

Data Mining and Machine Learning

Lecture 3 [Assignment Project Exam Help](https://powcoder.com)

Stopping, Stemming & TF-IDF <https://powcoder.com>

Similarity [Add WeChat powcoder](https://powcoder.com)

Peter Jančovič

Objectives

- Understand definition and use of **Stop Lists**
- Understand motivation and methods of **Stemming**
- Understand how to calculate the TF-IDF Similarity between two documents

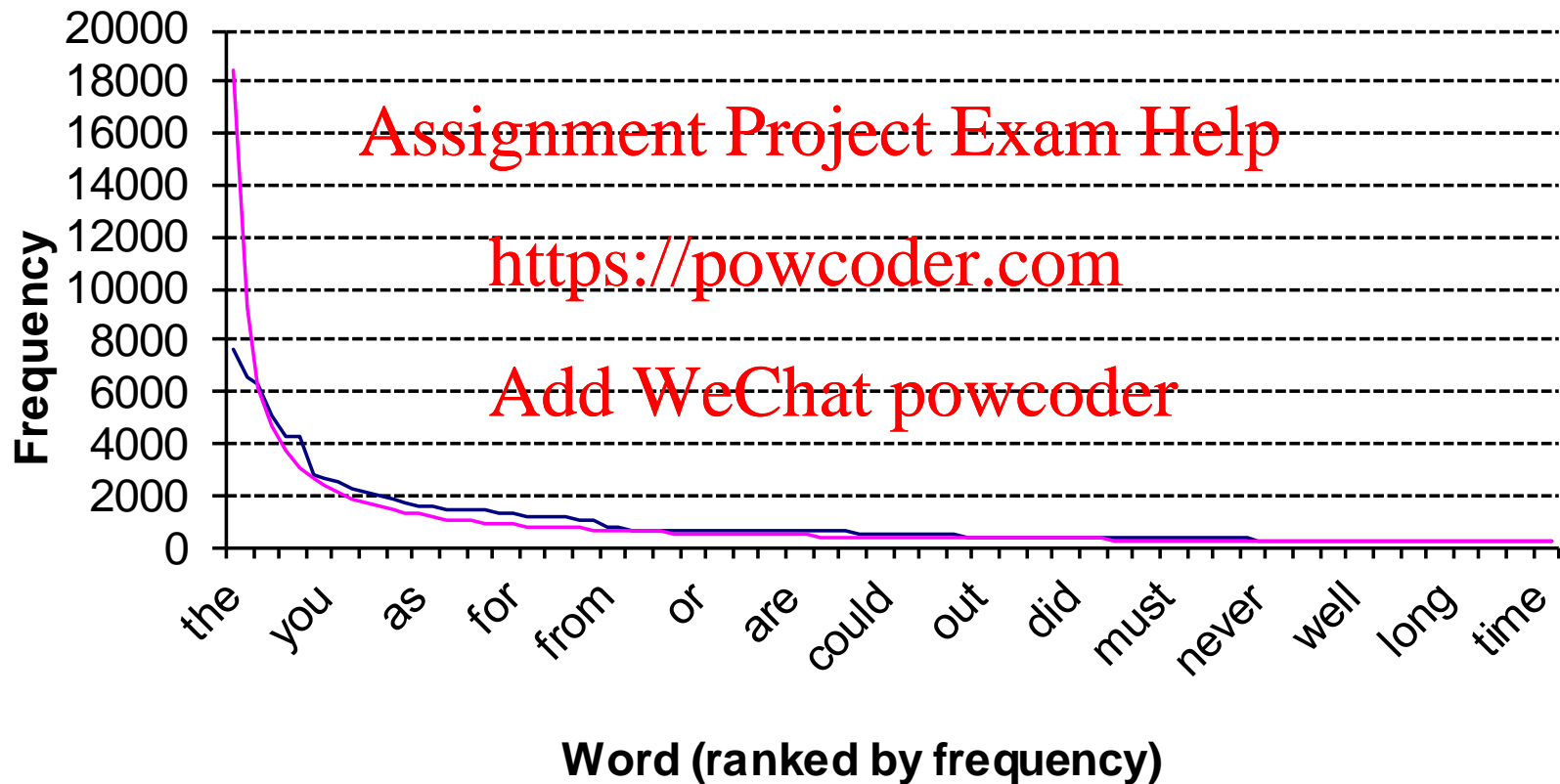
Assignment Project Exam Help

<https://powcoder.com>

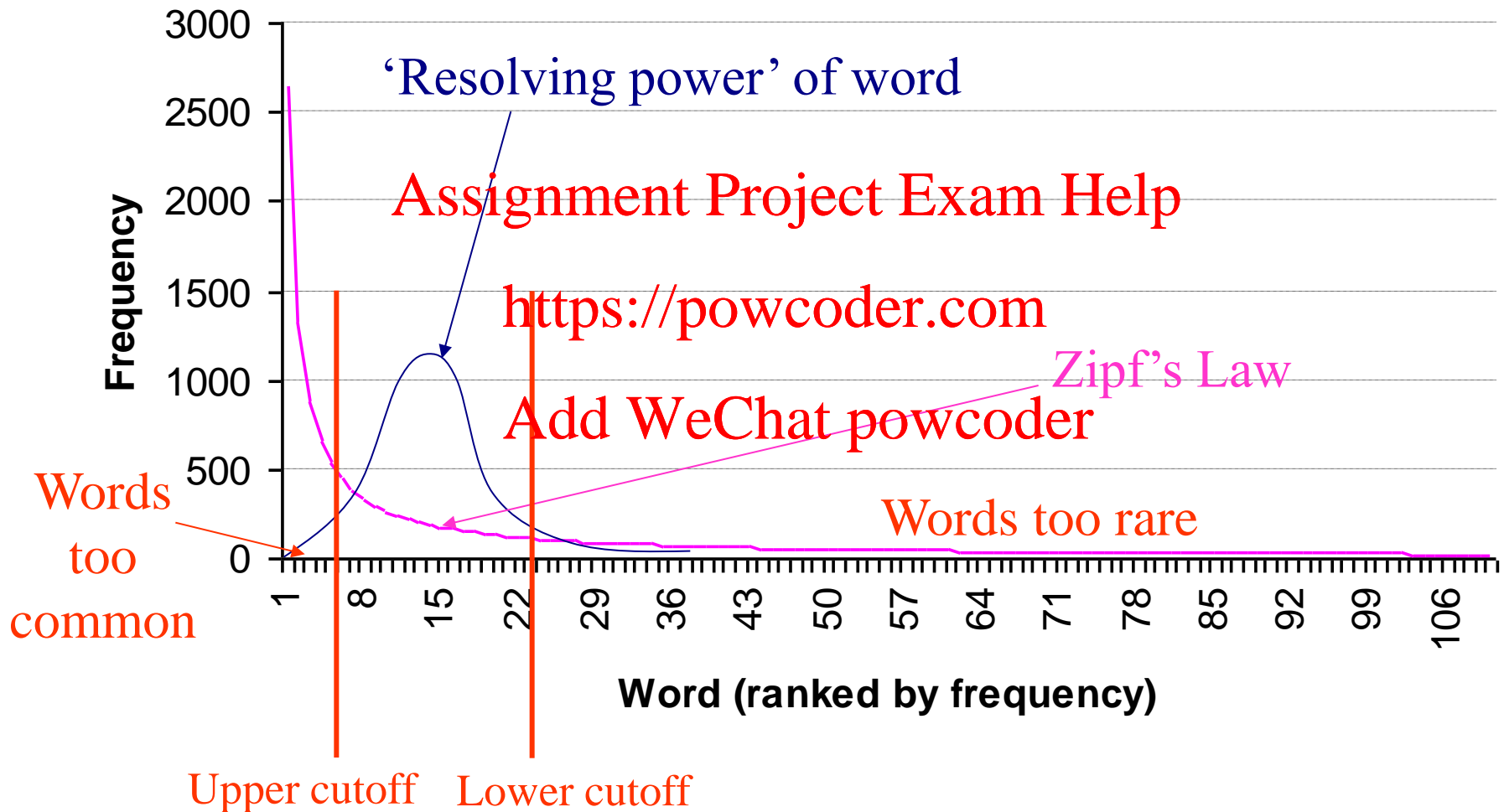
Add WeChat powcoder

Zipf's Law

Zipf's law ——— Actual statistics from “Jane Eyre” ———



‘Resolving Power’ of words



Text Pre-Processing

- Stop Word Removal: Simple techniques to remove ‘noise words’ from texts
 - Remove common ‘noise’ words which contribute no information to the IR process (e.g. “the”)
- Stemming: Remove irrelevant differences from different ‘versions’ of the same word
 - Identify different forms of the same word (e.g. “run” and “ran”) identify them with a common stem
- (Later) Exploit semantic relationships between words
 - If two words have the same meaning, treat them as the same word

Assignment Project Exam Help

<https://powcoder.com>

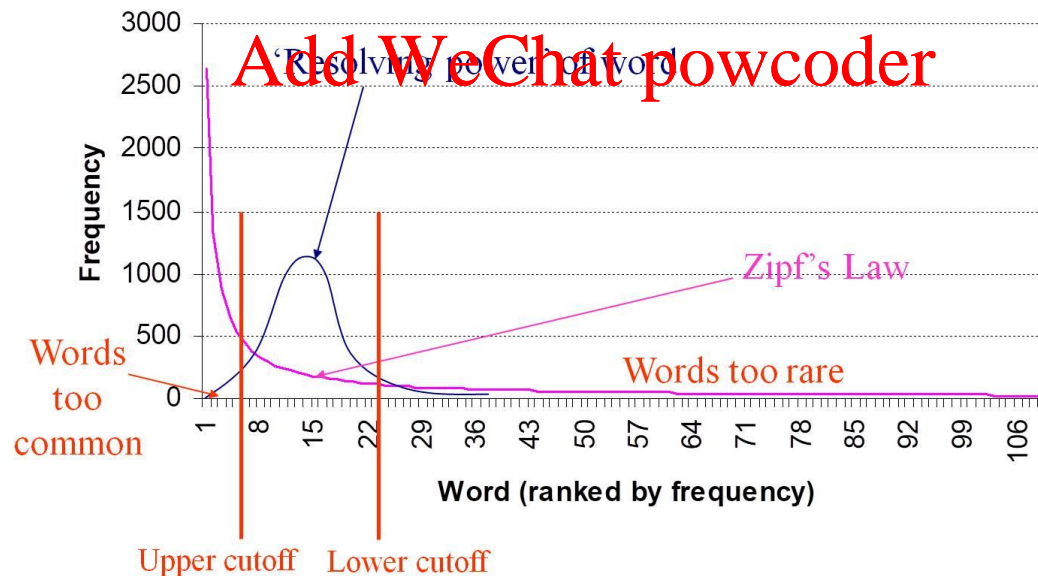
Add WeChat powcoder

Stemming (morphology)

- Basic idea: If a query and document contain different forms of the same word, then they are related
- Remove surface markings from words to reveal their basic form: <https://powcoder.com>
 - formsu → form, formingu → form
 - formedu → form, formeru → form
- “form” is the stem of forms, forming, formed, former

Stemming (morphology)

- Stemming replaces tokens (words) with equivalence classes of tokens (words)
- Equivalence classes are stems
 - Reduces the number of different words in a corpus
 - Increases the number of instances of each token



Stemming

- Of course, not all words obey simple, regular rules:

- running → run

- runs → run

- women → woman

- leaves → leaf

- ferries → ferry

- alumnus → alumni

- datum → data

- crisis → crises

[Belew, chapter 2]

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Stemming

- Linguists distinguish between different types of morphology:
 - Minor changes, such as plurals, tense
 - Major changes, e.g. *incentive* → *incentivize*, which change the grammatical category of a word
- Common solution is to identify sub-pattern of letters within words and devise rules for dealing with these patterns

Stemming

- Example rules [Belew, p 45]

