# Data Mining and Machine Learning

## Lecture 3 Stopping, TF-IDF Similarity

Peter Jančovič

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

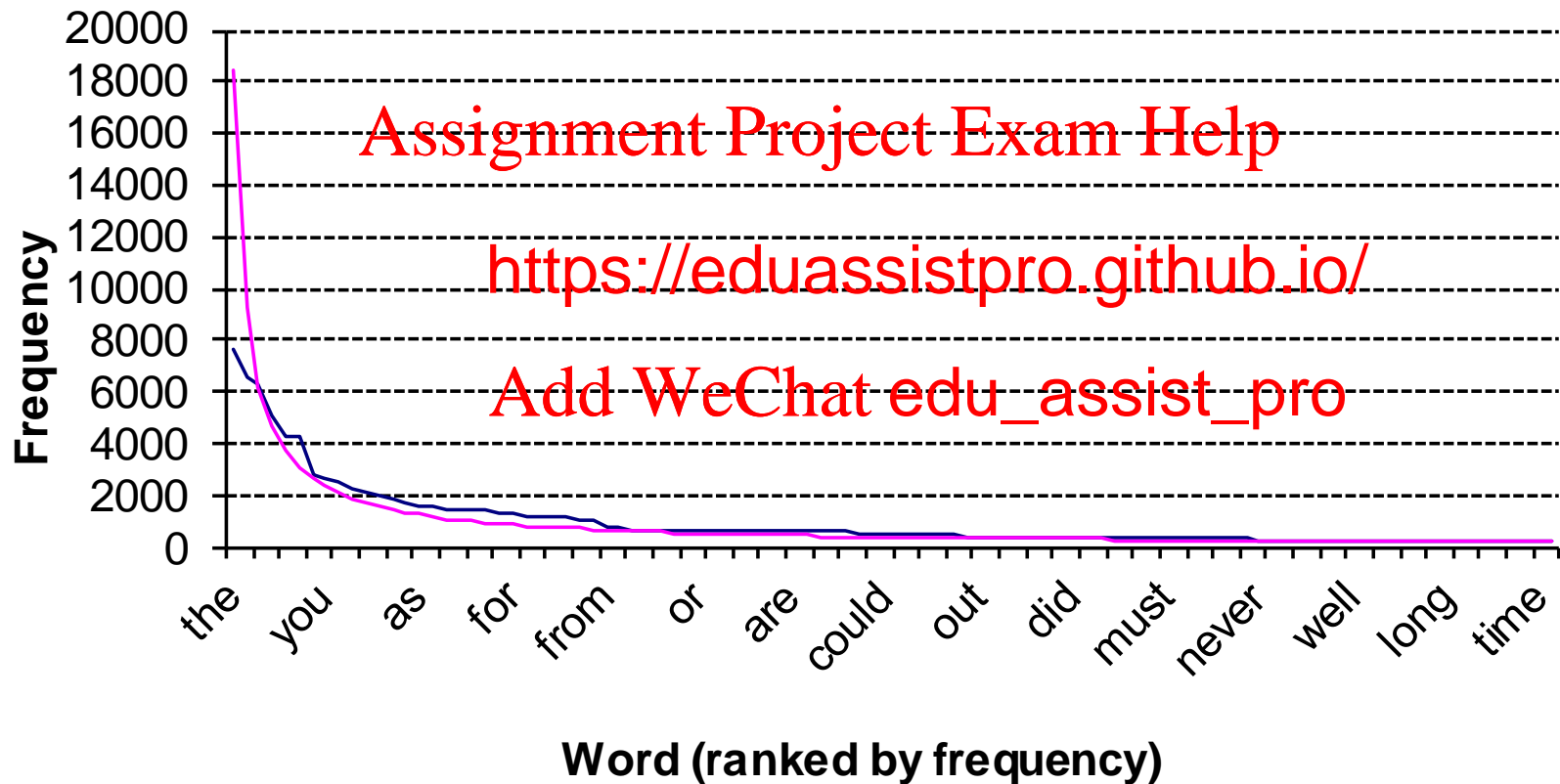# Objectives

- Understand definition and use of **Stop Lists**

- Understand motivation and methods of **Stemming**

Assignment Project Exam Help

- Understand                                -IDF Similarity

between two https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Zipf's Law

Zipf's law —————— Actual statistics from "Jane Eyre" ————



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 'Resolving Power' of words



'Resolving power' of word

Assignment Project Exam Help

https://eduassistpro.github.io/

pf's Law

Add WeChat edu_assist_pro

Words too rare

Words too common

Upper cutoff    Lower cutoff

Word (ranked by frequency)

Frequency

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Text Pre-Processing

- <u>Stop Word Removal</u>:  Simple techniques to remove 'noise words' from texts
  - Remove common 'noise' words which contribute no information to the IR process (e.g. "the")

- <u>Stemming</u>:  Re https://eduassistpro.github.io/fferent 'versions' of the same word
  - Identify different forms of t            d (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

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY<sup>OF</sup>
BIRMINGHAM

# Stemming (morphology)

- <u>Basic idea</u>: 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:

  - form<u>s</u> → fo ___
  - form<u>ed</u> → form, form<u>er</u> →

- "form" is the <u>stem</u> of forms, forming, formed, former

Data Mining and Machine Learning

UNIVERSITY<sup>OF</sup>
BIRMINGHAM

# Stemming (morphology)

- Stemming replaces tokens (words) with <u>equivalence classes</u> of tokens (words)

- Equivalence classes are <u>stems</u>

  – Reduces the number of different words in a corpus

  – Increases th ach token

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Stemming

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

  - running → run
  - runs → run
  - women →
  - leaves → le
  - ferries → ferry
  - alumnus → alumni
  - datum → data
  - crisis → crises

[Belew, chapter 2]

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Stemming

- Linguists distinguish between different types of morphology:
  - Minor changes, such as plurals, tense
  - Major chan                        ntivize, which change the                        word
- Common solution is to ide          attern of letters within words and devise <u>rules</u> for dealing with these patterns

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Stemming

- Example rules [Belew, p 45]
  - $(.*)SSES \rightarrow /1SS$
    - Any string ending SSES is stemmed by replacing SSES w
    - E.G: "cl
  - $(.[AEIOU].*)ED \rightarrow /1$
    - Any string containing a vowel and ending in ED is stemmed by removing the ED
    - E.G. "classed" $\rightarrow$ "class"
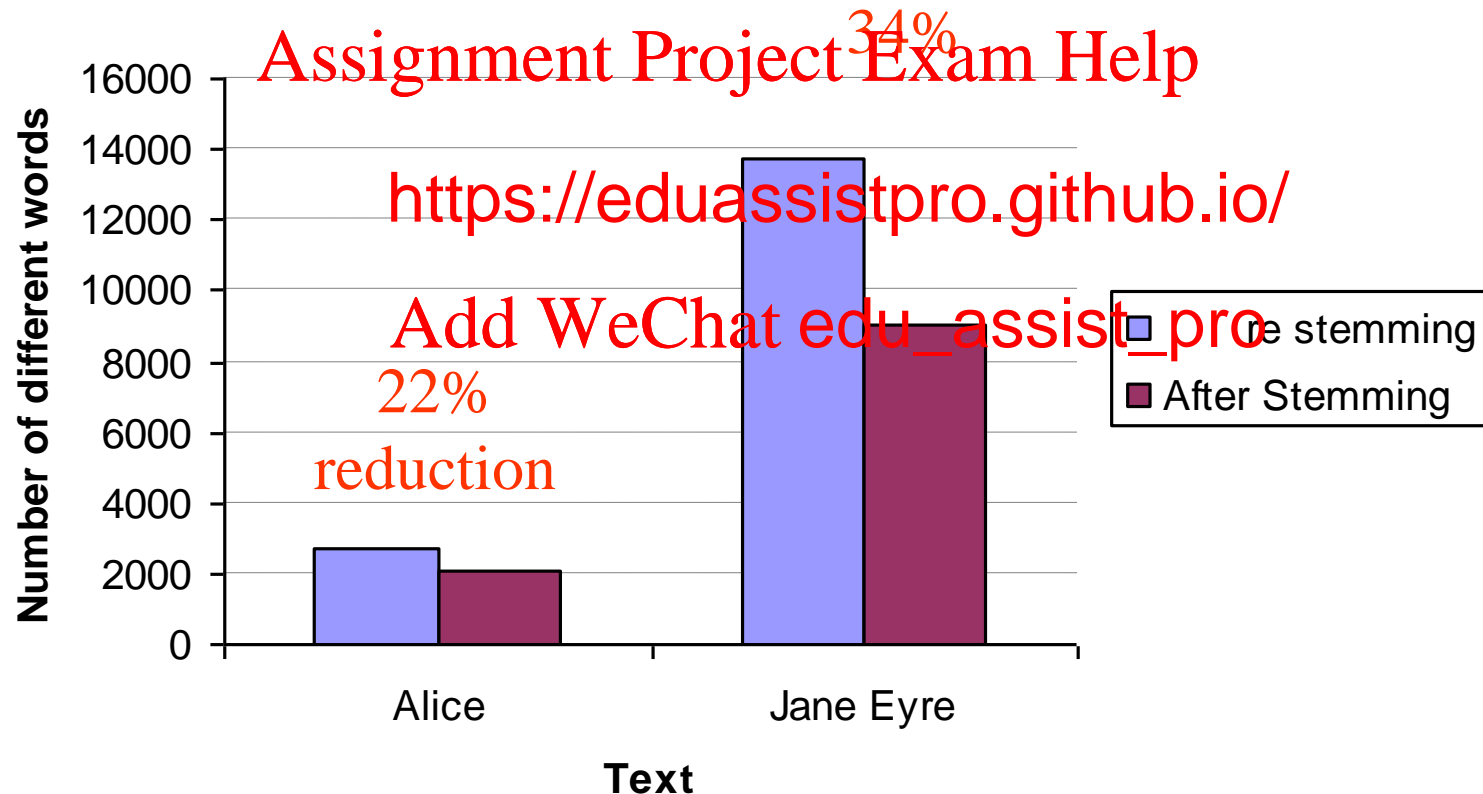
Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Stemmers

- A <u>stemmer</u> is a piece of software which implements a stemming algorithm

- The <u>Porter stemmer</u> is a standard stemmer which is available as                                   anvas)

- The Porter st                                    et of about 60 rules

- Use of a stemmer typically reduces vocabulary size by 10% to 50%

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example

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

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example

- Examples of results of Porter stemmer:
  - form → form
  - former → former
  - formed → form
  - forming → form
  - formal → form
  - formality → fo
  - formalism → formal
  - formica → formica
  - formic → formic
  - formant → formant
  - format → format
  - formation → format

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 twic
peeped into the bo
was reading, but it had
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
p into the book her
ad, but it had
rs in it, 'and what
is the us of a book,' thought alic
'without pictur or convers?'

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# 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 sombr , that further out-door exercise wa

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- 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"

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Stop Lists

- In Information Retrieval, these words are often referred to as <u>Stop Words</u>

- Rather than detecting stop words using rules, stop words are si stem in a text file: the <u>Stop</u>

- Stop Lists typically consis st common words from some large corpus

- There are lots of candidate stop lists online

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 | | | |
| in | be | | | |
| that | at | an | wat | |
| is | by | they | him | |
| was | i | which | been | |
| he | this | you | has | |
| for | had | were | when | |

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 | | seems | became |
| is | not | would | | behind | federal |
| was | are | their | | cannot | moment |
| he | but | we | about | this | 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 |

300 most common words from Brown Corpus

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 | | | | will | me |
| she | a | all | by | | | | more | |
| it | in | with | i | is | her | | | |
| of | that | had | this | one | all | | | |
| said | is | but | Had | down | she | | | |
| 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 | | |

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 fi https://eduassistpro.github.io/
  - Compares e                                    ord
  - Prints out words not in stop

- `stop stopListFile textFile > opFile`

Assignment Project Exam Help

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Examples

**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

**Original first paragraph**

Alice was beginning to get very tired of sitting by her sister on the bank, and of having noth
once or twice she had pe
the book her sister was r
it had no pictures or conversations
in it, `and what is the use of a book,'
thought Alice `without pictures or conversation?'

**wn removed**

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

conversation

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Matching

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

$$Sim(q,d)$$

which define $q$ and $d$

- Given the qu $q$ rn 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

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Similarity

- The <u>similarity</u> between *q* and *d* will depend on the number of <u>terms</u> which are common to *q* and *d*

- But we also need to know how <u>useful</u> each common term is for different documents.

- For example,

  - It is probably not significant if *q* and *d* share "*the*"

  - But it probably is significant if they share "*magnesium*"

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 $t$

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

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

UNIVERSITY OF BIRMINGHAM

# Why is IDF weighting useful?

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

- Case 1: *t* occurs equally often in all documents
  - $ND = ND_t,$
  - hence *IDF(*_____

- Case 2: *t* occurs in just a f            ents
  - $ND > ND_t$
  - hence *IDF(t) > 0*

- Note that *IDF(t)* ignores how often term *t* occurs in a document

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 'p

- Suppose doc $_2$       s, including $t$

  - Number of shared terms is
  - But in this case co-occurrence of $t$ appears less significant
- Intuitively the similarity measure $Sim(q,d)$ needs to include <u>normalisation</u> by some <u>function</u> of $N$ and $M$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# TF-IDF weight

- Let $t$ be a term and $d$ a document

- TF-IDF – Term Frequency – Inverse Document Frequency

- The TF-IDF w document $d$ is:

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

where:

$f_{td}$ = <u>term frequency</u> – the number of times $t$ occurs in $d$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# TF-IDF weight (continued)

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

- For $w_{td}$ to be
  – $f_{td}$ must be l                                           n in $d$
  – $IDF(t)$ must be large, t in relatively few documents

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Query weights

- Now suppose $t$ is a term and $q$ is a query.
- If $q$ is a <u>long</u> query, can treat $q$ as a document:

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

where $f_{tq}$ is the (query), i.e. the number of times the term $t$ occurs in the query $q$

- If $q$ is a <u>short</u> query, define the TF-IDF weight as

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# TF-IDF Similarity

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

Sum over all terms in both *q* and *d*

'Length' of query *q*

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

'Length' of document *d*

Data Mining and Machine Learning

UNIVERSITYOF BIRMINGHAM

# Document length

- Suppose $d$ is a document

- For each term $t$ in $d$ we can define the TF-IDF weight $w_{td}$

- The length o <span style="color:red">https://eduassistpro.github.io/</span> d by:

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">Add WeChat edu_assist_pro</span>

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

UNIVERSITY OF BIRMINGHAM

# 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 representatio̶n̶ ̶a̶n̶d̶ ̶L̶a̶t̶e̶n̶t̶ ̶S̶emantic Indexing (L

- For now, just remember th $x = (x_1, x_2, x_3)$ is a vector in 3 dimensional space, then the length of ***x*** is given by: $\|x\| = \sqrt{\phantom{x_1^2} \quad ^2 \quad ^2 \quad ^2}$

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Summary

- Understand definition and use of **Stop Lists**

- Understand motivation and methods of **Stemming**

  Assignment Project Exam Help

- Understand                         -IDF Similarity

  between two https://eduassistpro.github.io/

  Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM