

Intelligent Data Analysis 2020

## **Lecture 3**

# **Stopping, Stemming & TF-IDF Similarity**

Martin Russell

# 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

# Text Pre-Processing

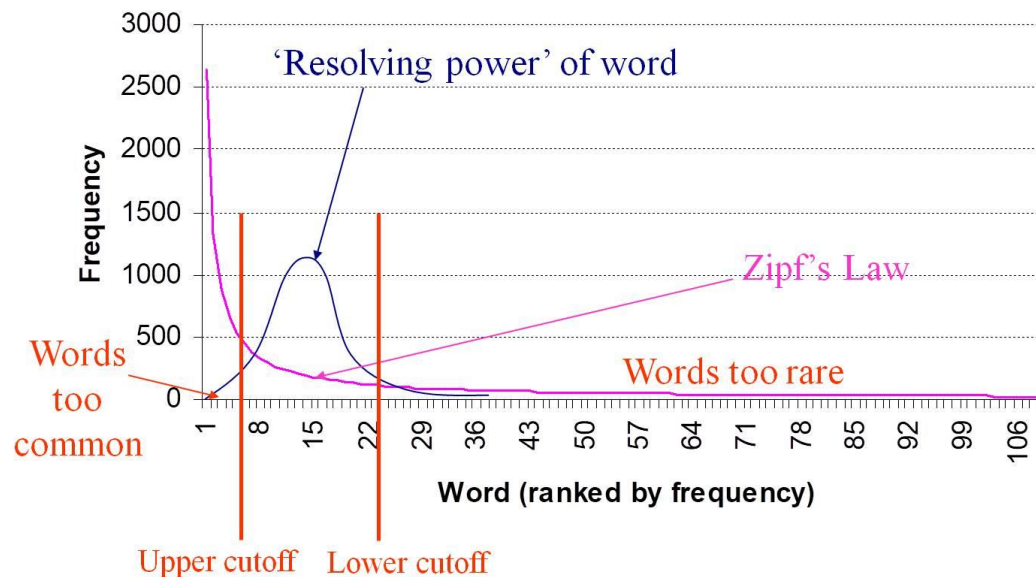
- **Stop Word Removal:** Simple techniques to remove ‘noise words’ from texts
  - Remove common ‘noise’ words which contribute no information to the IR process (e.g. “the”)
- **Stemming:** Remove irrelevant differences from different ‘versions’ of the same word
  - Identify different forms of the same 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

# 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:
  - forms → form, forming → form
  - formed → form, former → 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
- Common solution is to identify **sub-pattern of letters** within words and devise **rules** for dealing with these patterns

[Belew, chapter 2]

# 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”

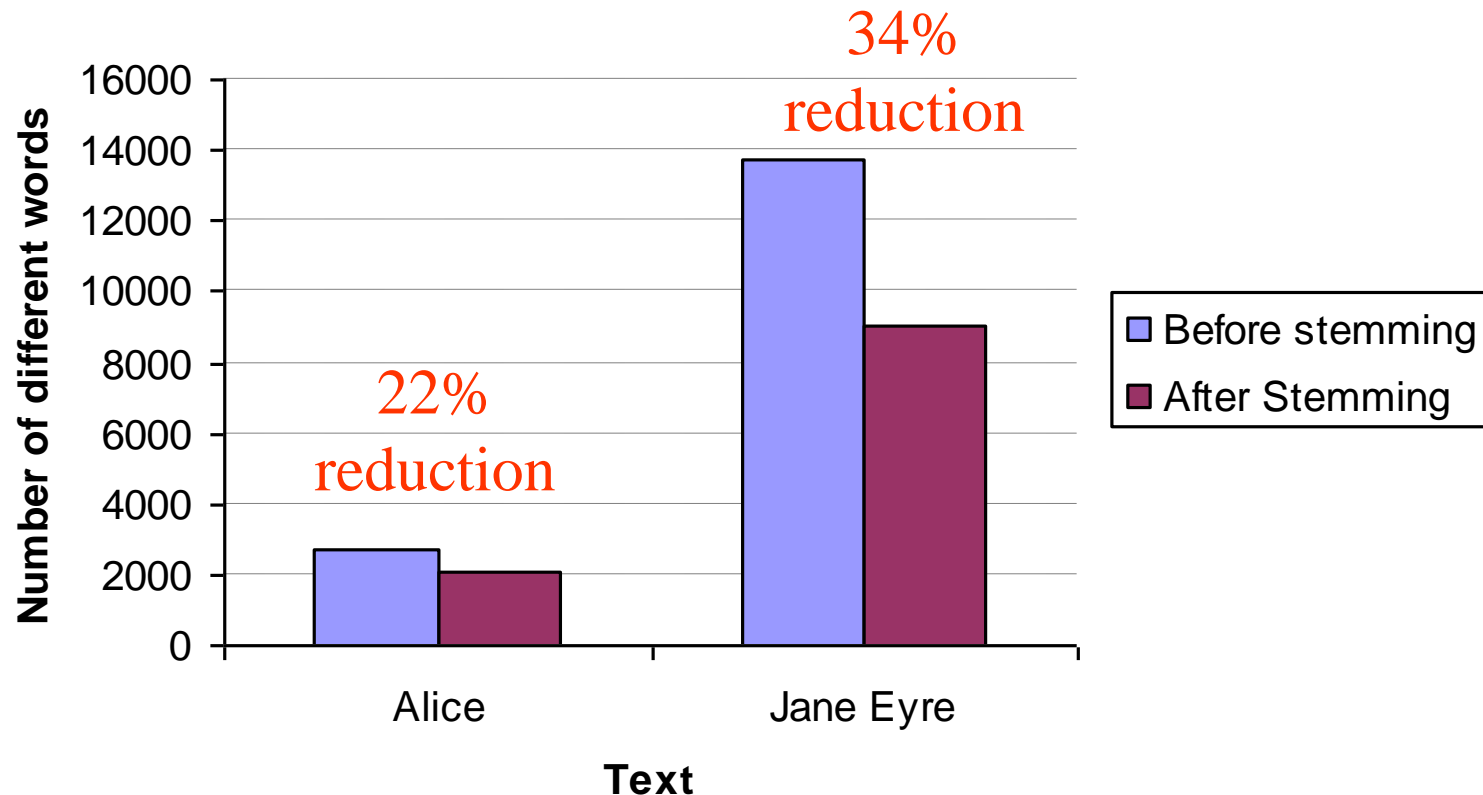
# Stemmers

- 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
- The Porter stemmer implements a set of about 60 rules
- Use of a stemmer typically reduces vocabulary size by 10% to 50%



# Example

- Apply Porter stemmer to *Jane Eyre* and *Alice in Wonderland*



# 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

# Example: First paragraph from 'Alice in Wonderland'

## Before

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?'

## After

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?'

# 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 out-door exercise was now out of the question

- Noise words in red
  - 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

# Example: 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	

# The text matters

## Alice vs Brown: Most Frequent Words

the	the	as	his	this	an	know	has	though
and	of	her	on	they	they	them	when	t
to	and	at	be	little	which	like	who	off
a	to	on	at	he	you	were	will	how
she	a	all	by	out	were	again	more	me
it	in	with	i	is	her	herself	if	
of	that	had	this	one	all	went	out	
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	wha	have	no	been	there		
		t						

# stop.c

- C program on course website
  - 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`



# 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

# Simple text retrieval

- Given two documents, d1 and d2
- Let's assume that:
  - All of the stopwords have been removed
  - Stemming has been applied
- We want to know if d1 and d2 are 'about' the same thing – we want to calculate the similarity between d1 and d2
- The simplest approach is to calculate the **TF-IDF similarity**

# Matching

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

$$Sim(q, d)$$

which defines the **similarity** between  $q$  and  $d$

- Given the query  $q$  we return 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*”

# IDF weighting

- Popular 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 total number of documents in the corpus
- $ND_t$  is number of those documents that include  $t$
- In this context,  $\log$  is  $\log_e$  – **natural logarithm**

# 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  $IDF(t) = 0$
- Case 2:  $t$  occurs in just a few documents
  - $ND > ND_t$
  - hence  $IDF(t) > 0$
- Note  $IDF(t)$  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  $t$  is less significant

# 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)

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

- 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

# 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:

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

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

- 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:

The diagram illustrates the TF-IDF similarity formula with four callout boxes explaining its components:

- Sum over all terms in both  $q$  and  $d$** : Points to the summation term  $\sum_{t \in q \cap d} w_{td} \cdot w_{tq}$ .
- 'Length' of document  $d$** : Points to the denominator term  $\|d\|$ .
- 'Length' of query  $q$** : Points to the denominator term  $\|q\|$ .

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

# 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  $x = (x_1, x_2, x_3)$  is a vector in 3 dimensional space, then the length of  $x$  is given by: 
$$\|x\| = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

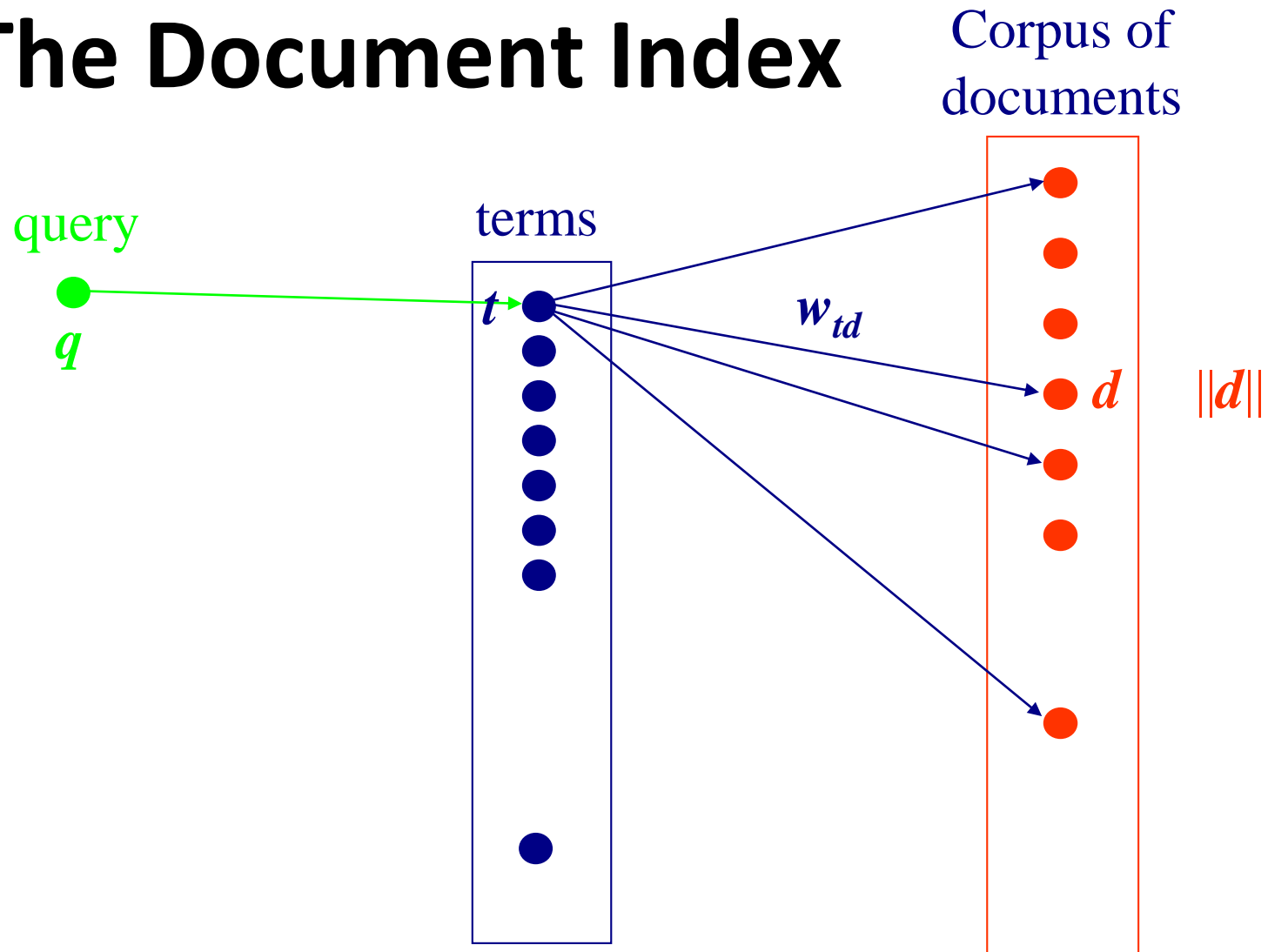
# Practical Considerations (1)

- Given a query  $q$ :
  - Calculate  $\|q\|$  and  $w_{tq}$  for each term  $t$  in  $q$
  - Not too much computation!
- For each document  $d$ 
  - $\|d\|$  can be computed in advance
  - $w_{td}$  can be computed in advance for each term  $t$  in  $d$
- Potential number of documents is **huge**
- Potential time to compute all  $Sim(q, d)$  is huge!

# Practical Considerations (2)

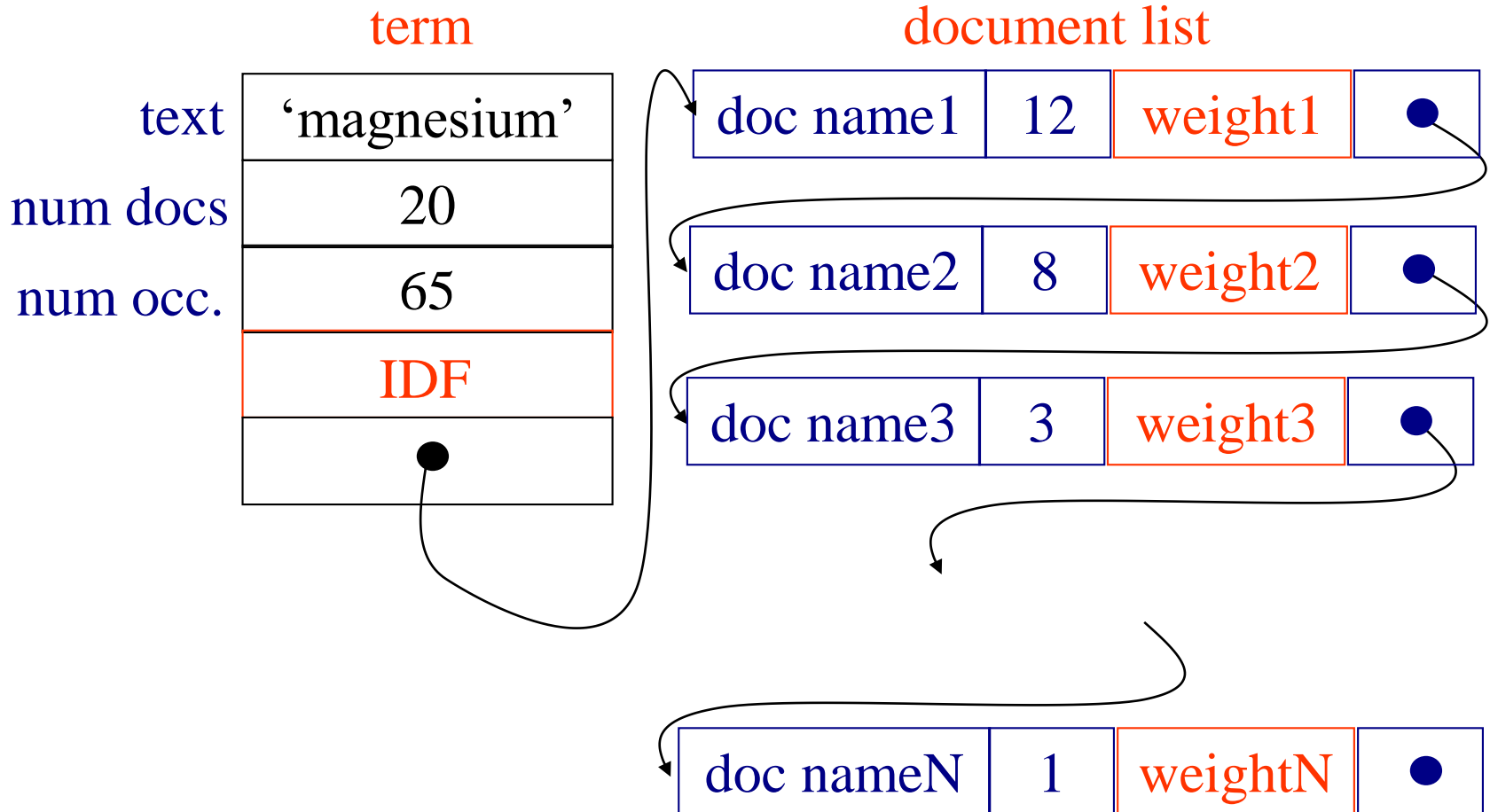
- Suppose the query  $q$  contains a term  $t$
- If  $t$  didn't already occur in the corpus it's of no use
- Need to identify all documents  $d$  which include  $t$   
(so that we can calculate  $Sim(q, d)$  for these  $d$ )
- This will take too long if the number of documents is very large (as it will be in real applications)
- To speed up this computation, we compute a data structure, called the Document Index, in advance

# The Document Index





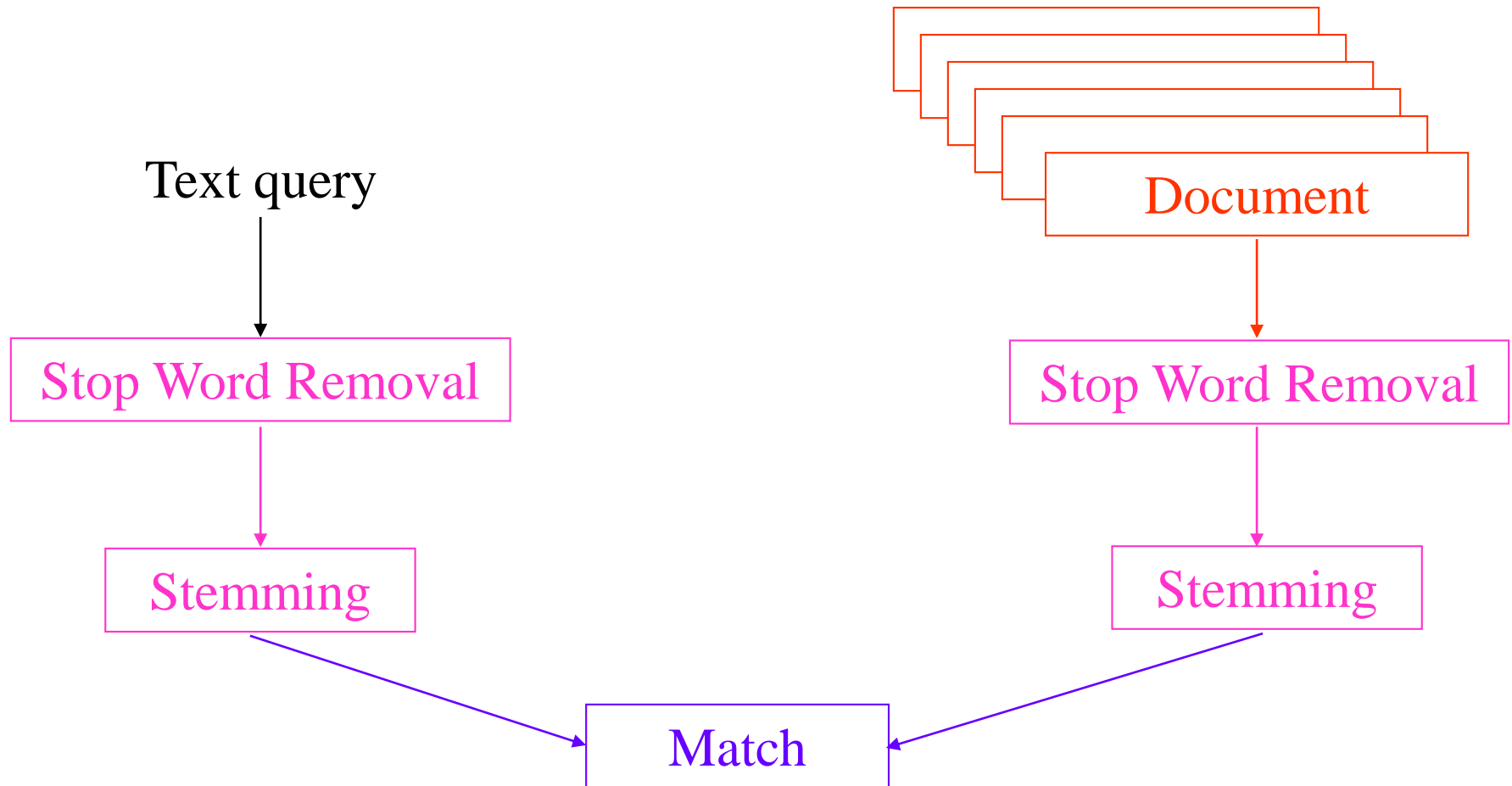
# The Document Index



# Practical considerations

- Order **terms** according to decreasing IDF
- For each term, order **documents** according to decreasing weight
- For each term in the query
  - Identify term in index
  - Increment similarity scores for documents in the list for this term
  - Option to stop when weight falls below some threshold

# Summary of the IR process



# Homework

- Download Porter Stemmer from the web
  - See URL on course web page
  - Compile and run it under your favourite OS
  - Try it out on some words and text corpora
- Download `stop.c` from the web
  - Download some stop lists
  - Compile and run `stop.c` under your favourite OS
  - Try it out on some stop lists and text corpora
- How can you make `stop.c` run on stemmed text?

# Summary

- Stemming and stop-word removal
- TF-IDF similarity
- Practical considerations – the document index