

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 16: Web ICS

Assignment Project Exam Help

Brief (non-technical) Add WeChat edu_assist_pro

- Early keyword-based engines ca. 1995-1997
 - Altavista, Excite, Infoseek, Inktomi, Lycos
- Paid search ed into
Overture.co <https://eduassistpro.github.io/>
 - Your search ranking depended on how much you paid
 - Auction for keywords: **casino** was expensive!

Assignment Project Exam Help

Brief (non-technical) Add WeChat edu_assist_pro

- 1998+: Link-based ranking pioneered by Google
 - Blew away all early engines save Inktomi
 - Great user experience in search of a business model
 - Meanwhile Got **Assignment Project Exam Help** were nearing \$1 billion
- Result: Google a <https://eduassistpro.github.io/>, independent of search results
 - Yahoo followed suit, acquiring Ov (for news placement) and Inktomi (for search)
- 2005+: Google gains search share, dominating in Europe and very strong in North America
 - 2009: Yahoo! and Microsoft propose combined paid search offering

Assignment Project Exam Help

Add WeChat edu_assist_pro

Assignment Project Exam Help

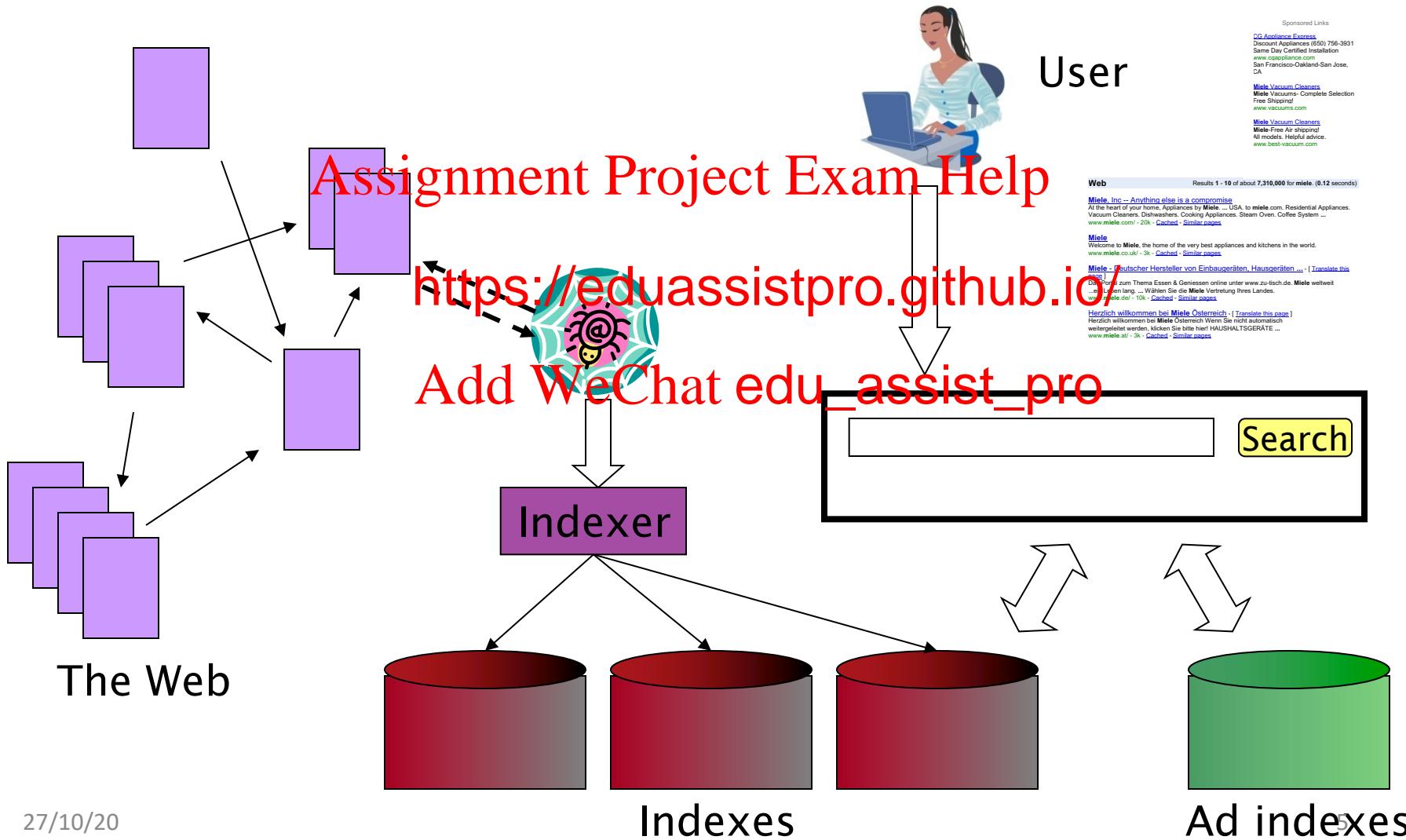
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Algorithmic results.

Assignment Project Exam Help

Web search basics



Assignment Project Exam Help

User Needs

Need [Brod02, RL04]

- Informational – want to learn about something (~40% / 65%)

Low hemoglobin

- Navigational – went to go to that page (~25% / 15%)

Airlines

- Transactional – <https://eduassistpro.github.io/> (~35% / 20%)

- Access a service

Seattle

- Downloads

M

- Shop

images

Canon S410

Gray areas

- Find a good hub

Car rental Brasil

- Exploratory search “see what’s there”

Assignment Project Exam Help

How far do people go? Add WeChat edu_assist_pro

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

(Source: [iprospect.com WhitePaper_2006_SearchEngineUserBehavior.pdf](#))

Assignment Project Exam Help

Users' empirical eval results

- Quality of pages varies widely
 - Relevance is not enough
 - Other desirable qualities (non IR!!)
 - Content: Trustworthy, diverse, non-duplicated, well maintained
 - Web readability: display correctly & fast
 - No annoyances
- Precision vs. r <https://eduassistpro.github.io/>
- What matters
 - Precision at 1? Precision above the fold?
 - Comprehensiveness – must be able to deal with obscure queries
 - Recall matters when the number of matches is very small
- User perceptions may be unscientific, but are significant over a large aggregate

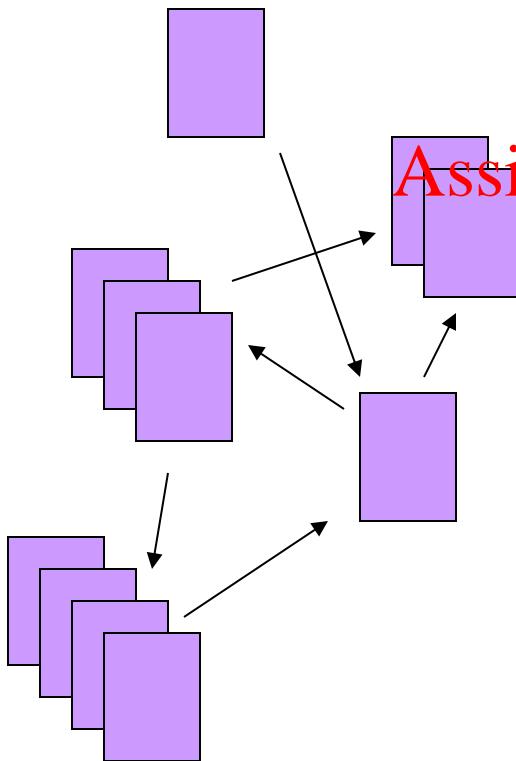
Assignment Project Exam Help

Users' empirical eval engines

- Relevance and validity of results
- UI – Simple, no clutter, error tolerant
- Trust – Results are objective
- Coverage of top
- Pre/Post proces <https://eduassistpro.github.io/>
 - Mitigate user errors (auto spell ch ist,...)
 - Explicit: Search within results, mor ne ...
 - Anticipative: related searches
- Deal with idiosyncrasies
 - Web specific vocabulary
 - Impact on stemming, spell-check, etc
 - Web addresses typed in the search box
- “The first, the last, the best and the worst ...”

Assignment Project Exam Help

The Web document ~~Add WeChat edu_assist pro~~



- No design/co-ordination
- Distributed content creation, linking, democratization of publishing
- Content includes truth, lies, obsolete diction ...
<https://eduassistpro.github.io/>
html, ...), semi-
structure
structure
structure
...)
- Scale much larger than previous text collections ... but corporate records are catching up
- Growth – slowed down from initial “volume doubling every few months” but still expanding
- Content can be *dynamically generated*

Assignment Project Exam Help

Spam

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

- (Search Engine Optimization)

Assignment Project Exam Help

Size of the web

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Assignment Project Exam Help

What is the size of Add WeChat edu_assist_pro?

- Issues
 - The web is really infinite
 - Dynamic content, e.g., calendar
 - Soft 404: w valid page
 - Static web c <https://eduassistpro.github.io/> due to mirroring (~3)
 - Some servers are seldom co
- Who cares?
 - Media, and consequently the user
 - Engine design
 - Engine crawl policy. Impact on recall.

Assignment Project Exam Help

What can we achieve?

- The relative sizes of search engines
 - The notion of a page being indexed is still *reasonably* well defined.
 - Already there <https://eduassistpro.github.io/>
 - Document extension: e.g. engine not yet crawled, by indexing anchor text.
 - Document restriction: All engines restrict what is indexed (first n words, only relevant words, etc.)
- The coverage of a search engine relative to another particular crawling process.

Assignment Project Exam Help

New definition? Add WeChat edu_assist_pro

(IQ is whatever the IQ tests measure.)

- The statically indexable web is whatever search engines in Assignment Project Exam Help index
- Different engines <https://eduassistpro.github.io/> index different things
 - max url depth, max count/hour, rules, priority rules, etc.
- Different engines index different things under the same URL:
 - frames, meta-keywords, document restrictions, document extensions, ...

Assignment Project Exam Help

Relative Size from O

Given two engines

$A \cap B$

Assignment

Sample URLs randomly from A

Check if contained in B and vice versa

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

* Size A

Add WeChat

$A \cap B$

* Size B

$$(1/2) * \text{Size A} = (1/6) * \text{Size B}$$

$$\therefore \text{Size A} / \text{Size B} =$$

$$(1/6) / (1/2) = 1/3$$

Assignment Project Exam Help

Sampling URLs

- Ideal strategy: Generate a random URL and check for containment in each index.

- Problem: Ran ~~Assignment Project Exam Help~~ find! Enough to generate a ran ~~https://eduassistpro.github.io/~~ a given Engine.

- Approach 1: ~~Generate a ran~~ contained in a given engine

- Suffices for the estimation of relative size

- Approach 2: Random walks / IP addresses

- In theory: might give us a true estimate of the size of the web (as opposed to just relative sizes of indexes)

Assignment Project Exam Help

Statistical methods

- Approach 1
 - Random queries
 - Random searches
- Approach 2 <https://eduassistpro.github.io/>
 - Random IP addresses
 - Random walks

Assignment Project Exam Help

Random URLs from ~~Add WeChat edu_assist_pro~~ queries

- Generate random query: how?
 - **Lexicon**: 400,000+ words from a web crawl
 - **Conjunctive Queries**: $w_1 \text{ and } w_2$
e.g., *vocalists A*<https://eduassistpro.github.io/>
- Get 100 result URLs from engine ~~Add WeChat edu_assist_pro~~
- Choose a random URL as the current page check for presence in engine B
- This distribution induces a probability weight $W(p)$ for each page.
- Conjecture: $W(SE_A) / W(SE_B) \sim |SE_A| / |SE_B|$

Not an English dictionary

Assignment Project Exam Help

Query Based Che

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- *Strong Query* to check whether an engine B has a document D :
 - Download D . Get list of words
 - Use 8 low freq words
 - Check if D is present at <https://eduassistpro.github.io/>
- Problems:
 - Near duplicates
 - Frames
 - Redirects
 - Engine time-outs
 - Is 8-word query good enough?

Assignment Project Exam Help

Advantages & disadvantages

Add WeChat edu_assist_pro

- Statistically sound under the induced weight.
- Biases induced by random query
 - Query Bias: Favors content-rich pages in the language(s) of the lexicon
 - Ranking Bias: So s & fetch all
 - Checking Bias: Duplicates, impoverished
 - Document or query restriction light not deal properly with 8 words conjunctive query
 - Malicious Bias: Sabotage by engine
 - Operational Problems: Time-outs, failures, engine inconsistencies, index modification.

Assignment Project Exam Help

Random searches

- Choose random searches extracted from a local log [Lawrence & Giles 97] or build “random searches”
[Notess] Assignment Project Exam Help
 - Use only que <https://eduassistpro.github.io/>
 - Count norma
 - Use ratio statistics

Assignment Project Exam Help

Advantages & disadvantages

- Advantage
 - Might be a better reflection of the human perception of coverage
- Issues <https://eduassistpro.github.io/>
 - Samples are correlated with log
 - Duplicates
 - Technical statistical problems (must have non-zero results, ratio average not statistically sound)

Assignment Project Exam Help

Random searches

- 575 & 1050 queries from the NEC RI employee logs
- 6 Engines in 1998, 11 in 1999
- Implemented Assignment Project Exam Help
 - Restricted
 - Counted URLs
 - Match
 - Computed size ratio & overlap
 - Estimated index size ratio & overlap by averaging over all queries

Assignment Project Exam Help

Queries from Lawrence

- *adaptive access control*
- *neighborhood preservation*
- *topographic hamiltonian structures*
- *right linear gramm pulse width modulation*
- *unbalanced prior probabilities ranked assignment method*
- *internet explorer favourites importing karvel thornber zili liu*
- *softmax activation function bose multidimensional system theory gamma mfp sis*
- *exploring neural marking counterpropagation network fat shattering dimension abelson amorphous computing*

Assignment Project Exam Help

Random IP address

- Generate random IP addresses
- Find a web server at the given address
 - If there's one
- Collect all pages at <https://eduassistpro.github.io/>
 - From this, choose a page at

Assignment Project Exam Help

Random IP addresses

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- HTTP requests to random IP addresses
 - Ignored: empty or authorization required or excluded
 - [Lawr99] Estimated 2.8 million IP addresses running crawlable web from observing 2500 servers
 - OCLC using IP sampling found 3.5 billion hosts in 2001
 - Netcraft [Netc02] accessed 37.2 million hosts in July 2002
- [Lawr99] exhaustively crawled 2500 servers and extrapolated
 - Estimated size of the web to be 800 million pages
 - Estimated use of metadata descriptors:
 - Meta tags (keywords, description) in 34% of home pages, Dublin core metadata in 0.3%

Assignment Project Exam Help

Advantages & disadvantages

- Advantages
 - Clean statistics
 - Independent of crawling strategies
- Disadvantages
 - Doesn't deal w/ <https://eduassistpro.github.io/>
 - Many hosts might not accept requests
 - No guarantee all pages are linked
 - Eg: employee pages
 - Power law for # pages/hosts generates bias towards sites with few pages.
 - But bias can be accurately quantified IF underlying distribution understood
 - Potentially influenced by spamming (multiple IP's for same server to avoid IP block)

Assignment Project Exam Help

Random walks

- View the Web as a directed graph
- Build a random walk on this graph
 - Includes various “jump” rules back to visited sites
 - Does not get stuck
 - Can follow all links
 - Converges to a stationary distribution
 - Must assume graph is finite and irreducible for the walk.
 - Conditions are not satisfied (cookie crumbs, flooding)
 - Time to convergence not really known
 - Sample from stationary distribution of walk
 - Use the “strong query” method to check coverage by SE

Assignment Project Exam Help

Advantages & disadvantages

- Advantages
 - “Statistically clean” method at least in theory!
 - Could work even for infinite web (assuming convergence) under certain conditions
- Disadvantage
 - List of seeds is a problem.
 - Practical approximation might not be valid.
 - Non-uniform distribution
 - Subject to link spamming

Assignment Project Exam Help

Conclusions

- No sampling solution is perfect.
- Lots of new ideas ...
-but the pro
- Quantitative <https://eduassistpro.github.io/> and a good research problem

Assignment Project Exam Help

Duplicate detection

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Assignment Project Exam Help

Duplicate document

- The web is full of duplicated content
- Strict duplicate detection = exact match
 - Not as common
- But many, many duplicates
 - E.g., Last modified date difference between two copies of a page

Assignment Project Exam Help

Duplicate/Near-Duplicate Detection

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- *Duplication*: Exact match can be detected with fingerprints
- *Near-Duplication*: Approximate match
 - Overview <https://eduassistpro.github.io/>
 - Compute syntactic similarity measure
 - Use similarity threshold to detect near-duplicates
 - E.g., Similarity > 80% => Documents are “near duplicates”
 - Not transitive though sometimes used transitively

Assignment Project Exam Help

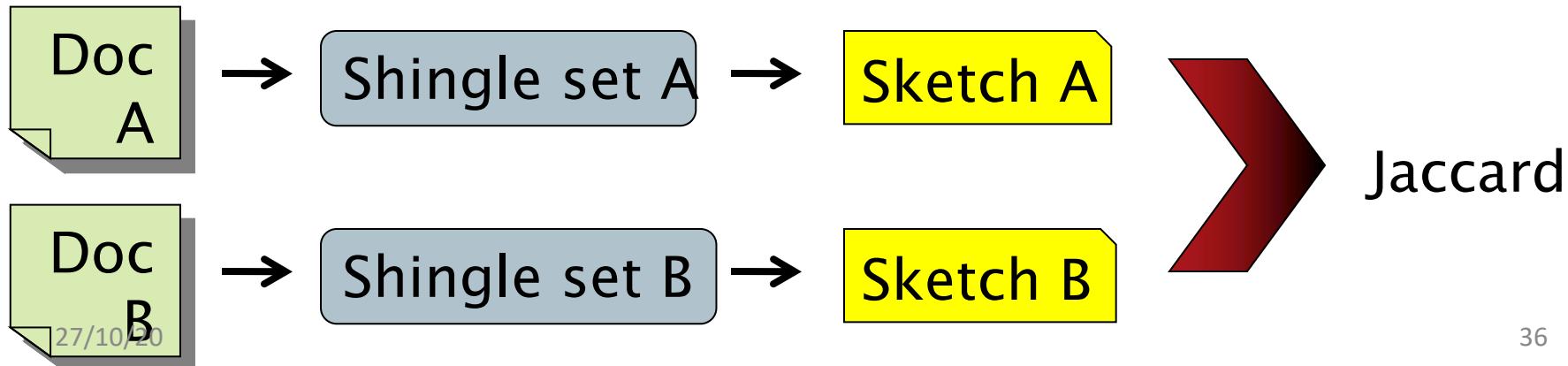
Computing Similarity

- Features:
 - Segments of a document (natural or artificial breakpoints)
 - Shingles (~~Assignment Project Exam Help~~)
 - *a rose is a ros* [https://eduassistpro.github.io/
a_rose_is_a](https://eduassistpro.github.io/a_rose_is_a)
rose_is_a_rose [https://eduassistpro.github.io/
rose_is_a_rose](https://eduassistpro.github.io/rose_is_a_rose)
is_a_rose_is [https://eduassistpro.github.io/
is_a_rose_is](https://eduassistpro.github.io/is_a_rose_is)
a_rose_is_a [https://eduassistpro.github.io/
a_rose_is_a](https://eduassistpro.github.io/a_rose_is_a)
- Similarity Measure between two docs (= sets of shingles)
 - Set intersection
 - Specifically (Size_of_Intersection / Size_of_Union)

Assignment Project Exam Help

Shingles +~~Add WeChat~~ Set Inter

- Computing exact set intersection of shingles between all pairs of documents is expensive/intractable
- Approximate <https://eduassistpro.github.io/shingles> from each (a *sketch*)
- Estimate $(\frac{\text{size_of_intersection}}{\text{size_of_union}})$ based on a short sketch



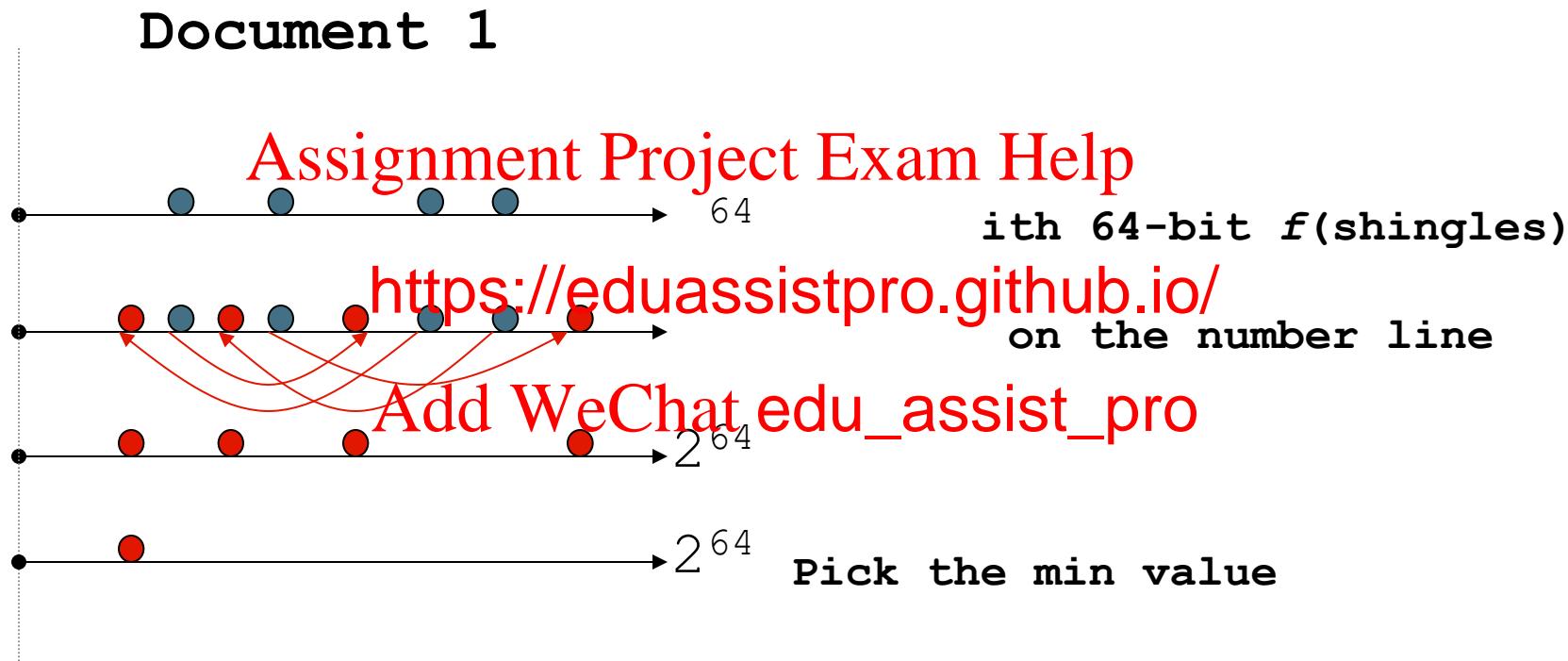
Assignment Project Exam Help

Sketch of a document

- Create a “sketch vector” (of size ~ 200) for each document
 - Documents that share $\geq t$ (say 80%) corresponding vector elements <https://eduassistpro.github.io/>
 - For doc D , $\text{sketch}_D[i]$ is
 - Let f map all shingles in t to \mathbb{F}_{2^m} (e.g., $f = \text{fingerprinting}$)
 - Let π_i be a *random permutation* on $[0, 2^m-1]$
 - Pick $\text{MIN} \{\pi_i(f(s))\}$ over all shingles s in D

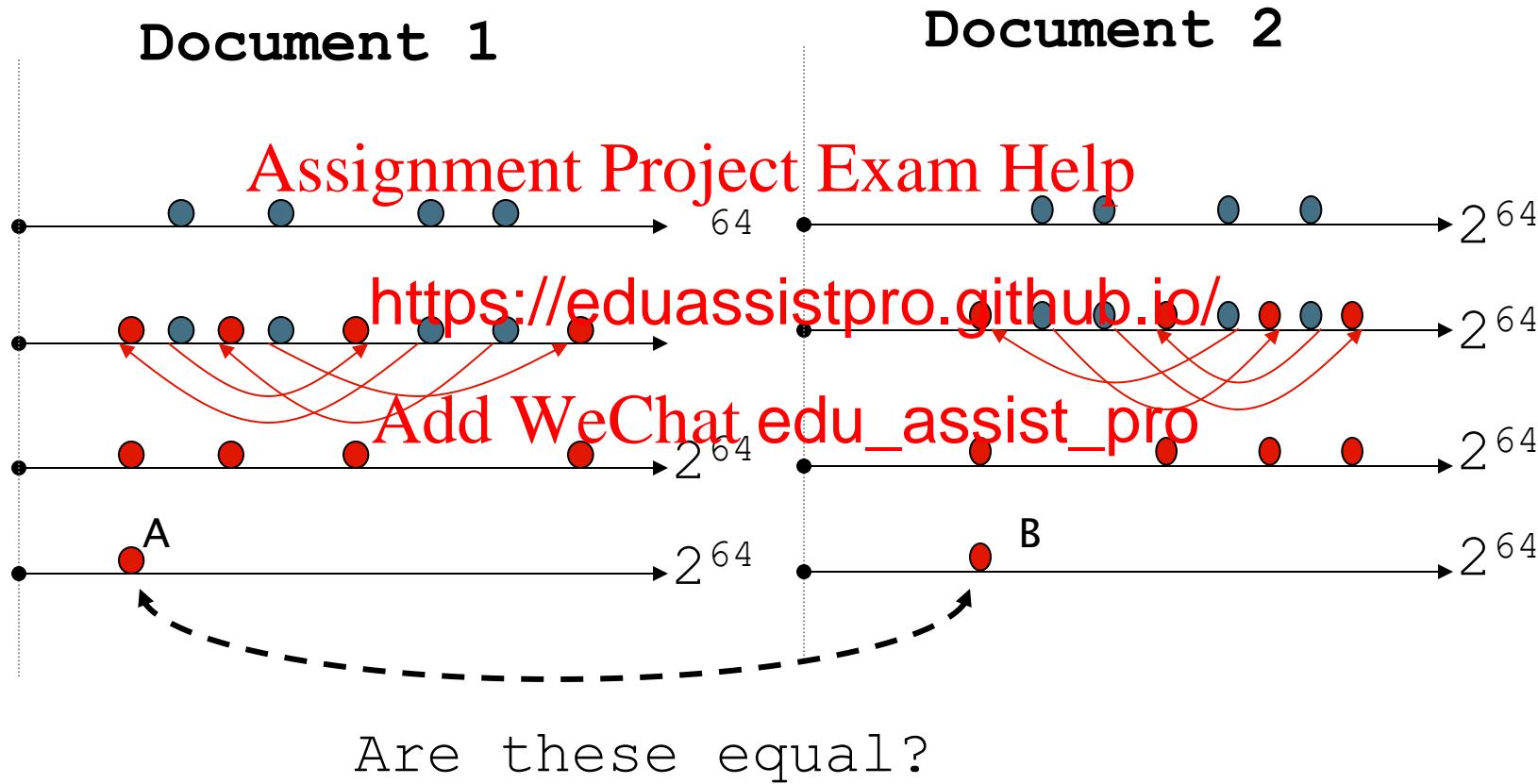
Assignment Project Exam Help

Computing Sketch[Add WeChat edu_assist_pro oo1]



Assignment Project Exam Help

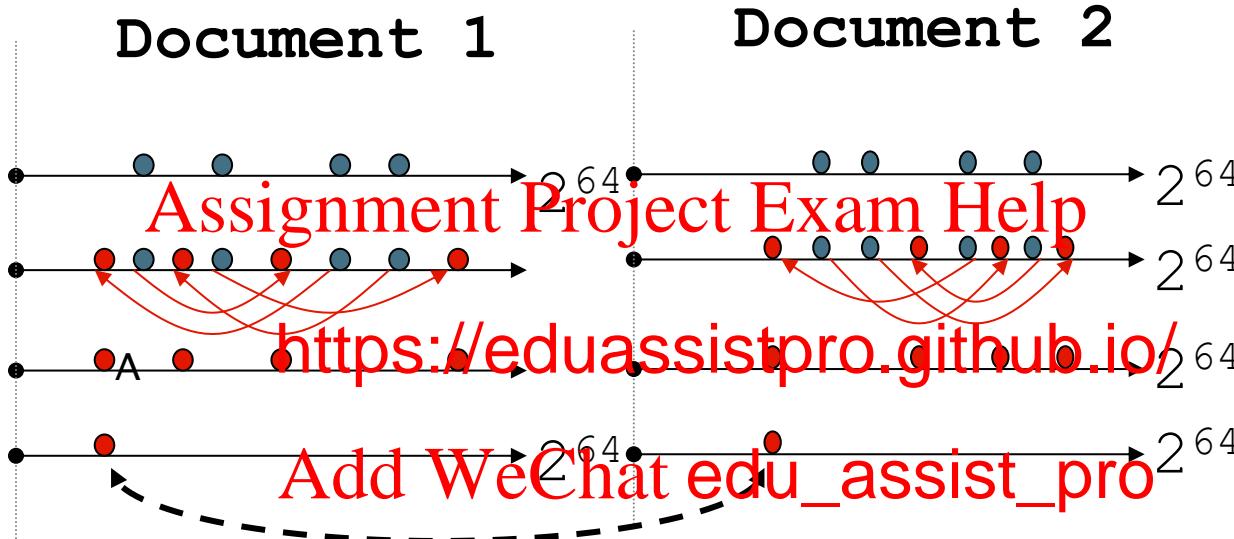
Test if $\text{Doc1.Sketch}[i] = \text{Doc2.Sketch}[i]$



Test for 200 random permutations: $\pi_1, \pi_2, \dots, \pi_{200}$

Assignment Project Exam Help

However... Add WeChat edu_assist_pro



$A = B$ iff the shingle with the MIN value in the union of Doc1 and Doc2 is common to both (i.e., lies in the intersection)

Claim: This happens with probability

$\text{Size_of_intersection} / \text{Size_of_union}$

Why?
↓

Assignment Project Exam Help

Set Similarity of C_i and C_j

$$\text{Jaccard}(C_i, C_j) = \frac{|C_i \cap C_j|}{|C_i \cup C_j|}$$

Assignment Project Exam Help

- View sets as columns
e row for each element in the <https://eduassistpro.github.io/>
tes presence of item i in set j
- Example

$C_1 \quad C_2$

0	1
1	0
1	1
0	0
1	1
0	1

$\text{Jaccard}(C_1, C_2) = 2/5 = 0.4$

Assignment Project Exam Help

Key Observation

- For columns C_i, C_j , four types of rows

	C_i	C_j	
type A	1	1	Assignment Project Exam Help
type B	1	0	https://eduassistpro.github.io/
type C	0	1	
type D	0	0	Add WeChat edu_assist_pro

- Overload notation: $A = \# \text{ of rows of type A}$
- Claim**

$$\text{Jaccard}(C_i, C_j) = \frac{A}{A + B + C}$$

Assignment Project Exam Help

“Min” Hashing

- Randomly permute rows
- Hash $h(C_i) = \text{index of first row with 1 in column } C_i$
- Surprising Pro <https://eduassistpro.github.io/>
 $\Pr [h(C_i) = h(C_j)] = \frac{1}{J}$
- Why?
 - Both are $A/(A+B+C)$
 - Look down columns C_i, C_j until first non-Type-D row
 - $h(C_i) = h(C_j) \leftrightarrow$ type A row

Assignment Project Exam Help

Min-Hash sketches

- Pick P random row permutations
- MinHash sketch

$\text{sketch}(C) = \text{list of } K \text{ indexes of first rows with 1 in column } C$

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

dom variabl

- Similarity of signatures
- Let $\text{sim}[\text{sketch}(C_i), \text{sketch}(C_j)] = \text{fraction of permutations where MinHash values agree}$
- Observe $E[\text{sim}(\text{sig}(C_i), \text{sig}(C_j))] = \text{Jaccard}(C_i, C_j)$

Assignment Project Exam Help

Practical Implementation

- Random permutation is hard to obtain; simulate them using universal hashing instead
 - $h: \{0, 1, 2, \dots, U\} \rightarrow \{0, 1, 2, \dots, M\}$
 - $h(x) = ((a^*x + b) \mod P) \mod M$
 - where
 - $P \gg U$ and is a prime number
 - a, b are two randomly chosen integers modulo P and $a \neq 0$
 - $\text{sketch}(C) = \{ \text{argmin}_{e \in C} \{ h_i(e) \} \mid 1 \leq i \leq k \}$

Assignment Project Exam Help

Example Add WeChat edu_assist_pro

			Signatures		
			S ₁	S ₂	S ₃
	C ₁	C ₂	1	2	1
	C ₂	C ₃	0	5	4
	C ₃		4	3	5
Assignment Project Exam Help			Perm 1 = (12345) Perm 2 = (54321)		
R ₁	1	0	1	2	2
R ₂	0	1	1		
R ₃	1	0	0	Add WeChat edu_assist_pro	
R ₄	1	0	1		
R ₅	0	1	0		

Similarities			
	1-2	1-3	2-3
Col-Col	0.00	0.50	0.25
Sig-Sig	0.00	0.67	0.00

Assignment Project Exam Help

Example Using the ~~Add WeChat edu_assist pro~~ Hashing

$$h(x) = (7x+1 \bmod 31) \bmod 9$$

$$g(x) = (17x+8 \bmod 31) \bmod 9$$

	C_1	C_2	C_3
R_1	1	0	1
R_2	0	1	1
R_3	1	0	0
R_4	1	0	1
R_5	0	1	0

Note: this example results in different sketches from the previous slide

$$S_1 = \{R_1, R_3, R_4\}$$

$$h(e) = \{8, 4, 2\} \rightarrow \min_elem = R_4$$

$$g(e) = \{7, 1, 5\} \rightarrow \min_elem = R_3$$

$$\text{sketch}(S1) = \{R_4, R_3\}$$

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

$$h(e) = \{8, 4, 2\} \rightarrow \min_elem = R_5$$

$$g(e) = \{7, 1, 5\} \rightarrow \min_elem = R_5$$

$$\text{sketch}(S1) = \{R_5, R_5\}$$



Therefore, estimated similarity between S_1 and S_2 is $0/2 = 0.0$

Assignment Project Exam Help

All signature pairs

- Now we have an extremely efficient method for estimating a Jaccard coefficient for a single pair of documents.
- But we still have <https://eduassistpro.github.io/> where N is the number
 - Still slow
- One solution: locality sensitive hashing (LSH)
- Another solution: Sorting (Henzinger 2006)

Assignment Project Exam Help

SimHash Add WeChat edu_assist_pro

- Generalization of LSH to other similarity measures
[Charikar, STOC 02]

- $\theta(\mathbf{x}, \mathbf{y})$: relates to cosine similarity
- $h_{\mathbf{u}}(\mathbf{x}) = \text{sig}(\mathbf{x}^T \mathbf{u})$ where \mathbf{u} is a unit vector
- then $\Pr[h_{\mathbf{u}}(\mathbf{x}) = h_{\mathbf{u}}(\mathbf{y})] = \frac{1}{\pi}$

Assignment Project Exam Help

Practical Implementation

- Near duplicate Web page detection from google [Henzinger, SIGIR06] [Manku et al, WWW07]
 - Document D → set of tokens with idf weighting → form a set of “features” $v(D)$
 - Each feature is a \mathbb{R}^n -dimensional binary vector of [-1,1]
 - Sum up the weighted projection vectors $v(D) \rightarrow r(D)$
 - a f -bit signature $\text{sig}(D) \leftarrow \text{sign}(r(D))$
- Results (in comparison with Shingling)
 - Fairly accurate and stable
 - Does not capture order among tokens

Assignment Project Exam Help

Simhash Add WeChat `edu_assist_pro`

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

Simhash Example

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

count-min sketch,
approximately
preserves the
inner product of
the raw feature
vectors



$$\begin{array}{r} 2 - 2 + 1 \\ - 1 - 1 + \\ 1 - 1 + 1 \\ + 1 - 1 + \\ 1 - 1 + 1 \\ = 1 \end{array}$$

Assignment Project Exam Help

More resources

- IIR Chapter 19

Assignment Project Exam Help

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

Add WeChat edu_assist_pro