# Natural Language Processing and Information Retrieval

## VSM and Optmization

### Alessandro Moschitti

Department of Computer Science and Information
Engineering
University of Trento
Email: moschitti@disi.unitn.it

# Summary: weighting

- Term Weighting

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

- The idf (inverse document frequency) of *t* by

$$\mathrm{idf}_t = \log_{10} (N/\mathrm{df}_t)$$

# Summary: tf-idf weighting

- The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$w_{t,d} = (1 + \log_{10} \text{tf}_{t,d}) \times \log_{10}(N / \text{df}_t)$$

- Best known weighting scheme in information retrieval

- Increases with the number of occurrences within a document

- Increases with the rarity of the term in the collection

# Recap: Queries as vectors

- **Key idea 1:** Do the same for queries: represent them as vectors in the space

- **Key idea 2:** Rank documents according to their proximity to the query in this space

- proximity = similarity of vectors

# Summary – vector space ranking

- Represent the query as a weighted tf-idf vector

- Represent each document as a weighted tf-idf vector

- Compute the cosine similarity score for the query vector and each document vector

- Rank documents with respect to the query by score

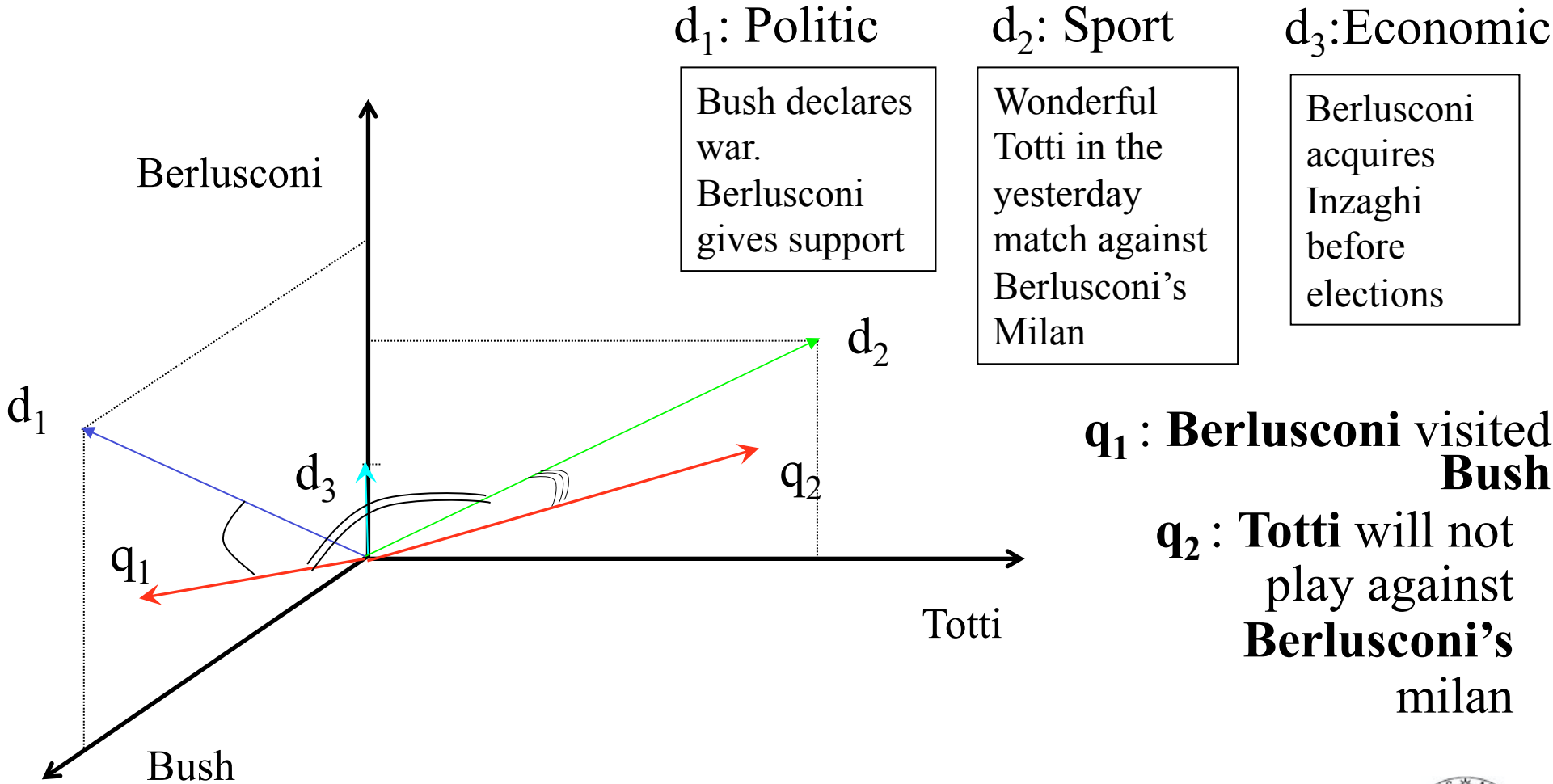- Return the top $K$ (e.g., $K = 10$) to the user

# VSM: formal definition (see Salton 89')

- Features are dimensions of a Vector Space

- Documents and Queries are vectors of feature weights

- A set of documents is retrieved based on $\vec{d} \cdot \vec{q}$,

- where $\vec{d}$, $\vec{q}$ are the vectors representing documents and query

# The Vector Space Model

d₁: Politic

Bush declares war. Berlusconi gives support

d₂: Sport

Wonderful Totti in the yesterday match against Berlusconi's Milan

d₃:Economic

Berlusconi acquires Inzaghi before elections

q₁ : **Berlusconi** visited **Bush**

q₂ : **Totti** will not play against **Berlusconi's** milan

Berlusconi

d₁

d₃

d₂

q₂

q₁

Totti

Bush

# tf-idf weighting has many variants

| Term frequency | | Document frequency | | Normalization | |
|---|---|---|---|---|---|
| n (natural) | $\text{tf}_{t,d}$ | n (no) | 1 | n (none) | 1 |
| l (logarithm) | $1 + \log(\text{tf}_{t,d})$ | t (idf) | $\log \frac{N}{\text{df}_t}$ | c (cosine) | $\frac{1}{\sqrt{w_1^2 + w_2^2 + \ldots + w_M^2}}$ |
| a (augmented) | $0.5 + \frac{0.5 \times \text{tf}_{t,d}}{\max_t(\text{tf}_{t,d})}$ | p (prob idf) | $\max\{0, \log \frac{N - \text{df}_t}{\text{df}_t}\}$ | u (pivoted unique) | $1/u$ |
| b (boolean) | $\begin{cases} 1 & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$ | | | b (byte size) | $1/CharLength^{\alpha}$, $\alpha < 1$ |
| L (log ave) | $\frac{1 + \log(\text{tf}_{t,d})}{1 + \log(\text{ave}_{t \in d}(\text{tf}_{t,d}))}$ | | | | |

Columns headed 'n' are acronyms for weight schemes.

Why is the base of the log in idf immaterial?

# Weighting may differ in queries vs documents

- Many search engines allow for different weightings for queries vs. documents

- SMART Notation: denotes the combination in use in an engine, with the notation *ddd.qqq,* using the acronyms from the previous table

- A very standard weighting scheme is: lnc.ltc

- Document: logarithmic tf (l as first character), no idf and cosine normalization

  A bad idea?

- Query: logarithmic tf (l in leftmost column), idf (t in second column), no normalization …

# tf-idf example: lnc.ltc

Document: *car insurance auto insurance*
Query: *best car insurance*

| Term | Query | | | | | | Document | | | | Prod |
|------|-------|------|------|------|------|---------|--------|-------|------|---------|------|
|      | tf-raw | tf-wt | df | idf | wt | n'lize | tf-raw | tf-wt | wt | n'lize |      |
| auto | 0 | 0 | 5000 | 2.3 | 0 | 0 | 1 | 1 | 1 | 0.52 | 0 |
| best | 1 | 1 | 50000 | 1.3 | 1.3 | 0.34 | 0 | 0 | 0 | 0 | 0 |
| car | 1 | 1 | 10000 | 2.0 | 2.0 | 0.52 | 1 | 1 | 1 | 0.52 | 0.27 |
| insurance | 1 | 1 | 1000 | 3.0 | 3.0 | 0.78 | 2 | 1.3 | 1.3 | 0.68 | 0.53 |

$$\text{Doc length} = \sqrt{1^2 + 0^2 + 1^2 + 1.3^2} \approx 1.92$$

Score = 0+0+0.27+0.53 = 0.8

# Computing cosine scores

$\textsc{CosineScore}(q)$
1   $\textit{float Scores}[N] = 0$
2   $\textit{float Length}[N]$
3   **for each** query term $t$
4   **do** calculate $w_{t,q}$ and fetch postings list for $t$
5      **for each** $\text{pair}(d, \text{tf}_{t,d})$ in postings list
6      **do** $\textit{Scores}[d] += w_{t,d} \times w_{t,q}$
7   Read the array $\textit{Length}$
8   **for each** $d$
9   **do** $\textit{Scores}[d] = \textit{Scores}[d]/\textit{Length}[d]$
10   **return** Top $K$ components of $\textit{Scores}[]$

# Efficient cosine ranking

- Find the $K$ docs in the collection "nearest" to the query $\Rightarrow K$ largest query-doc cosines.

- Efficient ranking:

  - Computing a single cosine efficiently.

  - Choosing the $K$ largest cosine values efficiently.

    - Can we do this without computing all $N$ cosines?

# Efficient cosine ranking

- What we're doing in effect: solving the *K*-nearest neighbor problem for a query vector

- In general, we do not know how to do this efficiently for high-dimensional spaces

- But it is solvable for short queries, and standard indexes support this well

# Special case – unweighted queries

- No weighting on query terms
  - Assume each query term occurs only once
- Then for ranking, don't need to normalize query vector
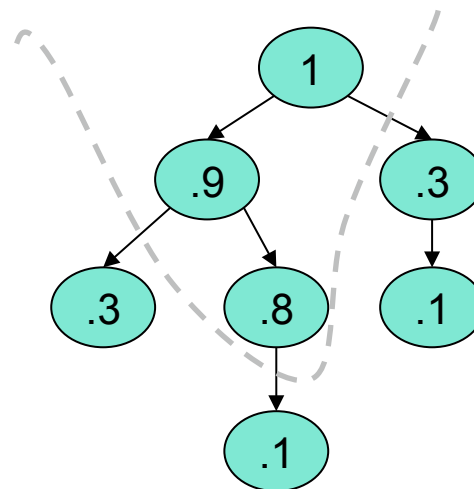  - Slight simplification of algorithm

# Computing the *K* largest cosines: selection vs. sorting

- Typically we want to retrieve the top *K* docs (in the cosine ranking for the query)
  - not to totally order all docs in the collection

- Can we pick off docs with *K* highest cosines?

- Let *J* = number of docs with nonzero cosines
  - We seek the *K* best of these *J*

# Use heap for selecting top *K*

- Binary tree in which each node's value > the values of children

- Takes *2J* operations to construct, then each of *K* "winners" read off in 2log *J* steps.

- For *J*=1M, *K*=100, this is about 10% of the cost of sorting.

# Bottlenecks

- Primary computational bottleneck in scoring: <u>cosine computation</u>

- Can we avoid all this computation?

- Yes, but may sometimes get it wrong
  - a doc *not* in the top $K$ may creep into the list of $K$ output docs
  - Is this such a bad thing?

# Cosine similarity is only a proxy

- User has a task and a query formulation

- Cosine matches docs to query

- Thus cosine is anyway a proxy for user happiness

- If we get a list of $K$ docs "close" to the top $K$ by cosine measure, should be ok

# Generic approach

- Find a set $A$ of *contenders*, with $K < |A| << N$
  - $A$ does not necessarily contain the top $K$, but has many docs from among the top $K$
  - Return the top $K$ docs in $A$

- Think of $A$ as <u>pruning</u> non-contenders

- The same approach is also used for other (non-cosine) scoring functions

- Will look at several schemes following this approach

# Index elimination

- Basic algorithm cosine computation algorithm only considers docs containing at least one query term

- Take this further:
  - Only consider high-idf query terms
  - Only consider docs containing many query terms

# High-idf query terms only

- For a query such as *catcher in the rye*

- Only accumulate scores from *catcher* and *rye*

- Intuition: **in** and **the** contribute little to the scores and so <u>don't alter rank-ordering much</u>

- Benefit:
  - Postings of low-idf terms have many docs → these (many) docs get eliminated from set *A* of contenders

# Docs containing many query terms

- Any doc with at least one query term is a candidate for the top *K* output list

- For multi-term queries, only compute scores for docs containing several of the query terms
  - Say, at least 3 out of 4
  - Imposes a "soft conjunction" on queries seen on web search engines (early Google)

- Easy to implement in postings traversal

# 3 of 4 query terms

| Antony | | 3 | 4 | 8 | 16 | 32 | 64 | 128 | |

| Brutus | | 2 | 4 | 8 | 16 | 32 | 64 | 128 | |

| Caesar | | 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 |

| Calpurnia | | 13 | 16 | 32 | | | | | |

Scores only computed for docs 8, 16 and 32.

