

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

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

Add WeChat `edu_assist_pro`  
Lecture 7: Scoring an assembly

# Assignment Project Exam Help

## Recap: tf-idf weight

---

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

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

- 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

# Assignment Project Exam Help

## Recap: Queries as Add WeChat edu\_assist\_pro

---

- 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 Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

Recap:  $\text{cosine}(\text{query}, \text{document})$

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

Add WeChat edu\_assist\_pro

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

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

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

# Assignment Project Exam Help

This lecture

Add WeChat edu\_assist\_pro

---

- Speeding up vector space ranking
- Putting together a complete search system <https://eduassistpro.github.io/>
  - Will require learning number of miscellaneous topics

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

# Assignment Project Exam Help

Term-at-a-time

Computing cosine similarity  
Add WeChat `edu_assist_pro`

---

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

# Assignment Project Exam Help

## Efficient cosine ran

---

- Find the  $K$  docs in the collection “nearest” to the query  $\Rightarrow K$  largest query-doc cosines.
- Efficient ran
  - Computing <https://eduassistpro.github.io/>
  - Choosing the  $K$  largest cosines efficiently.
    - Can we do this without computing all  $N$  cosines?

# Assignment Project Exam Help

Curse of dimensionality

## Efficient cosine ran

- 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-dime <https://eduassistpro.github.io/>
- But it is solvable for short q d standard indexes support this well

# Assignment Project Exam Help

## Special case—unw

- No weighting on query terms
  - Assume each query term occurs only once
- Then for rank vector <https://eduassistpro.github.io/>
  - Slight simplification of al

# Assignment Project Exam Help

Faster cosine:  
Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

---

Assignment Project Exam Help

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

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)



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

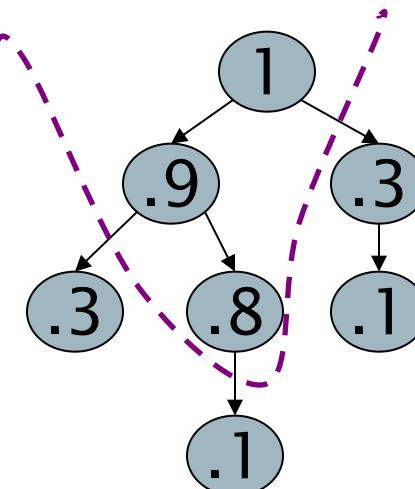
- Typically we want to retrieve the top  $K$  docs (in the cosine ranking for the query)
  - not to totally order all docs in the collection
- Can we pick o <https://eduassistpro.github.io/>?
- Let  $n$  of docs with nonzero
  - We seek the  $K$  best of th

# Assignment Project Exam Help

[http://en.wikipedia.org/wiki/Binary\\_heap](http://en.wikipedia.org/wiki/Binary_heap)

## Use heap for select Add WeChat edu\_assist\_kpro

- Max-heap:
  - Binary tree in which each node's value > the values of children
- Takes  $2n$  operations to sort Add WeChat edu\_assist\_kpro  $n$  elements into  $K$  “winners” read off in  $2\log n$
- Total time is  $O(n + K \log(n))$
- For  $n=1M$ ,  $K=100$ , this is about 1/1000 cost of sorting.

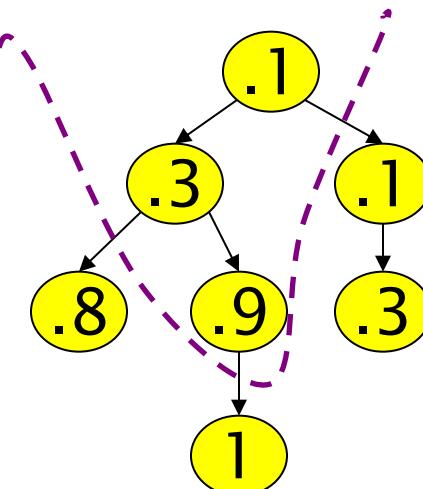


# Assignment Project Exam Help

[http://en.wikipedia.org/wiki/Binary\\_heap](http://en.wikipedia.org/wiki/Binary_heap)

