

Introduction to Assignment Project Exam Help **Informa** ai <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro Lecture 7: Scoring an ssembly

Recap: tf-idf weighting

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

$$w_{t,d} = (1 + \log(N/d_f_t)) \cdot (N/d_f_t)$$

- Best known weighting scheme for 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 <https://eduassistpro.github.io/>
- proximity = similarity of vectors

Recap: cosine(query,document)

$$\cos(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{\|\vec{q}\| \|\vec{d}\|}$$

Dot product **Unit vectors**

$\sum_{i=1}^{|V|} q_i d_i$

$\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}$

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$\cos(\vec{q}, \vec{d})$ is the cosine similarity of \vec{q} and \vec{d} ... or,
equivalently, the cosine of the angle between \vec{q} and \vec{d} .

This lecture

- Speeding up vector space ranking
- Putting together a complete search system
 - Assignment Project Exam Help
 - <https://eduassistpro.github.io/>
 - Will require learning a number of miscellaneous topics
 - Add WeChat `edu_assist_pro`

Question: Why don't we just use the query processing methods for Boolean queries?

Term-at-a-time

Computing cosine scores

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Efficient cosine ranking

- Find the K docs in the collection “nearest” to the query $\Rightarrow K$ largest query-doc cosines.
- Efficient ranking:
 - Computing <https://eduassistpro.github.io/>
 - Choosing the K largest cosines 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, ~~Assignment Project Exam Help~~ we do not know how to do this efficiently for high-dime <https://eduassistpro.github.io/>
- But it is solvable for short q d standard indexes support this well
~~Add WeChat edu_assist_pro~~

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 <https://eduassistpro.github.io/>
 - Slight simplification of all formulas

Faster cosine: unweighted query

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



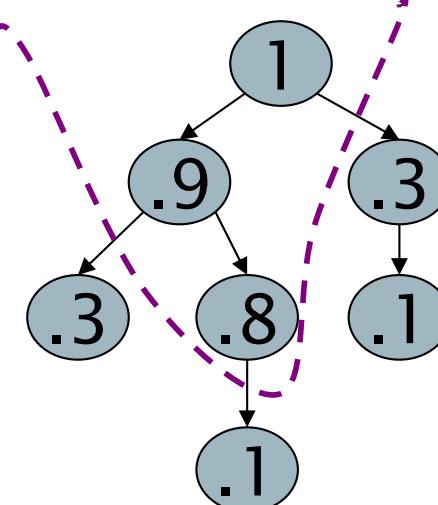
Computing the K largest cosines: selection vs. sorting

- 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 <https://eduassistpro.github.io/cosines/>?
- Let n of docs with nonzero [Add WeChat edu_assist_pro](#)
 - We seek the K best of these n

http://en.wikipedia.org/wiki/Binary_heap

Use heap for selecting top K

- Binary tree in which each node's value > the values of children
- Takes $2n$ operations to construct, then each of K “winners” read <https://eduassistpro.github.io/>
- Total time is $O(n + K * \log(n))$
- For $n=1M$, $K=100$, this is about $\frac{1}{2}$ the cost of sorting.



Quick Select

- QuickSelect is similar to QuickSort to find the top-K elements from an array
 - Takes $O(n \log(K))$ time
 - Sorting the to <https://eduassistpro.github.io/> time
 - Total time is $O(n + K \log(K))$
- Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Query Processing

- Document-at-a-time
 - Calculates complete scores for documents by processing all term lists
- Term-at-a-time
 - Accumulates one at a time
- Both approaches have optimization techniques that significantly reduce time required to generate scores
 - Distinguish between safe and heuristic optimizations

<https://eduassistpro.github.io/>

processing term lists

[Add WeChat edu_assist_pro](#)

Document-At-A-Time

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Document-At-A-Time

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Term-At-A-Time

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Term-At-A-Time

// accumulators

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist // A_d contains partial score

Optimization Techniques

- Term-at-a-time uses more memory for accumulators, but accesses disk more efficiently
- Two classes ~~Assignment Project Exam Help~~
 - Read less data <https://eduassistpro.github.io/>
 - e.g., skip lists
 - better for simple feature funct
 - Calculate scores for fewer documents
 - e.g., conjunctive processing
 - better for complex feature functions

Conjunctive Processing

- Requires the result document containing all the query terms (i.e., conjunctive Boolean queries)
 - More efficient
 - Can also be used for more complex queries <https://eduassistpro.github.io/>
 - Default for most engines
- Can be combined with both AND and NOT

Conjunctive Term-at-a-Time

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Conjunctive Document-at-a-Time

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Threshold Methods

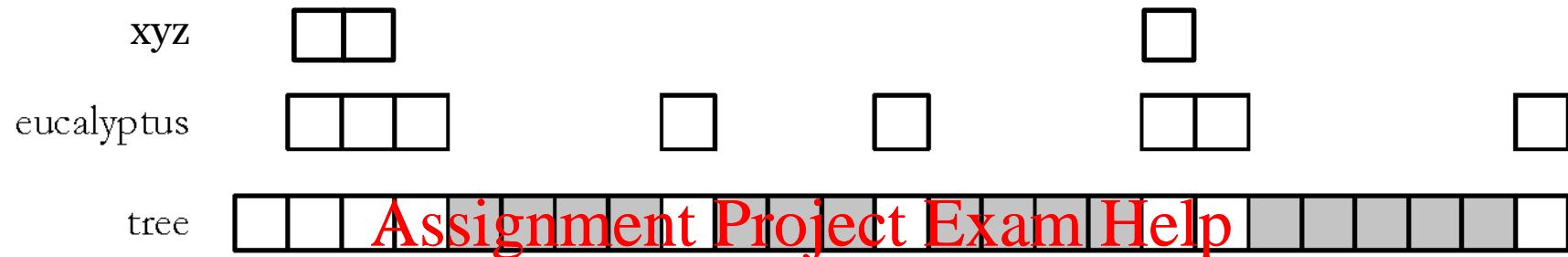
- Threshold methods use number of top-ranked documents needed (k) to optimize query processing
 - for most applications, k is small
- For any query, <https://eduassistpro.github.io/> core that each document needs can be shown to the user
 - score of the k th-highest scoring document
 - gives *threshold* τ
 - optimization methods estimate τ' to ignore documents

Threshold Methods

- For document-at-a-time processing, use score of lowest-ranked document so far for τ'
 - for term-at-a-time, have to use k^{th} -largest score in the accumulator t
- MaxScore met maximum score that remaining documents add to τ'
 - *safe* optimization in that ranking will be the same without optimization

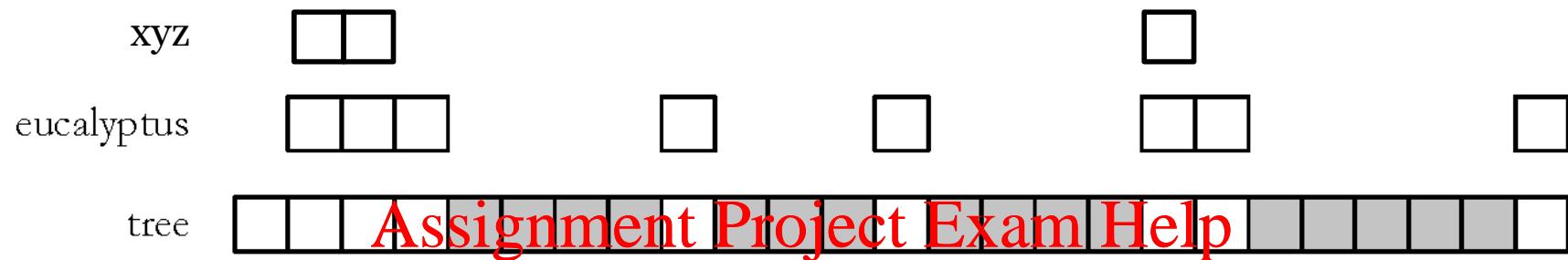
Better than the example in the textbook. See my Note 2 too.

MaxScore Example



- Compute max <https://eduassistpro.github.io/> and sort them in decreasing order (query processing)
- Assume $k = 3$, τ' is lowest score of the current top- k documents
- If $\mu_{tree} < \tau'$ → any doc that scores higher than τ' must contain at least one of the first two keywords (aka *required term set*)
 - Use postings lists of required term set to “drive” the query processing
 - Will only check some of the white postings in the list of “tree” to compute score → at least all gray postings are skipped.

MaxScore



<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Other Approaches

- Early termination of query processing
 - ignore high-frequency word lists in term-at-a-time
 - ignore documents at end of lists in doc-at-a-time
 - *unsafe* optimi <https://eduassistpro.github.io/>
- List ordering
 - order inverted lists by quality (., PageRank) or by partial score
 - makes unsafe (and fast) optimizations more likely to produce good documents

Bottlenecks

- Primary computational bottleneck in scoring: cosine computation
- Can we avoid ~~Assignment Project Exam Help~~ this computation?
- Yes, but may sometimes <https://eduassistpro.github.io/>
 - a doc *not* in the top K ~~Add WeChat edu_assist_pro~~ to the list of K output docs
 - Is this such a bad thing?

Cosine similarity is only a proxy

- Justifications
 - User has a task and a query formulation
 - Cosine matching does not query
 - Thus cosine is a proxy for happiness
- Approximate
 - If we get a list of K docs “close” to query, rank K by cosine measure, should be ok

Generic approach

- Find a set A of *contenders*, with $K < |A| \ll N$
 - A does not necessarily contain the top K , but has many docs among the top K
 - Return the top <https://eduassistpro.github.io/>
- Think of A as pruning non-candidates
- Add WeChat edu_assist_pro
- The same approach is also used for other (non-cosine) scoring functions
- Will look at several schemes following this approach

Index elimination

- Basic algorithm FastCosineScore of Fig 7.1 only considers docs containing at least one query term
- Take this further:
 - Only consider <https://eduassistpro.github.io/>
 - Only consider docs containing query terms
Add WeChat edu_assist_pro

High-idf query terms only

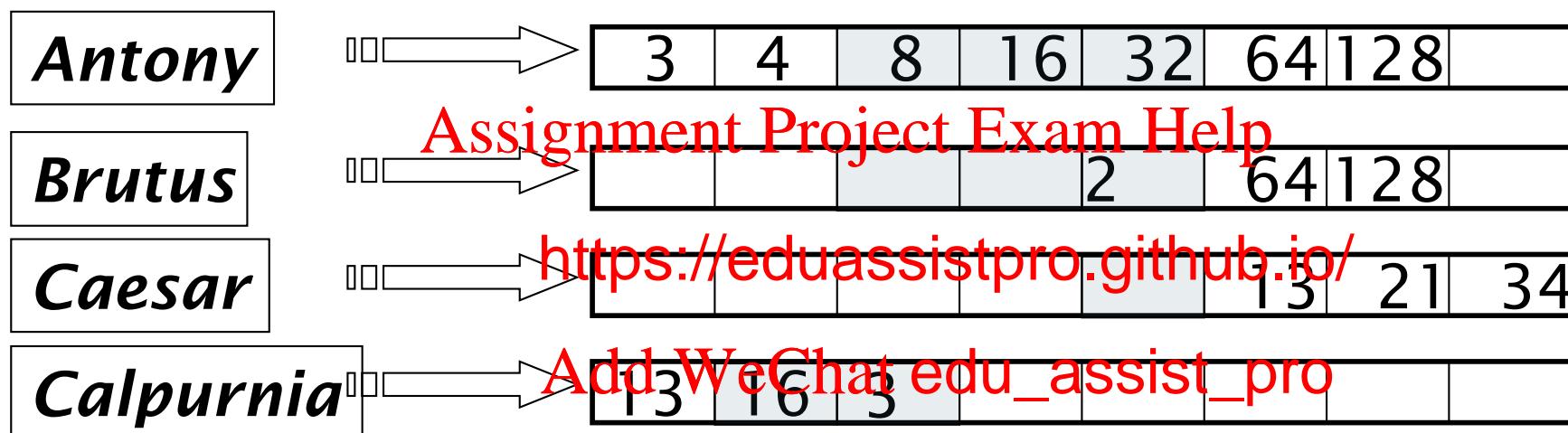
- For a query such as *catcher in the rye*
- Only accumulate scores from *catcher* and *rye*
- Intuition: *in* and *the* contribute little to the scores
and so don't a <https://eduassistpro.github.io/>
- Benefit: Add WeChat edu_assist_pro
 - Postings of low-idf terms have many docs → these (many) docs 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 terms
 - Say, at least 3 out of 4
 - Imposes a “soft conjunction” seen on web search engines (early Google)
- Easy to implement in postings traversal

Can generalize to WAND method (safe)

3 of 4 query terms



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 champion list for t
 - (aka fancy list) <https://eduassistpro.github.io/>
- Note that r has ex build time
 - Add WeChat edu_assist_pro
 - Thus, it's possible that $r < K$
- At query time, only compute scores for docs in $A = \bigcup_{t \in Q} \text{ChampionList}(t)$
 - Pick the K top-scoring docs from amongst these

Inspired by “fancy lists” of Google:

<http://infolab.stanford.edu/~backrub/google.html>

Exercises

- How do Champion Lists relate to Index Elimination?
Can they be used together?
- How can Champion Lists be implemented in an inverted index <https://eduassistpro.github.io/>
 - Note that the champion list to do with small docIDs Add WeChat edu_assist_pro

Static quality scores

- We want top-ranking documents to be both *relevant* and *authoritative*
- *Relevance* is being modeled by cosine scores
- *Authority* is ty <https://eduassistpro.github.io/> property of a document
- Examples of authority signs
 - Wikipedia among websites
 - Articles in certain newspapers
 - A paper with many citations
 - Many diggs, Y!buzzes or del.icio.us marks
 - (Pagerank)

Quantitative

Modeling authority

- Assign to each document a *query-independent quality score* in $[0,1]$ to each document d
 - Denote this by $g(d)$
- Thus, a quantity $\frac{\text{citations}}{\text{https://eduassistpro.github.io/}}$ is scaled into $[0,1]$
 - Add WeChat `edu_assist_pro`
 - Exercise: suggest a formula f

Net score

- Consider a simple total score combining cosine relevance and authority
 - $\text{net-score}(q, d) = \alpha \cosine(q, d) + \beta \text{authority}(d)$
 - Can use some weighting other than an equal weighting
 - Indeed, any function of the two terms of user happiness – more later
- Now we seek the top K docs by net score

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 <https://eduassistpro.github.io/>
 - Postings intersection [Add WeChat edu_assist_pro](#)
 - 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 https://eduassistpro.github.io/0.ms), this allows us to stop processing postings
 - Short of computing scores for

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 Assignment Project Exam Help $g(d) \cdot \text{tf-idf}_{td}$
- Seek top- K res <https://eduassistpro.github.io/> in these champion lists
Add WeChat edu_assist_pro

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 <https://eduassistpro.github.io/>, only traverse all the *high* lists first
 - If we get more than K docs, stop
 - Only union the high lists
 - 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 tiers

Impact-ordered postings

- We only want to compute scores for docs for which $wf_{t,d}$ is high enough
- We sort each ~~posting list~~ by $wf_{t,d}$
- Now: not all p <https://eduassistpro.github.io/>
- How do we compute score 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 <https://eduassistpro.github.io/>
 - One from the postings of each document
- Compute only the scores for 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 <https://eduassistpro.github.io/> query term
 - Add WeChat edu_assist_pro
 - Stop if doc scores relatively
- Can apply to cosine or some other net scores

Why $N^{1/2}$ leaders?

Cluster pruning: preprocessing

- Pick \sqrt{N} docs at random: call these *leaders*
- For every other doc, pre-compute nearest leader
 - Docs attached to each leader have $O(N^{1/2})$ followers;
 - Likely: each leader has $O(N^{1/2})$ followers.

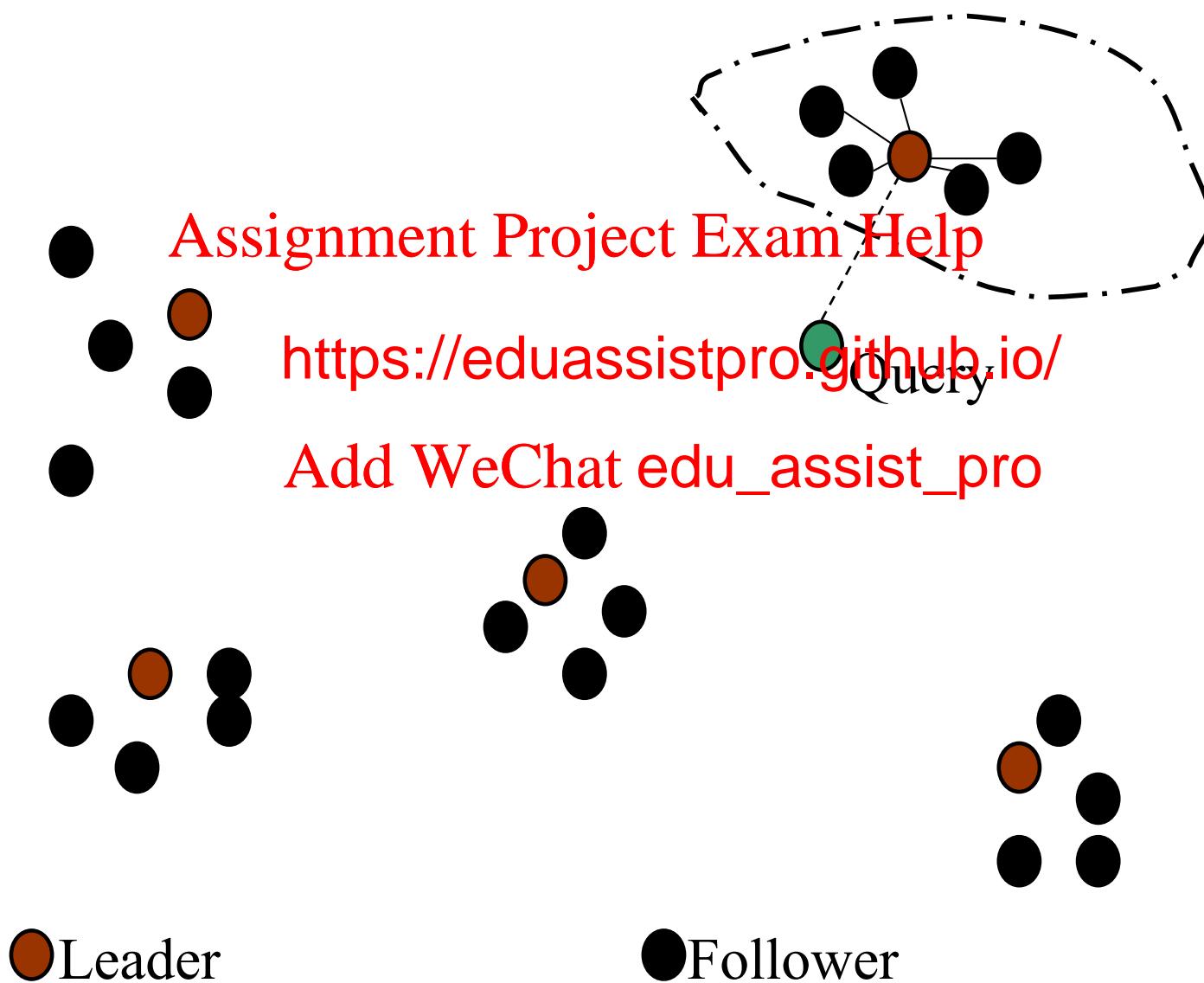
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Cluster pruning: query processing

- Process a query as follows:
 - Given query Q , find its nearest *leader* L .
Assignment Project Exam Help
 - Seek K nearest followers. ong L 's
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Visualization



Why use random sampling

- Fast
- Leaders reflect data distribution

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

General variants

- Have each follower attached to $b1=3$ (say) nearest leaders.
- From query [Assignment 4](#) ([say](#)) nearest leaders and their followers. <https://eduassistpro.github.io/>
- Can recur on leader/follow [Add WeChat edu_assist_pro](#) ction.

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?
 - Hint: write down the cost, and minimize the <https://eduassistpro.github.io/>
- What is the effect of the $\text{Add WeChat } \text{edu_assist_pro}$, $-b2$ 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 <https://eduassistpro.github.io/>
 - Title [Add WeChat edu_assist_pro](#)
 - Date of publication
 - Language
 - Format
 - etc.
- These constitute the metadata about a document

Fields

- We sometimes wish to search by these metadata
 - E.g., find docs authored by William Shakespeare in the year 1601
- Year = 1601 is an example
- Also, author is a field
- Field or parametric index: parameter 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 zone is a region of the doc that can contain an arbitrary amount of text e.g.,
 - Title Assignment Project Exam Help
 - Abstract <https://eduassistpro.github.io/>
 - References ...
- Build inverted indexes on ~~Add WeChat~~ ~~edu_assist_pro~~ ~~elt~~ to permit querying
- E.g., “find docs with *merchant* in the title zone and matching the query *gentle rain*”

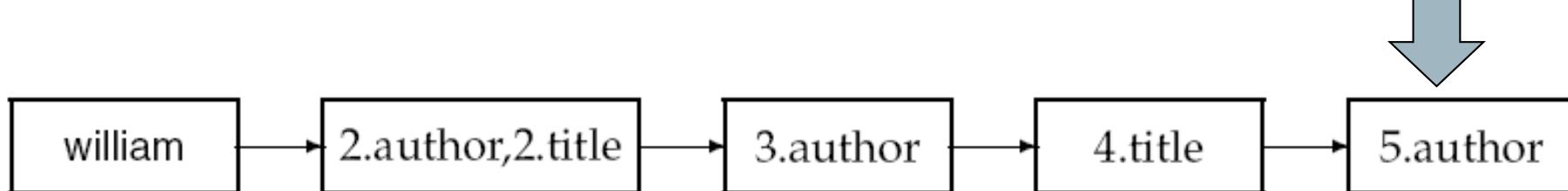
Example zone indexes

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Encode zones in dictionary vs. postings.



Tiered indexes

- Break postings up into a hierarchy of lists
 - Most important
 - ... **Assignment Project Exam Help**
 - Least important <https://eduassistpro.github.io/>
- Can be done by **Add WeChat edu_assist_pro** **sure**
- Inverted index thus broken **rs** 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

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Query term proximity

- Free text queries: just a set of terms typed into the query box – common on the web
- Users prefer ~~Assignment Project Exam Help~~ which query terms occur within close proximity <https://eduassistpro.github.io/>
- Let w be the smallest window 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 4 (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*
Assignment Project Exam Help
 - Run the quer <https://eduassistpro.github.io/>
 - If $<K$ docs con *erest ra*, run the two phrase query *Add WeChat edu_assist_pro*
 - 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 query parser

Aggregate scores

- We've seen that score functions can combine cosine, static quality, proximity, etc.
- How do we know the best combination?
- Some applications:
<https://eduassistpro.github.io/>
- Increasingly common:
mac
WeChat
edu_assist_ed_pro
 - See later lecture

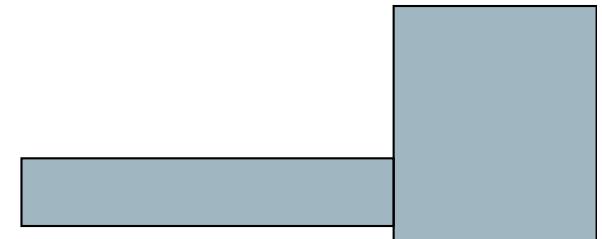
Putting it all together



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Resources

- IIR 7, 6.1

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro