

Web Information Retrieval Notes

Giuliano Abruzzo

Basato sugli appunti di Edoardo Puglisi

<https://github.com/machine1104/masterNotes>

July 13, 2020

Contents

1 Cap 1: Inverted Index & Boolean Retrieval	4
2 Cap 2: Document Ingestion	7
2.1 Tokenization	7
2.2 Stopwords	7
2.3 Normalization	8
2.4 Stemming and Lemmatization	8
2.5 Faster Postings List	8
3 Cap 3: Dictionaries and Tolerant Retrieval	9
3.1 Dictionaries	9
3.2 Wild-Card Queries	9
3.2.1 Spelling Correction	11
3.2.2 Phonetic Correction	12
4 Cap 4: Index Construction Algorithm	12
4.1 BSBI Algorithm	12
4.2 SPIMI Algorithm	13
4.3 Distributed Indexing	14
4.4 Dynamic Indexing	15
5 Index Compression Algorithms	17
5.1 Heaps' Law	17
5.2 Zipf's Law	17
5.3 Dictionary Compression	17
5.3.1 Dictionary as an array of fixed-width entry	18
5.3.2 Dictionary as string	18
5.3.3 Dictionary as string with blocking	18
5.4 Postings Compression	19
5.4.1 Variable Bytecode	20
5.4.2 Gamma Codes	20
6 Cap 6: Ranked Retrieval	21
6.1 Vector space model for scoring	22
6.1.1 Dot Products	22
7 Cap 7: Computing Scores	24
7.1 Efficient Scoring and Ranking schemes	24
7.1.1 Heap Tree	24
7.1.2 Inexact top-K Retrieval	24
7.1.3 Index Elimination	24
7.1.4 Champion Lists	25
7.1.5 Static Quality Score	25
7.1.6 Cluster Pruning	26
7.2 Combining components	27

7.2.1	Parametric and Zone Indexes	27
7.2.2	Tiered Indexes	27
7.2.3	Query-term Proximity	28
7.2.4	Query Parser	28
7.3	Putting all together	28
8	Cap 8: Evaluation of search results	29
8.1	Unranked Sets Measures	29
8.2	Evaluation of ranked retrieval results	30
8.3	Relevance Assessment	32
9	Cap 9: Text classification/categorization	34
9.1	Bayesian Methods	35
9.2	KNN, K-Nearest Neighbors Classification	38
9.3	SVM - Support Vector Machines	38
10	Cap 10: Link Analysis	40
10.1	Page Rank	41
10.1.1	Markov Chains	42
10.2	Topic-Specific Page Rank	45
10.3	Hubs and Authorities	46
11	Cap 11: Recommendations	48
11.1	Content-Based	49
11.2	Collaborative Filtering	49
11.3	Local & Global Effects	51
11.4	Latent Factor Model	51
12	Cap 12: Web Advertising	52
12.1	Performance-based Advertising	53
13	Cap 12: Exercises	54

1 Cap 1: Inverted Index & Boolean Retrieval

Information retrieval means finding *material* (usually *docs*) of an **unstructured** (a *data* which doesn't have an easy-for-a-computer *structure*) nature (usually *text*) that satisfies an information need from within large collections. The field of *information retrieval* also covers supporting *users* in *browsing* or *filtering* document collections. Given a set of documents, **clustering** is the task of coming up with a *good grouping* of the *documents* based on their contents. Given a set of topics, **classification** is the task of deciding which *class*, if any, each of a set of *documents* belongs to. *Information retrieval* systems can also be distinguished by the **scale** at which they operate, in *web search*, the system has to provide search over billions of *documents* stored on millions of computers instead at the other extreme there is *personal information retrieval*.

Information retrieval in the Web is different because:

- *Infos* can be **linked**;
- We have huge *sources* with different kinds of files;
- *Sources* can be not *homogeneous*;

Links can be used for suggestion, from a *page* to another one, every *web page* has a **rank** that depends on how many *links* point to it, and this *rank* is used by Google for **relevance criteria**. We can have also a **quality** difference between *infos* on the same *source*. *Web* in past was different, so now we can distinguish:

- **Web 1.0:**

- In 90s, *infos* on the *web* were almost always generated from a *source* instead *users* were quite passive, in fact web was used for reading, owning, consuming, so very few interactions of the *user*. The *web* was structured in a *hierarchical* way, called **Taxonomy**.

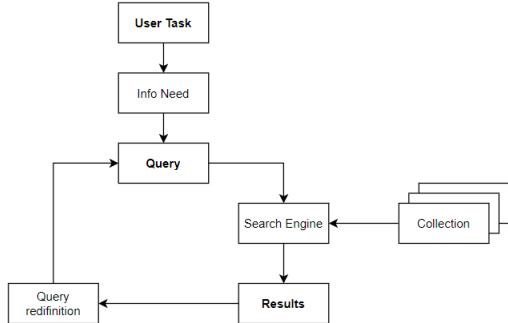
- **Web 2.0:**

- Now, a lot of *web content* is **user-generated** and there are a lot of interactions between *users* and *information* that we can retrieve. The structure is not anymore *hierarchical*, we have *tags* and every *tag* is a *category*, this structure is called **Folksonomy**.

If we consider the *web* as a **graph**, where the **nodes** are the *pages* and **edges** are the *hyperlinks* between them:

- **Crawling the web** means that by visiting this *graph* we will collect a **web corpus** (*multiple documents*), starting from a *page* and through this *links* we iterate on another *pages*. Often we don't want to *crawl* some *pages*, cause we don't need them, or we want to add some kind of *restrictions*, and we need to be sure to visit all *relevant pages* in order to avoid *information missing*, so this problem requires an huge amount of *computing power* and *competitiveness*;
- We also need **Indexing the web**, so we need to give a *structure* to the *web corpus* that we crawled;
- Last, we need some **Search Algorithms**, in order to solve *user's information need* using *queries*, they need to be based on **relevance** (results are coherent with the query) and **authority** (pages must give user real information);

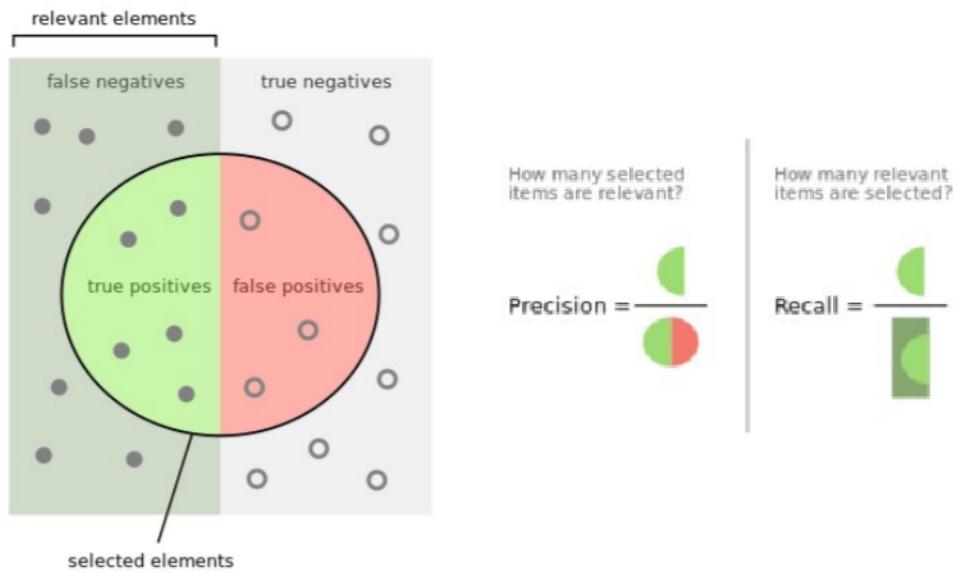
The **Meta-Search** is when we research using combined results from multiple *search engines*. A **collection** is a set of *documents*, not always *static*, and the **goal** is to collect *documents* that are relevant to *user's information* need. The **classic Search Model** is:



For example, the **user task** is: '*how to get rid of a mice in a politically correct way*', so the **info need** is '*remove mouse without killing it*', the **query** will be '*how trap mouse alive*'

Between **user task** and **query** there could be some *misperception* or *misperformulation* by the *search engine*. In order to measure a **search engine performance**, we have to compare the *ground truth* (what we expect from a search), with the *result* we actually obtain, and this is done through **Precision** and **Recall**:

- **Precision:** $\frac{\text{relevant docs}}{\text{returned docs}}$
- **Recall:** how many documents over the relevant ones are shown;



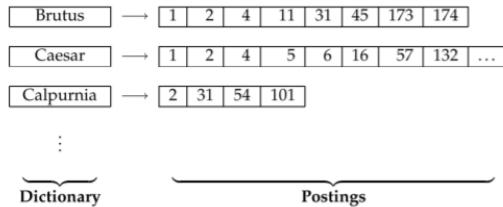
Usually we have a **trade-off** between them, so we need to obtain either a good *precision* or *recall* depending on the problem's type.

The first type of *query* we are going to see is the **Boolean IR query**, a *query* in which we research *keywords* in *documents*. For example we want to know which plays of Shakespeare contains ‘*brutus*’ and ‘*caesar*’ **and not** ‘*calpurnia*’. A way could be **grepping** (search a *string* in a *document*) all the *plays*, but this is **strongly inefficient** for large *docs* and it’s not flexible for other *matching operations*. Another option is using a **binary term-document incidence matrix**, where the *words* are the *rows*, and the *plays* are the *column*, and we will have a 1 in the cell if that *word* appear in that *play*:

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	...
Antony	1	1	0	0	0	1	
Brutus	1	1	0	1	0	0	
Caesar	1	1	0	1	1	1	
Calpurnia	0	1	0	0	0	0	
Cleopatra	1	0	0	0	0	0	
mercy	1	0	1	1	1	1	
worser	1	0	1	1	1	0	
...	...						

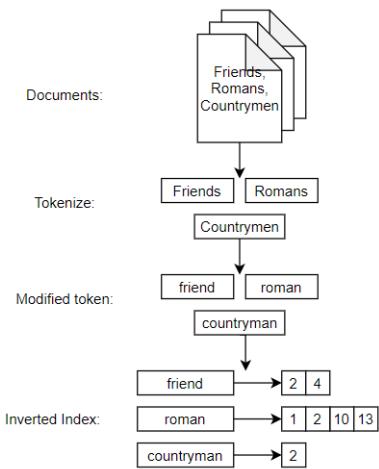
In order to *answer* the **query** we will take the vectors of *brutus*, *caesar* and *calpurnia* and we do a *bitwise and*:
 $101000 \text{ AND } 110111 \text{ AND } (\text{NOT } 010000)$
 $= 100100$, so the answer will be *Antony* and *Cleopatra* and *Hamlet*.

But if we have **big collection**, we will have a *matrix* too big to fit in a *computer memory*, and also it will be full of *zeroes*, with few *non-zero entries*. A much better representation is to record only the *1’s occurrences*, by using an **Inverted Index**:



This means that for each *term* we will have a *linked list* that records which *documents* the *term* occurs in. Every *document* must have a **docID** that identify it and these are *sorted*, so we avoid to use *fixed length array*.

In order to build an **inverted index** we need to:



- Collect the *documents* to be *indexed*;
- *Tokenize* the *text*, turning each *document* into a list of *tokens*;
- Do *linguistic preprocessing*, producing a list of *normalized tokens*, which are the *indexing terms*;
- Index the *documents* that each *term* occurs in by creating an *inverted index*, consisting of a *dictionary* and *postings*.

In *linguistic modules* we **stem** the text and also remove *stopwords*, not only *lowercase*. The core **indexing step** is *ordering* both *words* and *docID*, and often also the *frequency* of every *word* is added. So the *cost* for a *query* $word_1 \text{ AND } word_2$ is $O(X + Y)$ where X and Y are the *postings length* of the two *words*. The *intersect* of the two *posting lists* is the crucial operation, and it's also called **merging**. The **query optimization** is the process of selecting how to organize the work of answering a *query*:

- (... AND ...) AND (... AND ...): we will start from *intersect* the two smallest *posting lists*, so the next *intersection* will be more efficient;
- (... OR ...) AND (... OR ...): we use the *frequency* to estimate which *OR* we will use, usually the one which *term's frequency sum* is lower;

Sometimes we want a *query* to match a *phrase*, for example we want to answer queries with ‘*stanford university*’, we call these **Phrase Queries**. These *words* have to appear in these *precise order*, so a *posting list* is not good anymore, we could index every consecutive pair of *terms* in the *text* as a *phrase*, also called **Biword Indexes**, for example ‘*Friends, Romans, Countryman*’ generate two biwords ‘*friends romans*’ and ‘*romans countryman*’, where each of these is a *dictionary term*, but we will have too many *dictionaries* and *false positive*. So we need to use a **Positional Index**, in which we consider the **position** of the *term* in the *docs*:

- -Stanford: 2: 1,12,24; 4: 8,16,190; ...
- -University: 1: 17,19; 4: 17,191; ...

Where at the left we have *docID*, at the right the *positions* of the *word*. Using position we can even do **proximity search** and we can also combine *Biword* and *Positional*.

2 Cap 2: Document Ingestion

2.1 Tokenization

We said that in order to construct an *inverted index*, we need to **tokenize** the *text*. **Tokenization** is the process of chopping *character streams* into *tokens*, a **token** is a sequence of chars grouped together as a *semantic unit* for *processing*. **Tokenization** depends on how we want to process the *index* (by *words*, by *biwords*, ...), the most simple strategy is to split on *whitespaces*, but there could be the need of grouping special sets of *characters*. The problem of *tokenization* is finding the correct *token* to use, for ‘*Finland's*’ which is the best *tokenization*? ‘*Finland AND s*’, ‘*Finlands*’ or ‘*Finland's*’, these issues are **language-specific**, so we require to know the *language* of the document, we can have also problems with *dates* and *numbers*.

2.2 Stopwords

Stopwords are *words* that are extremely common and their *semantic content* is almost useless, they are generated by *collection frequency* (number of times each *term* appears in *documents*). Using a *stop list* significantly reduces the number of *postings* that a *system* has to store, and a lot

of the time not indexing *stop words* does little harm: keyword searches with terms like *the* and *by* don't seem very useful. However, this is not true every time, for example, the meaning of *flights to London* is likely to be lost if the word *to* is stopped out. For most modern *IR systems*, the additional cost of including *stop words* is not that big, neither in terms of *index size* nor in terms of *query processing time*.

2.3 Normalization

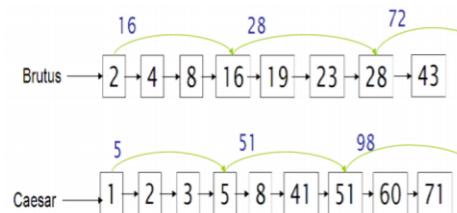
Normalization is the process of canonize *tokens* so that matches occur despite superficial differences in the character sequences of the *tokens*. The most standard way to *normalize* is to implicitly create **equivalence classes**, for example we delete dots (*U.S.A* → *USA*) or hyphens (*anti-discriminatory* → *antidiscriminatory*). It can also be applied on accents or variants and there could misunderstanding too, but even on *synonyms*. During *normalization* all *token* characters are reduced to *lower case*.

2.4 Stemming and Lemmatization

For grammatical reasons, *documents* are going to use different forms of a *word*, such as *organize*, *organizes*, and *organizing*. The goal of both **stemming** and **lemmatization** is to reduce *inflectional forms*, for instance: *am, are, is* → *be*, or *car, cars, car's, cars'* → *car*. **Stemming** usually refers to a heuristic process that *chops off* the ends of *words* in the hope of achieving this goal correctly most of the time, includes the removal of *derivational affixes*. **Lemmatization** usually refers to doing things properly with the use of a *vocabulary* and *morphological analysis of words*, normally aiming to remove inflectional endings only and to return the base or *dictionary form* of a *word*, which is known as the **lemma**. The most common algorithm for *stemming English* is **Porter's Algorithm**, which consists of 5 phases of *word reductions*, applied sequentially, some typical rules are: substitution like *SSES* → *SS* or *IES* → *I*, weight of word sensitive rules, ($m > 1$) Ement: *replacement* → *replac* or *cement* → *cement*.

2.5 Faster Postings List

In the basic **postings list intersection**, if the list lengths are m and n , the *intersection* takes $O(m + n)$ operations. One way to do better than this is to use a **skip list** by augmenting *postings lists* with **skip pointers** (at indexing time), which allow us to avoid *processing parts* of the *postings list* that will not figure in the *search results*. We have a *trade-off*, more *skips* means shorter *skip spans* (space of a skip), but also means lot of *comparisons* and lots of space storing *skip pointers*. Building effective *skip pointers* is easy if an *index* is relatively static, it is harder if a *postings list* keeps changing because of updates.