## Use heap for select Add WeChat edu\_assist\_kpro

- What about using a min-heap?
- Use the min-heap to maintain the top k scores so far.
- For each new assignment, update the min-heap
  - H.push (s)
  - H.pop()
- Total time is  $O(n * \log(k) + k * \log(k))$ ; space complexity is  $O(k)$



# Quick Select

---

- QuickSelect is similar to QuickSort to find the top-K elements from an array
  - Takes  $O(n)$  time (in expectation)
- Sorting the to <https://eduassistpro.github.io/> takes  $O(n \log(K))$  time
- Total time is  $O(n + K * \log(K))$

# Query Processing

- Document-at-a-time
  - Calculates complete scores for documents by processing all term lists one document at a time
- Term-at-a-time <https://eduassistpro.github.io/>
  - Accumulates processing term lists one at a time
- Both approaches have optimization techniques that significantly reduce time required to generate scores
  - Distinguish between safe and heuristic optimizations

# Assignment Project Exam Help

Document At-A-Ti  
Add WeChat [edu\\_assist\\_pro](#)

---

Assignment Project Exam Help

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

Add WeChat [edu\\_assist\\_pro](#)

# Assignment Project Exam Help

Document At-A-Ti  
Add WeChat edu\_assist\_pro

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

Term-At-A-Time  
Add WeChat edu\_assist\_pro

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

Term-At-A-Time Add WeChat edu\_assist\_pro

---

// accumulators

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro // Add contains partial score

# Assignment Project Exam Help

## Optimization Tech

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

# Assignment Project Exam Help

## Conjunctive Processes

---

- Requires the result document containing all the query terms (i.e., conjunctive Boolean queries)
  - More efficient
  - Can also be found at <https://eduassistpro.github.io/>
  - Default for TAAT
- Can be combined with both TAAT

# Assignment Project Exam Help

Add WeChat **edu\_assist\_pro**  
Conjunctive  
Term-at-a-Time

Assignment Project Exam Help

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

Add WeChat **edu\_assist\_pro**

# Assignment Project Exam Help

Add WeChat **edu\_assist\_pro**  
Conjunctive  
Document-at-a-Time

Assignment Project Exam Help

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

Add WeChat **edu\_assist\_pro**

# Assignment Project Exam Help

## Threshold Method

- 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*  $\tau$
  - optimization methods estimate  $\tau'$  to ignore documents

# Assignment Project Exam Help

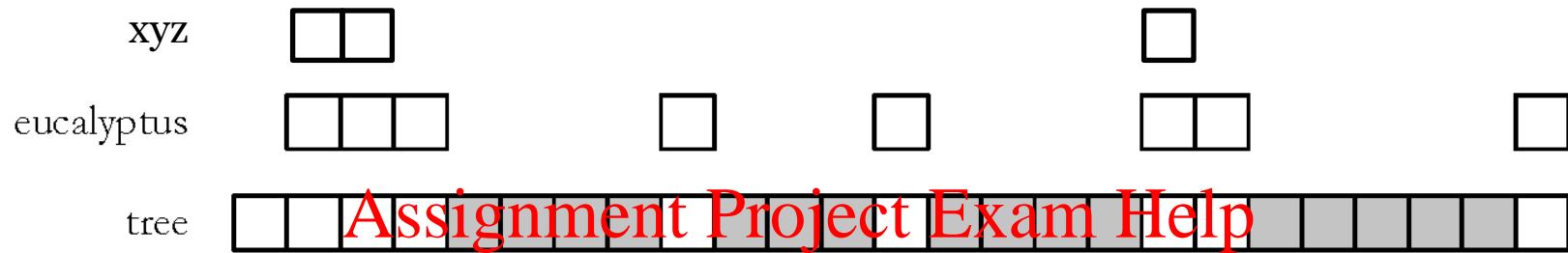
## Threshold Method

- For document-at-a-time processing, use score of lowest-ranked document so far for  $\tau'$ 
  - for term-at-a-time, have to use  $k^{th}$ -largest score in the accumulator
- MaxScore me ximum score that remaining documents to  $\tau'$ 
  - safe* optimization in that ranking will be the same without optimization

