{sim}(A, B) > \text{sim}(A, C)$ : we are interested in high rated items and similar so green are not similar and red is not important cause A rated it 1.

Jaccard similarity:  $\text{sim}(A, B) = \frac{1}{5}$ ,  $\text{sim}(A, C) = \frac{2}{4} \rightarrow \text{sim}(A, B) < \text{sim}(A, C) \rightarrow \text{NOT TRUE}$

Cosine similarity:  $\text{sim}(A, B) = 0.386$ ,  $\text{sim}(A, C) = 0.322 \rightarrow \text{sim}(A, B) > \text{sim}(A, C) \rightarrow \text{TRUE}$  but considers missing ratings "bad ratings" (too similar)

Solution:

	HP1	HP2	HP3	TW	SW1	SW2	SW3
A	2/3			5/3	-7/3		
B	1/3	1/3	-2/3				
C				-5/3	1/3	4/3	
D		0				0	

The main idea of Pearson Correlation coefficient is to calculate cosine similarity subtracting average ratings

e.g. user A's avg  $= \frac{(4+5+1)}{3} = \frac{10}{3} \rightarrow \text{HP1}_A = \frac{4 - 10}{3} = \frac{-2}{3}$

$\rightarrow \text{sim}(A, B) = 0.082$ ,  $\text{sim}(A, C) = -0.559 \rightarrow \text{sim}(A, B) > \text{sim}(A, C)$  and they are very different now!

Note: 0's in second matrix don't bother anyone

Now we have  $\mathbf{r}_x$  and  $N = k$  users most similar to  $x$  who rated item  $i$

Prediction for item  $i$  of user  $x$ :

$$r_{xi} = \frac{1}{k} \sum_{y \in N} r_{yi} \quad \text{or} \quad r_{xi} = \frac{\sum_{y \in N} S_{xy} \cdot r_{yi}}{\sum_{y \in N} S_{xy}} \quad (S_{xy} = \text{sim}(x, y)) \quad \text{or many other tricks...}$$

Some concept (user-user collaborative filtering) can be applied to item-item scenario:

for item  $i$  find similar items and estimate rating for item  $i$  based on ratings for similar items

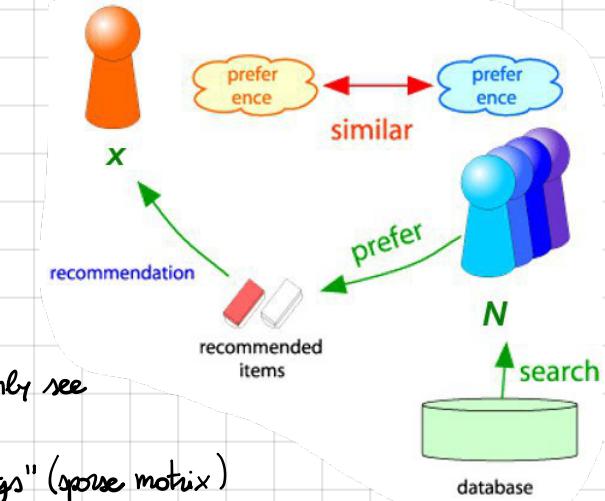
$$r_{xi} = \frac{\sum_{j \in N(i, x)} S_{ij} \cdot r_{xj}}{\sum_{j \in N(i, x)} S_{ij}}$$

$\left\{ \begin{array}{l} S_{ij} = \text{similarity of items } i \text{ and } j \\ r_{xj} = \text{rating of user } x \text{ on item } j \\ N(i, x) = \text{set of items rated by } x \text{ similar to } i \end{array} \right.$

It has been observed that item-item is better than user-user cause items are simpler while users have multiple tastes.

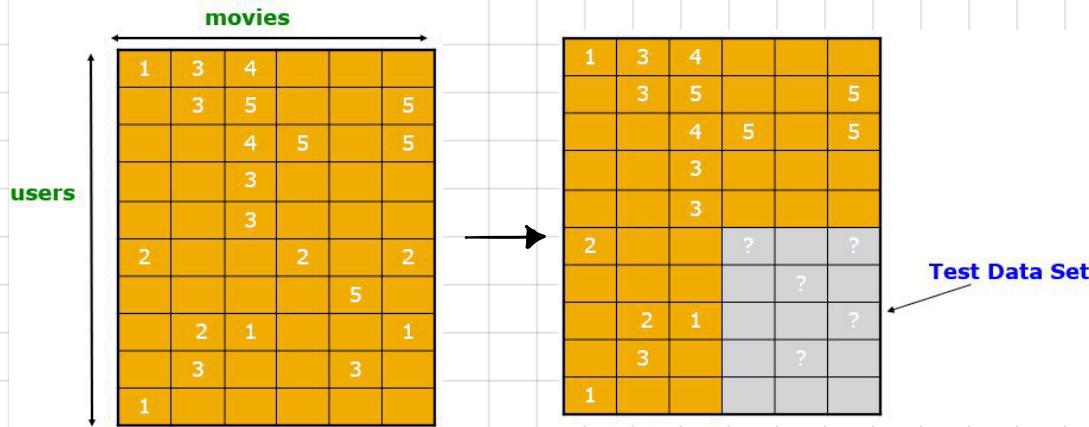
PRO: works for any kind of item

CONS: cold start (need data to find match), sparsity, cannot recommend items not already rated, tend to recommend popular items



In many case we can combine multiple methods.

But how do we evaluate performance? 3



Basically we compute predictions for test data set with known ratings: root-mean-square error (RMSE)

$$\sqrt{\sum_{x_i} (r_{xi} - \hat{r}_{xi})^2}$$
 where  $\hat{r}_{xi}$  is predicted and  $r_{xi}^*$  is the true rating of  $x$  on  $i$

Then we consider precision at top "10" cause errors on high ratings are more relevant than ones on low ratings.

### LOCAL & GLOBAL EFFECTS

In practice we get better estimates if we model deviations:  $\hat{r}_{xi} = b_{xi} + \frac{\sum_{j \neq i} s_{ij} (r_{xj} - b_{xj})}{\sum_{j \neq i} s_{ij}}$

$\mu$  = overall mean rating

$b_x$  = rating deviation of user  $x$  = (avg rating of user  $x$ ) -  $\mu$

$b_j$  = (avg rating of movie  $j$ ) -  $\mu$

e.g.

Mean movie rating = 3.7 , "The sixth sense" is 4.2 stars

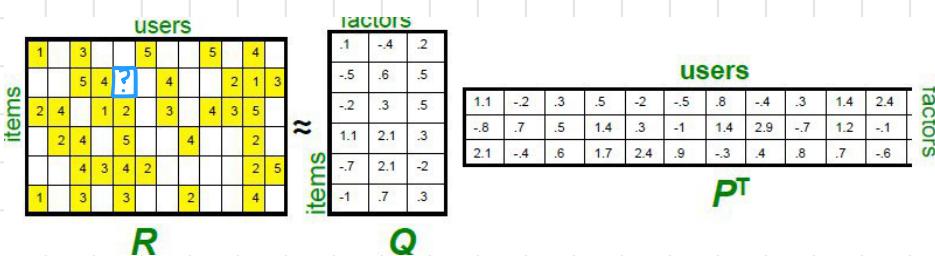
Joe usually rates 0.2 stars below avg.  $\rightarrow 4.2 - 0.2 = 4$  = baseline estimation

Joe didn't like related movie "Signs"  $\rightarrow 4 - 0.2 = 3.8$  stars = final estimation

$\hookrightarrow$  we are combining local (user  $x$ ) and global (all users) effects

To improve recommendations we want lower RMSE and try to build a system that works well on known ratings and hope it predict well the unknown ratings.  $\downarrow$

### LATENT FACTOR MODELS



Suppose we can do this:  $R = Q \cdot P^T$

The results is a 3D space in which users and items are 3D points.

$$\hat{r}_{xi} = q_i \cdot p_x = \sum_f q_{if} \cdot p_{xf} \quad q_i = \text{row } i \text{ of } Q \quad p_x = \text{column } x \text{ of } P^T \rightarrow \text{even missing ratings}$$

$$? = [-0.5, 0.6, 0.5] \cdot \begin{bmatrix} -2 \\ 0.3 \\ 2.4 \end{bmatrix} = 1 + 0.18 + 1.2 \approx 2.4$$

How do we find P and Q?

We want to minimize error  $\rightarrow \min_{P, Q} \sum_{(x_i, r_i) \in R} (r_{xi} - q_i \cdot p_x)^2$

To avoid overfitting we introduce **regularization**: allow rich model where there are sufficient data and shrink aggressively where data are scarce

$$\min_{P, Q} \underbrace{\sum_{\text{training}} (r_{xi} - q_i \cdot p_x)^2}_{\text{error}} + \left[ \lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right] \quad \text{minimize error and complexity}$$

$\lambda_1$  and  $\lambda_2$  are user set regularization parameters

$\Rightarrow$  **Stochastic Gradient Descent**:

initialize P and Q then iterate over the ratings and update factors:

for each  $r_{xi} \rightarrow \epsilon_{xi} = 2(r_{xi} - q_i \cdot p_x) \rightarrow$  derivative of the error

$q_i = q_i + \mu_1 (\epsilon_{xi} p_x - \lambda_2 q_i) \rightarrow$  update

$p_x = p_x + \mu_2 (\epsilon_{xi} q_i - \lambda_1 p_x) \rightarrow$  update

for until convergence:

for each  $r_{xi}$

compute gradient

Not very clear!

$\rightarrow$  all together:  $r_{xi} = \mu + b_x + b_i + q_i p_x$

to complicate this a bit more we can make  $b_x$  and  $b_i$  time dependent:  $r_{xi} = \mu + b_x(t) + b_i(t) + q_i p_x$

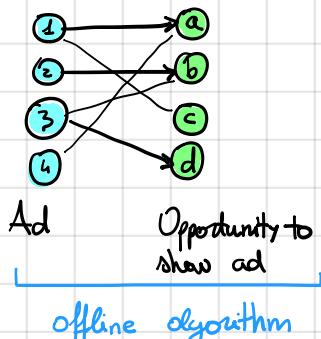
## WEB ADVERTISING

First of all we have to make a difference between offline and online algorithms. In the second ones the input isn't known at prior but one piece at time (like a continuous data stream) and "decisions" depends on its history.

The goal: show right ad depending on query (on its history)

Note: we don't know anything about future

e.g.



We talk about perfect matching when all vertices are matched  
Maximum matching when we have the max number of connections

Online version consist of graph not know upfront. For example we want to pair boys and girls according to their preferences ([tinder](#)?). Initially we only have boys set and in each round one girl make choices: she reveals her edges → we have to decide to pair her with someone or not at that time

Greedy algorithm: pair the new girl with only eligible boy, if there is none do not pair

How good is it? Compare it with an offline matching

→ competitive ratio =  $\min_{\text{all possible inputs } I} (|M_{\text{greedy}}| / |M_{\text{optimal}}|)$  where  $M_{\text{greedy}}$  is the matching of greedy algorithm and  $M_{\text{optimal}}$  is the one of offline optimal algorithm

Competitive ratio gives the greedy's worst performance over all possible inputs I

Consider  $M_g \neq M_{\text{opt}}$  and set G of girls matched in  $M_{\text{opt}}$  but not in  $M_g$

$$\bullet |M_{\text{opt}}| \leq |M_g| + |G|$$

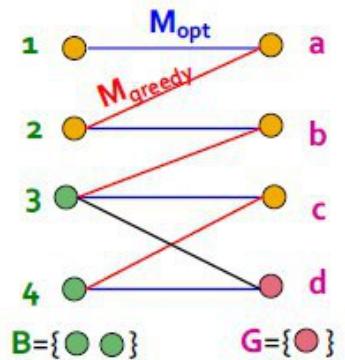
$$\bullet B = \text{boys linked to girls in } G \rightarrow |M_g| \geq |B|$$

(boys in B are already matched in  $M_g$  because a non matched boy linked to a non matched girl are matched by  $M_g$ )

Opt matches all girls in G to some boy in B →  $|G| \leq |B|$

$$\Rightarrow |G| \leq |B| \leq |M_g| \quad \text{worst is when } |G| = |B| = |M_g|$$

$$\Rightarrow |M_{\text{opt}}| \leq |M_g| + |G| \rightarrow |M_g| / |M_{\text{opt}}| \geq \frac{1}{2}$$



Performance-based Advertising: advertisers bid on search keywords called [adwords](#)

A stream of query arrives at search engine  $q_1, \dots, q_n$ ; several advertisers bid on each query; when  $q_i$  arrives search engine select a subset of ads to show

Goal: maximize profit → simple solution  $BID \cdot CTR$  (click rate = probability ad is clicked)

Limitation: CTR unknown and advertisers have limited budgets and bid on multiple query.

CTR is predicted or measured historically

Search engine has to be sure that the size of ads set is not larger than the number of ads per query, each advertiser has bid on the search query and each advertiser has enough budget to pay if ad is clicked.

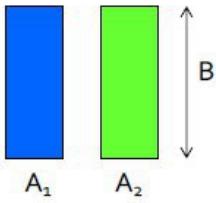
In the simplified scenario all budget and CTR are the same, we accept 1 ad for query and ad price = 1  
 $c.\text{ratio} = \frac{1}{2}$  (greedy algorithm)

Bad scenario: A bids on X, B bids on X and Y. Budgets = 4\$, query stream xxxx yyyy  
 → worst case greedy choice = BBBB--- cause of Y B finished budget and no one can bid for it → c ratio =  $\frac{1}{2}$   
 → optimal = AAAA BBBB

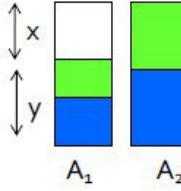
→ Balance algorithm by Mehta, Saberi, Vazirani and Vazirani (MSVV): for each query, pick the advertiser with largest unspent budget  
 → balance choice ABA BBBB c. ratio =  $\frac{M_{\text{opt}}}{R_{\text{opt}}} = \frac{6}{8} = \frac{3}{4}$

In a simple case with 2 advertisers  $A_1$  and  $A_2$  with budgets  $B \geq 1$  optimal solution exhaust both budgets.  
 Balance must exhaust at least one budget (if not, allocates more queries) → Balance profit = opt. profit -  $x$   
 $= 2B - x$ . ( $x$  = queries not allocated)

e.g.



Blue = queries allocated to  $A_1$  in optimal solution  
 Green = " to  $A_2$  "

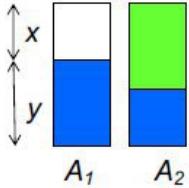


Suppose we exhaust  $A_2$  budget: balance profit =  $2B - x = B + y$   
 $B$  is  $A_2$  exhausted budget and  $y$  is the allocated profit on  $A_1$ .

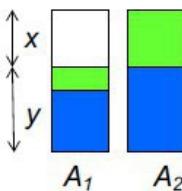
Goal: show  $y \geq B/2$

• Case 1: Balance assigns  $\geq B/2$  blue queries to  $A_1$

→  $y \geq B/2$



• Case 2: Balance assigns  $\geq B/2$  blue queries to  $A_2$   
 Consider the last blue query assigned to  $A_2$ : at this time  $A_2$ 's unspent budget  $\geq A_1$ 's because balance assigns queries to the one with greatest budget  
 This means  $A_1$  and  $A_2$  have the same number of queries.  
 $\rightarrow x \leq B/2$  and  $x+y=B \rightarrow y \geq B/2$



→ Profit =  $B+y$  but  $y \geq B/2 \rightarrow \text{Profit} = 3B/2$  competitive ratio =  $\frac{3B/2}{2B} = \frac{3}{4}$

In general case, worst competitive ratio is  $1 - \frac{1}{e} \approx 0.63$  THE BEST RESULT POSSIBLE

THE END

## EXERCISES

- Prove that for any graph the pagerank of each node is at least  $\frac{d}{N}$ .

$$\forall v: \pi_v \geq \frac{d}{N}$$

PageRank  $\rightarrow$  Ergodic MC  $\Rightarrow \pi^T = \pi^T P$

$$\pi_v = \frac{d}{N} + (1-d) \sum_{u: u \rightarrow v} \frac{\pi_u}{d_u}$$

def of pagerank

$\pi_v$   $\xrightarrow{\text{Teleport}}$   $\xrightarrow{\text{following links}} > \phi$

Probability

$$\Rightarrow \pi_v \geq \frac{d}{N} \quad \forall v$$

- An algorithm gives 30 docs as result in the following order

RRNRN NRRNN NRNNN NRNNR NNNNN NRNRN

- ① What is the precision on top 20?

$$P@20 = (\# \text{relevant docs retrieved}) / (\# \text{total retrieved}) = \frac{8}{20} = \frac{2}{5}$$

- ② What is the recall on top 20?

$$R@20 = (\# \text{relevant docs retrieved}) / (\# \text{total relevant docs}) = \frac{8}{10} = \frac{4}{5}$$

- ③ What is the F1 measure on top 20?

$$F1@20 = 2 \cdot \left[ \frac{P@20 \cdot R@20}{P@20 + R@20} \right] = \frac{8}{15}$$

- ④ What is MAP (mean average precision)?

$$\forall k: P@k \quad \forall q = \text{query}$$

$$\boxed{MAP(q)} = \frac{1}{m} \sum_i^m P@k \quad (\Rightarrow q \in Q \Rightarrow MAP(Q) = \frac{1}{|Q|} \sum_{q \in Q} MAP(q))$$

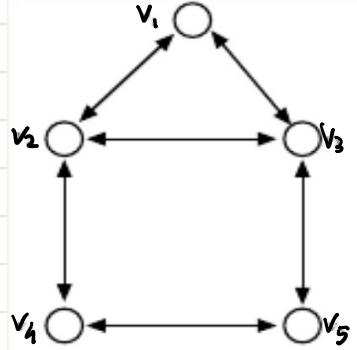
- Compute PR (pageRank scores) for  $\alpha=0$  and  $\alpha=\frac{1}{2}$  ( $\alpha = \text{prob to teleport}$ )

- $\pi_1 = \frac{\alpha}{5} + \frac{1-\alpha}{3} \pi_2 + \frac{1-\alpha}{3} \pi_3$
- $\pi_2 = \frac{\alpha}{5} + \frac{1-\alpha}{2} \pi_1 + \frac{1-\alpha}{2} \pi_4 + \frac{1-\alpha}{3} \pi_3$
- $\pi_4 = \frac{\alpha}{5} + \frac{1-\alpha}{2} \pi_5 + \frac{1-\alpha}{3} \pi_2$

$\pi_3 = \pi_2$   $\xrightarrow{\text{symmetry}}$

$$\pi_5 = \pi_4$$

$$\sum_i \pi_i = 1$$



- ①  $\alpha=0$ , random walk on undirected graph (bidirectional)

$$\forall v_i: \pi_i = \frac{d_i}{2m} \quad m = \# \text{undirected edges}$$

$$\pi_1 = \frac{2}{2 \cdot 6} = \frac{1}{6} = \pi_4 = \pi_5$$

$$\pi_2 = \pi_3 = \frac{1}{4}$$

[What if?  $\alpha=1 \rightarrow \pi_i = \frac{1}{N} \quad \forall i$ ]

- ② Just compute  $\pi_1, \pi_2$  and  $\pi_4$  with  $\alpha=\frac{1}{2}$

- True or false

- In a Boolean retrieval system, stemming never lowers precision. **false**  
e.g. { house } → stem. → house  $\Rightarrow$  house's postings list will contain both argument  $\rightarrow$  lower precision  
housing
- In a Boolean retrieval system, stemming never lowers recall **true**  
Stemming a term more documents are matched increasing the total number of relevant docs
- Stemming increases the size of the vocabulary. **false**  
Because more terms are "unified" with the same stem.

