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

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

This lecture; IIR Sections 6.2-6.4.3

- 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 retrieval

- 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.

Problem with Boolean search: feast or famine

- Boolean queries often result in either too few (=0) or too many (1000s) results.
- Query 1: 'Standard User dink 550m H200,000 hits
- Query 2: "sta https://eduassistpro.gitl@bdidpund": 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 retrieval models

- Rather than a set of documents satisfying a query expression, in ranked retrieval models, the system returns an ordering of vertice (top) and telements 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 or famine: not a problem in ranked retrieval

- When a system produces a ranked result set, large result sets are not an issue
 - Indeed, the sign ment 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

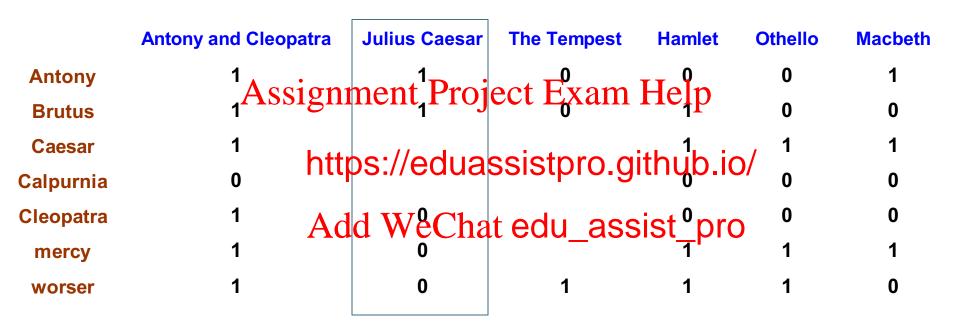
Scoring as the basis of ranked retrieval

- 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 to each doc Add WeChat edu_assist_pro
 This score measures how w ent an
- This score measures how w ent and query "match".

Query-document matching scores

- We need a way of assigning a score to a query/document pair
- Let's start with a one term query Help
- If the query tehttps://eduassistpro.githubgoument: 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.

Boolean Model: Binary termdocument incidence matrix



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

Take 1: Jaccard coefficient

- 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) = Add Awe 色雨 edu_assist_pro
- A and B don't have to be th
- Always assigns a number between 0 and 1.

Jaccard coefficient: Scoring example

- What is the query-document match score that the Jaccard coefficient computes for each of the two documents setument 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

Issues with Jaccard for scoring

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

Improved Modelling: Term-document count matrices

- Consider the number of occurrences of a term in a document:
 - Each dochregums to universe Finance: Helpumn below

https://eduassistpro.github.io/

	Antony and Cleopath	quyye@edaa	tedu_as	sistmpro	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 words model

- Vector representation doesn't consider the ordering of words in a document
- John is quicker than Mary and Wary 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 frequency tf

- 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.

Idea 2: Normalized tf is an important scoring factor

Log-frequency weighting

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 frequency

- 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 frequency, continued

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., mon, increase, Final) 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 weight

- df_t is the <u>document</u> frequency of t: the number of documents that contain t
 - df, is an Aveignmentuler of the Fram attemess of t
 - df_t ≤ N https://eduassistpro.github.io/
- - We use $\log (N/df_t)$ instead of N/df_t to "dampen" the effect of idf.

Idea 3: Normalized idf is an important scoring factor

Will turn out the base of the log is immaterial.

idf example, suppose N = 1 million

term	df_t	idf _t
calpurnia	1 · · · · · · · · · · · · · · · · · · ·	TT 1
animal ASS18	nment Project Exam	Help
sunday	ttps://eduassistpro.gi	ithub.io/
fly	Add WeChat edu_ass	ist nro
under	du weenat euu_ass	οιδι_μισ
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 idf on ranking

- 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 vs. Document frequency

 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

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

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

- Best known weightiMgeschet 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

Binary \rightarrow count \rightarrow weight matrix

	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)/ o	0
Cleopatra	2.85	0		. 0	0	0
mercy	1.51 AC	ld WeChat	t edu_ass	sist _{.1} 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 as vectors

- 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 we empt edu_assiste brearch engine
- These are very sparse vectors most entries are zero.

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 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 vector space proximity

- 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 distat edu_assist for different lengths.

Why distance is a bad idea

```
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 \overrightarrow{A} e d d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e d e e
distribution of
terms in the
document \overrightarrow{d_2} are
very similar.
```

Use angle instead of distance

- Thought experiment: take a document d and append it to itself. Call this document d'.
- "Semantically gument Project Fram Helpontent
- 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 angles to cosines

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

From angles to cosines

Assignment Project Exam Help

https://eduassistpro.github.io/

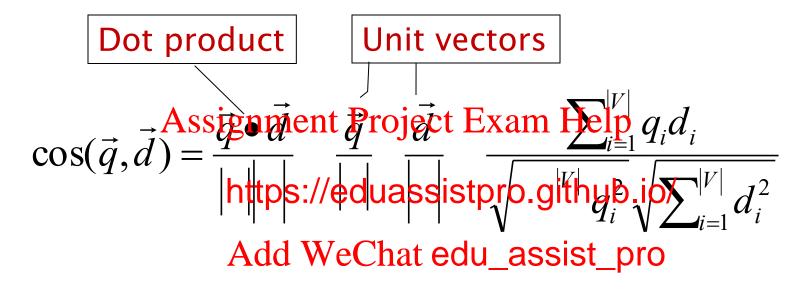
Add WeChat edu_assist_pro

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

Length normalization

- 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 L₂ 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(query,document)



 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

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

Cosine for length-normalized vectors

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

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

for q, d length-normalized.

Cosine similarity illustrated

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Cosine similarity amongst 3 documents

How similar are

the novels	i on m	term	SaS	PaP John	WH
SaS: Sense and	ıgımı	affection	115	1eip 58	20
Sensibility	https	s://eduas	sistpro.git	hub.io/ 7	11
PaP: Pride and	Add	we'Chat	edu_assi	st_pro 0	6
<i>Prejudice,</i> 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 example contd.

Log frequency weighting

After length normalization

term	SaS	Assigni	nelltPr	01	ecternan	n Help	PaP	WH
affection	3.06			J		0.789	0.832	0.524
jealous	2.00	http	os://edu	ıa	ssistpro.	githats.	<mark>O</mark> /0.555	0.465
gossip	1.30	0	1.78		g it edu_as	.335	0	0.405
wuthering	0	Aa	a wet 1	na	it edu_as	ssist_pi	0	0.588

```
cos(SaS,PaP) ≈ 0.789 × 0.832 + 0.515 × 0.555 + 0.335 × 0.0 + 0.0 × 0.0 ≈ 0.94
```

 $cos(SaS,WH) \approx 0.79$ $cos(PaP,WH) \approx 0.69$

Computing cosine scores

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

tf-idf weighting has many variants

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 in queries vs documents

- 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: Inc.ltc

Document: car insurance auto insurance

Query: best car insurance

Term	Assignment Project Exam Helpument								Pro d		
	tf- raw	tf-wt	http	s://e	edua	ssist	oro.gi	thrento.	i <mark>O</mark> /wt	n'liz e	
auto	0	0	5 000 0	1 2/86	eClo a	ıt edu	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/feature perspective

- 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 – vector space ranking

- 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 Respirated assistue property by score
- Return the top K (e.g., K = 10) to the user

Resources for today's lecture

- 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