# Assignment Project Exam Help

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

# MaxScore Example

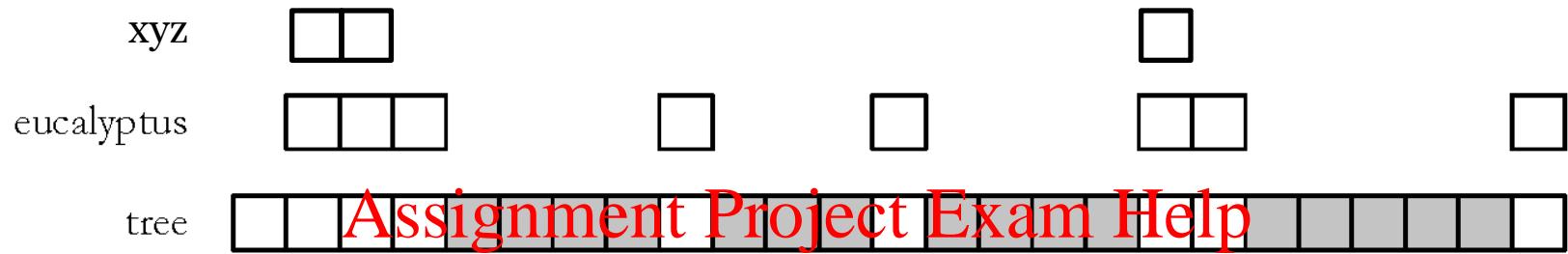


- Compute main search list and sort them in decreasing query processing) **Add WeChat edu\_assist\_pro**
  - Assume  $k = 3$ ,  $\tau'$  is lowest score of the current top- $k$  documents
  - If  $\mu_{tree} < \tau' \rightarrow$  any doc that scores higher than  $\tau'$  must contains *at least one of* 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  $\rightarrow$  at least all gray postings are skipped.

# Assignment Project Exam Help

MaxScore Add WeChat edu\_assist\_pro

---



<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* optim <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

# Assignment Project Exam Help

## Bottlenecks

Add WeChat edu\_assist\_pro

---

- Primary computational bottleneck in scoring: cosine computation
- Can we avoid all this computation?
- Yes, but may <https://eduassistpro.github.io/>
  - a doc *not* in the top  $K$  m to the list of  $K$  output docs
  - Is this such a bad thing?

# Assignment Project Exam Help

## Cosine similarity is ~~Add WeChat~~ ~~edu\_assist\_pro~~ ~~proxy~~

---

- Justifications
  - User has a task and a query formulation
  - Cosine matches ~~Assignment Project Exam Help~~ ~~does to query~~
  - Thus cosine is <https://eduassistpro.github.io/> ~~happiness~~
- Approximate
  - If we get a list of  $K$  docs "closely related" by cosine measure, should be ok

# Assignment Project Exam Help

## 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 from among the top  $K$
  - Return the t <https://eduassistpro.github.io/>
- Think of  $A$  as pruning non-c
- The same approach is also her (non-cosine) scoring functions
- Will look at several schemes following this approach

# Assignment Project Exam Help

## 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

# Assignment Project Exam Help

## High-idf query terms

- For a query such as *catcher in the rye*
- Only accumulate scores from *catcher* and *rye*
- Intuition: *in a* contributes little to the scores  
and so don't consider <https://eduassistpro.github.io/>
- Benefit:
  - Postings of low-idf terms have many docs → these (many) docs get eliminated from set *A* of contenders

# Assignment Project Exam Help

## Docs containing ~~Add WeChat~~ ~~edu\_assist\_pro~~ terms

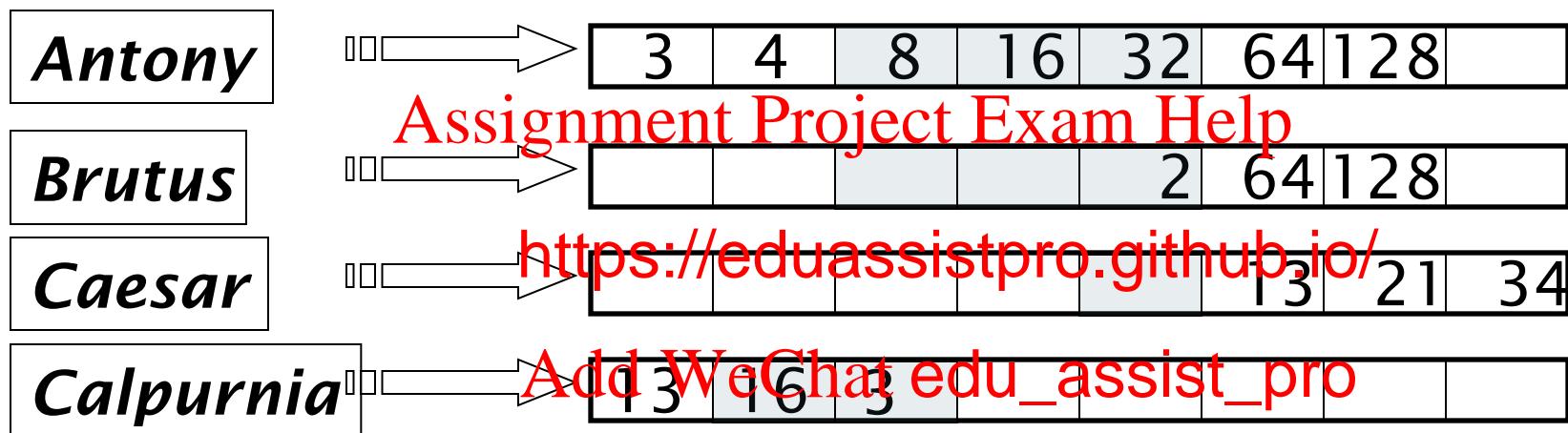
---

- 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

# Assignment Project Exam Help

Can generalize to WAND method (safe)

3 of 4 query terms



Scores only computed for docs 8, 16 and 32.

# Assignment Project Exam Help

## 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  $\propto$  build time
  - 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>

# Assignment Project Exam Help

## Exercises Add WeChat edu\_assist\_pro

---

- 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

# Assignment Project Exam Help

## Static quality score

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

Quantitative

# Assignment Project Exam Help

## 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]
  - Exercise: suggest a formula  $f$

# Assignment Project Exam Help

## Net score

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

# Assignment Project Exam Help

## Top $K$ by net score

- First idea: Order all postings by  $g(d)$
- Key: this is a common ordering for all postings
- Thus, can consider terms' postings for <https://eduassistpro.github.io/>
  - Postings intersection
  - Cosine score computation
- Exercise: write pseudocode for cosine score computation if postings are ordered by  $g(d)$

# Assignment Project Exam Help

## Why order posting

- Under  $g(d)$ -ordering, top-scoring docs likely to appear early in postings traversal
- In time-bound applications (say, we have to return whatever seen in <https://eduassistpro.github.io/> 0 ms), this allows us to stop pos
  - Short of computing scores for  $\text{postings}$

# Assignment Project Exam Help

Champion lists in <https://eduassistpro.com>

---

- Can combine champion lists with  $g(d)$ -ordering
- Maintain for each term a champion list of the  $r$  docs with highest  $g(d) + \text{tf-idf}_{td}$
- Seek top- $K$  re <https://eduassistpro.github.io/> champion lists

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## 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 if
  - If we get more than  $K$  docs, stop at  $K$  and 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

# Assignment Project Exam Help

## Impact-ordered post

- 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 phttps://eduassistpro.github.io/ order!
- How do we compute score <sup>o pick off top  $K$ ?</sup>
  - Two ideas follow

# Assignment Project Exam Help

## 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 in the union
- Compute only the scores for documents in the union

# Assignment Project Exam Help

## 2. idf-ordered term

---

- 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
  - Stop if doc scores relatively
- Can apply to cosine or some other net scores

# Assignment Project Exam Help

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

## Cluster pruning: ~~Add WeChat edu\_assist\_pro~~

- Pick  $\sqrt{N}$  docs at random: call these *leaders*
- For every other doc, pre-compute nearest leader
  - Docs attached to leaders have fewer followers.
  - Likely: each leader has  $\sqrt{N}$  followers.

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

# Assignment Project Exam Help

## Cluster pruning: Add WeChat edu\_assist\_pro

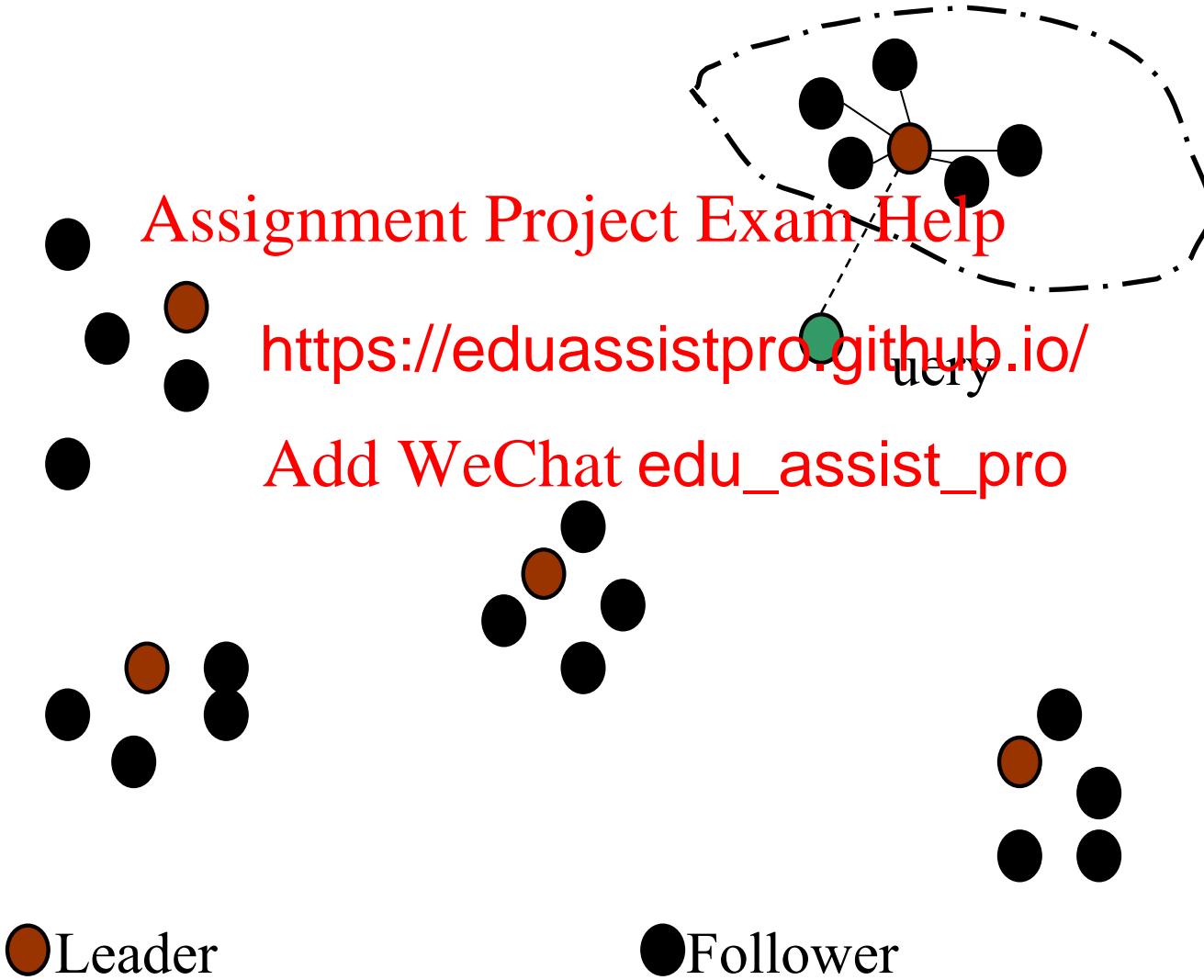
---

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

# Assignment Project Exam Help

## Visualization

Add WeChat edu\_assist\_pro



# Assignment Project Exam Help

## Why use random's Add WeChat edu\_assist\_pro

---

- Fast
- Leaders reflect data distribution

### Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## General variants

- Have each follower attached to  $b1=3$  (say) nearest leaders.
- From query, find  $b2=4$  (say) nearest leaders and their followers. <https://eduassistpro.github.io/>
- Can recur on leader/follower relation.

# Assignment Project Exam Help

## Exercises Add WeChat edu\_assist\_pro

---

- To find the nearest leader in step 1, how many cosine computations do we do?
  - Why did we have VN in the first place?
  - Hint: write  $d$  | its cost, and minimize the <https://eduassistpro.github.io/>
- What is the effect of the  $\text{Add WeChat edu\_assist\_pro}$  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.

# Assignment Project Exam Help

## Parametric and ~~Add WeChat edu\_assist\_pro~~ 20

---

- Thus far, a doc has been a sequence of terms
- In fact documents have multiple parts, some with ~~Assignment Project Exam Help~~ special semantics.
  - Author <https://eduassistpro.github.io/>
  - Title
  - Date of publication
  - Language
  - Format
  - etc.
- These constitute the metadata about a document

# Assignment Project Exam Help

## Fields

Add WeChat edu\_assist\_pro

---

- We sometimes wish to search by these metadata
  - E.g., find docs authored by William Shakespeare in the year 1601, containing *Hamlet*
- Year = 1601 is an example
- Also, author | date, etc
- Field or parametric index: provides range search over each field value
  - Sometimes build range trees (e.g., for dates)
- Field query typically treated as conjunction
  - (doc *must* be authored by shakespeare)

# Assignment Project Exam Help

## Zone

Add WeChat `edu_assist_pro`

---

- 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 `z`      Add WeChat `edu_assist_pro`      H to permit querying
- E.g., “find docs with *merchant* in the title zone and matching the query *gentle rain*”

# Assignment Project Exam Help

Example zone index  
~~Add WeChat~~ <https://eduassistpro.github.io/>

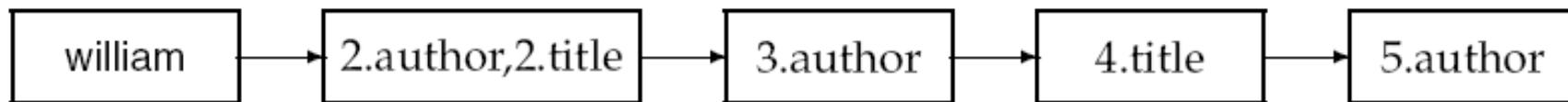
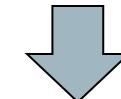
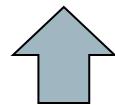
---

Assignment Project Exam Help

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

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

Encode zones in dictionary vs. postings.



# Assignment Project Exam Help

## Tiered indexes

- Break postings up into a hierarchy of lists
  - Most important
  - ...
- Least important <https://eduassistpro.github.io/>
- Can be done
  - Add WeChat edu\_assist\_pro
- Inverted index thus broken
  - sure
  - rs of decreasing importance
- At query time use top tier unless it fails to yield  $K$  docs
  - If so drop to lower tiers

# Assignment Project Exam Help

Example tiered ind   
 Add WeChat edu\_assist\_pro

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Query term ~~Add WeChat~~ proxim

- Free text queries: 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 <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?

# Assignment Project Exam Help

## 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 <https://eduassistpro.github.io/>
  - If  $<K$  docs come back, run the two phrase query *rising 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 query parser

# Assignment Project Exam Help

## 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
  - See later lecture

## Assignment Project Exam Help

Putting it all together  
Add WeChat [edu\\_assist\\_pro](#)

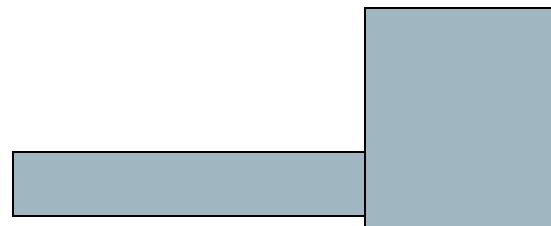
---



Assignment Project Exam Help

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

Add WeChat [edu\\_assist\\_pro](#)



# Assignment Project Exam Help

## Resources

Add WeChat `edu_assist_pro`

---

- IIR 7, 6.1

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`