

Introduction to Assignment Project Exam Help

Informa

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

Add WeChat edu_assist_pro

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
- Collection sta <https://eduassistpro.github.io/>
- Weighting schemes [Add WeChat 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.
Add WeChat edu_assist_pro
- 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: “~~Assignment Project Exam Help~~” → 200,000 hits
- Query 2: “~~sta~~ ~~https://eduassistpro.github.io/~~ ~~card found~~”: 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 ~~Assignment Project Exam Help~~ ordering over the ~~(top)~~ documents in the collection with <https://eduassistpro.github.io/>
- Free text queries: Rather than the language of operators and expressions, the 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 size of the result set is not an issue
 - We just show <https://eduassistpro.github.io/>
 - We don't overdo it
- 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 with <https://eduassistpro.github.io/>
- Assign a score – say in [0, 1] document
- This score measures how well the document 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
- If the query term is <https://eduassistpro.github.io/>:
score should be 0
- The more frequent the query term in the document,
the higher the score (should be)
- We will look at a number of alternatives for this.

Take 1: Jaccard coefficient

- Recall from Lecture 3: A commonly used measure of overlap of two sets A and B
- $jaccard(A,B) = \frac{|A \cap B|}{|A \cup B|}$
- $jaccard(A,A) = 1$
- $jaccard(A,B) = 0$ if $A \cap B = \emptyset$
- A and B don't have to be the same.
- 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 ~~Assignment 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 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 more informative than frequent terms. <https://eduassistpro.github.io/> Add WeChat edu_assist_pro
- We need a more sophisticated weighting for length

Recall (Lecture 1): Binary term-document incidence matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1		0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1			1	1	1
Calpurnia	0	https://eduassistpro.github.io/		0	0	0
Cleopatra	1	0	Add WeChat edu_assist_pro	0	0	0
mercy	1	0		1	1	1
worser	1	0	1	1	1	0

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

Term-document count matrices

- Consider the number of occurrences of a term in a document:
 - Each document is a count vector in \mathbb{N}^T : a column below

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

	Antony and Cleopatra	Julius Caesar	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0
Brutus	4	157	0	1	0
Caesar	232	227	0	2	1
Calpurnia	0	10	0	0	0
Cleopatra	57	0	0	0	0
mercy	2	0	3	5	5
worser	2	0	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 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 [ditional 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 ~~Assignment Project Exam Help~~ ~~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 occurrences of a term 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 weighting

- The log frequency weight of term t in d is

$$w_{t,d} = \begin{cases} 1 + \log_{10} \text{tf}_{t,d}, & \text{if } \text{tf}_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}$$

- $0 \rightarrow 0, 1 \rightarrow 1$, <https://eduassistpro.github.io/>
 - Score for a document query q over terms t in both q and d :
$$\text{score} = \sum_{t \in q \cap d} (1 + \log \text{tf}_{t,d})$$
 - 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
- Consider a term *Assignment Project Exam Help* are in the collection (e.g., <https://eduassistpro.github.io/>)
- A document containing this ~~Add WeChat~~ *edu_assist_pro* is 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., *high*, *increase*, *line*)
Assignment Project Exam Help
- A document containing <https://eduassistpro.github.io/> is more likely to be relevant than a document containing *Add WeChat edu_assist_pro*
- But it's not a sure indicator
- → 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 document frequency of t : the number of documents that contain t
 - df_t is an inverse measure of the informativeness of t
 - $df_t \leq N$
- We define the t frequency) of t by
 - 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 example, suppose $N = 1$ million

term	df_t	idf_t
calpurnia		1
animal	Assignment Project Exam Help	100
sunday	https://eduassistpro.github.io/	
fly	Add WeChat edu_assist_pro	
under		
the		1,000,000

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

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.io/>
 - idf affects the ranking of doc Add WeChat edu_assist_pro queries with at least two terms
 - 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.
- Example: <https://eduassistpro.github.io/>

Word	Collection frequency	Document frequency
Add WeChat edu_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.

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

- Best known weightings for information 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 documents for a query

$$\text{Score}(q, d) = \sum_{\substack{\text{Assignment} \\ \text{Project} \\ t \in q}} \text{tf.idf}_{t,d}$$

Exam Help

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

Add WeChat edu_assist_pro

Binary → count → weight matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	Assignment Project Exam Help	0	1	0	0
Caesar	8.59			1.51	0.25	0
Calpurnia	0	https://eduassistpro.github.io/	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	Add WeChat edu_assist_pro	0.12	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
- Documents are points or vectors in this space
- Very high-dim dimensions when you apply web search engine
 - Assignment Project Exam Help
 - <https://eduassistpro.github.io/>
 - Add WeChat `edu_assist_pro`
- 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 t <https://eduassistpro.github.io/>
- proximity = similarity of v_e
[Add WeChat](#) [edu_assist_pro](#)
- proximity \approx inverse of distance
- 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 distance? Assignment Project Exam Help
- Euclidean dist <https://eduassistpro.github.io/>
- . . . because Euclidean dist e for vectors of different lengths.

Why distance is a bad idea

The Euclidean

distance between \vec{q}

and \vec{d}_2 is large even though the

distribution of terms in the query \vec{q} and the document \vec{d}_2 are

similar. Even though the distribution of terms in the query \vec{q} and the document \vec{d}_2 are

similar.

Even though the distribution of terms in the query \vec{q} and the document \vec{d}_2 are

similar.

Use angle instead of distance

- Thought experiment: take a document d and append it to itself. Call this document d' .
- “Semantically” d and d' have the same content
- The Euclidean <https://eduassistpro.github.io/two.doc> documents can be quite large
[Add WeChat edu_assist_pro](#)
- The angle between the two vectors 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 decreasing order of the angle between query and document
 - Rank documents by cosine(query, document)
- Cosine is a monotonically decreasing function for the interval $[0^\circ, 180^\circ]$

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 L_2 norm:
$$\|\vec{v}\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$$

<https://eduassistpro.github.io/>
- Dividing a vector by its L_2 norm makes it a unit (length) vector (on surface of sphere)
- 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)

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

Assignment Project Exam Help
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Dot product Unit vectors

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

$\cos(\vec{q}, \vec{d}_i)$ <https://eduassistpro.github.io/>
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

SaS: *Sense and*

Sensibility

PaP: *Pride and
Prejudice, and*

WH: *Wuthering
Heights?*

	term	SaS	PaP	WH
	affection	115	58	20
	https://eduassistpro.github.io/	7		11
	gossip	0		6
	wuthering	0	0	38

Term frequencies (counts)

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	PaP	WH	term	SaS	PaP	WH
affection	3.06				0.789	0.832	0.524
jealous	2.00				0.515	0.555	0.465
gossip	1.30	0	1.78	g	0.335	0	0.405
wuthering	0	0	2.58	Add WeChat edu_assist_pro	0	0	0.588

$$\cos(\text{SaS}, \text{PaP}) \approx$$

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

$$\approx 0.94$$

$$\cos(\text{SaS}, \text{WH}) \approx 0.79$$

$$\cos(\text{PaP}, \text{WH}) \approx 0.69$$

Why do we have $\cos(\text{SaS}, \text{PaP}) > \cos(\text{SaS}, \text{WH})$?

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: denotes the combination in use in an engine, with <https://eduassistpro.github.io/> using the acronyms from [Add WeChat edu_assist_pro Inc.Itc](#)
- A very standard weighting
 - 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 ...

A bad idea?

tf-idf example: Inc.ltc

Document: *car insurance auto insurance*

Query: *best car insurance*

Term	Query						Document			Pro d
	Assignment	Project	Exam	Help	Document	Score	tf-wt	wt	n'liz e	
	tf-raw	tf-wt								
auto	0	0	50000	2.3	0					0
best	1	1	50000	1.3	1.3	0.34	0	0	0	0
car	1	1	10000	2.0	2.0	0.52	1	1	1	0.52
insurance	1	1	1000	3.0	3.0	0.78	2	1.3	1.3	0.53

Exercise: what is N , the number of docs?

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

$$\text{Score} = 0 + 0 + 0.27 + 0.53 = 0.8$$

Representation/feature perspective

- Inc.ltc
 - doc vector:
 - tf-vector **Assignment Project Exam Help**
 - normalized
 - query vector **<https://eduassistpro.github.io/>**
 - tf-idf-vector **Add 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
 - Assignment
 - Project
 - Exam
 - Help
- Compute the vector and each for the query
- Rank documents 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.miislita.com/information-retrieval-tutorial/cosin.html)
<http://www.miislita.com/information-retrieval-tutorial/cosin.html> <https://eduassistpro.github.io/>
 - Term weighting and cosine similarity tutorial for SEO folk!

Add WeChat **edu_assist_pro**