- $(.*)SSES \rightarrow /1SS$

- Any string ending SSES is stemmed by replacing SSES with SS

- E.G: “classes” \rightarrow “class”

- $(.[AEIOU].*)ED \rightarrow /1$

- Any string containing a vowel and ending in ED is stemmed by removing the ED

- E.G. “classed” \rightarrow “class”

Assignment Project Exam Help

<https://powcoder.com>

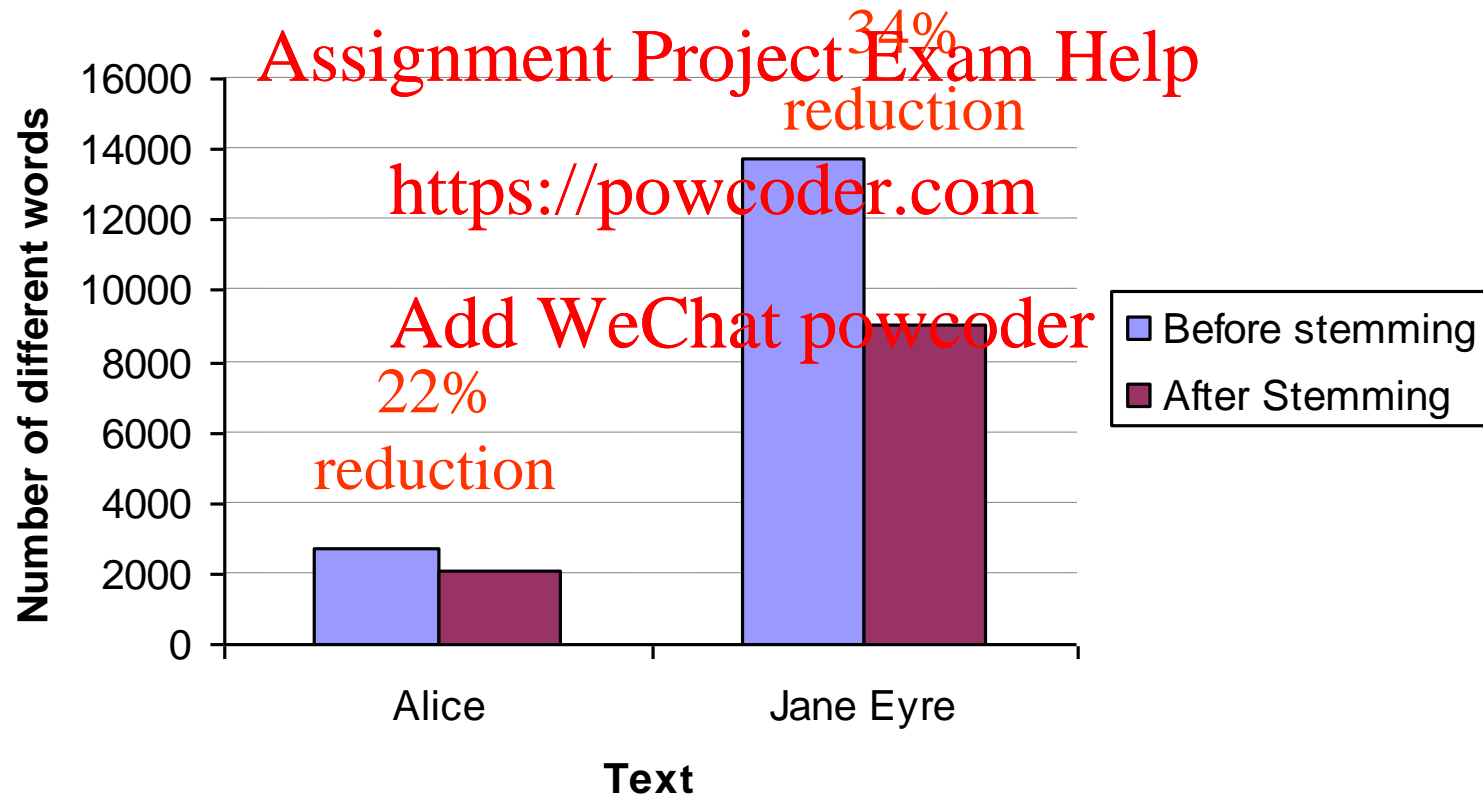
Add WeChat powcoder

Stemmers

- A stemmer is a piece of software which implements a stemming algorithm
- The Porter stemmer is a standard stemmer which is available as a free download (see Canvas)
<https://powcoder.com>
- The Porter stemmer implements a set of about 60 rules
Add WeChat powcoder
- Use of a stemmer typically reduces vocabulary size by 10% to 50%