3 Cap 3: Dictionaries and Tolerant Retrieval

3.1 Dictionaries

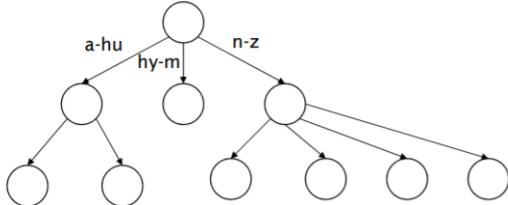
Data Structure for *posting lists* are usually *dictionaries*, but we need to store in memory efficiently, and also we need to quickly look up elements at *query time*. There are two main choices: **Hash Tables** and **Trees**:

- **Hash Tables:**

- Each *vocabulary term* (the **key**) is **hashed** into an *integer* over a large enough *space* (in order to avoid *hash collisions*), at *query time* we **hash** each *query term* separately and following a *pointer* to the corresponding *postings*. So the *lookup* is really fast $O(1)$, the problem is that it's not easy to find *similarity* (minor variants of a *query term*), and there is no *prefix search*. If the size of the *vocabulary* keeps growing we may need to *rehash* everything to avoid *collisions*.

- **Trees:**

- The best-known *search tree* is the **binary tree**: the search for a *term* begins at the *root* of the *tree*, each *internal node* (including the *root*) represents a *binary test*. A **lookup** costs $O(\log N)$ if the *tree* is **balanced** (number of *terms* under the two *sub-trees* of any *node* is equal or differ by one), else is $O(N)$. It solves the *prefix search* problem, but we can have problems of *re-balancing*, in fact when we insert new *terms* or we delete old ones the *tree* needs to be re-balanced. One approach to avoid it is to allow the number of *sub-trees* under an *internal node* to vary in a fixed interval, these are called **B-Trees**, in which every *internal node* has a number of children in the interval $[a, b]$. Unlike *hashing*, **search trees** demand that the characters used in the document collection have a **prescribed ordering**.



3.2 Wild-Card Queries

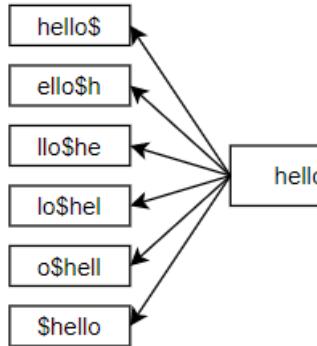
Wildcard Queries are used in any of the following situations:

- The user is uncertain of the spelling of a *query term*, like *Sydney* vs *Sidney*, so the wildcard query is S^*dney ;
- The user is aware of multiple variants of spelling a term, and seeks documents containing any of the variants, like *color* vs *colour*;

- The user seeks documents containing variants of a term that would be caught by *stemming*, but is unsure whether the search engine performs stemming like *judicial* vs. *judiciary*, the wildcard query *judicia**;
- The user is uncertain of the correct rendition of a foreign word or phrase, like *Universit** *Stuttgart*;

A *query* such as *mon** is **trailing wildcard query**, because the *** symbol occurs only once, at the end of the search *string*. This query is easy to find in a B-tree, but the query **mon* is very harder, called **leading wildcard queries**, we can use a reverse *B-Tree* in which each *root-to-leaf* path corresponds to a term written backwards, for example *lemon* will be represented by the path *root-n-o-m-e-l*. Using a regular *B-Tree* together with a *reverse B-Tree* we can handle a general case *se*mon* where the regular B-tree is used to find the prefix *se*, and the reverse B-tree is used to find the suffix *mon*, and then we will **intersect** the two *sets* obtained, but this solution is pretty expensive.

A more efficient way is using **Permutation Index**, a form of *inverted index*, we will use a special symbol \$ used to mark the end of a term.



If we have to search *hel*o*, the key is to rotate in such a way the *** symbol is at the end of the string, we set $X = hel$ and $Y = o$, we search for $Y\$X*$ so $o\$hel^*$. we search in the regular b-tree, all the words that start with *hel* and ends with *o* (? probably?).

If we have a query like *fi*mo*er* first we ignore the *mo*, and we search for *fi*er* that permuted is *er\$fi**, then we will filter the words that doesn't contain *mo*. The problem with permutation index is that its dictionary becomes quite large.

In order to solve this problem, we can use the **Bigram Index** (or **K-gram Index** where *k-gram* is a sequence of *k* characters) for a single word, for example in ‘*April is the cruelest month*’, we get the bigrams: ‘\$a,ap,pr,ri,il,l\$,i,is,s\$,t,th,he,e\$,c,cr,ru,ue,el,le,es,st,t\$, \$m,mo,on,nt,h\$’. We maintain a *second inverted index* from *bigrams* to *dictionary terms* that match each *bigram*. For example *\$m* will point to *mace*, *madden*, *mo* will point to *among*, *amortize*, *on* to *along*, *among*. The query *mon** can now be run as *\$m AND mo AND on*, but these can return also the word *moon*, so we need another filter at the end. This method is *fast* and *space efficient* (compared to *permutation*). The processing of a *wildcard query* can be quite expensive because of the added *lookup* in the *special index*, *filtering* and finally the *standard inverted index*. A **search engine** may support such rich functionality, but most commonly, the capability is hidden behind an interface that most users never use.

3.2.1 Spelling Correction

There are two methods to solve the problem of **correcting spelling** errors in *queries*: the first based on **edit distance** and the second based on **k-gram overlap**. There are two basic principles underlying most spelling correction algorithms:

- Correct spellings for a *misspelled query*, so we choose the **nearest one**;
- When two correctly spelled *queries* are tied, select the one that is **more common**;

We will focus on two types of *spelling correction*: **isolated-term correction** and **context-sensitive correction**:

- **Isolated-Term correction:**
 - In **isolated-term correction**, we attempt to correct a single *query term* at a time, but such *isolated-term correction* would fail to detect, for instance, that the query *flew form Heathrow* contains a misspelling of the term *from*, because each term in the *query* is correctly spelled in isolation.
 - **Edit Distance:**
 - * Given two strings s_1 and s_2 , the **edit distance** between them is the minimum number of *edit operation* to convert one string to the other. *Edit distance* is also called **Levenshtein distance**, for example the edit distance of *cat* and *dog* is 3. There is also the **weighted edit distance** in which we have a *weight* based on the *character* involved (char m mis-typed as n more than z , so replacing m with n is a smaller edit than z), but these requires a *weight matrix* as input. In order to compute the *weighted edit distance* we need $O(|s_1| \times |s_2|)$ where $|s_1|$ denotes the length of the string. In order to correct the *spelling*, given a *query*, first we enumerate all *char sequences* within a preset *weighted edit distance*, then we intersect this set with list of *correct words*, then we show these *words* to the user as a suggestion.
 - **N-Gram Overlap:**
 - * We can use the **n-gram index** to retrieve *vocabulary terms* that have many *n-grams* in common with the *query*, and the *retrieval process* is essentially that of a single scan through the *postings* for the *n-grams* in the *query string*. For example if the text is *november* (trigrams are *nov*, *ove*, *vem*, *emb*, *mbe*, *ber*) and the query is *december* (trigrams are *dec*, *ece*, *cem*, *emb*, *mbe*, *ber*) so we have 3 *trigrams* that overlap. The measure of *overlapping* is given by the **Jaccard Coefficient** by two sets X and Y and it is: $|X \cap Y| / |X \cup Y|$, this will be a value between 0 and 1, so we will choose the *terms* over a certain *threshold*.
- **Context-Sensitive Correction:**
 - *Isolated-term correction* would fail to correct typographical errors where all *query terms* are correctly spelled, for example ‘*I flew form Milan*’, ‘*form*’ is an error (*from*). The simplest way to correct these errors is to enumerate *corrections* of each of the *query terms* even though each *query term* is correctly spelled, then try substitutions of each *correction* in the phrase. For each such *substitute phrase*, the *search engine* runs the *query* and determines the number of *matching results*. This enumeration can be expensive, several

heuristics are used to reduce the *search space*. We wanna rank alternatives probabilistically: $\text{argmax}_{\text{corr}} P(\text{corr}|\text{query})$, that with bayes rules is: $\text{argmax}_{\text{corr}} P(\text{query}|\text{corr}) * P(\text{corr})$, where *query* is the noisy channel, and *corr* the language model.

3.2.2 Phonetic Correction

Misspellings that arise because the user types a *query* that sounds like the *term target*. The main idea here is to generate, for each term, a **phonetic hash** so that *similar-sounding terms hash* to the same value. Algorithms for such *phonetic hashing* are commonly collectively known as **Soundex algorithms**, common steps are:

- We turn every *term* to be indexed into a *4-character reduced form*, we build an **inverted index** from these reduced forms to the original *terms* called *soundex index*;
- We do the same with the *query terms*;
- When the *query* calls for a *soundex match*, search this *soundex index*;

4 Cap 4: Index Construction Algorithm

In order to answer a *query*, we need to build an *inverted index* for a set of terms, this construction process is called **Index Construction** or **Indexing**, that takes advantage of *secondary storage*. It's important to note that: *main memory >> secondary storage*, it's way faster. There are some consideration:

- To optimize **transfer time** we have a big *chunk* of data and not several small *chunks*;
- **Disk I/O** is *block-based* (of fixed lights);
- **Fault tolerance** is handled with *replication* (several instead of a single fault-tolerant machine);
- The **main memory** is the better;

We need larger average *word token size* to handle *words*, especially if we strip out *stopwords*. The goal is construct the *inverted index*, but we can't do the whole sorting in *main memory*, we need intermediate steps.

4.1 BSBI Algorithm

BSBI Algorithm or **Blocked Sort-Based Indexing Algorithm**, is an **external sorting algorithm** which try to minimize the number of *random disk seeks* during *sorting*, for example to sort *100M postings* (made of pairs *term-docID*) we define *10 blocks* of *10M postings* each, in such a way that each *block* can fit in the *memory* and:

- For each *block*:
 - Accumulate *postings*;
 - Sort in *memory*;
 - Write to *disk*;
- Then merge all the *blocks sorted* in order to obtain the *final index*;

Algorithm 2 Blocked sort-based indexing. The algorithm stores inverted blocks in files f_1, \dots, f_n and the merged index in f_{merged} .

```

1: function BSBI_CONSTRUCTION
2:    $n \leftarrow 0$ 
3:   while all documents have not been processed do
4:      $n \leftarrow n + 1$ 
5:      $block \leftarrow PARSE\_NEXT\_BLOCK()$ 
6:      $BSBI\_INVERT(block)$ 
7:      $WRITE\_BLOCK\_TO\_DISK(block, f_n)$ 
8:   end while
9:    $MERGE\_BLOCKS(f_1, \dots, f_n, f_{merged})$ 
10: end function

```

The algorithm parses *documents* into *termID–docID* (*termID* is a unique mapped in a dictionary from the term) pairs and accumulates the pairs in *memory* until a *block* of a fixed size is full, the *block* is then inverted and written to *disk*. **Inversion** involves two steps, first we sort the *termID–docID pairs*, next, we collect all *termID–docID* pairs with the same *termID* into a *postings list*, where a *posting* is simply a *docID*, and the result is then written to *disk*. In the final step, the algorithm simultaneously merges the *blocks* into one large *merged index*. It's time complexity is $O(T \log T)$, cause the step with the highest time complexity is *sorting*, and T is an upper-bound for the number of *items* we must sort. The actual *indexing time* is usually dominated by the time it takes to parse the document (*parse block function*), and the final merge (*merge blocks function*). The key decision is the **block size** that need to be *optimized*.

4.2 SPIMI Algorithm

Blocked sort-based indexing has excellent *scaling properties*, but it needs a *data structure* for mapping *terms* to *termIDs*. For very large collections, this *data structure* does not fit into *memory*. An alternative is **Single-Pass In-Memory Indexing or SPIMI**, which uses separate *dictionaries* for each *block* (so we don't need to maintain *term-termID* mapping across *blocks* like in *BSBI*), and the other idea is **don't sort**, so we accumulate *postings* in *postings lists* as they occur, so we can generate a complete *inverted index* for each *block* that we will merge into one *big index*. Each *postings list* is dynamic so it's immediately available to collect *postings*, and this has two advantages, it's *faster* (no sorting) and saves *memory* cause we keep track of the term a *posting lists* belongs to, so the *termID* of *postings* doesn't need to be stored. *SPIMI* has also an important component, the **compression**, in fact both *postings* and the *dictionary terms* can be store compactly on *disk*. The *time complexity* is $O(T)$ since no sorting is required and all operations are linear in the size of the collection. *SPIMI_INVERT* is called repeatedly on the *token stream* until the entire *collection* has been processed. *Tokens* are processed one by one during each successive call of *SPIMI_INVERT*. When a *term* occurs for the first time, it is added to the *dictionary*, and a new *postings list* is created, at the end it returns this *postings list* for subsequent occurrences of the term. (NOT SURE we don't sort *postings* cause they are already sorted by the *docsID* that is incremental when received)

Algorithm 3 Inversion of a block in single-pass in-memory indexing.

```
1: function SPIMI_INVERT(token_stream)
2:   output_file = NEW_FILE()
3:   dictionary = NEW_HASH()
4:   while free memory available do
5:     token  $\leftarrow$  next(token_stream)
6:     if term(token)  $\notin$  dictionary then
7:       postings_list = ADD_TO_DICTIONARY(dictionary, term(token))
8:     else
9:       postings_list = GET_POSTINGS_LIST(dictionary, term(token))
10:    p2  $\leftarrow$  next(p1)
11:   end if
12:   if full(postings_list) then
13:     postings_list = DOUBLE_POSTINGS_LIST(dictionary, term(token))
14:   end if
15:   ADD_TOPOSTINGS_LIST(postings_list, docID(token))
16: end while
17: sorted_terms  $\leftarrow$  SORT_TERMS(dictionary)
18: WRITE_BLOCK_TO_DISK(sorted_terms, dictionary, output_file)
19: return output_file
20: end function
```

4.3 Distributed Indexing

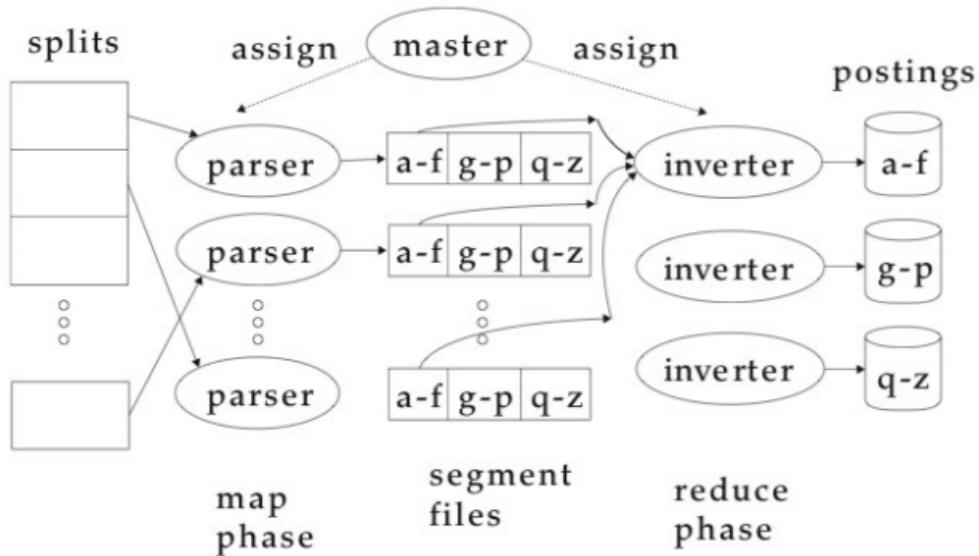
Collections are often so large that we cannot perform **index construction** efficiently on a *single machine*, in this case we will use **distributed indexing algorithms** for *index construction*. The result of the construction process is a *distributed index* that is partitioned across several machines. We have a **Master Machine** which is a *fault-tolerant machine*, and **Idle Machines**, that are part of a *pool* and they have *indexing tasks* in parallel assigned by the *master*.

The distributed index construction method is an application of **MapReduce** a general architecture for *distributed computing*. A robust *distributed indexing* needs the work to be divided in re-assignable chunks, a **master node** directs the process of assigning and reassigning *tasks* to individual **worker nodes** (*idle machine*).

First, the input data, are split into *n splits* where the size of the *split* is chosen to ensure that the work can be distributed evenly and efficiently, these *splits* are not preassigned to *machines*, but are instead assigned by the *master node*: as a *machine* finishes processing one *split*, it is assigned the next one, if a machine dies, the *split* it is working on is simply reassigned to another machine. In general, *MapReduce* breaks a large computing problem into smaller parts by *recasting* it in terms of manipulation of *key-value pairs*, in the form of $\langle \text{termID}, \text{docID} \rangle$.

The **map phase** consists of *mapping splits* of the input data to *key-value pairs*, the machines that execute the map phase are called **parser**. Each *parser* writes its output to local intermediate

files, the *segment files*. For the **reduce phase**, we want all values for a *given key* to be stored close together, this is achieved by *partitioning the keys* into j term partitions (range of letters like *a-f*) and having the *parsers* write *key-value pairs* for each term partition into a separate *segment file*. Collecting all values for a *given key* into one list is the task of the **inverters**, the *master* assigns each *term partition* to a different *inverter*. Finally, the list of *values* is sorted for each *key* and written to the final *sorted postings list*. *Parsers* and *inverters* are not separate sets of machines, the *master* identifies *idle machines* and assigns *tasks* to them, the same machine can be a *parser* in the *map phase* and an *inverter* in the *reduce phase*.



4.4 Dynamic Indexing

Most collections are *modified frequently* with *documents* being *added*, *deleted*, and *updated*, new *terms* need to be added to the *dictionary*, and *postings lists* need to be updated. The simplest way to achieve this is to *periodically reconstruct* the **index** from scratch, this is a good solution if the number of changes over time is small and if enough resources are available to construct a new *index* while the old one is still available for *querying*.

If there is a requirement that new *documents* be included quickly, the simplest approach is to use *two indexes*, a large **main index** on *disk*, and a small **auxiliary index** that store new documents in *memory*. *Searches* are run across both *indexes*, *deletions* are stored in an **invalidation bit vector**, then we filter out *deleted documents* before returning the *search result*. When the *auxiliary index* becomes too large, we merge it with the *main index*. Unfortunately, this scheme is infeasible because most *file systems* cannot efficiently handle very large numbers of files and merges are computationally expensive. The simplest alternative is to store the *index* as one large file as a concatenation of all *postings list*, but this would generate a lot of files.

An alternative is the **Logarithmic Merge**, that reduces the cost of *merging indexes* over time. We maintain *several indexes, each twice larger than the previous one*, and the smallest Z_0 is maintained in memory , and the others (I_0, I_1, \dots) are on disk, when Z_0 becomes too large, we will write it to disk as I_0 , or we merge it with I_0 if it already exist and write the merger to I_1 and so on. It use a *binary number* to save which index is full for example 1011 (positions 3 2 1 0), means that I_2 has space, instead the others are full. For example if we have a as size of a *memory index*, then on the *disk* we will have *indexes* of size $2a$, $4a$ and so on, so *logarithmic scale*. The **time complexity** of *index construction* in the worst case is $O(T \log T)$, because $\log T$ is the number of indexes with T number of postings, so for a query we need to merge $O(\log T)$ indexes and the worst case is $O(T \log T)$ because each of T posting is merged $O(\log T)$ times. This is more efficient than the $O(n^2)$ of the previous alternative. (*supponiamo che una parola ha T postings nella sua posting list, se devo fare il merge su tutte le I_i allora avrò $T \log T$ dato che i vari I_i sono ognuno il doppio del precedente*)

```

LMERGEADDTOKEN(indexes, Z0, token)
1    $Z_0 \leftarrow \text{MERGE}(Z_0, \{\text{token}\})$ 
2   if  $|Z_0| = n$ 
3     then for  $i \leftarrow 0$  to  $\infty$ 
4       do if  $I_i \in \text{indexes}$ 
5         then  $Z_{i+1} \leftarrow \text{MERGE}(I_i, Z_i)$ 
6           ( $Z_{i+1}$  is a temporary index on disk.)
7            $\text{indexes} \leftarrow \text{indexes} - \{I_i\}$ 
8         else  $I_i \leftarrow Z_i$  ( $Z_i$  becomes the permanent index  $I_i$ .)
9            $\text{indexes} \leftarrow \text{indexes} \cup \{I_i\}$ 
10          BREAK
11       $Z_0 \leftarrow \emptyset$ 

LOGARITHMICMERGE()
1    $Z_0 \leftarrow \emptyset$  ( $Z_0$  is the in-memory index.)
2    $\text{indexes} \leftarrow \emptyset$ 
3   while true
4   do LMERGEADDTOKEN(indexes, Z0, GETNEXTTOKEN())

```

Because of this *complexity of dynamic indexing*, some large search engines adopt a **reconstruction from-scratch strategy**. They do not construct *indexes dynamically*, instead, a new *index* is built from scratch periodically. *Query processing* is then switched from the *new index* and the *old index* is deleted.

5 Index Compression Algorithms

Compression techniques for *dictionary* and *inverted index* are essential for efficient *IR systems*, one benefit of *compression* is straightforward: we need less *disk space*. Compression can be **lossy** (discard some *infos*, like *stopwords* or *lowering*), or **lossless** (all *infos* are preserved). We will use some variables in statistics:

- **N**: documents;
- **L**: average number of tokens per document;
- **T**: average number of bytes per term, non-positional postings;
- **M**: word types, average number of bytes per token;

5.1 Heaps' Law

In order to get the number of *distinct terms M* in a collection is to use the **Heaps's Law** which estimates *vocabulary size* as a function of collection size:

$$M = kT^b$$

Where T is the number of *tokens* in the collection, k and b are two parameters usually: $30 \leq k \leq 100$ and $b \cong 0.5$, this law's suggests that the *dictionary size* continues to increase with more *documents* in the *collection*, and the size of the *dictionary* is quite large for large collections. If a *term* is *frequent*, it will not characterize a document so much, and the contrary, so we can *rank* terms by their frequency (their *relevance*)

5.2 Zipf's Law

We also want to understand how *terms* are *distributed* across documents, a commonly used model of the distribution of *terms* in a collection is **Zipf's Law**. It states that, if t_1 is the most common *term* in the collection, t_2 is the next most common, and so on, then the **collection frequency** cf_i of the i -th most common *term* is proportional to $1/i$:

$$cf_i \propto \frac{1}{i}$$

So if the most *frequent term* occurs cf_1 times, then the second most *frequent term* has half as many occurrences, the third most frequent *term* a third as many occurrences, and so on. The *frequency* decreases very rapidly with *rank*. The goal is to encode most *frequent terms* to smaller size *encoding*. So the most frequent *word* will appear cf_1 the second one will appear $\frac{1}{2}cf_2$ times and so on.

5.3 Dictionary Compression

One of the *primary factors* in determining the *response time* of an *IR system* is the number of *disk seeks* necessary to process a *query*. If parts of the **dictionary** are on *disk*, then many more *disk seeks* are necessary in *query evaluation*, so, the main goal of compressing the *dictionary* is to fit it in *main memory*.

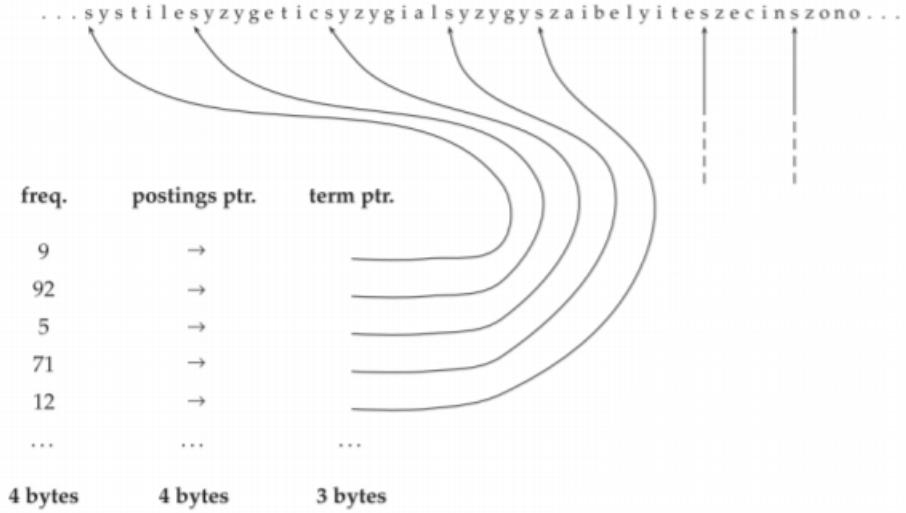
5.3.1 Dictionary as an array of fixed-width entry

The simplest *data structure* for the *dictionary* is to sort the *vocabulary lexicographically* and store it in an *array of fixed-width entries*. We allocate 20 bytes for the *term* itself, 4 bytes for its *document frequency*, and 4 bytes for the *pointer* to its *postings list*. For example for *Reuters-RCV1* (a dataset with 400.000 elements) we need $M \times (20 + 4 + 4) = 400.000 \times 28 = 11.2MB$ for storing the *dictionary* in this scheme.

space needed:	20 bytes	4 bytes	4 bytes
	term	document frequency	pointer to postings list
a	656,265	→	
aachen	65	→	
...	
zulu	221	→	

5.3.2 Dictionary as string

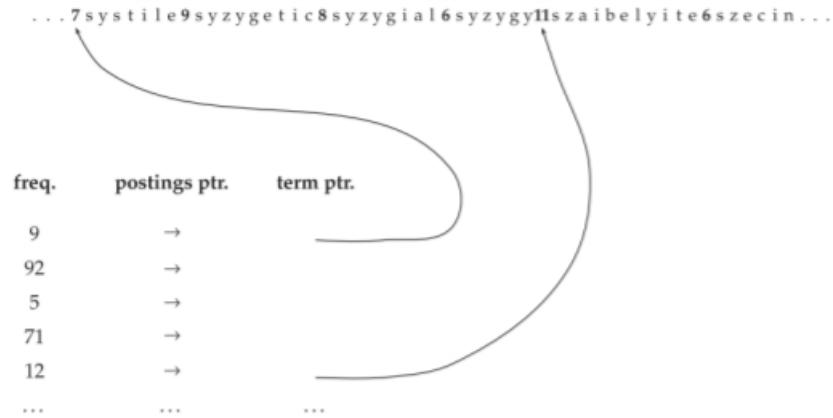
Using *fixed-width entries* for *terms* is wasteful, in fact average length of a *term* in English is about eight characters and we also we need a way of storing *terms* with more than twenty characters. We can overcome these problems by storing the *dictionary terms* as one **long string of characters**, so we use a pointer to demarcate the end of the current term. For *Reuters-RCV1* we will use $(4 + 4 + 3 + 8) \times 400.000 = 7.6 MB$, where 8 bytes are the average of the *term*.



5.3.3 Dictionary as string with blocking

We can also compress the *dictionary* by grouping *terms* in the *string* into **blocks** of size k and keeping a *term pointer* only for the first *term* of each *block*. We store the *length* of the *term* in the *string* as an *additional byte* at the beginning of the *term*, we thus eliminate $k - 1$ *term pointers*,

but we need an additional k bytes for storing the *length*. Dictionary of *Reuters-RCV1* is reduced by 0.5 MB, to 7.1 MB. (*se prima avevamo un puntatore da tre byte per ogni termine ora abbiamo un puntatore ogni blocco e un byte per ogni per la lunghezza e risparmiamo circa 5 byte per ogni blocco*)



There is a **tradeoff** between *compression* and *lookup time*, if we increase *block size* k we have better *compression*, but we have more *lookup time* for a *word* in that *block*. Consecutive *entries* in an *alphabetically sorted* list share common prefixes, this observation leads to **front coding**. A *common prefix* is identified for a *sub-sequence* of the *term list* and then referred to with a *special character*. In the case of Reuters, front coding saves another 1.2 MB.

One block in blocked compression ($k = 4$) ...
8automata8automate9automatic10automation

↓

...further compressed with front coding.
8automat*a1◦e2◦ic3◦ion

A sequence of terms with identical prefix ("automat") is encoded by marking the end of the prefix with * and replacing it with ◦ in subsequent terms. As before, the first byte of each entry encodes the number of characters.

5.4 Postings Compression

Posting lists total size is about 10 times larger than the *total dictionary size*, so we need to compress their *size*. One idea is to store **gaps** (*difference between the two document id index*) instead of *docsID*, for example the *word computer* → 28154, 28159, 28160,... we can store is like *computer* → 28154, 5, 43, With the *original one* we use 20 bits for each *docID*, instead now with *gaps* that are usually shorter we use way less than 20 bits. For an economical representation of this distribution of *gaps*, we need a *variable encoding method* that uses fewer bits for *short gaps*. To

encode small numbers in less space than large numbers, we look at two types of methods: **bytewise compression** and **bitwise compression**.

5.4.1 Variable Bytecode

Variable Bytecode or **VB encoding** uses an *integral number* of bytes to encode a *gap*, and last 7 bits of a byte are “*payload*” and encode part of the *gap*. It dedicate 1 bit to be a **continuation bit** c , if the gap G fits within 7 bit, binary encode it in the 7 available bits and we set $c = 1$ else we use more than a block (*at the beginning of each byte we can have 0 or 1, if 0 then is a number which requires more than a block, if it's 1 it fit in the 7 bit*). For example we have **docID** 824 and 829, so the gap is 5, 824 in binary is 1100111000 so it doesn't fit in 7 bit, we will write it as: **0000110 10111000** where the black numbers are *continuation bits* (so we use two blocks), instead *gap* 5 in binary is **10000101** cause it fit in 7 bits. The length of *7+1 bits* can be changed depending on how the gaps are, for small *gaps 4 bits block* are usually better.

5.4.2 Gamma Codes

VB codes use an adaptive number of *bytes* depending on the size of the *gap*, **Bit-level codes** (also called **Gamma Codes**) adapt the length of the *code* on the *finer grained bit level*. The simplest bit-level code is **unary code**, the *unary code* of n is a string of n 1s followed by a 0 (for example $3 = 1110$, $4 = 11110$), but this is not a very efficient code. **Gamma code** uses *length* and *offset* of a *gap G*. **The offset is the gap in binary without the leading 1**. For example 13: binary is 1101 and offset is 101, length encodes the length of offset in unary code, for 13 the length of offset is 3 bit that is 1110 in unary (*ao prendi l'offset, 101, vedi quant'è lungo, 3, lo schiaffi in unario, 1110*), so we have that the **gamma code** of 13 is: 1110 concat 101 = 1110101.

The *offset length* is: $\lfloor \log_2 G \rfloor$ bits, instead the length of *length part* is: $\lfloor \log_2 G + 1 \rfloor$ bits, so the length of the *entire code* is $2 \times \lfloor \log_2 G + 1 \rfloor$, this means that *gamma codes* are always of *odd length* and they are within a factor of 2 of what we claimed to be the optimal *encoding length* $\log_2 G$. In fact assuming the 2^n *gaps G* with $1 \leq G \leq 2^n$ are all equally likely, the *optimal encoding* uses n bits for each G , so some *gaps* cannot be encoded with fewer than $\log_2 G$ bits and our goal is to get as close to this **lower bound** as possible.

Gamma codes have also some important *properties*:

- **Prefix-Free:** no *gamma codes* is the *prefix* of the other, this means that there is always a unique *decoding* of a sequence of gamma codes, and we don't need *delimiters* between them.
- **Universal:** we can use it far any data distribution;
- **Parameter-Free:** there are no parameters in this procedure;

However machines have *word boundaries* so compressing at a *bit level* can be *expensive (slow)*, for this reason the **VB code** can be a better solution.

6 Cap 6: Ranked Retrieval

Thus far, our queries have all been **Boolean**, *documents* either match or don't, this is good for *expert users* with precise understanding of their needs and of the *collection*, but it's not good for the majority of *users*. In fact for these *queries* we have:

- *Docs* that either match or don't (*too strict*);
- Thousands of results (*inefficient to present all of them in a web page*);
- Too few or too many results;

Users need a **ranked series** of results and, let's say, the first 10 ones, so not a complete result of thousands of elements. **Ranking** is done with respect to specific criteria of **relevance**, and is measured through a **score** in $[0, 1]$ assigned to each *query-document pair*. For example in the **1-Word Query**: the most the *query term* appears in a *doc*, the higher the *ranking* will be, if it doesn't appear *score* is 0. There are also some alternatives:

- **Jaccard Coefficient**:
 - A commonly used measure of overlap of two sets A and B , $JACCARD(A, B) = \frac{|A \cap B|}{|A \cup B|}$, A and B don't have to be the same size. The problem is that this method doesn't take *frequency* into consideration, and doesn't use *rare terms* (that are more informative than *frequent terms*).
- **Bags of words**:
 - A method in which we represent each *document* as a **count vector**, with *frequency* for every *term*, the problem is that we don't consider *order* so we will not use it;

The **term frequency** $tf_{t,d}$ of term t in document d , is defined as the number of times that t occurs in d . We need to use it when we compute *query-document match scores*, but we don't want that *relevance* is increased proportionally with the *term frequency*. We use a **log frequency weight** for t in d (can also be called $wf_{t,d}$):

$$w_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d} & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

So the **score** for a *document-query* pair is: $matching-score(q, d) = \sum_{k \in (q \cap d)} (1 + \log_{10} tf_{t,d})$, which means the sum of the **weighted frequencies** of that *term* in both *query* and *document*.

We also want to use the **frequency of the term in the collection** for *weighting* and *ranking*, in fact let's consider a *term* that is really *rare* in a *collection*, it appears in only one *document*, this probably means that such *document* has a great relevance for the *collection*. So we want *high weights* for **rare terms** and *low weights* for *frequent words*. We will use **document frequency** to factor this into computing the *matching score*, that is the number of *documents* in the *collection* that the *term* occurs in.

We call df_t the number of *docs* in which *term t* occurs, it is an *inverse measure* of the **informativeness** of term *t*, so we define $idf_t = \log_{10} \frac{N}{df_t}$ the measure of the *informativeness* of a term (where *N* is the number of *docs* in the *collection*). The best known weighting scheme is called **tf-idf weight** that is the product of *tf weight* (*term frequency* in a document) and *idf weight* (*rarity measure* of a *word* in the collection):

$$tf - idf_{t,d} = tf_{t,d} \times idf_t = (1 + \log_{10} tf_{t,d}) \times \log_{10} \frac{N}{df_t}$$

This scheme assign to term *t* a weight that is:

- Highest when *t* occurs many times within a small number of *documents*;
- Lower when *t* occurs fewer times in a document, or occurs in many *documents*;
- Lowest when the *term* occurs in all *documents*;

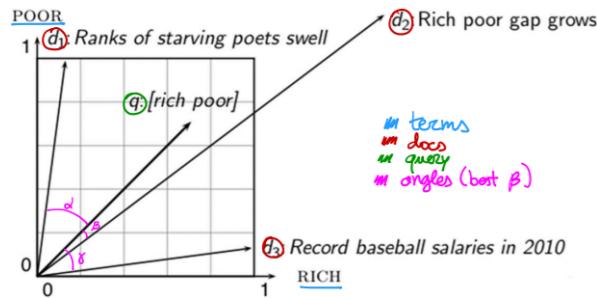
It's important to note that this formula uses the $wf_{t,d}$ (so the *log weighted frequency*), not just the *term frequency* itself.

6.1 Vector space model for scoring

The representation of a set of *documents* as *vectors* in a common *vector space* is known as the **vector space model**.

6.1.1 Dot Products

We denote by $\vec{V}(d)$ the **vector** derived from *document d*, with one component in the *vector* for each *dictionary term*. The set of *documents* in a collection then may be viewed as a set of *vectors* in a *vector space*, in which there is one *axis* for each *term*. This *representation* loses the relative *ordering* of the terms in each document. We need to quantify the *similarity* between two documents in the *vector space*, we can consider the **magnitude** (*modulo*) of the *vector difference* between two *document vectors*, but two *documents* with very *similarity* content can have a significance vector difference if one document is much longer than the other.



Another way to quantify the similarity is to compute the **cosine similarity** of the *vector representations* of two documents:

$$sim(d_1, d_2) = \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| \cdot |\vec{V}(d_2)|}$$

Where the numerator represents the dot product of the two vectors, while the denominator is the product of the **Euclidean Length** (that will length normalize the two vectors). This measure is the *cosine* of the angle θ between the two vectors. We can also rewrite as: $\vec{v}(d_1) = \frac{\vec{V}(d_1)}{|\vec{V}(d_1)|}$ and $\vec{v}(d_2) = \frac{\vec{V}(d_2)}{|\vec{V}(d_2)|}$ so the equation will be:

$$sim(d_1, d_2) = \vec{v}(d_1) \cdot \vec{v}(d_2)$$

We can also view the *query* as a *vector*, we consider a *query* q , this *query* into the *unit vector* $\vec{v}(q)$, so the idea is to assign to each document d a score equal to:

$$\cos(\vec{q}, \vec{d}) = sim(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| \cdot |\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i \cdot d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \cdot \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

Where:

- q_i is the **tf-idf weight** of the term i in the **query**;
- d_i is the **tf-idf weight** of the term i in the **document**;
- $|\vec{q}|$ and $|\vec{d}|$ **length** of \vec{q} and \vec{d} ;

Algorithm 4 The basic algorithm for computing vector space scores.

```

1: function COSINE_SCORE( $q$ )
2:   float  $Scores[N] = 0$ 
3:   Initialize  $Length[N]$ 
4:   for all query term  $t$  do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5:     for all pair( $d, tf_{t,d}$ ) in postings list do  $Scores[d] += wf_{t,d} \times w_{t,q}$ 
6:     end for
7:   end for
8:   Read the array  $Length[d]$ 
9:   for all  $d$  do  $Scores[d] = Scores[d] / Length[d]$ 
10:  end for
11:  return Top  $K$  components of  $Scores[]$ 
12: end function

```

In a typical setting we have a collection of *documents* each represented by a *vector*, a *text query* represented by a *vector*, and a positive integer K . We seek the K *documents* of the collection with the **highest vector space scores** on the given *query*. Typically, we seek these ordered by

decreasing score (usually $K = 10$). The algorithm computes the *vector space scores*, and for each term t of the *query* it will update the *score* of the *document* by adding in the contribution from term t . This process is also known as **term-at-a-time scoring** or *accumulation*, and the N documents of the *array scores* are therefore known as *accumulators*. It is wasteful to store the *weight* of term t in document d since storing this *weight* may require a floating point number. The most complex and expensive operation is the *extraction* of the top K scores, this requires a *priority queue data structure*, often implemented using a **heap**. Each of the K top scores can be extracted from the *heap* at a cost of $O(\log N)$ comparisons.

7 Cap 7: Computing Scores

We saw in the *Cosine Score Algorithm*, that we return the first K elements of *scores array*, so the most relevant ones, but now we want to do this without calculating all the *cosines*. So we are doing the **K-nearest neighbors problem** for a *query vector*. In general we don't know how to do this for *high dimensional spaces*, but it is solvable for *short queries*. We assume that each *query term* will occur only once, this means that for *ranking* we don't need to *normalize* query vector.

7.1 Efficient Scoring and Ranking schemes

7.1.1 Heap Tree

Let J be the total number of *non-zero cosine documents*, we need to find the K best of these J documents. The solution is to use a **heap**: in fact it takes $2J$ for constructing and $2 \log J$ for reading each *winner*. The *bottleneck* is that the *cosine computation* in this case has to be done for all the *tree elements* in order to build the *tree*.

7.1.2 Inexact top-K Retrieval

Another method is called **inexact top-K retrieval**, that is not from the user's perspective a bad thing, in fact we avoid all these calculations, even if we lose some *accuracy*: we find a set of **contenders** A with $K < |A| \ll N$, so A doesn't contain the top K but it has many *docs* from among the top K , and in such a way we return the top K docs in A , so we can think A as a **pruning non-contenders**, and this scheme is also used for *non-cosine based functions of scoring*.

7.1.3 Index Elimination

With **Index Elimination** scheme we use *two pruning*, we select only *query terms* with high *idf weight* and we only select *documents* which contains several *query terms*, for a **multi-term query** q :

- We only consider *documents* containing *terms* whose *idf weight* exceeds a preset threshold. In the *postings traversal* (*scorrere le varie posting list tipo un ciclo for*), we only traverse the *postings* for *terms* with high *idf*. The *postings lists* of *low-idf terms* are generally long, thus the set of documents for which we compute *cosines* is greatly reduced. *Low-idf terms* are treated as *stop words* and do not contribute to *scoring*.
- We only consider *documents* that contain many of the *query terms*. A danger of this *scheme* is that by requiring all (or even many) *query terms* to be present in a *document* before

considering it for *cosine computation*, we may end up with fewer than K candidate *documents* in the *output*.

7.1.4 Champion Lists

The idea of **champion lists** is to pre-compute, for each term t in the *dictionary*, the set r of documents with the *highest weights* for t . For *tf-idf weighting*, these would be the r documents with the highest *tf* values for term t . We call this set the **champion list** for term t .

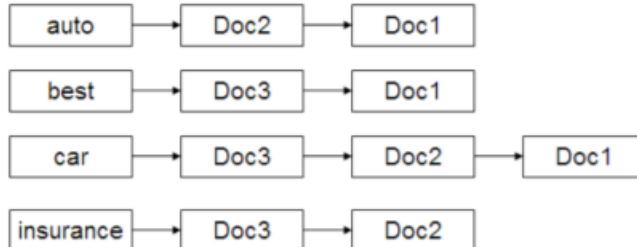
Given a query q we create a set A as follows: we take the union of the *champion lists* for each of the *terms* comprising q . We restrict *cosine computation* to only the *documents* in A . One issue here is that the value r is set at the time of **index construction**, whereas K is *application dependent* and may not be available until the *query* is received, so as a result we may find ourselves with a set A that has fewer than K documents.

7.1.5 Static Quality Score

In order to **top-ranking documents** we need to guarantee two important *properties*:

- **Relevance:** modeled by *cosine scores*;
- **Authority:** is *query-independent property* of a document that indicates its *validity*;

In many *search engines*, we have available a measure of **quality** $g(d) \in [0, 1]$ for each document d that is *query-independent* and thus **static**:



A static quality-ordered index. In this example we assume that Doc1, Doc2 and Doc3 respectively have static quality scores $g(1) = 0.25$, $g(2) = 0.5$, $g(3) = 1$.

The **Net Score** for a document d is a combination of $g(d)$ with the *query-dependent score*:

$$NetScore(q, d) = g(d) + \text{cosine}(q, d) = g(d) + \frac{\vec{V}(q) \cdot \vec{V}(d)}{|\vec{V}(q)| \cdot |\vec{V}(d)|}$$

In this form, the static quality and the query score have equal contribution. Now we seek the top k docs by the net-score by using two parameters: α and β in order to give more importance to a term instead of the other:

$$NetScore(q, d) = \alpha \cdot g(d) + \beta \cdot \text{cosine}(q, d)$$

In order to get the first K documents, we need to **order** the documents with respect to $g(d)$, because this raises the probability to find the most relevant *docs* early in *posting traversal*, so we have better *performance*. There are three ways of acting:

- **Global Champion List:**

- We maintain for each term t a **global champion list** of the r documents with the highest values for $g(d) + tf - idf_{t,d}$ score, that will be sorted by a common order, so at *query time* we only compute the *net scores* for documents in the union of these *global champion lists*. We seek the *top-K results* from only the *docs* in these *champions lists*;

- **High and Low List:**

- We maintain for each term t an **high list**, that contains the *documents* with the highest values for t , and a **low list** which contains the other *documents* containing t , we first scan only the *high list* of the *query terms*, if we obtain scores for K documents in the process we terminate, if not we continue in the *low list*;

- **Impact Ordering:**

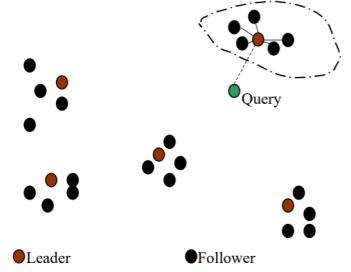
- In all the *postings lists* we order the documents by some *common ordering*. With **Impact order** we compute only *documents* with $wf_{t,d}$ high enough (should be *weighted log frequency* $1 + \log_{10} tf_{t,d}$, *così de botto senza senso, sulle slide ogni tanto con tf intende wf e i due termini sembrano interscambiabili*), so we will sort each *postings list* by $wf_{t,d}$, this means that all the *postings lists* are in a different *order* (a *document* can be very relevant for a term, and useless for another). In order to compute the *score* for picking off top K there are two ways:
 - * **Early termination:** when traversing t 's *postings*, we stop after a fixed r docs or when $wf_{t,d}$ is too low, so we take the **union** of the resulting sets of *docs*, one for each *query term*, and we compute the *score* only for *documents* in this union;
 - * **Idf-Ordered terms:** when considering *postings* of *query terms*, we look at them in order of decreasing **idf** (since high *idf* terms contribute most to the score). As we update *score contribution* for each *query term*, we stop if *document* scores relatively unchanged (?).

7.1.6 Cluster Pruning

Cluster Pruning is a technique in which we have a *pre processing step* during which we **cluster** the *document vectors*, then at *query time*, we consider only *documents* in a small number of *clusters* as *candidates* for which we compute *cosine scores*. In the *pre processing step* we take \sqrt{N} *documents* randomly, and we call them **leaders**, for all the other, called **followers**, we compute the *nearest leader*.

The expected number of *followers* for each *leader* is $\approx \frac{N}{\sqrt{N}} = \sqrt{N}$. The *query processing* proceeds as follows:

- Given a **query** q , find the **leader** L closest to q , this means computing *cosine similarity* between q and the \sqrt{N} *leaders*;
- We seek for K nearest documents from the *leader's* \sqrt{N} *followers*;



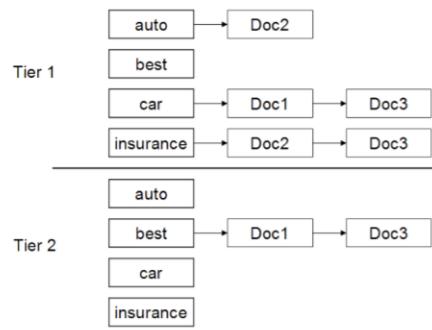
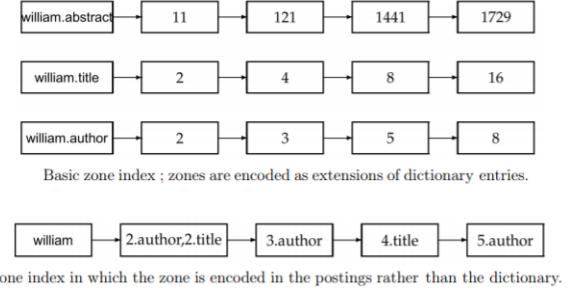
7.2 Combining components

7.2.1 Parametric and Zone Indexes

Some *docs* have specific *term* with *special semantics* called **metadata** which includes **fields** such as *format*, *author*, *year*, ... , and sometimes we want to search by these *fields* so we can build some **parametric indexes**, one for each *field*, in such a way we can select only the *documents* matching a specified *field* (like match a *date* in a *query*). There are also **zones**, similar to *fields*, but the contents of a *zone* can be arbitrary *free text*, like the *document title*. Whereas the *dictionary* for a *parametric index* comes from a *fixed vocabulary*, the *dictionary* for a *zone index* must structure whatever *vocabulary stems* from the text of that *zone*.

7.2.2 Tiered Indexes

We may occasionally find ourselves with a set of *contenders* that has fewer than K *documents*. A common solution in this case is the use of **Tiered Indexes**, which break *postings* up into a *hierarchy* of lists, from the most important to the last, this order can be done by $g(d)$ (*measure of quality of a document*) or another measure. At *query time* we seek *docs* in *top tier*, if it's not enough we go below.



Tiered indexes. If we fail to get K results from tier 1, query processing “falls back” to tier 2, and so on. Within each tier, postings are ordered by document ID.

7.2.3 Query-term Proximity

Especially for *free text queries*, users prefer a *document* in which most or all of the *query terms* appear **close** to each other, w is the smallest *window* in a *doc* which contains all the *query terms* (the *distance* between the *query terms* in the doc). The smallest w is, the better the *document* matches the query. Such **proximity-weighted scoring** functions are a departure from pure *cosine similarity* and closer to the *soft conjunctive semantics*.

7.2.4 Query Parser

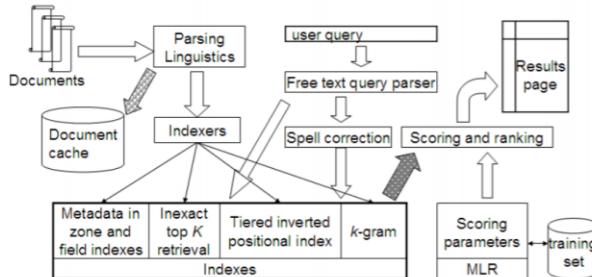
Common *search interfaces* tend to mask *query operators* from the end user, inviting **free text queries**. Typically, a **query parser** is used to translate the *user-specified keywords* into a *query* with various *operators* that is executed against the underlying *indexes*. This execution can entail *multiple queries* against the underlying *indexes*, for example, the *query parser* may issue a stream of *queries*:

- Run the *user-generated query string* as a **phrase query**, rank them by *vector space scoring* using as *query* the *vector* consisting of , for example, the 3 terms *rising interest rates*;
- If fewer than K *documents* contain the phrase *rising interest rates*, run the two **2-term phrase queries** *rising interest* and *interest rates*, rank these using *vector space scoring*;
- If we still have fewer than K results, run the *vector space query* consisting of the three individual *query terms*;

Each of these steps will compute a *score*, so we must combine these using an **aggregate scoring function** that accumulates evidence of a *document's relevance* from multiple sources.

7.3 Putting all together

Document will be parsed and will be applied a **language processing** (*tokenization, stemming, ...*). The resulting **tokens** will feed into two modules: a copy of each parsed *document* will go in a **document cache**, instead a second copy is fed to the **indexers**, that will create *indexes* like **zone and field indexes**, **tiered positional index**, **indexes for spelling correction** and other **tolerant retrieval**, and structures for accelerating **inexact top-K retrieval**. A **free text user query** is sent down to the *indexes* both directly and through a module for generating *spelling-correction candidates*. **Retrieved documents** (dark arrow) are passed to a **scoring module** that computes scores based on **machine-learned ranking (MLR)**. Finally, these **ranked documents** are rendered as a **results page**.



8 Cap 8: Evaluation of search results

8.1 Unranked Sets Measures

A *search engine* can be easily **evaluated** by measurable factors like:

- How fast does it **index** (*docs/hour*);
- How fast does it **search** (*latency as a function of index size*);
- Expressiveness of **query language**;

But the *key measure* is **User Happiness**, this cannot be quantified as the above factor, but it's something crucial for the *validity* of a *search engine*, it includes *speed of response*, *user-friendly UI*, and of course **relevance**. In order to measure ad hoc IR (*information retrieval*) we need a **test collection** composed by three things:

- A benchmark **document collection**;
- A benchmark **suite of queries**;
- A set of relevance *judgments*, usually binary assessment of either **relevant** or **non-relevant** for each *query-document pair*;

In fact for each *document*, with respect to a *user information need*, will be assigned a *binary classification* as *relevant* or *non-relevant*, and this decision is referred as the **ground truth** judgment of *relevance*. A *document* is **relevant** if it addresses the *stated information need*, not because it just happens to contain all the *words* in the *query*.

The two most *frequent* and basic measures for *information retrieval* effectiveness are **precision** and **recall**:

- **Precision, P**, is the fraction of *retrieved documents* that are *relevant*:

$$\text{Precision} = \frac{\#\text{(relevant items retrieved)}}{\#\text{(retrieved items)}} = P(\text{relevant} \mid \text{retrieved})$$

- **Recall, R**, is the fraction of *relevant documents* that are *retrieved*:

$$\text{Recall} = \frac{\#\text{(relevant items retrieved)}}{\#\text{(relevant items)}} = P(\text{retrieved} \mid \text{relevant})$$

These notions can be summarized in this table:

	Relevant	Nonrelevant
Retrieved	true positives (tp)	false positives (fp)
Non retrieved	false negatives (fn)	true negatives (tn)

We can also write *precision* and *recall* as:

- $P = tp / (tp + fp)$ $F = tp / (tp + fn)$

An alternative is using **accuracy** = $(tp + tn) / (tp + fn + gn + tn)$, but this is useless for *IR*, because the *true negative documents* set is huge, and a system tuned to maximize *accuracy* will give a lot of *false positive documents*. *Precision* and *recall* will trade off against one another, for example you can always get a *recall* of 1 but with very low *precision* by retrieving all *documents* for all *queries*, in general we want to get some amount of *recall* while tolerating only a certain percentage of *false positive*.

Another measure of **relevance** that trades off *precision* versus *recall* is called **F-Measure**, which is the weighted harmonic mean of *precision* and *recall*:

$$F = \frac{1}{\alpha \cdot \frac{1}{P} + (1 - \alpha) \cdot \frac{1}{R}} = \frac{(\beta^2 + 1) \cdot P \cdot R}{\beta^2 \cdot P + R} \text{ where } \beta^2 = \frac{1 - \alpha}{\alpha}$$

Where $\alpha \in [0, 1]$ and $\beta \in [0, \infty]$, the most used is the so called **Balanced F-Measure** which uses $\alpha = 1/2$ and $\beta = 1$ commonly written as F_1 so the formula is simplified as:

$$F_{\beta=1} = \frac{2 \cdot P \cdot R}{P + R}$$

Using values $\beta < 1$ will emphasize *precision*, instead values $\beta > 1$ will emphasize *recall*. We use **harmonic mean** because 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*:

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.

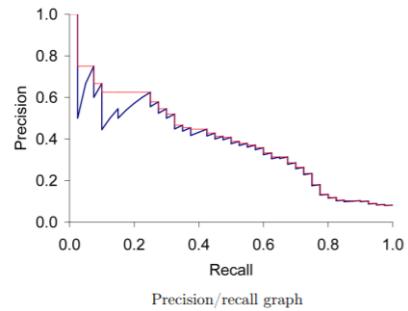
However if you drive for the same distance at each speed -- say 10 miles -- then the average speed over 20 miles is the harmonic mean of 30 and 40, about 34.3mph.

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.

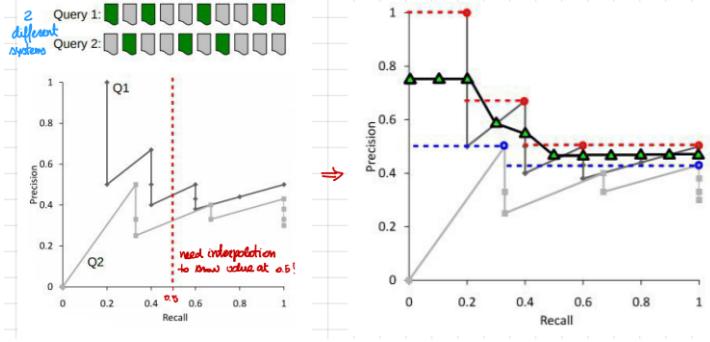
8.2 Evaluation of ranked retrieval results

Precision, *recall*, and *F*, are measures made for **unranked sets**, in fact they are computed using *unordered sets of documents*. In a **ranked retrieval context**, the set of *retrieved documents* are naturally given by the top K retrieved *documents*, so *precision* and *recall* values can be plotted to give a **precision-recall curve**: if the $(k + 1)$ -th document retrieved is *non-relevant*, then *recall* is the same of the top K documents, but *precision* has dropped, instead if it *relevant* then both *precision* and *recall* increase, and the curve jags to the right.



In order to remove the *jiggles* the standard way is using the **interpolation precision** p_{interp} . The idea is, if locally *precision* increases with increasing *recall*, then you take the max of *precision* to right of value. At certain *recall level* r , the *interpolation* is defined as the **highest precision** found for any recall level $r' \geq r$:

$$p_{interp}(r) = \max_{r' \geq r} p(r')$$



Praticamente trovi il miglior bilancio tra precision e recall trovando i punti ottimali (i triangoli) dove ottieni il massimo tra le due

Examining the entire *precision-recall curve* is very informative, but there is often a desire to reduce this *information* down to a few numbers. The traditional way of doing this is the **11-point interpolated average precision**, which for each *information need*, the *interpolated precision* is measured at the **11 recall levels** of 0.0, 0.1, 0.2, ..., 1.0, then for each *recall level* we will calculate the arithmetic *mean* of the *interpolated precision* for each *information need* in the test collection.

There are also other measures that have become more common, like:

- **MAP:**

- The **Mean Average Precision**, for a *single information need*, is the average of the *precision* obtained for the set of top K documents, each time a *relevant document* is retrieved. It avoid to use *interpolation*, since it use **fixed recall levels**.

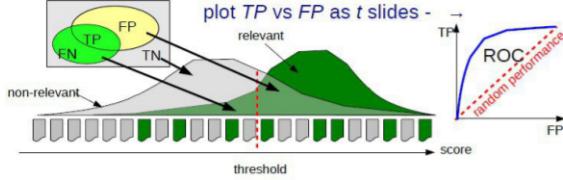
- **DCG:**

- The **Discounted Cumulative Gain**, is a measure of *ranking quality*, given n documents, each has a rating r_1, r_2, \dots, r_n :

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

8.3 Relevance Assessment

- 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):



In order to properly evaluate a *system*, the *test information* need to be sure that the *documents* are **relevant** for the usage of the *system*, this is a *time-consuming* and *expensive process* involving human beings. For large modern collections, it is usual for *relevance* to be assessed only for a subset of the *documents* for each *query*. The most standard approach is **pooling**, where *relevance* is assessed over a subset of the *collection* that is formed from the top K *documents* returned by a number of different *IR systems*, but the problem is a human is not a device, all have different *judgments system*. In the social sciences, a common measure for **agreement** between judges is the **Kappa Measure**, designed for *categorical judgments* and a simple *agreement rate for the rate of chance agreement*:

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

Where $P(A)$ is the proportion of times the judges agreed, $P(E)$ is the the proportion of times they would be expected to agree by chance. The *Kappa value* will be 1 if two *judges* always agree, 0 if the agree only at the rate given by chance, and negative if they are worse than random, and if there are more than two judges, it is normal to calculate an average *pairwise kappa value*. Usually a Kappa vale $\in [\frac{2}{3}, 1]$ is acceptable. An example:

$$\begin{aligned} P(A) &= 370/400 = 0.925 \\ P(\text{non-relevant}) &= (10+20+70+70)/800 = 0.2125 \\ P(\text{relevant}) &= (10 + 20 + 300 + 300)/800 = 0.7878 \\ P(E) &= P(\text{non-relevant})^2 + P(\text{relevant})^2 = \\ &0.2125^2 + 0.7878^2 = 0.665 \\ \kappa &= (0.925 - 0.665)/(1 - 0.665) = 0.776 \end{aligned}$$

Number of docs	Judge 1	Judge 2
300	Relevant	Relevant
70	Nonrelevant	Nonrelevant
20	Relevant	Nonrelevant
10	Nonrelevant	Relevant

Another problem with the *relevance-based assessment* is the distinction between *relevance* and **marginal relevance**: whether a *document* still has distinctive *usefulness* after the user has looked at certain other *documents*, even if a *document* is highly relevant, its information can be completely **redundant** with other *documents* which have already been examined. Maximizing *marginal relevance* requires returning *documents* that exhibit diversity and novelty.

- **Interactive Relevance Feedback:**

- We make a *query*, we take the first results set and ask user to select **relevant** and **non-relevant documents**, so the user will see the result summaries and will click the *docs* that he thinks are *relevant*, so we improve the *overall performance*. The best method is the *Rocchio Feedback*.

- **Query Expansion:**

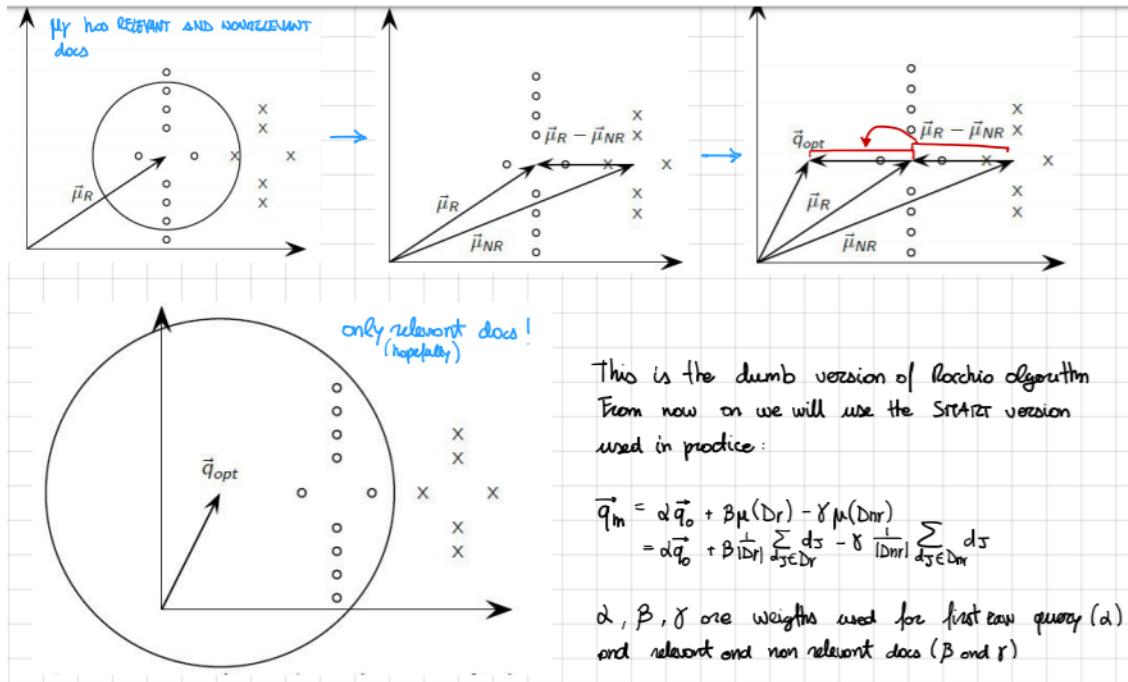
- With this method we improve the *retrieval results* by adding *synonymous* or *related term* to the *query*.

The main idea is to do a *local on-demand analysis* for an *user query* with **relevance feedback** and *global one* (e.g. of collection) to produce *thesaurus* (*dizionario dei sinonimi*) 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*. The search engine computes a new representation of information need, hopefully better, and runs a new query with new results.
- **Centroid** is the center of mass of a set of points that summarize their characteristics (we represent *docs* as *vector/points* in a high dimensional space): $\vec{\mu}(D) = \frac{1}{|D|} \sum_{d \in D} \vec{v}(d)$.
- **Rocchio Algorithms** chooses the *query* \vec{q}_{opt} defined as:

$$\vec{q}_{opt} = argmax_{\vec{q}} [sim(\vec{q}, \mu(D_R)) - sim(\vec{q}, \mu(D_{NR}))]$$



Positive feedback β is more valuable than negative one γ and many systems only allow positive feedback so $\gamma = 0$

- **Relevance Feedback** has 2 *assumptions*:
 - * User knows the *terms* in collection for first *query*, violation: mismatch of *user's vocabulary* and *collection* one (e.g. cosmonaut-astronaut);
 - * *Relevant docs* contain *similar terms*, violation: docs with different terms but some *query concept* are not helped by the *feedback*.
- **Edited query** as better than the raw one, a *fair evaluation* must on docs not yet judged by user, an alternative can be *user* revises and resubmits *query*, but this is expensive to process, since users are reluctant to give *feedback*.
- **Pseudo-Relevance Feedback**:
 - * Retrieve a *ranked list* of hints for the *user's query*;
 - * Assume top k docs are *relevant*;
 - * Do *relevance feedback*;
- It works very well on *average* but several iteration can cause a **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 a *database* that collect *synonymous* is called **thesaurus** (e.g. hospital, medical).
- Create a *manual thesaurus* is expensive and usually used in specialized search engine for science and engineering. Instead *automatic thesaurus* has a fundamental notion: **similarity between two words**:
 - * Two words are *similar* if the **co-occur with similar words**, e.g. *car,motorbike*, cause they both occur with *road, gas, ..., robust* rule;
 - * Two words are *similar* if they **occur in a given grammatical relation with the same words**, *accurate* rule;
- Anyway *query log* is the main source of **query expansion**: e.g. user searching *flower pics* and user searching for *flower clipart* both click on some link, these two are *potential expansions* of each other, *herbal medicine* is a potential expansion of *herb*.

9 Cap 9: Text classification/categorization

Classification means distinguish *relevant doc* within *collection*, that is different from *ranking*. A **Standing Query** is a *query* used for *classification* that is run periodically to retrieve new *relevant documents* according to an *information need*, for example the *google alert*, if the *query* is bike and yamaha, will notify new docs relevant for yamaha bikes.

Classification is used for **spam filtering** too, in fact when you mark an *email* as *spam*, the system tries to create a *query* which describes it and will use it to find new *spam* (so it adds elements to *training set* for *ml*).

- **Bayes Theorem** says that: $P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A|B) \cdot P(B)}{P(A)}$
- $P(A_1, A_2, A_3, A_4) = P(A_1|A_2, A_3, A_4) \cdot P(A_2|A_3, A_4) \cdot P(A_3|A_4) \cdot P(A_4)$

Given a representation of a doc, we want to identify its category, so we need a function called classifier. There are several methods of classification:

- Manual: that is accurate when is done by expert, but it's difficult and expansive to scale;
- Hand-coded rule-based: it assign category if document contains a given boolean combination of words;
- Supervised Learning: what we said for spam, dataset can be built by amateurs;

Keeping the last method in mind (*relevance feedback*) we define the **Probabilistic relevance feedback**: if *user* has told us some *relevant* and some *non-relevant docs* we can build a **probabilistic identifier**. For example the probability that a *relevant doc* contains the word *dog*:

$$P(\text{dog}|\text{Relevant}) = \frac{N_{r\ dog}}{N_r} = \frac{\text{relevant docs with dog}}{\text{total relevant docs}}$$

Instead the probability that a *non relevant doc* contains the word *dog* is:

$$P(\text{dog}|\text{NR}) = \frac{N_{nr\ dog}}{N_{nr}}$$

9.1 Bayesian Methods

Bayesian Methods are *learning* and *classification methods* based on *probability theory*, where the *Bayes theorem* plays a critical role , it builds a generative model that approximates how data is produced. Has *prior probability* of each *category* given no information about an item. *Naïve Bayes* methods use a **bag of words** as the item description, where a text is represented as the set of its *words* keeping only track of *multiplicity*. The **Bayes rule for text classification** for a document d and a class c is:

$$P(c, d) = \frac{P(d|c) \cdot P(c)}{P(d)}$$

Where $P(c)$ is the probability that the doc that I pick is class c , it can be estimated from the *frequency of classes* in the training examples. Assuming that $d = \langle x_1, x_2, x_3, \dots \rangle$ in position order then:

$$P(d|c) = P(x_1, x_2, x_3, \dots | c)$$

But this is very complicated, so to simplify we assume that **terms probability are independent**: $P(d|c) = \prod_{i=1}^n p(x_i|c)$ but this is still an huge number because we still keep track of position. Then: $\forall t P(x_i = t|c) = P(x_j = t|c) \forall i, j \neq i$. In this way *I go home*, and *go I home*, are computed in the same way because we don't keep track of the position:

$$\frac{P(d|c) \cdot P(c)}{P(d)} = \frac{\prod_{i=1}^n p(x_i|c) \cdot P(c)}{\prod_{i=1}^n p(x_i)}$$

Now we want to **maximize this probability**:

$$\gamma(d) = \operatorname{argmax}_{c \in C} \frac{\prod_{i=1}^n p(x_i|c) \cdot P(c)}{\prod_{i=1}^n p(x_i)}$$

But since the *denominator* doesn't depend on *classes* we have that:

$$\gamma(d) = \operatorname{argmax}_{c \in C} \prod_{i=1}^n p(x_i|c) \cdot P(c)$$

So we have that $\forall c \hat{P}(c) = \frac{N_c}{N}$ and that $\forall t, \forall c, \hat{P}(t|c)$ this probability represents the occurrences of term t over the total number of terms in all docs of class c , that can be easily done during *dictionary* creation. We have a problem if $P(x_i|c) = 0$ that happened when we have a **misspell**, in this case:

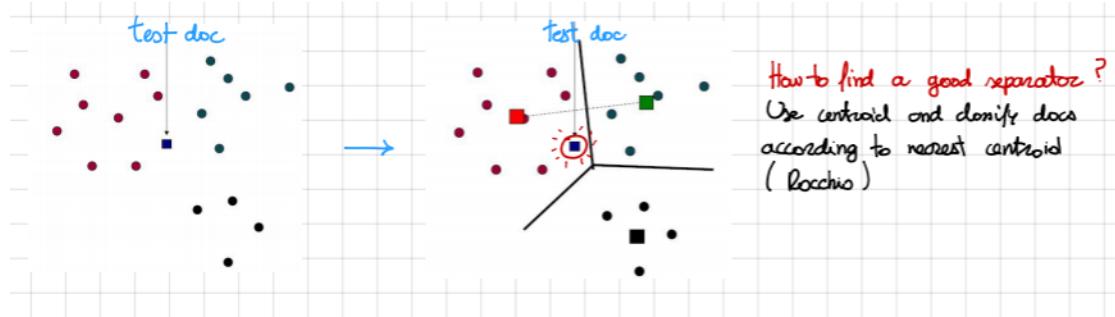
$$\hat{P}(t|c) = \frac{N_{t,c}}{\prod_{i=1}^n N_{t_i,c}} \rightarrow \hat{P}(t|c) = \frac{N_{t,c} + 1}{\prod_{i=1}^n (N_{t_i,c} + 1)}$$

Adding 1 doesn't make any problem even with a *zero probability*.

We can also use the **Multivariate Bernoulli Model** in which: $P(d|c) = P(x_1, x_2, x_3, \dots | c)$ and every *doc* is a **binary vector** with length $n = |V|$ (*dictionary*), and like before we can make some assumption: $P(d|c) = \prod_{i=1}^n P(x_i|c)$ and this represents how many *documents* over the total have x_i in them (**different from before**).

Naive Bayes is the edited version of *Bayes Rules*, and it is widely used for *spam filtering* and it's not so weak: in fact it's very fast and robust and has low storage requirements. **Evaluation** must be done on *test data* that are independent of *training data*, sometimes its possible to do a *cross-validation* (averaging results between multiple training and test sets). This *evaluation* is done according to *precision*, *recall*, *F1 score* and *classification accuracy*.

We are working with *documents* represented as *vector* in *high dimensional space*, and we need to do this **classification** with two premise: *docs* in some *class* form a *contiguous region* in *space*, and *docs* from *different classes* don't tend to *overlap*. **Learning a classifier** means build *surface* to delineate *classes* in the space:



Centroid is calculated as: $\vec{\mu}(c) = \frac{1}{|D_c|} \sum_{d \in D_c} \vec{v}(d)$ where D_c is the set of all documents that belong to class c and $v(d)$ is the **vector space representation** of d (they in general not represent

a *unit vector* even when the inputs are *unit vectors*, **unit vectors** are the vector scaled down so they don't represent the *absolute length*). **Rocchio** defines line or hyperplan as: $\prod_{i=1}^n w_i \cdot d_i = \theta$

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

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

With **Rocchio** we map each document to the nearest *centroid*, is little used outside *text classification* and in general is worse than *Naive Bayes* but it's cheap to *train* and *test document*. It's is also considered a **Linear Classifier** in fact: $x^T \cdot \vec{w} = \frac{1}{2}\theta$.

Also the **Naive Bayes** is a **linear classifier**, in fact we want to estimate:

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

Where n_d is the number of terms in document d , so given the *class* what is the probability of finding that *token* t_k in that *class* and this probability is the same for every same *token* regardless of *position* (only in models not in real life). We assume that we are in a *bi-class model* so we have also another class \bar{c} so $\hat{P}(\bar{c}|d) \propto \hat{P}(\bar{c}) \cdot \prod_{k=1}^{n_d} \hat{P}(t_k|\bar{c})$.

When we try to **classify** a new doc we are doing:

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

We decide by checking this *add ratio* if this formula gives a value ≥ 1 then the doc is class c else is class \bar{c} . For numerical reason we are going to use the logarithms:

$$\log \frac{\hat{P}(c|d)}{\hat{P}(\bar{c}|d)} = \log \frac{\hat{P}(c)}{\hat{P}(\bar{c})} \sum_{k=1}^{n_d} \log \frac{\hat{P}(t_k|c)}{\hat{P}(t_k|\bar{c})}$$

We replace the *produttoria* with the *sommatoria* since $\log(A \cdot B) = \log(A) + \log(B)$.

We can also write this sum no longer with respect the **position** (so to *the number of tokens in the document* n_d) but with respect to the **terms** that appear in the *document*, for example in the phrase: '*a rose is a rose*' we have $n_d = 5$ but the vocabulary V has only three terms so we can write:

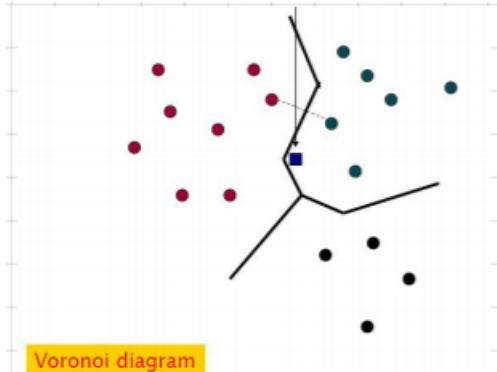
$$\log \frac{\hat{P}(c|d)}{\hat{P}(\bar{c}|d)} = \log \frac{\hat{P}(c)}{\hat{P}(\bar{c})} \sum_{i \in V} x_i \log \frac{\hat{P}(t_i|c)}{\hat{P}(t_i|\bar{c})}$$

Where x_i represents the number of times a term i appears in d , so the category is c when $\log \frac{\hat{P}(c)}{\hat{P}(\bar{c})} = \theta$ is larger than 0: $\theta \geq 0$ and $\log \frac{\hat{P}(t_i|c)}{\hat{P}(t_i|\bar{c})} = w_i$ so we have that: $\sum x_i \cdot w_i \geq -\theta$

9.2 KNN, K-Nearest Neighbors Classification

In order to classify document d we define its **K-Neighborhood** as the k nearest neighbors and pick the *majority class label* in this *neighborhood* for larger k : $P(c|d) \approx \frac{\#C}{K}$, so the probability that the class of a document d is c is given by the number of classes divided by the number of k nearest neighbors. It is a **non-linear classifier** in fact it separate the *classes* not based on *hyper plane* (*linee rette*) but on hyper surfaces (*linee segmentate*), or polytops in *highdimensional spaces*.

Usually k is an *odd* number to avoid *ties*. Finding **nearest neighbors** requires a *linear search* through $|D|$ documents in collection but determining **k nearest neighbors** is the same as determining the **k best retrievals** using the *test document* as a *query* to a *database of training documents* and in order to that we use an *inverted index*. No *training* is necessary and it *scales* well, and it can also be more accurate than *Naive Bayes* and *Rocchio*, but this is expensive at *test-time* and small changes can have *ripple effect*.

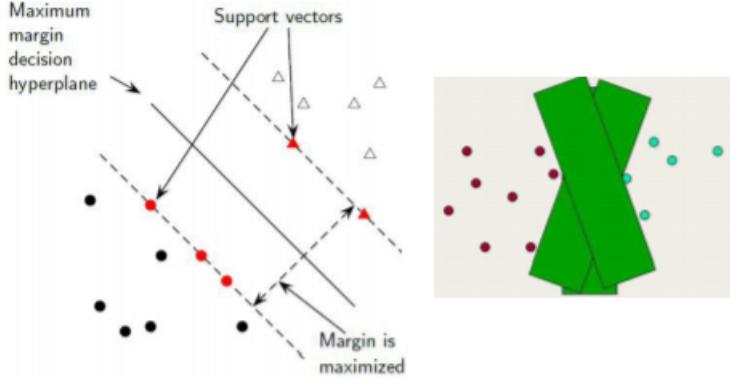


Sometimes we can also have that a *doc* that belongs to **multiple classes**, if the classes are for example $C = \{c_1, c_2, c_3\}$ then we will first classify $c = c_1$ and $\hat{c} = \{c_2, c_3\}$ then $c = c_2$ and so on, so we obtain a *vector* of results like $d: \{0.001, 0.2, 0.1\}$ that represents the score for which the document is that class ($c_1 = 0.001, \dots$) and in this case we can use two different *classifier*:

- **one-of classifier:** we get only the most *accurate class* (in the example $c_2 = 0.2$);
- **any-of classifier:** we get a *vector* composed by *scores* given to each class based on this *accuracy*;

9.3 SVM - Support Vector Machines

Learning algorithms for **vector space classification** can be divided in two types: **simple** (like *Rocchio* or *KNN*) and **iterative** (like *SVM* and *Perceptron*) that are usually the best choice. With **SVM** we are still in *vector space* but we want a **large margin classifier**: we want to increase the *division region* between *classes* and aim to find a separating *hyperplane* that is maximally far from any point in the *training set*. In case of *non-linear-separability* we accept some points as *noise* (errors). We want to maximize the **margin**, cause points near the *decision boundary* are less confident, *uncertain classification decision* and at the same time gives a *classification safety margin* respect to *errors* and random variation.



The line from which we define the *margin* is also called **decision boundary** and is a *linear separator*. Vectors on *margin lines* are called **support vectors**. We can image this *region* as a fat version of the *hyperplane*, and of course the further a *point* is classified away from this *region*, the more precise the *classification* will be.

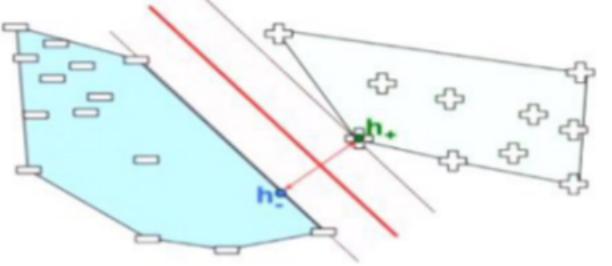
All points \vec{x} on the hyperplane satisfy the equation: $\vec{w}^T \vec{x} + b = 0$ where b is what we called θ

- Find boundary vector **w** that satisfies conditions:

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

- Convex hull: polytope around all positives / negatives

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



So we want to classify **negative** or **positive documents**, where d is the document, so all *positive documents* will get the score ≥ 1 , the *negatives* instead have a score ≤ 1 . We use these *alpha* values in the *weight definition* because some *documents* are useless for the definition of the *hyperplane* (we only need the *points* near the *border*) so for a lot of documents *alpha* will be set to 0. **Polytopes** are the lines around positive and negative docs (*the regions blue and white in the figure*), any point in the border of the *polytope* can be defined as a *linear combination* of the other points in the *border*. We search for the two **closest points** h_+, h_- , and this points are the *linear combination* of the other points (for positive: $\sum_{d \in +} \alpha_d \vec{d} \dots$). We will search to **minimize** this vector (*distance*) $\vec{w} = \vec{h}_+ - \vec{h}_-$ because we are searching for the *closest points* in order to **maximize** the *region* between the two *polytopes*, and this vector will be *perpendicular* to the **hyperplane** that

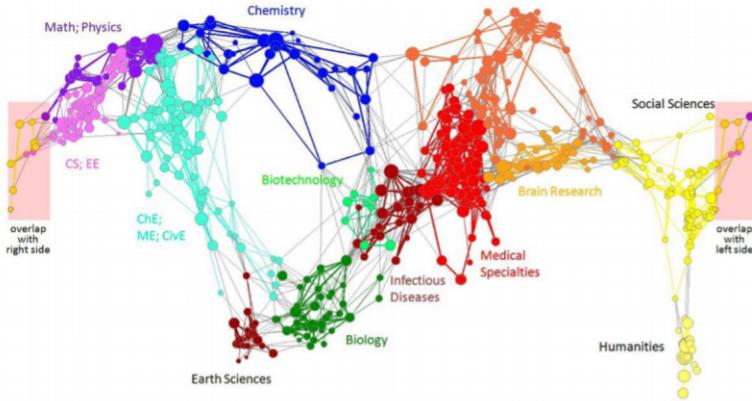
we are searching. If a new *document* is classified within the **margin** we can return a *don't know* or use a *probability function* to assign to this new *point* a *class*.

There are some cases in which a *dataset* is **not linearly separable**, in this case we can choose to accept to make **errors** in order to increase the *margin* and reduce *overfitting*. We introduce a **Slack Variable** ε_i a non zero value that takes in consideration these *errors* in such a way the sum of ε_i gives an *upper bound* on the number of errors.

We can also choose between different ways to use **SVM** for **more classes**:

- Train and run k *classifier* and then select the class with the highest confidence (**one-of classification**);
- We build $\frac{k \cdot (k-1)}{2}$ **one-versus-one classifiers**, and choose the *class* selected by the most of the *classifiers*, less time of training due the *dataset* is much smaller (you are not considering all the *dataset* but only a *given class*);
- **Structured prediction;**

10 Cap 10: Link Analysis



Since the web contains many *sources* we have to establish what pages are **trusted**, and to do that we can use two hints: usually *trustworthy pages* may point to each other, and if we are searching for example to the query *newspaper* pages that actually know about *newspaper* might be pointing to many *newspaper*. We will **rank** the *graph's nodes* by the **link structure**. There are several **link analysis approach** for computing the *importance of nodes* in a *graph*:

- **Page rank;**
- **Hubs and authorities (HITS);**
- **Topic-specified (Personalized) page rank;**
- **Web Spam detection Algorithms;**

10.1 Page Rank

The idea is to use the *links* as **votes**, so more *links* more points, and *links* to/from **important pages** count more. For example page j with importance r_j has n *out-links*, each links get $\frac{r_j}{n}$ votes, the page j own importance is the sum of the votes of its *in-links*. A page is important if it's pointed to by other important pages to the rank r_j is defined as:

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

Where d_i represent the number of *out-links* from a node. But this leads to not a *unique solution*, so we need to introduce an **additional constraint** like $r_y + r_y + r_a = 1$ (the sum of all the *ranks* is equal to one). But since this **Gaussian Elimination method** works only for *small examples* we need a better method for *large web-size graph*.

We need to introduce the **Stochastic adjacency matrix**, if the page i has d_i out-links, then if $i \rightarrow j$ then $A_{i,j} = 1$ else $A_{i,j} = 0$ so we can rewrite the rank as:

$$r_j = \sum_{i \rightarrow j} \frac{A_{i,j} \cdot r_i}{d_i} = \sum_{i \rightarrow j} M_{i,j} \cdot r_i$$

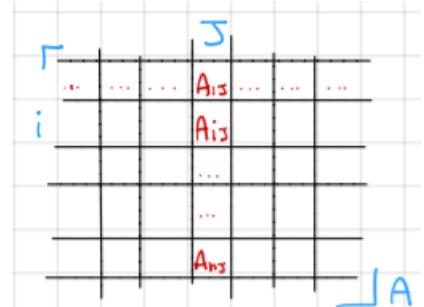
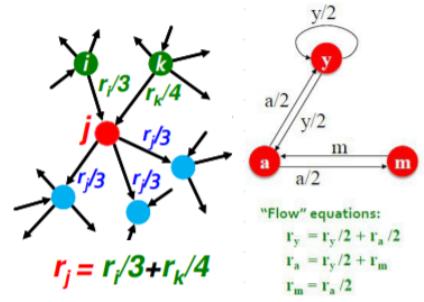
And the **flow equation** can be written as (*scriviamo r^T perché de solito il vettore sono colonne e non righe quindi scrivo trasposta*): $r^T = M \cdot r^T$ Where r^T is an *eigenvector* with *eigenvalue* 1 and that the **stochastic property** is:

$$\forall i \sum_{j=1}^n M_{i,j} = 1$$

In order to solve $r^T = M \cdot r^T$ we need to use the **Power Iteration Method**, and to simplify the formula we start from $r = M^T \cdot r = A \cdot r$ (A now is a column stochastic matrix different from before and it's $A = M^T$). We initialize $r^{(0)} = x$ and iterate $r^{(t+1)} = A \cdot r^{(t)}$ (so for example for $r^{(1)} = A \cdot x$ then $r^{(2)} = A \cdot r^{(1)} = A^2 \cdot x$ and so on) so we obtain $r^{(t)} = A^t \cdot x$. And usually the x vector is $[\frac{1}{N}, \dots, \frac{1}{N}]$

$(w_1, \lambda_1), \dots (w_n, \lambda_n)$ are eigenvector w_i and eigenvalues λ_i of A , then in a **stochastic matrix** (*no demonstration required just math formula here*) $\lambda_1 = 1$ and all *eigenvalues* are real $\leq \lambda_1$. Since this is system of *eigenvector* every vector can expressed a i of these ones: $x = \sum_{i=1}^n \alpha_i \cdot w_i$.

Then we can write $A^t \cdot x = A^t \cdot (\sum_{i=1}^n \alpha_i \cdot w_i) = \sum_{i=1}^n \alpha_i \cdot A^t \cdot w_i$, then since $A^t \cdot w_i = \lambda_i \cdot w_i$ we obtain that $A^t \cdot x = \sum_{i=1}^n \alpha_i \cdot \lambda_i^t \cdot w_i$. When you compute the **power method** actually you are

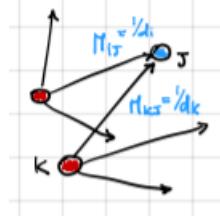


computing the *powers* of a *linear combination* of the *eigenvectors* and this is important because in *stochastic matrix* when you compute the module of the *eigenvalues*: $|\lambda_1| = \lambda_1 = 1$ and that $|\lambda_i| < 1, i > 1$. This is *fundamental* cause, when t (number of iterations of the *power method*) tends to *infinite* then all the *summatory* tends to 0 except for $\lambda_1 \cdot w_i$ so the main *eigenvalue* is 1 and r is the **main eigenvector**: $\lambda_1 \cdot w_1 = \lambda_1 \cdot r$. Whenever we get a vector like this, we *rescale* the *vector* so the sum of the components is equal to 1 and we can get r and this will always be **not negative**. So to solve $r^t = r^t \cdot M$ we do $(r^{(t+1)})^T = (r^{(t)})^T \cdot M$.

$\lambda_1 \cdot w_1 + \sum_{i=2}^n \alpha_i \cdot \lambda_i^t \cdot w_i$ is going to converge faster based on how small $|\lambda_i|$ are, they depend on $\max_{i \geq 2} |\lambda_i|$, this means that the number t that i need before it becomes as small as i like is *logarithmic* in 1 over as small as i like. From this equation we can see that this term (*the second term*) is going down **exponentially** fast with respect to maximum value of λ_i

Let's assume now that $r^{(0)} = \begin{pmatrix} 1/n \\ \dots \\ 1/n \end{pmatrix}$, this gives all probability to be at time $t = 0$ in any given point. We have that: $(r_j^{(1)}) = \sum_{i \rightarrow j} M_{i,j} \cdot r_i^{(0)}$ and $(r^{(1)})^T = (r^{(0)})^T \cdot M$

The probability to be in node j at time $t = 1$ ($r_j^{(1)}$) is equal to the probability to be in any incoming neighbors ($r_i^{(0)}$) times the probability to go to j , ($M_{i,j}$), at time $t = 0$, so the **probability distribution** is given by: $(r^{t+1})^T = (r^t)^T \cdot M$, that is called **random walk** and r its **stationary distribution**.



Qual'è la probabilità che al tempo 1 mi trovo nel nodo j ? Sarà uguale alla probabilità che mi trovo in uno dei nodi vicini con ramo che entra in j moltiplicato per la probabilità che andrà effettivamente a j , tipo se metrovo su k a tempo 0 allora la probabilità di andare a j sarà M_{kj} moltiplicato per la probabilità che me trovo effettivamente in k e indovina mpo random walk è un tipo particolare di markov chain quindi mo te becchi pure quello

10.1.1 Markov Chains

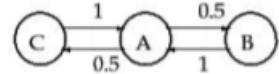
A **Markov chain** is a *discrete-time stochastic process*: a process that occurs in a series of *time-steps* in each of which a **random choice** is made and it consists of N states. Each *web page* will correspond to a state in the *Markov Chain*. It's characterized by an $N \times N$ **transition probability matrix** P each of whose entries is in the interval $[0, 1]$. The *chain* can be in one of the N states at any given *time-step*; then the entry $P_{i,j}$ tell us the *probability* that the state at next time-step is j , conditioned on the current state being i .

Each entry $P_{i,j}$ is known as a *transition probability*, and depends only on the *current state* i , and this is known as the **Markov property**: $\forall_{i,j} P_{i,j} \in [0, 1]$ and $\forall_i \sum_{j=1}^n P_{i,j} = 1$ (*that is the second property which is not always true*), and this second property is the same of the *stochastic matrix*.

e.g.

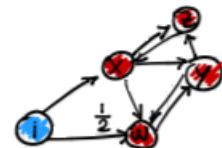
$$\begin{matrix} A & B & C \\ \begin{pmatrix} 0 & 0.5 & 0.5 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} \end{matrix}$$

In a **Markov Chain**, the *probability distribution* of the *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 reach *B* or *C*, but from them we can only reach *A* with probability 1.



Se ti chiedi perché stiamo trattando sti random walk con le markov chain? La risposta è immagina il web 15 anni fa, senza google, allora tu sei un surfista di internet che clicca randomicamente sui vari link di ogni pagina in modo da vedere a che pagina finirai, bene se dopo un milione di link cliccati sarai andato molte più volte su una pagina rispetto ad un'altra, allora è scontato che quella pagina sarà più importante (rank). Queste r_i^{t+1} rappresenta la probabilità che ci troviamo al nodo (la pagina web) i al tempo $t+1$, uso lo stesso simbolo perché dimostra che il rank e la probabilità corrispondono.

Back to the **page rank**, we have that: $r_j^{t+1} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$ or $r^T = M \cdot r^T$ and this can be rank or probability, and $(r^{(t+1)})_j^T = (r^{(t)}) \cdot M$ it's equal to r_j^{t+1} , and we have that $M_i = i^{\text{th}} \text{ row} = (\frac{1}{d_1}, 0, 0, \frac{1}{d_2}, \dots)$ for example: $(i, x, y, z, w) = (0, \frac{1}{2}, 0, 0, \frac{1}{2})$



For the **page rank formula** we have three questions:

- Does this converge?
- Does it converge to what we want?
- Are result reasonable?

The difference with a **markov chain** (if not *stochastic*) is that all *non-zero entries* in the same row are equal.

$$\forall i, \sum_j p_{i,j} = 1$$

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

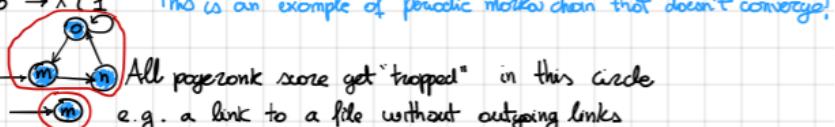
$$\hookrightarrow Z_j^{(t)} = P(X^{(t)} = j) \longrightarrow M = \begin{pmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} \quad [(Z^{(t)})^T = (Z^{(t)})^T M]$$

$$\begin{array}{ccccccccc} & & & & & & & & \\ \begin{array}{c} a \\ \xrightarrow{\hspace{1cm}} \\ b \end{array} & & & & & & & & \\ \left. \begin{array}{cccccc} V_a = 1 & 0 & 1 & \dots \\ V_b = 0 & 1 & 0 & \dots \\ t = 0 & 1 & 2 & 3 \dots \end{array} \right\} & \rightarrow & a \begin{pmatrix} a & b \\ 0 & 1 \\ b & 0 \end{pmatrix} = A & \rightarrow & \det(A - \lambda I) = 0 = P(\lambda) & & & \\ & & & & & & & & \end{array}$$

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

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

Other problem: Spider traps
Similar problem is: DeadEnd



So, we can have 3 kinds of problems, **Flip Flop**, in which we are stacked in two *nodes* (*periodic markov chain*), **Dead End**, no *outgoing links* so this *page* obtains all the *rank*, and **Spider Traps**, it happens (look the figure) when from x we enter in a *circle* where we can exit anymore, so x will get a *rank* of 0 because we always continue to be inside this **group** of nodes.

In order to fix **Dead End** we adjust the *matrix* (so the *graph*) in order to add a *link* from the *dead node* to all the *nodes* of the *graph* with probability $\frac{1}{N}$, adding these links will not affect the **page rank** at all (*it equally treats all other pages including itself, and since $1/N$ is a very low number*).

In order to fix **Spider Traps** instead they created the **Teleports**, where at each time step the *random walker* has two options:

- With probability β we follow a link at random (like usual);
- With probability $1 - \beta$ we jump to a *random node*;

This β is usually in range $[0.8, 0.9]$, with this solution we remove the *spider traps* and we remove the **periodicity** so the *chain* cannot be *periodic* anymore.

When any nodes is reachable from any other nodes then the *Markov Chain* is called **irreducible**, and when a *Markov Chain* is *irreducible* and *aperiodic* is called **Ergodic**: $\lim_{t \rightarrow \infty} r^{(t)} = r$ and r is unique. This means that there is only one *vector* that solves the *equation* (*quella principale* $r^T \cdot M = r^T$).

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

with a_T = state at end of round T

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

We ensure 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 different probabilities. $p^{(t)}$ corresponds to a probability distribution, an n -dimensional vector such that its j^{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 $\sum p_j^{(t)} = 1 \quad \forall t$

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

The steps for **random walk**:

- Remove Dead Ends:

- A **Dead End** corresponds to a row of zeros, so we replace it with a row of $\frac{1}{N}$ which corresponds to adding links to all nodes even itself. From now on we consider $A = \hat{P} + (\frac{1}{n} \cdot a)[1]^T$, where a is a *column vector* that is equal to 1 if it's a *dead end*, else is 0, and the $[1]^T$ is a row of all 1 for example:

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

$$v = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

basically just replace 0's with $\frac{1}{n}$

- **Teleportation:**

- The *random walk* is modified when at a generic node i with probability β we follow one of i outgoing links, with probability $1 - \beta$ we jump to any node at random (even i), so $r_j = \sum_{i \rightarrow j} \beta \cdot \frac{r_i}{d_i} + (1 - \beta) \cdot \frac{1}{n}$,

longer β is (< 1) than more similar the graph is to the original one but without periodicity and irreducible

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

$$\sum \mathcal{Z}_j^{(t-1)} = 1 \quad \mathbf{1} = \beta [\mathcal{Z}^{(t-1)}]^T A + \frac{1-\beta}{n} [\mathcal{Z}^{(t-1)}]^T \mathbf{1} \quad = [\mathcal{Z}^{(t-1)}]^T (\beta A + \frac{1-\beta}{n} \mathbf{1} \mathbf{1}^T) \quad (\rightarrow \mathcal{Z}^{(t)} = \mathcal{Z}^{(t-1)} M)$$

M stochastic

Smaller β means faster convergence but lose original topology → find a compromise $\alpha \sim [0.8, 0.9]$ usually

There are also some problems with this kind of **Page Rank**:

- 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-links created in order to boost page rank → **Trust Rank**;

10.2 Topic-Specific Page Rank

The goal of **Topic-Specific Page Rank** is to 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. The probability with *teleportation* is $r_j = \sum_{i \rightarrow j} \beta \cdot \frac{r_i}{d_i} + (1 - \beta) \cdot \frac{1}{n}$, and let's assume we want to rank pages according to topic *sports*, then we will update *teleportation* in such a way:

$$\beta \cdot M_{i,j} + \frac{1 - \beta}{|S|}$$

If $i \in S$ where S is sport pages, this means that when *teleporting*, jump only to *sports pages*, or $\beta \cdot M_{i,j} + 0$ otherwise:

$$\mathcal{Z}(t)^T = \mathcal{Z}(t-1)^T (\beta A + \frac{1-\beta}{|S|} \mathbf{1} \mathbf{1}^T) \quad \mathbf{1} = \begin{cases} 1 & i \in S \\ 0 & \text{otherwise} \end{cases} \quad (\text{none or } \alpha \text{ in dead ends elimination})$$

P^T is called **Personalization Vector**, and usually you don't compute it cause there are *open source repositories* with million pages handmade categorized called *DMOZ*. P can be composed by *multiple topics* with different percentage of **usage** by the user for example:

e.g. $P_{\text{User}} = (0.3, \dots, 0, \dots, 0.2, \dots, 0, \dots, 0.5, \dots, 0, \dots)$ (not a personalization vector)

$\begin{matrix} \uparrow \\ \text{Books} \end{matrix} \quad \begin{matrix} \uparrow \\ \text{Movies} \end{matrix} \quad \begin{matrix} \uparrow \\ \text{Music} \end{matrix}$

$$\rightarrow \begin{bmatrix} P_{\text{Music}} \\ P_{\text{Movies}} \\ P_{\text{Books}} \end{bmatrix} \rightarrow P = 0.3P_{\text{Books}} + 0.2P_{\text{Movies}} + 0.5P_{\text{Music}}$$

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

$$Z_1^T = Z_1^T (PA + \frac{1-\beta}{15} [1] P_1^T) = Z_1^T (\beta A + (1-\beta)[1] P_1^T) \text{ with } \sum P_{1,i} = 1$$

$$Z_2^T = Z_2^T (PA + \frac{1-\beta}{15} [1] P_2^T) = Z_2^T (\beta A + (1-\beta)[1] P_2^T) \text{ with } \sum P_{2,i} = 1$$

Note: Personalization vectors are stored somewhere

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

Can we compute it without applying power method again?

P_1 and P_2 of course are two different MC.

$$\begin{cases} Z_1^T = Z_1^T P_1 \\ Z_2^T = Z_2^T P_2 \end{cases} \xrightarrow{\text{L}} Z^T = Z^T (\beta A + (1-\beta)[1] P^T) ?$$

$\xrightarrow{\text{yes}} Z^T = (\beta Z_1^T + (1-\beta)Z_2^T) PA + (1-\beta)(\alpha P_1^T + \beta P_2^T) \xrightarrow{\text{P}^T} (Z^T \cdot [1] = 1)$

So, with this new method, we don't have to use the **power method** again, we just use the **linear combination**, and this is not done for every user, since users are classified into *groups*. It's important to note that $\sum_{i=1}^n \pi_i = 1$ (the page rank of i) but if there are *dead ends*: $\sum_{i=1}^n n\pi_i \leq 1$

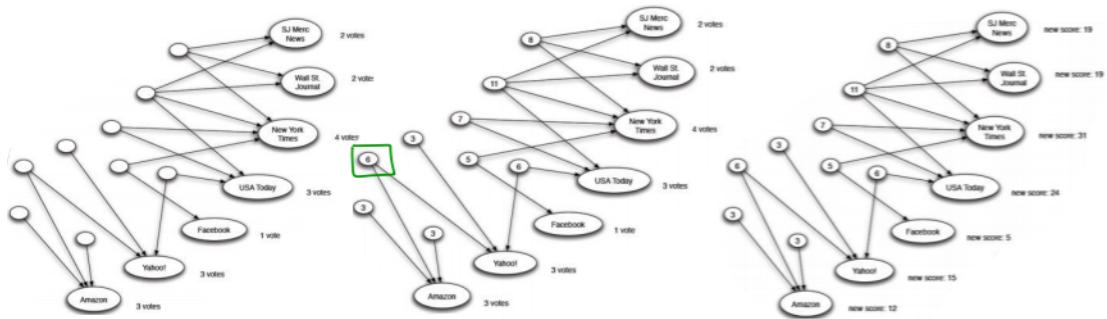
10.3 Hubs and Authorities

HITS is the **Hypertext-Induced Topic Selection**, a measure of *importance* similar to the *Page Rank*, but here for example we want to find good *newspaper*, we just don't find only *newspaper*, we will find also some **experts** who link in a coordinated way to good *newspapers*. **Votes** are always **links**, but this time also **outgoing links** count. So here a *page* is good, if *trusted pages* link to this *page*, but also if this *page* links to *trusted pages*. So we have two scores:

- As **Hub**: total sum of *votes of authorities* pointed to;
- As **Authority**: total sum of *votes coming from experts*;

How does it work:

- Each *page* starts with an **Hub score** of 1, so the *authorities* collects their votes (e.g. hubs points to Amazon \rightarrow vote 3 as authority);
- *Hubs* collect **Authority Score** (e.g. green hub score is now $6 = 3(\text{Amazon}) + 3(\text{Yahoo})$);
- **Authorities** again collect *Hub score* to compute a *final score*;



This is just an example, in real life graph is not bipartite and each page has both the scores. And now other boring math:

- good hub links to many good authorities } two vectors h and a
 - good authority is linked from many good hubs

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

$$A_{xy} = \begin{cases} 1, & x \rightarrow y \\ 0, & \text{otherwise} \end{cases}$$

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

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

$$A = \begin{pmatrix} x & | & y \end{pmatrix}$$

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

and so on

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

for $t \geq 1$:

$$\underline{a}(t) = A^T h(t-1)$$

$$\underline{h}(t) = A \underline{a}(t)$$

$$\Rightarrow \begin{cases} \underline{a}(t) = A^T A \underline{a}(t-1) \\ \underline{h}(t) = A A^T \underline{h}(t-1) \end{cases}$$

$[\underline{e}(t) = M \cdot \underline{e}(t-1)]$ power method

Entries in h and a tend to grow, no:

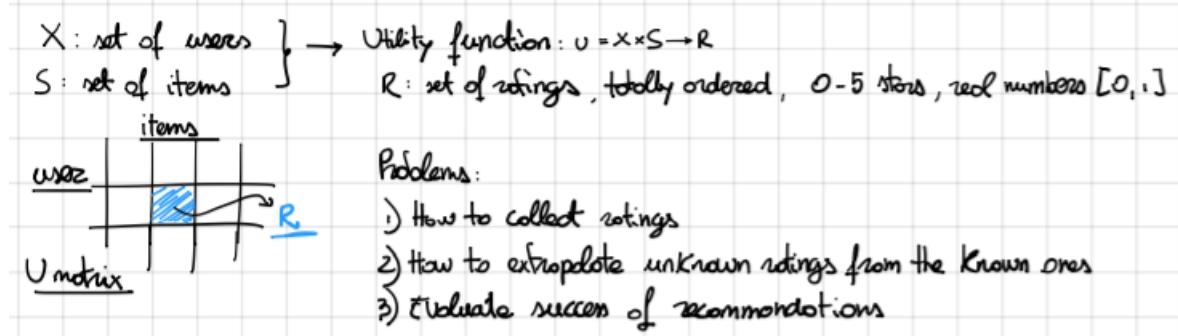
$$\underline{a}(t) = \frac{\underline{a}(t)}{\|\underline{a}(t)\|_2}$$

$$\underline{h}(t) = \frac{\underline{h}(t)}{\|\underline{h}(t)\|_2}$$

11 Cap 11: Recommendations

When an *user* search for something on the *web* usually he get **recommendations** from different sites. We came 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. There are three different types:

- **Editorial and hand curated:** like a list of *favorites*;
- **Simple Aggregates:** like *top10, most popular, recent, ...* ;
- **Tailored to individual user:** like *Amazon or Netflix*, that is the category we will talk about;



- **Explicit Way:** ask user to rate items (usually not so good);
- **Implicit Way:** learn from past, if purchase, then high rate, but low rating uncertain;
- **U Matrix** is sparse, so *new items* have *no rating*, *new users* have *no history*, so we have 3 approaches:
 - Content Based;
 - Collaborative;
 - Latent Factor Based;

Content Based and *Collaborative* are the ones used today.

11.1 Content-Based

Recommend items to costumer x similar to previous items rated highly by x :

- **Item Profiles:**

- These are represented as a vector of **features** (for example *Movie* as *author*, *title*, *actors*) through *TF-IDF*, so the *doc profile*, made of a set of words with highest *TF-IDF scores*.

- **User Profiles:**

- *User* has related items with profiles $i_1, \dots, i_n \rightarrow i_K$ vector. The *User Profile 1*: average of related items, *User Profile 2*: weighted average using rates as weights.

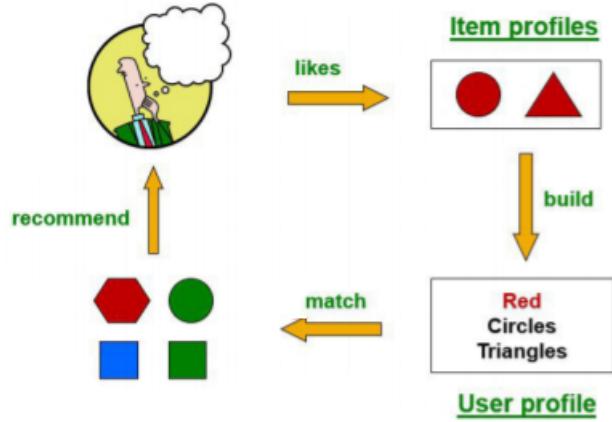
- **Match:**

- Given a user profile x and an item profile i , we estimate it as:

$$u(x, i) = \cos(x, i) = \frac{x \cdot i}{\|x\| \cdot \|i\|}$$
 the **cosine similarity**.

There are pro and cons:

- **Pro:** we don't have any *cold start*, able to recommend to *users* with unique tastes, able to recommend new and unpopular *items*, able to provide *explanations* (like a list of common features);
- **Cons:** hard to find *appropriate features*, no recommendations for *new users*, too *specific*, unable to exploit *quality judgments* of other users;

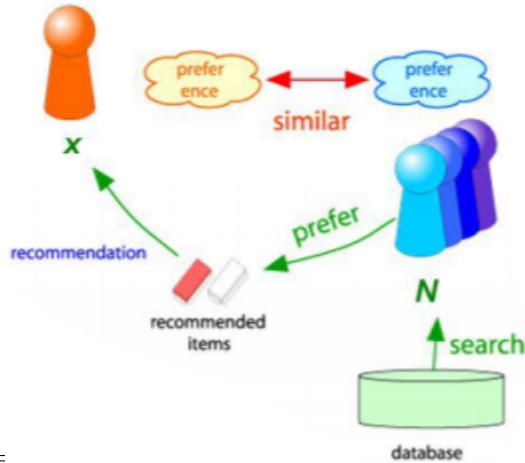


11.2 Collaborative Filtering

Consider *customer x*, we need to find a set of N other users whose *ratings* are **similar** to x ratings. So we will estimate the x ratings based on *rating* of users in N . Let's call r_x the vector of user x ratings, there are three ways:

- **Jaccard Similarity:** no, cause it ignores the values of *ratings*, in fact it only see if both *users* rated the same item;
- **Cosine Similarity:** no, cause treats missing *ratings* as *bad ratings* (*sparse matrix*);
- **Pearson Correlation Coefficient:** It calls S_{xy} the items rated by both users x and y , and \bar{r}_i the average ratings of r_i . The similarity is calculated as:

$$\frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x) \cdot (r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \cdot \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$



		ITEMS						
		HP1	HP2	HP3	TW	SW1	SW2	SW3
USERS	A	4	OK	5	5			
	B	5		4		1		
C				2				
D			not similar		4	5		

Jaccard similarity: $\text{sim}(A, B) = \frac{1}{5}$

Cosine similarity: $\text{sim}(A, B) = 0.386$, $\text{sim}(A, C) = 0.322$

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

$\rightarrow \text{sim}(A, B) < \text{sim}(A, C) \rightarrow \text{NOT TRUE}$

missing ratings "bad ratings" (too similar)

Solution:

	HP1	HP2	HP3	TW	SW1	SW2	SW3
A	2/3	1/3	-2/3	5/3	-7/3		
B	1/3	1/3					
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 $= \frac{(4+5+1)}{3} = 1\frac{2}{3} \rightarrow \text{HP1}_A = \frac{4 - 1\frac{2}{3}}{3} = \frac{2}{3}$

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

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

Now we have V_x and $N = k$ users most similar to x which 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:

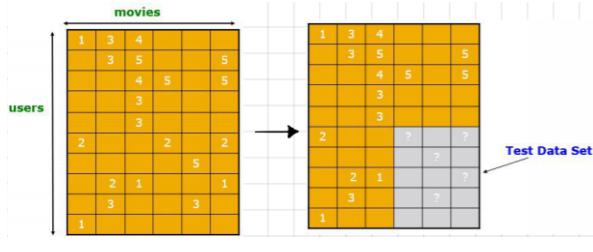
$$r_{xi} = \frac{\sum_{j \in N(i,x)} s_{ij} \cdot r_{xj}}{\sum_{j \in N(i,x)} s_{ij}}$$

Where s_{ij} is the similarity between the two *items* i and j , r_{xi} is the rating of user x on item i , and

$N(i, x)$ is the set of items rated by x similar to i . It has been observed that **item-item** is better than **user-user** cause *items* are simpler while *users* have multiple *tastes*. There are Pro and Cons:

- **Pro:** it 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**?



Basically we compare *predictions* for *test dataset* with known ratings, like **root-mean-square error** or **RMSE**:

$$\sqrt{\sum_{x_i} (r_{xi} - r_{xi}^*)^2}$$
 where r_{xi} is the predicted one 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*.

11.3 Local & Global Effects

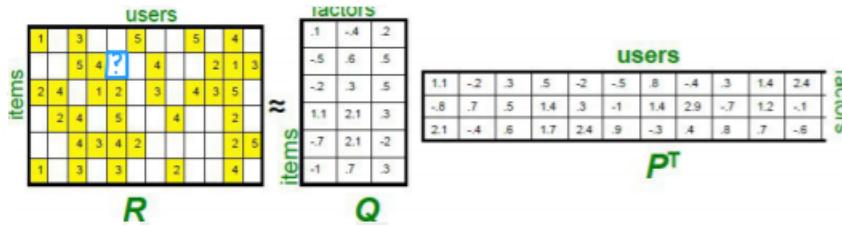
In practice we get better estimates if we **model deviations**:

$$\hat{r}_{xi} = b_{xi} + \frac{\sum_{j \in N(i,x)} s_{ij} \cdot (r_{xj} - b_{xj})}{\sum_{j \in N(i,x)} s_{ij}}$$

Where $b_{xi} = \mu + b_x + b_i$, where μ is the **overall mean rating**, b_x is the **rating deviation of the user** $x = (\text{avg rating of user } x) - \mu$, and $b_i = (\text{avg rating of item } i) - \mu$. For example, the *mean movie rating* is 3.7 the film *The sixth sense* has 4.2 starts, Joe usually rates 0.2 starts below the average: $4.2 - 0.2 = 4$ this is called the **baseline estimation** and Joe's didn't like recommendation film *Signs*: $4 - 0.2 = 3.8$ starts that is the **final estimation**.

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*.

11.4 Latent Factor Model



Suppose we can do $R = Q \cdot P^T$, the result is a **3D space** in which users and items are 3D points. $\hat{r}_{xi} = q_i \cdot p_x = \sum_f q_{if} \cdot p_{xf}$ where q_i is the row i of Q , P_x is the column x of P^T , even

missing rates. In the example if we are searching for the **box** ? inside, we are going to calculate:

$$\begin{bmatrix} -0.5 & 0.6 & 0.5 \end{bmatrix} \cdot \begin{bmatrix} -2 \\ 0.3 \\ 2.4 \end{bmatrix} = 1 + 0.18 + 1.2 = 2.4$$

In order to find P and Q we need to **minimize** the error: $\min_{P, Q} \sum_{(i, x) \in R} (r_{xi} - q_i \cdot p_x)^2$. To avoid **overfitting** we introduce **regularization**: allow *rich model* where there are sufficient data and *shrink aggressively* where data are scarce.

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

λ_1 and λ_2 are user set regularization parameters

⇒ **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}

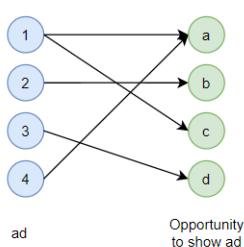
compute gradient ↴

Not very clear!

→ 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: $r_{xi} = \mu + b_x(t) + b_i(t) + q_i p_x$

12 Cap 12: Web Advertising



First of all we have to make a difference between *offline* and *online algorithms*. For **online algorithm**, the *input* isn't known at prior but one piece at time (like a *contiguous data stream*) and *decisions* depends on its history. The goal is to show right **ad** depending on the *query* (and its *history*), and of course we don't know anything about future. For **offline algorithms** the graph is something like that:

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* (like tinder), initially we have only *boys set*, and in each round one *girl* make choices: she reveals her *edges*, so we have to pair

her with someone or not at that time.

Greedy Algorithm: pair the new *girl* with any *eligible boy* if there is none don't pair, and we can compare it with an *offline matching*, using the **competitive ratio** that is equal to $= \min(\text{all possible inputs } i) \frac{|M_{\text{greedy}}|}{|M_{\text{optimal}}|}$ where M_{greedy} is the *matching* of *greedy algorithm* and M_{optimal} is the *matching* of the *offline 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_{opt} but not in M_g

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

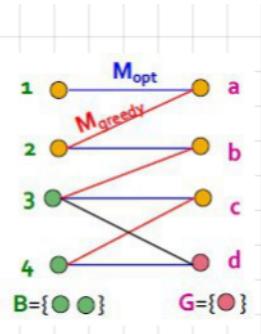
- $B = \text{boys linked to girls in } G \rightarrow |M_g| \geq |B|$

(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 $B \rightarrow |G| \leq |B|$

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

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



12.1 Performance-based Advertising

In this case, we use **advertisers** bid on search *keywords* called **adwords**. A *stream* of *query* arrives at search engine q_1, \dots, q_n , several *advertiser* bid on each *query*, and when q_i arrives *search engine* select a subset of *ads* to show. So the goal is to maximize the *profit*, and the simple solution is called **BID-CTR** (**click trough rate**, the probability that an *ad* is clicked), but this as some limitations, in fact *CTR* is *unknown* and *advertisers* have limited budget and bid on *multiple query*.

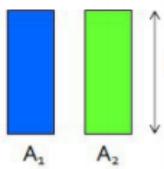
Search engine has to be sure that the **size of ads** set is not larger than the number of *ads* per *query*, each *advertiser* has bid on the *search query* and each *advertiser* has enough *budget* to pay if *ad* is clicked.

In the **simplified scenario** all budget and *CTR* are the same, we accept 1 *ad* for *query* and *ad price* is 1, then the *c ratio* is $\frac{1}{2}$ (greedy algorithm). We describe an example of a **bad scenario**: *A* bids on *x*, *B* bids on *x* and *y*, and the budgets is 4\$, the *query stream* is: *xxxxyyyy*.

- **Worst case greedy choice:** *BBBB* cause at *y* *B* finished budget and none can bid for it, so the *c ratio* is $\frac{1}{2}$;
- **Optimal case:** *AAAABBAA*;
- **Balance Algorithm MSVV:** for each *query* pick the *advertiser* with largest unspent *budget*: *ABABBB* so *c ratio* is $\frac{6}{8}$;

In the simple case with 2 *advertiser* A_1 and A_2 with budgets $B \geq 1$ **optimal solution** exhaust both budgets. *Balance* must exhaust at least one *budget* (if not, allocate more *queries*), so the $BP = \text{optimal profit} - x = 2B - x$, where BP is **Balance Profit**, and x are the *queries* not allocated.

e.g.

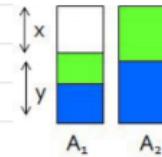
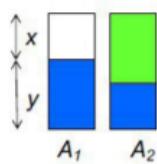


Blue = queries allocated to A_1 in optimal solution
Green = " to A_2 "

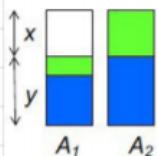
Suppose we exhaust A_2 budget: balance profit = $2B - x = B + y$
 B is A_2 exhausted budget and y is the allocated profit on A_1 .

Goal: show $y \geq B/2$

- Case 1: Balance assigns $\geq B/2$ blue queries to A_1
 $\rightarrow y \geq B/2$



- 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_1 's unpaid budget $\geq A_2$'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 $x + y = B \rightarrow y \geq B/2$



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

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

13 Cap 12: Exercises

- Prove that for any graph the pagerank of each node is at least α/N .

$\forall v: \pi_v \geq \alpha/N$
PageRank \rightarrow Ergodic MC \Rightarrow

$$\pi_v^T = \pi_v^T \rho$$

$$\pi_v^T = \frac{\alpha}{N} + (1-\alpha) \sum_{u \text{ with } v} \frac{\pi_u}{d_u}$$

def of pagerank

Probability \rightarrow $\pi_v^T = \frac{\alpha}{N} + (1-\alpha) \sum_{u \text{ with } v} \frac{\pi_u}{d_u}$
 Teleport \rightarrow following links $> \emptyset$

$$\Rightarrow \pi_v \geq \alpha/N \quad \forall v$$

- 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)?

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

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

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

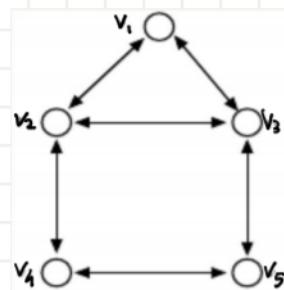
$$\bullet \pi_1 = \frac{\alpha}{5} + \frac{1-\alpha}{3} \pi_2 + \frac{1-\alpha}{3} \pi_3$$

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

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

$$\pi_3 = \pi_2 \quad \left. \begin{array}{l} \text{symmetry} \\ \pi_5 = \pi_4 \end{array} \right\}$$

$$\bullet \sum_i \pi_i = 1$$



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 π_1, π_2 and π_4 with $\alpha=\frac{1}{2}$

- True or false

- In a Boolean retrieval system, stemming never lowers precision. **false**
e.g. {house} → stem. → house ⇒ house's postings list will contain both argument → lower precision housing
- In a Boolean retrieval system, stemming never leaves 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.

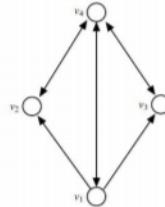
- Write down all the necessary equations needed to calculate the pagerank for a general teleporting probability α

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

∇_v

↓
 teleport
 ↓
 follow back

$$-\sum_v \pi_v = 1$$



- Given the following graph with teleport probability α

- Write all necessary equations needed to calculate personalized pagerank with respect to personalization vector $[1, 0, 0, 0]$

$$\begin{aligned}\pi_1 &= \alpha + (1-\alpha) \frac{\pi_4}{3} \\ \pi_2 &= \pi_3 = (1-\alpha) \frac{\pi_1}{3} + (1-\alpha) \frac{\pi_4}{3} \\ \pi_4 &= (1-\alpha) \left[\frac{\pi_1}{3} + \pi_2 + \pi_3 \right] \\ \sum \pi_i &= 1\end{aligned}$$