# Champion lists

- Precompute for each dictionary term $t$, the $r$ docs of highest weight in $t$'s postings
    - Call this the <u>champion list</u> for $t$
    - (aka <u>fancy list</u> or <u>top docs</u> for $t$)

- Note that $r$ has to be chosen at index build time
    - Thus, it's possible that $r < K$

- At query time, only compute scores for docs in the champion list of some query term
    - Pick the $K$ top-scoring docs from amongst these

# Exercises

- How do Champion Lists relate to Index Elimination? Can they be used together?

- How can Champion Lists be implemented in an inverted index?

  - Note that the champion list has nothing to do with small docIDs

# Static quality scores

- We want top-ranking documents to be both *relevant* and *authoritative*

- *Relevance* is being modeled by cosine scores

- *Authority* is typically a query-independent property of a document

- Examples of authority signals

    - Wikipedia among websites

    - Articles in certain newspapers

    - A paper with many citations

    - Account from website, e.g. delicious.com

    - Pagerank

Quantitative

# Modeling authority

- Assign to each document a *query-independent* <u>quality score</u> in [0,1] to each document *d*
  - Denote this by *g(d)*

- Thus, a quantity like the number of citations is scaled into [0,1]
  - Exercise: suggest a formula for this.

# Net score

- Consider a simple total score combining cosine relevance and authority

- net-score($q,d$) = $g(d)$ + cosine($q,d$)
  - Can use some other linear combination
  - Indeed, any function of the two "signals" of user happiness – more later

- Now we seek the top *K* docs by <u>net score</u>

# Top *K* by net score – fast methods

- First idea: Order all postings by *g(d)*

- Key: this is a common ordering for all postings

- Thus, can concurrently traverse query terms' postings for

  - Postings intersection

  - Cosine score computation

- Exercise: write pseudocode for cosine score computation if postings are ordered by *g(d)*
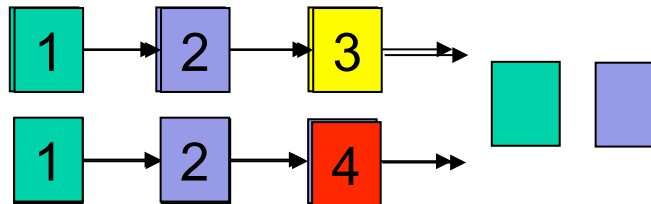
# Why order postings by *g(d)*?

- Under *g(d)*-ordering, top-scoring docs likely to appear early in postings traversal

- In time-bound applications (say, we have to return whatever search results we can in 50 ms), this allows us to stop postings traversal early

  - Short of computing scores for all docs in postings

# Re-ordering with respect to g(d)

1    0.1

2    0.9

g(d)

3    0.7

4    0.3

1
2
3
4

*Brutus*   1 → 2 → 3 →

*Caesar*   1 → 2 → 4 →

# Champion lists in *g(d)*-ordering

- Can combine champion lists with *g(d)*-ordering

- Maintain for each term a champion list of the *r* docs with highest $g(d) + $ tf-idf$_{td}$

- Seek top-*K* results from only the docs in these champion lists

# High and low lists

- For each term, we maintain two postings lists called *high* and *low*

  - Think of *high* as the champion list

- When traversing postings on a query, only traverse *high* lists first

  - If we get more than *K* docs, select the top *K* and stop

  - Else proceed to get docs from the *low* lists

- Can be used even for simple cosine scores, without global quality *g(d)*

- A means for segmenting index into two <u>tiers</u>

# Impact-ordered postings

- We only want to compute scores for docs for which $wf_{t,d}$ is high enough

- We sort each postings list by $wf_{t,d}$

- Now: not all postings in a common order!

- How do we compute scores in order to pick off top *K?*
  - Two ideas follow

# 1. Early termination

- When traversing $t$'s postings, stop early after either
  - a fixed number of $r$ docs
  - $wf_{t,d}$ drops below some threshold
- Take the union of the resulting sets of docs
  - One from the postings of each query term
- Compute only the scores for docs in this union

# 2. idf-ordered terms

- When considering the postings of query terms

- Look at them in order of decreasing idf

  - High idf terms likely to contribute most to score

- As we update score contribution from each query term

  - Stop if doc scores relatively unchanged

- Can apply to cosine or some other net scores

# Cluster pruning: preprocessing

- Pick √N *docs* at random: call these *leaders*

- For every other doc, pre-compute nearest leader

  - Docs attached to a leader: its *followers;*
  - <u>Likely</u>: each leader has ~ $\sqrt{N}$ followers.

# Cluster pruning: query processing

- Process a query as follows:
    - Given query $Q$, find its nearest *leader L.*
    - Seek $K$ nearest docs from among $L$'s followers.

# Visualization



Query

Leader    Follower

# Why use random sampling

- Fast

- Leaders reflect data distribution

# General variants

- Have each follower attached to *b1*=3 (say) nearest leaders.

- From query, find *b2*=4 (say) nearest leaders and their followers.

- Can recurse on leader/follower construction.

# Exercises

- To find the nearest leader in step 1, how many cosine computations do we do?
  - Why did we have √N in the first place?

- What is the effect of the constants *b1, b2* on the previous slide?

- Devise an example where this is *likely to* fail – i.e., we miss one of the *K* nearest docs.
  - *Likely* under random sampling.

# Parametric and zone indexes

- Thus far, a doc has been a sequence of terms

- In fact documents have multiple parts, some with special semantics:
  - Author
  - Title
  - Date of publication
  - Language
  - Format
  - etc.

- These constitute the <u>metadata</u> about a document

# Fields

- We sometimes wish to search by these metadata
  - E.g., find docs authored by William Shakespeare in the year 1601, containing *alas poor Yorick*

- Year = 1601 is an example of a <u>field</u>

- Also, author last name = shakespeare, etc.

- Field or parametric index: postings for each field value
  - Sometimes build range trees (e.g., for dates)

- Field query typically treated as conjunction
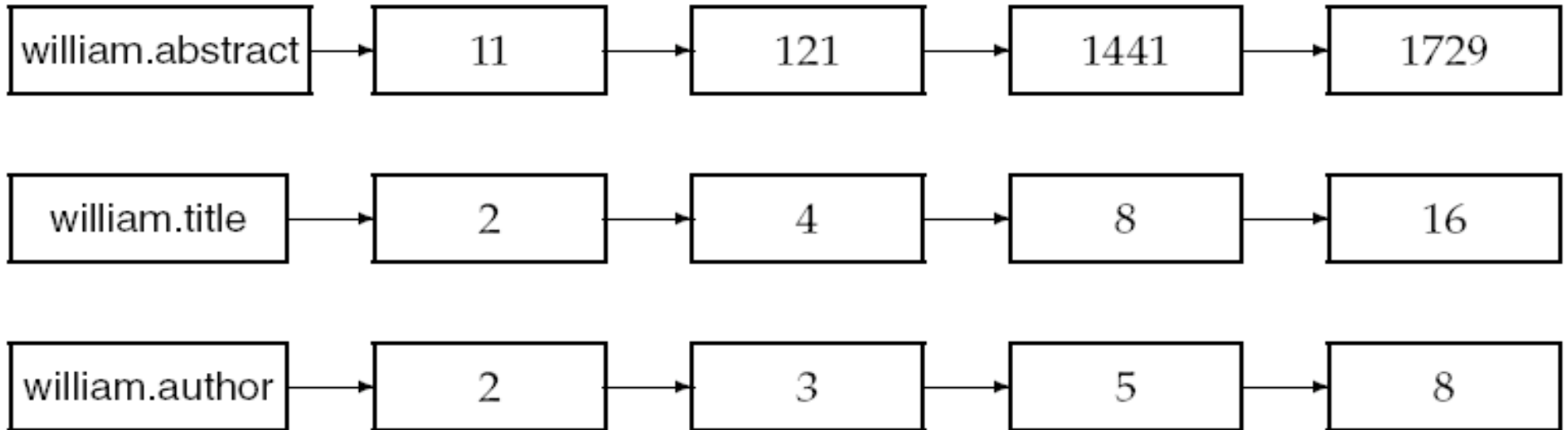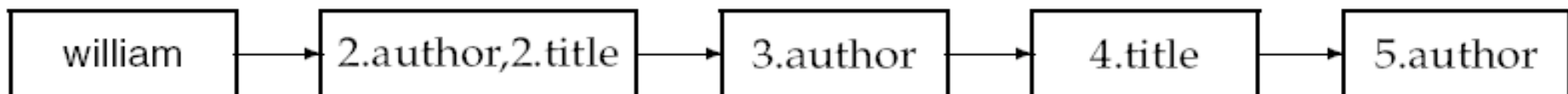  - (doc *must* be authored by shakespeare)

# Zone

- A <u>zone</u> is a region of the doc that can contain an arbitrary amount of text, e.g.,
  - Title
  - Abstract
  - References …

- Build inverted indexes on zones as well to permit querying

- E.g., "find docs with *merchant* in the title zone and matching the query *gentle rain*"

# Example zone indexes



| william.abstract | → | 11 | → | 121 | → | 1441 | → | 1729 |

| william.title | → | 2 | → | 4 | → | 8 | → | 16 |

| william.author | → | 2 | → | 3 | → | 5 | → | 8 |

Encode zones in dictionary vs. postings.

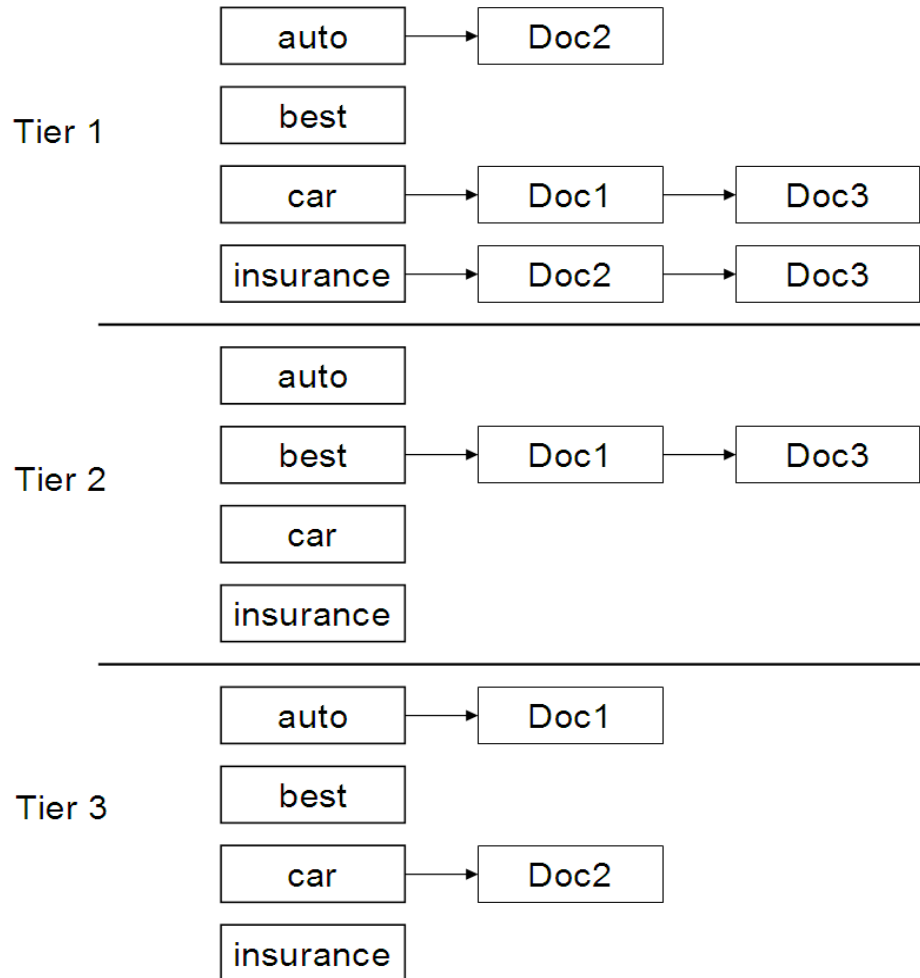| william | → | 2.author,2.title | → | 3.author | → | 4.title | → | 5.author |

# Tiered indexes

- Break postings up into a hierarchy of lists
  - Most important
  - …
  - Least important
- Can be done by *g(d)* or another measure
- Inverted index thus broken up into <u>tiers </u>of decreasing importance
- At query time use top tier unless it fails to yield *K* docs
  - If so drop to lower tiers

# Example tiered index



Tier 1
- auto → Doc2
- best
- car → Doc1 → Doc3
- insurance → Doc2 → Doc3

Tier 2
- auto
- best → Doc1 → Doc3
- car
- insurance

Tier 3
- auto → Doc1
- best
- car → Doc2
- insurance

# Query term proximity

- <u>Free text queries</u>: just a set of terms typed into the query box – common on the web

- Users prefer docs in which query terms occur within close proximity of each other

- Let *w* be the smallest window in a doc containing all query terms, e.g.,

  - For the query *strained mercy* the smallest window in the doc *The quality of mercy is not strained* is <u>4</u> (words)

- Would like scoring function to take this into account – how?

# Query parsers

- Free text query from user may in fact spawn one or more queries to the indexes, e.g., query *rising interest rates*

  - Run the query as a phrase query
  - If $<K$ docs contain the phrase *rising interest rates*, run the two phrase queries *rising interest* and *interest rates*
  - If we still have $<K$ docs, run the vector space query *rising interest rates*
  - Rank matching docs by vector space scoring
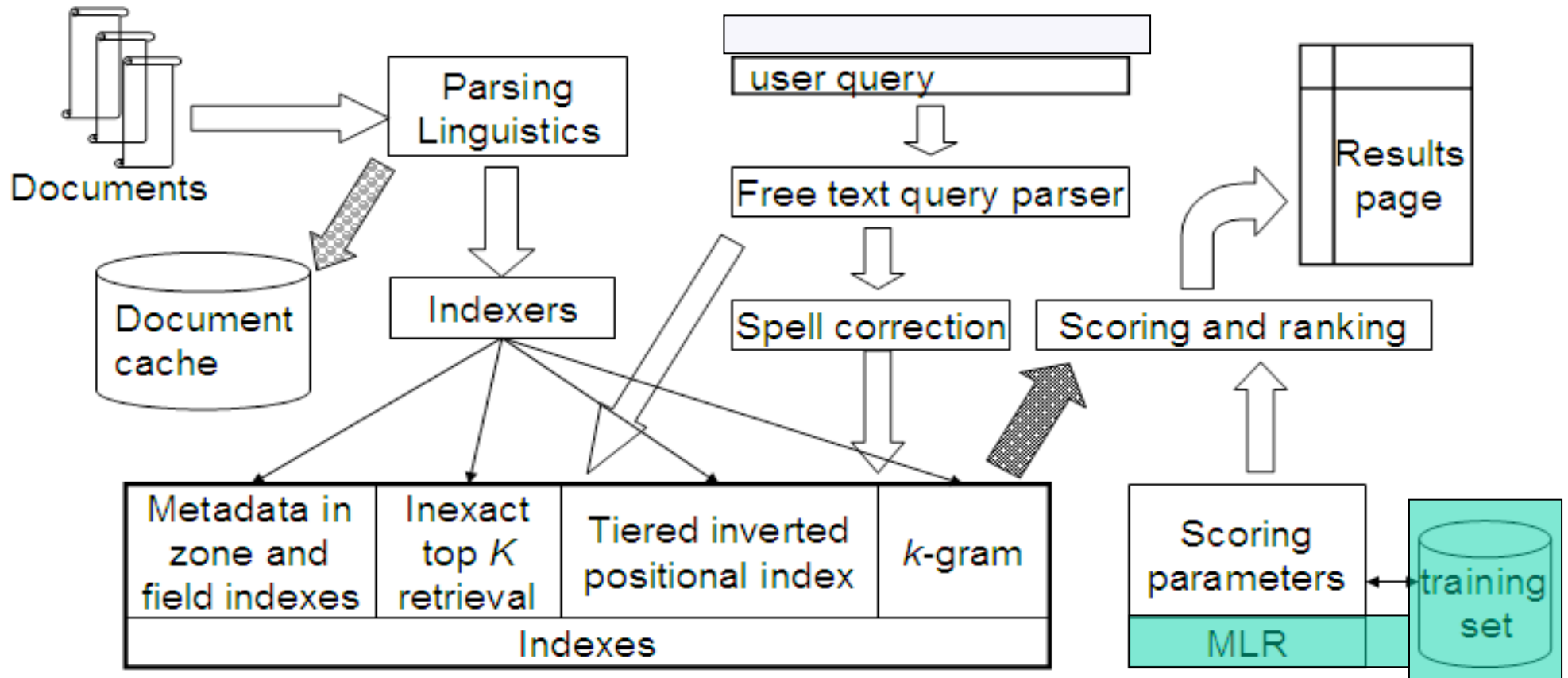
- This sequence is issued by a <u>query parser</u>

# Aggregate scores

- We've seen that score functions can combine cosine, static quality, proximity, etc.

- How do we know the best combination?

- Some applications – expert-tuned

- Increasingly common: machine-learned

# Putting it all together

# End Lecture

- ## Next time
  - Performance Measures for Retrieval Systems
  - Connected with Machine Learning and Natural Language Processing

- ## Introduction to ML if time allows

# Query Expansion

- N, the overall number of documents,

- $N_f$, the number of documents that contain the feature $f$

- $o_f^d$ the occurrences of the features $f$ in the document $d$

- The weight $f$ in a document is:

$$\omega_f^d = \left( \log \frac{N}{N_f} \right) \times o_f^d = IDF(f) \times o_f^d$$

- The weight can be normalized:

$$\omega'^d_f = \frac{\omega_f^d}{\sqrt{\sum_{t \in d} (\omega_t^d)^2}}$$

# Relevance Feeback and query expansion: the Rocchio's formula

- $\omega_f^d$ , the weight of $f$ in $d$
  - Several weighting schemes (e.g. TF * IDF, Salton 91')

- $\vec{q}_f$ , the profile weights of $f$ in $C_i$:

$$\vec{q}_f = \max\left\{0, \ \frac{\beta}{|T|}\sum_{d\in T}\omega_f^d \ - \ \frac{\gamma}{|\overline{T}|}\sum_{d\in\overline{T}}\omega_f^d\right\}$$

- $T_i$ , the training documents in $q$

# Similarity estimation between query and documents

- Given the document and the category representation

$$\vec{d} = \left\langle \omega_{f_1}^d, ..., \omega_{f_n}^d \right\rangle, \quad \vec{q} = \left\langle \Omega_{f_1}, ..., \Omega_{f_n} \right\rangle$$

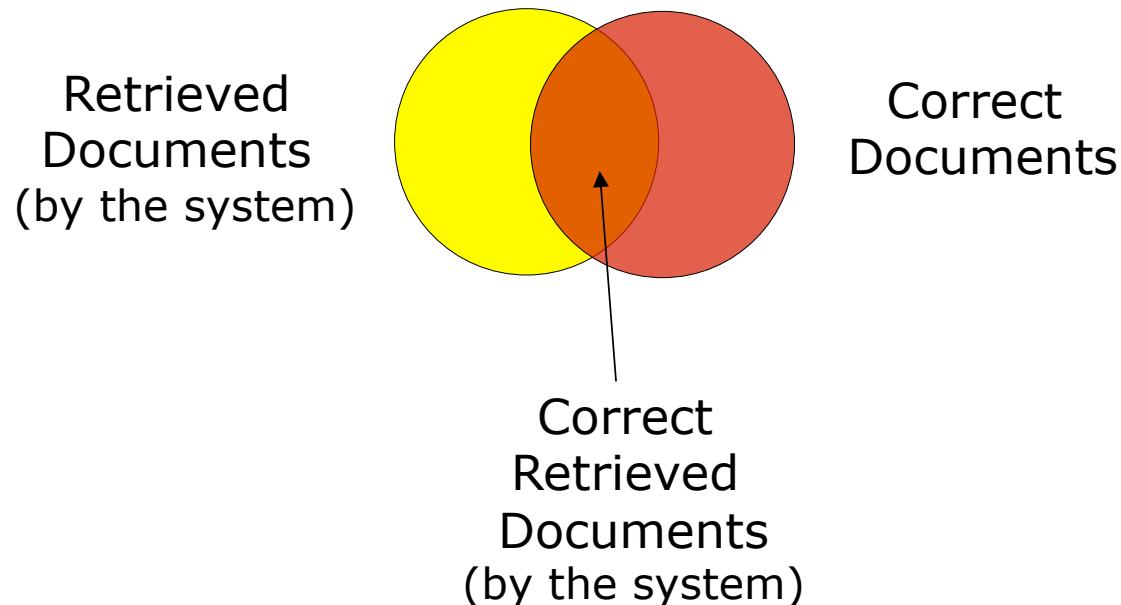- It can be defined the following similarity function (cosine measure

$$s_{d,i} = \cos(\vec{d}, \vec{q}) = \frac{\vec{d} \cdot \vec{q}}{\|\vec{d}\| \times \|\vec{q}\|} = \frac{\sum_f \omega_f^d \times \Omega_f^i}{\|\vec{d}\| \times \|\vec{q}\|}$$

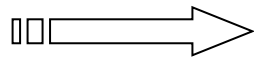- $d$ is assigned to $\vec{q}$ if $\vec{d} \cdot \vec{q} > \sigma$

# Performance Measurements

- Given a set of document *T*

- Precision = # Correct Retrieved Document / # Retrieved Documents

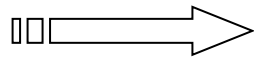- Recall = # Correct Retrieved Document/ # Correct Documents

Retrieved
Documents
(by the system)

Correct
Documents

Correct
Retrieved
Documents
(by the system)

| 1 | 1 | 2 | 2 | 3 | 4 | 5 | |

**Antony** →

| 3 | 4 | 8 | 16 | 32 | 64 | 128 | |

**Brutus** →

| 2 | 4 | 8 | 16 | 32 | 64 | 128 | |