Example

- Apply the Porter stemmer to the 'Jane Eyre' and 'Alice in Wonderland' texts



Example

- Examples of results of Porter stemmer:

- form → form
- former → former
- formed → form
- forming → form
- formal → formal
- formality → formal
- formalism → formal
- formica → formica
- formic → formic
- formant → formant
- format → format
- formation → format

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Example: First paragraph from 'Alice in Wonderland'

Before

After

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

alic wa begin to get veri tire of sit by her sister on the bank, and of have noth to do: onc or twice she had peep into the book her sister wa read, but it had no pictur or convers in it, 'and what is the us of a book,' thought alic 'without pictur or convers?'

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Noise Words – “Stop words”

There was no possibility of taking a walk that day. We had been wandering, indeed, in the leafless shrubbery an hour in the morning; but since dinner (Mrs. Reed, when there was no company, dined early) the cold winter wind had brought with it clouds so sombre, and a rain so penetrating, that further outdoor exercise was now out of the question

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- Noise words
 - Vital for the grammatical structure of a text
 - Of little use in the ‘bundle of words’ approach to identifying what a text is “about”

Stop Lists

- In Information Retrieval, these words are often referred to as Stop Words
- Rather than detecting stop words using rules, stop words are simply specified to the system in a text file: the Stop List
- Stop Lists typically consist of the most common words from some large corpus
- There are lots of candidate stop lists online

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Example 1: Short Stop List (50 wds)

the	it	not	her	who
of	with	are	all	will
and	as	but	she	more
to	his	from	there	if
a	on	or	would	out
in	be	have	their	so
that	at	an	we	
is	by	they	him	
was	i	which	been	
he	this	you	has	
for	had	were	when	

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Example 2: 300 Word Stop List

the	on	one	more	held	whose
of	be	you	no	keep	special
and	at	were	if	sure	heard
to	by	her	out	probably	major
a	i	all	so	free	problems
in	this	she	said	real	ago
that	had	there	what	seems	became
is	not	would	up	behind	federal
was	are	their	its	cannot	moment
he	but	we	about	miss	study
for	from	him	into	political	available
it	or	been	than	air	known
with	have	has	them	question	result
as	an	when	can	making	street
his	they	who	only	office	economic
	which	will	other	brought	boy

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

300 most common words from Brown Corpus

The text matters

Alice vs Brown: Most Frequent Words

the	the	as	his	this	an	know	has	thought
and	of	her	on	they	they	them	when	off
to	and	at	be	little	which	like	who	how
a	to	on	at	he	you	were	will	me
she	a	all	by	out	were	again	more	
it	in	with	i	is	her	herself	if	
of	that	had	this	one	all	when	it	
said	is	but	Had	down	she	would	so	
i	was	for	not	up	there	do		
alice	he	so	are	his	would	have		
in	for	be	but	if	their	when		
you	it	not	from	about	we	could		
was	with	very	or	then	him	or		
that	as	what	have	no	been	there		

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powder

stop.c

- C program on course Canvas page
 - Reads in a stop list file (text file, one word per line)
 - Stores stop words in char **stopList
 - Read text file one word at a time
 - Compares each word with each stop word
 - Prints out words not in stop list
- stop stopListFile textFile > opFile

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Examples

Original first paragraph

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

Stop list 50 removed

alice beginning get very tired
sitting sister bank having nothing
do once twice peeped into book
sister reading no pictures
conversations what use book
thought alice without pictures
conversation

Stop list Brown removed

alice beginning tired sitting sister
bank twice peeped book sister
reading pictures conversations
book alice pictures

conversation

Matching

- Given a query q and a document d we want to define a number:

Assignment Project Exam Help
 $Sim(q, d)$

which defines the similarity between q and d
<https://powcoder.com>

- Given the query q we will then return the documents $d_1 d_2 \dots d_N$ such that:
 - d_1 is the document for which $Sim(q, d)$ is biggest
 - d_2 has the next biggest value of $Sim(q, d)$,
 - etc

Similarity

- The similarity between q and d will depend on the number of terms which are common to q and d
- But we also need to know how useful each common term is for discriminating between different documents.
- For example,
 - It is probably not significant if q and d share “*the*”
 - But it probably is significant if they share “*magnesium*”

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

IDF weighting

- One commonly used measure of the significance of a term for discriminating between documents is the Inverse Document Frequency (IDF)

- For a token t define:

$$IDF(t) = \log \left(\frac{ND}{ND_t} \right)$$

- ND is the total number of documents in the corpus
- ND_t is the number of those documents that include t

Why is IDF weighting useful?

$$IDF(t) = \log\left(\frac{ND}{ND_t}\right)$$

- Case 1: *t* occurs equally often in all documents
 - $ND = ND_t$, <https://powcoder.com>
 - hence $IDF(t) = 0$
- Case 2: *t* occurs in just a few documents
 - $ND > ND_t$
 - hence $IDF(t) > 0$
- Note that $IDF(t)$ ignores how often term *t* occurs in a document

Effect of Document Length

- Suppose query q consists only of term t
- Suppose document d_1 also consists only of t
 - Number of shared terms is 1
 - Match is ‘perfect’
- Suppose document d_2 has 100 terms, including t
 - Number of shared terms is 1
 - But in this case co-occurrence of t appears less significant
- Intuitively the similarity measure $Sim(q, d)$ needs to include normalisation by some function of N and M

TF-IDF weight

- Let t be a term and d a document
- TF-IDF – Term Frequency – Inverse Document Frequency
- The TF-IDF weight w_{td} of term t for document d is:

$$w_{td} = f_{td} \cdot IDF(t)$$

where:

f_{td} = term frequency – the number of times t occurs in d

TF-IDF weight (continued)

$$w_{td} = f_{td} \cdot IDF(t)$$

Assignment Project Exam Help

- For w_{td} to be large:
 - f_{td} must be large, so t must occur often in d
 - $IDF(t)$ must be large, so t must only occur in relatively few documents

<https://powcoder.com>

Add WeChat powcoder

Query weights

- Now suppose t is a term and q is a query.
- If q is a long query, can treat q as a document:

Assignment Project Exam Help

$$w_{tq} = f_{tq} \cdot IDF(t)$$

<https://powecoder.com>

where f_{tq} is the (query) term frequency – the number of times the term t occurs in the query q

Add WeChat powecoder

- If q is a short query, define the TF-IDF weight as

$$w_{tq} = IDF(t)$$

TF-IDF Similarity

- Define the similarity between query q and document d as:

Assignment Project Exam Help

Sum over all
terms in both
 q and d

<https://powcoder.com>

Add WeChat powcoder

$$Sim(q, d) = \frac{\sum_{t \in q \cap d} w_{td} \cdot w_{tq}}{\|d\| \cdot \|q\|}$$

‘Length’ of
query q

‘Length’ of
document d

Document length

- Suppose d is a document
- For each term t in d we can define the TF-IDF weight w_{td}
- The length of document d is defined by:

$$Len(d) = \|d\| = \sqrt{\sum_{t \in d} w_{td}^2}$$

Comments on Document Length

- This definition of $Len(d)$ may not seem very intuitive at first
- It will become more intuitive when we study vector representations of documents and Latent Semantic Indexing (LSI)
- For now, just remember that if $\mathbf{x} = (x_1, x_2, x_3)$ is a vector in 3 dimensional space, then the length of \mathbf{x} is given by:

$$\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

Summary

- Understand definition and use of **Stop Lists**
- Understand motivation and methods of **Stemming**
- Understand how to calculate the TF-IDF Similarity between two documents

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder