Add WeChat edu_assist_pro

Introduction to

Assignment Project Exam Help
Informa
https://eduassistpro.github.io/

Lecture 6: Scoring, Ter g and the Vector Space Model

This lectured McSet edu_assist2pt6.4.3

- Ranked retrieval
- Scoring documents
- Term frequency Assignment Project Exam Help
- Collection sta https://eduassistpro.github.io/
- Weighting schemesWeChat edu_assist_pro
- Vector space scoring

Ranked retrieWaChat edu_assist_pro

- Thus far, our queries have all been Boolean.
 - Documents either match or don't.
- Good for expert users with precise understanding of their needs a https://eduassistpro.github.io/
 - Also good for applications: A can easily consume 1000s of results.
- Not good for the majority of users.
 - Most users incapable of writing Boolean queries (or they are, but they think it's too much work).
 - Most users don't want to wade through 1000s of results.
 - This is particularly true of web search.

Problems with Broke are arch: feast or family eChat edu_assist_pro

- Boolean queries often result in either too few (=0) or too many (1000s) results.
- Query 1: 'Standard User dink & 50m H200,000 hits
- Query 2: "sta https://eduassistpro.gitl@pdidpund": 0 hits
 Add WeChat edu_assist_pro
- It takes a lot of skill to com query that produces a manageable number of hits.
 - AND gives too few; OR gives too many

Ranked retrieWaClmq edu_assist_pro

- Rather than a set of documents satisfying a query expression, in ranked retrieval models, the system returns an ordering overtieet (top) additionents in the collection wit https://eduassistpro.github.io/
- Free text queries: Rather th Add WeChat edu_assist_pro operators and expressions, query is just one or more words in a human language
- In principle, there are two separate choices here, but in practice, ranked retrieval models have normally been associated with free text queries and vice versa

Feast of farmine! rolet axproblem in ranked retrieval Chat edu_assist_pro

- When a system produces a ranked result set, large result sets are not an issue
 - Indeed, the sign of the Project Exam Halpue
 - We just show t https://eduassistpro.github.io/
 - We don't ov

Add WeChat edu_assist_pro

Premise: the ranking algorithm works

Scoring as Atld Wbasit edu_assisted retrieval

- We wish to return in order the documents most likely to be useful to the searcher
- How can We rank-order the documents in the collection wit https://eduassistpro.github.io/
- Assign a score say in [0, 1 document Add WeChat edu_assist_pro
- This score measures how w ent and query "match".

Query-docath Yeachat edu_assist_scores

- We need a way of assigning a score to a query/document pair
- Let's start with a one term query ment Project Exam Help
- If the query tehttps://eduassistpro.githubqoument: score should be 0 Add WeChat edu_assist_pro
- The more frequent the que the document,
 the higher the score (should be)
- We will look at a number of alternatives for this.

Take 1: Jackel Mechat edu_asstst_pro

- Recall from Lecture 3: A commonly used measure of overlap of two sets A and B
- jaccard(A, B) signment | Project Exam Help
- jaccard(A,A) = https://eduassistpro.github.io/
- jaccard(A,B) = QiftAweethaeedu_assist_pro
- A and B don't have to be th
- Always assigns a number between 0 and 1.

Jaccard coefficient edu_assist_example

- What is the query-document match score that the Jaccard coefficient computes for each of the two documents setoment Project Exam Help
- Query: ides of https://eduassistpro.github.io/
- Document 1: caesar died in Add WeChat edu_assist_pro
- Document 2: the long marc

the term *Ides of March* is best known as the date that Julius Caesar was killed in 709 AUC or 44 B.C

Issues with dia wearn edu_assisting

- 1 It doesn't consider *term frequency* (how many times a term occurs in a document)
- 2 Rare terms in a collection are more informative than freque https://eduassistpro.github.ns/der this information
- We need a more sophisticated w
 We need a more sophisticated w

Recall (Leisture In) Birany Herm-document And Wellmtedu_assist_pro

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1 A seign	mant Drain	oot Evom	Halm	0	1
Brutus	ASSIGII	ment Proje	ect Exam	перр	0	0
Caesar	1	no://oduor	ociotoro a	ithulh ic	, 1	1
Calpurnia	0	ps://eduas	ssistpro.g	ուսար.ու	0	0
Cleopatra	1 A	ld WeCha	tedu ass	sist ⁰ nro	0	0
mercy	1		t caa_asc	1 1	1	1
worser	1	0	1	1	1	0

Each document is represented by a binary vector $\in \{0,1\}^{|V|}$

Term-docarder the edu_assistrices

- Consider the number of occurrences of a term in a document:
 - Each document in Each

https://eduassistpro.github.io/

	Antony and Cleopath	Glappe Gestals.	tedu_as	sisŧ <u>m</u> pro	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

Bag of workers Who Chat edu_assist_pro

- Vector representation doesn't consider the ordering of words in a document
- John is quicker than Mary and Mary is quicker than John have the https://eduassistpro.github.io/
- This is called the bag of wo Add WeChat edu_assist_pro
- In a sense, this is a step bac itional index was able to distinguish these two documents.
- We will look at "recovering" positional information later in this course.
- For now: bag of words model

Term frequently Othat edu_assist_pro

- 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 want to use the when computing query-document match scores. https://eduassistpro.github.io/
- Raw term frequency is not ant:
 Add WeChat edu_assist_pro
 A document with 10 occurre erm is
 - A document with 10 occurre erm is more relevant than a document with 1 occurrence of the term.
 - But not 10 times more relevant.
- Relevance does not increase proportionally with term frequency.

NB: frequency = count in IR

Log-frequency we tedu_assist_pro

The log frequency weight of term t in d is

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

- $0 \rightarrow 0, 1 \rightarrow 1$, https://eduassistpro.githulq.je/tc.
- Score for a documentequent edu_assistopen terms t in both q and d:
- score $=\sum_{t \in q \cap d} (1 + \log t f_{t,d})$
- The score is 0 if none of the query terms is present in the document.

Document frequent edu_assist_pro

- Rare terms are more informative than frequent terms
 - Recall stop words Assignment Project Exam Help
- Consider a ter are in the collection (e.g. https://eduassistpro.github.io/
- A document containing thist edu_assisy_likely to be relevant to the query arachn
- → We want a high weight for rare terms like arachnocentric.

Document of the cou_assisting code

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., mon, increase, Fine) Help
- A document chttps://eduassistpro.githmpre/likely to be relevant than a docume Add WeChat edu_assist_pro
- But it's not a sure indicator ce.
- → For frequent terms, we want high positive weights for words like high, increase, and line
- But lower weights than for rare terms.
- We will use document frequency (df) to capture this.

idf weightAdd WeChat edu_assist_pro

- df_t is the <u>document</u> frequency of t: the number of documents that contain t
 - df, is an Aveige measure of the fixer attitudeness of t
 - $df_t \le N$ https://eduassistpro.github.io/
- We define the by Add WeChat edu_assist_pro idf $t = log_{10} (N/df_t)$
 - We use $log(N/df_t)$ instead of N/df_t to "dampen" the effect of idf.

Will turn out the base of the log is immaterial.

idf exampled Wpp by edu_assist_prollion

term	df_t	idf_t
calpurnia	1 · · · · · · · · · · · · · · · · · · ·	TT 1
animal	signment Project Exam	Help
sunday	https://eduassistpro.g	github.io/
fly	Add WeChat edu_ass	sist nro
under	Add Weenat edd_as	3131_p10
the	1,000,000	

$$idf_t = log_{10} (N/df_t)$$

There is one idf value for each term t in a collection.

Effect of idel by Chat edu_assist_pro

- Does idf have an effect on ranking for one-term queries, like
 - iPhone Assignment Project Exam Help
- idf has no eff https://eduassistpro.github.rips
 - idf affects the ranking of doc least two terms dd WeChat edu_assist_pro
 - For the query capricious person, idf weighting makes occurrences of capricious count for much more in the final document ranking than occurrences of person.

Collection And Wook edu_assiste graency

 The collection frequency of t is the number of occurrences of t in the collection, counting multiple occurrences. Project Exam Help

Example: https://eduassistpro.github.io/

Word	Collection frequen Add WeChat ec	ent frequency lu_assist_pro
insurance	10440	3997
try	10422	8760

Which word is a better search term (and should get a higher weight)?

tf-idf weighting Chat edu_assist_pro

 The tf-idf weight of a term is the product of its tf weight and its idf weight. Assignment Project Exam Help

$$\mathbf{w}_{t,d} = (1 \frac{N/\mathrm{df}_t}{\text{https://eduassistpro.github.io/}^t})$$

- Best known weightingeschet edu_assistnation retrieval
 - Note: the "-" in tf-idf is a hyphen, not a minus sign!
 - Alternative names: tf.idf, tf x idf
- Increases with the number of occurrences within a document
- Increases with the rarity of the term in the collection

Final ranking of Chot edu_assist_for a query

$$Score(q,d) = \sum_{\text{Assignment Project} Example lp} tf.idf_{t,d}$$

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Binary -> cody technic edu_assist patrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.Assign	ment Proje	ect Exam	Help	0	0
Caesar	8.59			1.51	0.25	0
Calpurnia	o htt	ps://eduas	ssistpro.g	ithub.ic	0 /	0
Cleopatra	2.85	0		. 0	0	0
mercy	1.51 AC	ld WeChat	t edu_ass	sist _{.4} pro	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

Each document is now represented by a real-valued vector of tf-idf weights $\in \mathbb{R}^{|V|}$

Documents de Welcet edu_assist_pro

- So we have a |V|-dimensional vector space
- Terms are axes of the space
- Assignment Project Exam Help
 Documents ar this space
- Very high-dimhttps://eduassistpro.githubfio/ dimensions when you apple edu_assist_proearch engine
- These are very sparse vectors most entries are zero.

Queries as we'to hat 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 https://eduassistpro.github.io/
- proximity = similarity of ve Add WeChat edu_assist_pro
- proximity ≈ inverse of dista
- Recall: We do this because we want to get away from the you're-either-in-or-out Boolean model.
- Instead: rank more relevant documents higher than less relevant documents

Formalizing detector edu_assistroximity

- First cut: distance between two points
 - (= distance between the end points of the two vectors)
- Euclidean Assignment Project Exam Help
- Euclidean disthttps://eduassistpro.github.io/
- ... because Euclidean dist edu_assist for vectors of different lengths.

Why distance ishat edu_assistapro

```
The Euclidean
distance between q
and \overrightarrow{d_2} is large Avenignment Project Exam Help
though the
                       https://eduassistpro.github.io/
distribution of ter
in the query \overrightarrow{q} and t \underbrace{Aeld}_{Me} WeChat edu_assist_pro
distribution of
terms in the
document \overrightarrow{d_2} are
very similar.
```

Use angle Arts Weadat edu_assist gero

- Thought experiment: take a document d and append it to itself. Call this document d'.
- "Semantically guarant Project the same content
- The Euclideanhttps://eduassistpro.github.tbg/cuments can be quite large Add WeChat edu_assist_pro
 The angle between the two ts is 0,
- The angle between the two ts is 0, corresponding to maximal similarity.

 Key idea: Rank documents according to angle with query.

From angled twee bat edu_assist_pro

- The following two notions are equivalent.
 - Rank documents in <u>decreasing</u> order of the angle between query and stogument Project Exam Help
 - Rank docum cosine(query https://eduassistpro.github.io/
- Cosine is a mondth weally at edu_assisture on for the interval [0°, 180°]

From angled tweethat edu_assist_pro

Assignment Project Exam Help

https://eduassistpro.github.io/

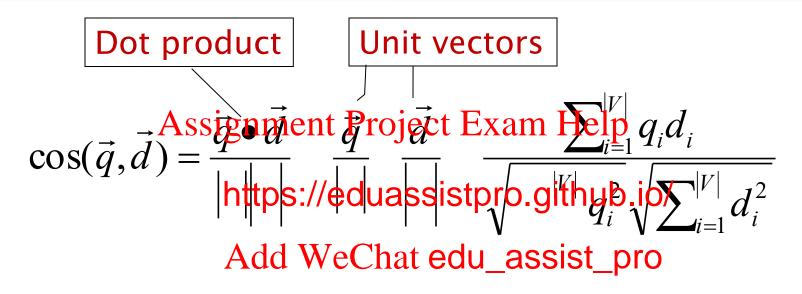
Add WeChat edu_assist_pro

But how – and why – should we be computing cosines?

Length notified Witte edu_assist_pro

- A vector can be (length-) normalized by dividing each of its components by its length for this we use the L₂ norm: Assignment Project Exam Help https://eduassistpro.github.io/
- Dividing a vector by its L, n it a unit Add WeChat edu_assist_pro (length) vector (on surface ersphere)
- Effect on the two documents d and d' (d appended to itself) from earlier slide: they have identical vectors after length-normalization.
 - Long and short documents now have comparable weights

cosine(quety, We that edu_assist_pro



 q_i is the tf-idf weight of term i in the query d_i is the tf-idf weight of term i in the document

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

Cosine for Alen gthanedu_assiste provectors

 For length-normalized vectors, cosine similarity is simply the dot product (or scalar product):
 Assignment Project Exam Help

 $\cos(\vec{q}, \frac{\text{https://eduassistprolgithub.io/}}{\text{Add WeChat edu_assist_pro}}$

for q, d length-normalized.

Cosine simila Wto hist edu_assist_pro

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Cosine similarity shat edu_assisteruments

How similar are

the novels	ionm	term	SaS	PaP John	WH
SaS: Sense and	igiiii	affection	115	1erp 58	20
Sensibility	https	s://eduas	sistpro.git	hub.io/ 7	11
PaP: Pride and	Add	anseichat	edu_assi	st_pro 0	6
Prejudice, and		wuthering	0	0	38

WH: Wuthering

Term frequencies (counts)

Heights?

Note: To simplify this example, we don't do idf weighting in this example.

3 documents lessist opro

Log frequency weighting

After length normalization

term	SaS	As <mark>sig</mark> nr	nehtPr	0 1	ecte <mark>Fi</mark> xan	n Help	PaP	WH
affection	3.06	C		J		0.789	0.832	0.524
jealous	2.00	http	s://edu	ıa	ssistpro.	githats.	O/0.555	0.465
gossip	1.30	0	1.78		g it edu_as	.335	0	0.405
wuthering	0	Agi	d Wet 1 2.58	na	it edu_as	ssist_pi	0	0.588

```
cos(SaS,PaP) \approx
```

$$0.789 \times 0.832 + 0.515 \times 0.555 + 0.335 \times 0.0 + 0.0 \times 0.0$$

 ≈ 0.94

 $cos(SaS,WH) \approx 0.79$

 $cos(PaP,WH) \approx 0.69$

Computing do Mir Chas edu_assist_pro

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

tf-idf weighting that edu_assistariants

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Columns headed 'n' are acronyms for weight schemes.

Why is the base of the log in idf immaterial?

Weighting may differ implieries vs documented WeChat edu_assist_pro

- Many search engines allow for different weightings for queries vs. documents
- SMART Notation. Genotes the combination in use in an engine, withtps://eduassistpro.ginubing the acronyms fro
- Add WeChat edu_assist_pro A very standard weighting Inc.ltc
 - Document: logarithmic tf (l as first character), no idf and cosine normalization
 - Query: logarithmic tf (l in leftmost column), idf (t in second column), cosine normalization ...

tf-idf example edu_assist_pro

Document: car insurance auto insurance

Query: best car insurance

Term	Assignment Project Exam Hebpument									Pro d	
	tf- raw	tf-wt	http	s://e	dua	ssist	pro.gi	thrends.	io/wt	n'liz e	
auto	0	0	5 00 0	1 386	Cha	ıt edi	ı_ass	ist_pr	O 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

Exercise: what is N, the number of docs?

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

Score =
$$0+0+0.27+0.53 = 0.8$$

Representation for the edu_assietrspective

- Inc.ltc
 - doc vector:
 - tf-vectorssignment Project Exam Help
 - normalized
 - query vector https://eduassistpro.github.io/
 - tf-idf-vectorAdd WeChat edu_assist_pro
 - normalized to unit length
 - score = similarity = inner product between the two vectors

Summary Add Wetchpt edu_assist rpking

- Represent the query as a weighted tf-idf vector
- Represent each document as a weighted tf-idf vector
- Compute the vector and eachttps://eduassistpro.github.io/
- Rank documents of the Respit edu_assistue proby score
- Return the top K (e.g., K = 10) to the user

Resources Afd Mediat edu_assistrero

- IIR 6.2 6.4.3
- Assignment Project Exam Help http://www. n-retrieval-tutorial/cosin https://eduassistpro.github.io/
 - Term weighting and cosine s Add WeChat edu_assist_pro