



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Introduction

Instructor:  
**Walid Magdy**

22-Sep-2021

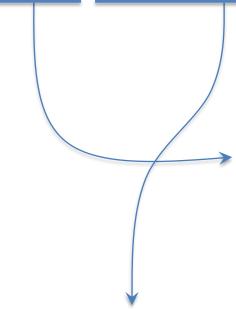
1

## Lecture Objectives

- Know about the course:
  - Topic
  - Objectives
  - Requirements
  - Format
  - Logistics
- Note:
  - No much technical content today
  - Don't assume next lectures would be the same!



## Text Technologies for Data Science



= documents, words, terms, ...  
≠ images, videos, music (*with no text*)

Information Retrieval  
Text Classification  
Text Analytics

## Search Engines Technologies



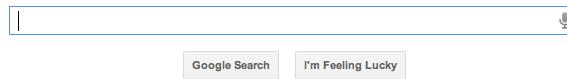
Walid Magdy, TTDS 2021/2022

3

## What is Information Retrieval (IR)?

IR is NOT just

Google



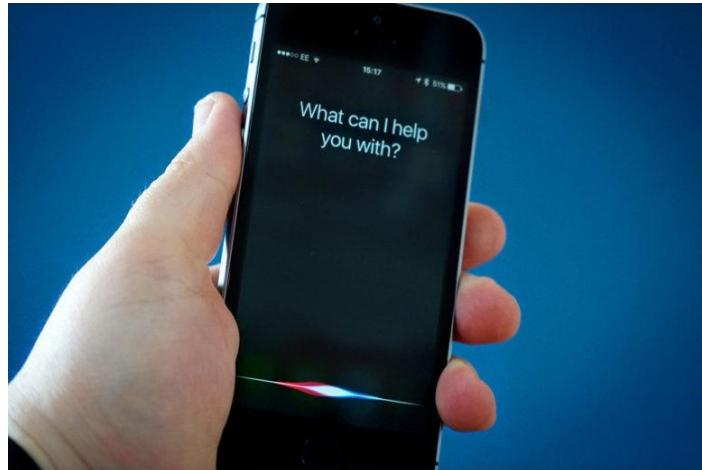
Web search

Walid Magdy, TTDS 2021/2022



4

# What is IR?



Speech - QA

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

5

# What is IR?

#Harvey

Top Latest People Photos Videos News Broadcasts

Search filters · Date

Who to follow · Refresh · View all

- harvey @missharvey Follow
- Harvey @Harveyofficial Follow
- Steve Harvey @IAmSteveHarvey Follow

Find people you know

Worldwide trends · Change

- #LOVE\_YOURSELF 1,000 Tweets
- #ReadABookDay 66,2K Tweets
- اللهم إنا نسألك العافية 41,6K Tweets
- Seif Marin 4,331 Tweets
- #WednesdayWisdom 806 Tweets

HSSawakening @HSSawakening · 48s Close the loopholes, stop tax evasion and use corporations fair tax share to care for our citizens after #Harvey and #Hurricaneirma2017

United Way Dallas @UnitedWayDallas · 1m We help those who want to do good, do great, which is why we are donating \$500K to the @SalArmyDFW for #Harvey evacuees. #letsdogreat

4 new results

Recommendation

Information Filtering

Social search

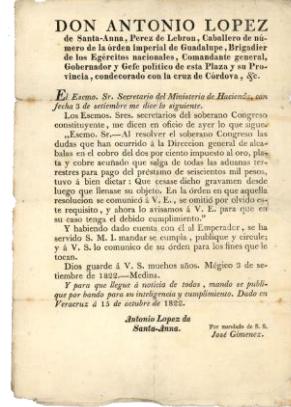
Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

6

## What is IR?



Library (book) search  
1950's

Walid Magady, TTDS 2021/2022



## What is IR?

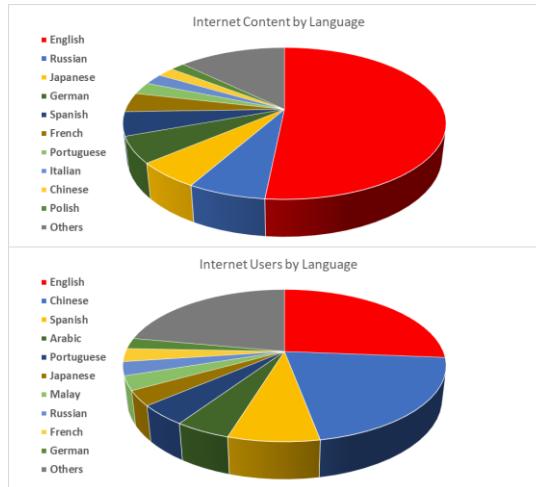


Legal search

Walid Magady, TTDS 2021/2022



## What is IR?



Cross-Language search

Walid Magdy, TTDS 2021/2022



9

## What is IR?



Content-based music search

Walid Magdy, TTDS 2021/2022



10

\*Source: Matt Lease (IR Course at U Texas)

# What is IR?



**Query suggestion / correction**

**Snippet selection / summarisation**

**Categorisation (search verticals)**

**Advertising**

Results 1 - 10 of about 132,000,000 for haiti (0.12s)

- Haiti - Wikipedia** [Sponsored Link]
- Haiti Earthquake Relief** [Sponsored Link]
- Latest News on Haiti**
- Haiti Earthquake**
- Earthquake in Haiti**
- Global Disasters Maps**
- Aid Haiti Quake Victims**
- Haiti News Summary**
- Haiti Earthquake Appeal**

THE UNIVERSITY of EDINBURGH

Walid Magdy, TTDS 2021/2022

11

# What is IR? Find?

File Edit View Window Help

Home Tools QU-RT.pdf x

terent enough from any of the pushed tweets), otherwise, the system does not consider pushing it to the user.

## 2.2 Push Notifications Scenario

The push notifications scenario simulates a recommender system that sends pop-up messages to users on their mobile phones after capturing tweets that match their interests. The task design restricts the number of pushed tweets per profile to a maximum of 10 tweets per day to avoid overwhelming the users. Having such constraint on the number of tweets to push, the system should wisely select the best candidate tweets to elect to the user in a timely fashion. We explain next how we used tweet freshness to nominate tweets to be pushed for an interest profile.

### 2.2.1 Tweets Nomination

While tracking all interest profiles simultaneously and monitoring the tweets stream, the system maintains, for each of the interest profiles, a list of *candidate tweets* that contains the tweets that were found relevant and novel so far. The RTS system periodically nominates a tweet to push to the

Alter scoring all terms, the top  $\vec{q}$  terms, are added to the top drift, the topic vector is reset to the vector) before each expansion, as  $\vec{q} = \vec{q} + \beta *$

where  $\vec{z}$  is the normalized vector of the  $k$  expansion terms, and  $\beta$  is a parameter used to restrict the influence of expansion terms on the new topic vector.

## 2.3 Periodic E-mail Digest Scenario

In this scenario, the RTS system is required to compile a daily list of a maximum of  $m$  tweets per interest profile and send it as an email digest to the user. For that, we adopted a similar but even simpler approach than the approach for push notification scenario. At the end of each day of the evaluation period, the system issues the title of the interest profile against the local tweet index that is incrementally updated over time. We experimented with three retrieval

<sup>1</sup><https://dev.twitter.com/rest/public/search>

**IR ≠ Find**

- Sequential
- Exact match

THE UNIVERSITY of EDINBURGH

Walid Magdy, TTDS 2021/2022

12

## What is IR?

- IR is finding material of an unstructured nature that satisfies an information need from within large collections
- Find → Task
- Unstructured → Nature
- Information need → Target
- Satisfies → Evaluation

Walid Magdy, TTDS 2021/2022



13

## Text classification

The screenshot shows the BBC News homepage. At the top, there is a navigation bar with links for BBC, Sign in, Home, News, Sport, Weather, iPlayer, and T. Below this is a large red banner with the word "NEWS" in white. Underneath the banner is another navigation bar with links for Home, UK, World, Business, Politics, Tech, Science, Health, Education, and Entertainment & Arts. The "UK" link is highlighted with a yellow box. Below these bars, there is a link for England, N. Ireland, Scotland, Alba, Wales, and Cymru. A news article titled "Second man held" is displayed, featuring a thumbnail image of a person. At the bottom of the page, there is a navigation bar with links for Home, Moments, Notifications, Messages, and a search bar labeled "Search Twitter". The "Today" link is highlighted with a red box. The footer of the page includes the text "Walid Magdy, TTDS 2021/2022" and the University of Edinburgh logo.

14

## Text classification

1–21 of 21 < >

Primary	Social Google+	Promotions Google Offers, Zagat	Updates Google Play
<input type="checkbox"/> James, me (2)	Hiking   Hiking trip on Saturday - Yay - so glad you can join. We should leave from I	3:14 pm	
<input type="checkbox"/> Hannah Cho	Thank you - Keri - so good that you and Steve were able to come over. Thank you :	3:05 pm	
<input type="checkbox"/> Ian Dudson	School   Incoming school conference dates... Hello everyone. A few people have	10:00 am	

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY of EDINBURGH

15

## Text classification



US008881191B2

(12) **United States Patent**  
Magdy et al.

(10) **Patent No.:** US 8,881,191 B2  
(45) **Date of Patent:** Nov. 4, 2014

(54) **PERSONALIZED EVENT NOTIFICATION  
USING REAL-TIME VIDEO ANALYSIS**

(51) **Int. Cl.**  
**H04H 60/65** (2008.01)  
**H04H 60/48** (2008.01)  
**G06F 17/30** (2006.01)

(75) Inventors: **Walid Magdy**, Giza (EG); **Motaz El-Saban**, Giza (EG)

(52) **U.S. Cl.**  
CPC ..... **H04H 60/48** (2013.01); **H04H 60/65** (2013.01); **G06F 17/30787** (2013.01); **G06F 17/30831** (2013.01)  
USPC ..... **725/32**; 725/43; 725/52; 382/181;  
348/460

(73) Assignee: **Microsoft Corporation**, Redmond, WA (US)

(\*) Notice: Subject to any disclaimer. the term of this

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY of EDINBURGH

16

## What is text classification?

- **Text classification** is the process of classifying documents into predefined categories based on their content.
  
- Input: Text (document, article, sentence)
- Task: Classify into one/multiple categories
- Categories:
  - Binary: relevant/irrelevant, spam .. etc.
  - Few: sports/politics/comedy/technology
  - Hierarchical: patents

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

17

## In this course, we will learn to

- How to build a search engine
  - which search results to rank at the top
  - how to do it fast and on a massive scale
- How to evaluate a search algorithm
  - is system A really better than system B
- How to work with text
  - two tweets talk about the same topic?
  - handle misspellings, morphology, synonyms
- How to classify text
  - into categories (sports, news, comedy, ...)
  - features to use
  - evaluate classification quality
- Apply text analytics
  - Find what makes a set of document different from others

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

18

## How this course is different from others?

- ANLP, FNLP
  - Some text processing
  - Text laws
  - No NLP (word/phrase level vs document level)
- ML practical
  - Text classification
  - No ML (using off-the-shelf ML tool)
- It does not overlap with others on:
  - Search engines
  - IR methods/models
  - IR evaluation
  - Text analysis
  - Processing large amount of textual data

Walid Magdy, TTDS 2021/2022



19

## Some terms you will learn about

- Inverted index
- Vector space model
- Retrieval models: TFIDF, BM25, LM
- Page rank
- Learning to rank (L2R)
- MAP, MRR, nDCG
- Mutual information, information gain, Chi-square
- binary/multiclass classification, ranking, regression

Walid Magdy, TTDS 2021/2022



20

## This Course is Highly Practical

- 70% of the mark is on practical work
- You will implement 50+% of what you learn
- By W5, you should have developed a basic working Search Engine from scratch
- Practical Lab every week
- Two coursework, mostly coding
- A course group project to develop a full system

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

21

## Pre-requests (1/3)

- Maths requirements:
  - Linear algebra: vectors/matrices (addition, multiplication, inverse, projections ... etc).
  - Probability theory: Discrete and continuous univariate random variables. Bayes rule. Expectation, variance. Univariate Gaussian distribution.
  - Calculus: Functions of several variables. Partial differentiation. Multivariate maxima and minima.
  - Special functions: Log, Exp, Ln.

$$\text{BM25}(D, Q) = \sum_{i=1}^n \log \frac{N - n(q_i) + 0.5}{n(q_i) + 0.5} \cdot \left[ \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdl}}\right)} + \delta \right]$$

Walid Magdy, TTDS 2021/2022



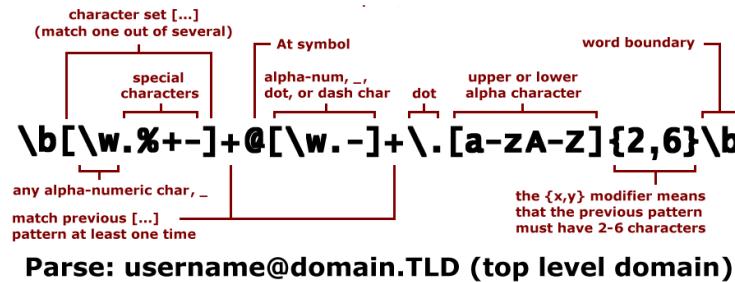
THE UNIVERSITY  
of EDINBURGH

22

## Pre-requests (2/3)



- Programming requirements:
  - Python
  - Knowledge in regular expressions
  - Shell commands (cat, sort, grep, uniq, sed, ...)
  - Data structures and software engineering for course project.



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

23

## Pre-requests (3/3)

- Team-work requirement:
  - Final course project would be in groups of 5-6 students.
  - Working in a team for the project is a requirement.
  - No exceptions will be allowed!



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

24

## Skills to be gained !!!

- Working with large text collections
- Few shell commands
- Some Python programming
- Software engineering skills
- Build text classifier in few mins
- TEAM WORK
  - Project management
  - Time management
  - Task assignment + system integration

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

25

## Course Structure

- 20 Lectures:
  - 2 lectures → Introduction (today)
  - 14 lectures → IR (50% practical lectures)
  - 4 lectures → Text Analytics/Classification
- 8-10 Labs:
  - Practice what you learn
- No Tutorials
- Some self-reading
- Lots of system implementation
- Few online videos

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

26

## Course Instructors



**Walid Magdy**

Reader

(14 lectures)

**Bjorn Ross**

Lecturer

(5 lectures)

+ 1 guest lecture

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

27

## Lecture Format

- 2 Lectures at a time
- Questions are allowed any time. Feel free to interrupt
- 5-10 mins break after L1
  - Feel free to go out and come back
  - Discuss 1<sup>st</sup> lecture with friends
  - Questions on L1 are allowed before starting L2
  - Mind teaser math problem (for fun)
- Some lectures are interactive. Please participate
- Some lectures will include demos (running code)

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

28

## Labs

- Online!
- How it will work?
  - Relevant lab will be announced with each lecture on Wednesday
  - You should implement lab directly after lecture
  - Any issues → ask on Piazza (tag question by lab number)
  - Produced output → Share on Piazza (publicly)
  - Demonstrators → answer questions + validate your output
  - DO NOT ask a question before checking if it was asked before
  - Tuesdays → Optional Teams meetings for those still require support
- Live lab times: Tuesday, 9am, 11am, 6pm
- Demonstrators:  
Zheng Zhao, Ibrahim Abu Farha, and Youssef Al Hariri

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

29

## Lab Zero (Lab 0)

- Please check Lab 0 before next week lectures
- Lab 0 is designed for one purpose:  
Help you decide to take TTDS or not
- Lab content:  
Read a text file word by word, lower-case letters, print
- If you Lab 0 challenging then:
  - Probably, TTDS would be very challenging to you
  - You will need much extra effort to implement labs and CW
  - Think wisely before you decide to take the course

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

30

## Assessments

- Coursework 1: **10%**  
The same as labs 1-3 → Build your first search engine
- Coursework 2: **20%**  
IR Evaluation, Text classification/analytics
- Group project: **40%**  
A full running search engine supported by text technologies
- Final Exam: **30%**

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

31

## Group Project

- The largest weight: 40% of the total mark
- Teamwork → Group 5-6 (you select your group)
- Design a full end-to-end search engine that searches a large collection of documents with many functionalities.
- $\text{Mark} = \text{Mark}_{\text{project}} \times \text{weight}_{\text{individual}}$ 
  - $\text{Mark}_{\text{project}}$  → the same for all team members
    - How complete/effective/fast/nice is your search engine?
  - $\text{weight}_{\text{individual}}$  → weight for individual contribution.
    - ranges from 0 to 1. It should be 1.0 by default but can be different for each member according to their contribution.
- Project prize → a prize will be awarded to best project

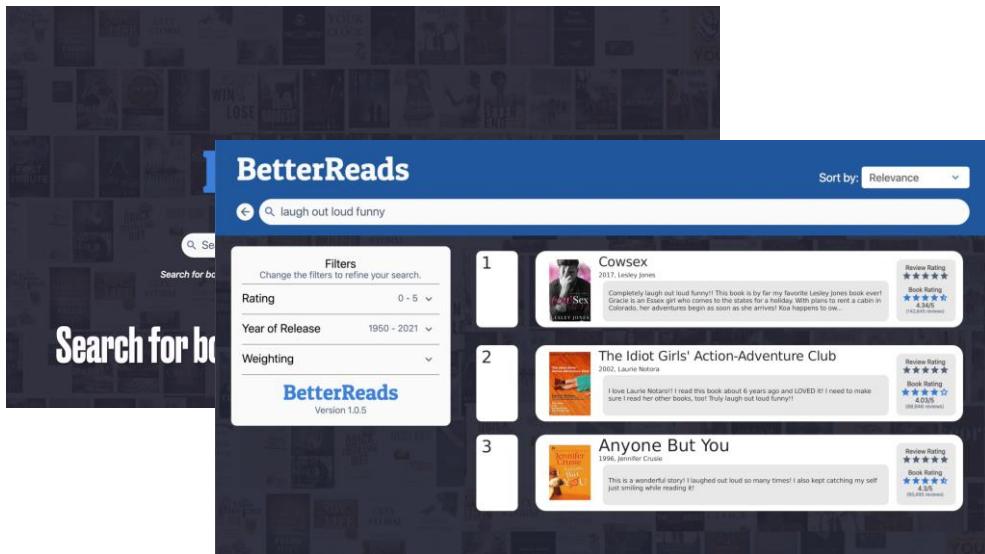
Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

32

## Example: BetterReads



Walid Magdy, TTDS 2021/2022



33

## Example: BetterReads

- 11.5M Book reviews from Good reads
- Average query time: 1.3 secs
- New reviews are crawled and indexed automatically every month
- Ranking: Relevance + Sentiment
- Engine hosted on Google cloud compute
- *Note: we will provide credit to Google cloud to host your engine*

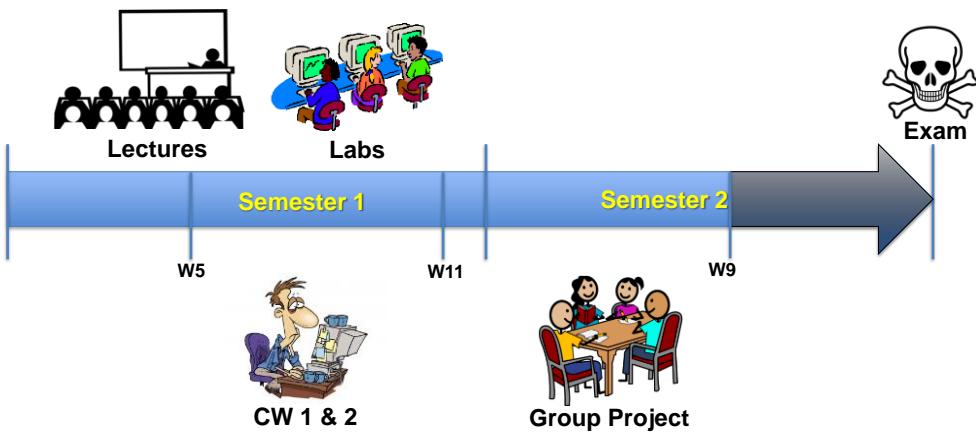
Walid Magdy, TTDS 2021/2022



34

## Timeline

- 2 Semesters (or one?)



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

35

## Logistics (1/2)

- Lectures:
  - Live on 2 Wednesdays, 12.00-14.30 (some exception might occur)
  - Recording will be available
  - Handouts to be posted on the day of the lecture
- Course webpage:
  - Link: <http://www.inf.ed.ac.uk/teaching/courses/tts/>
  - Handouts, Labs, CW details, link to recordings
- Learn:
  - Lecture recordings
  - Deadlines
- Note: all course materials are made public including recordings. Feel free to share with anyone interested

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

36

## Logistics (2/2)

- Pizza:
  - All communication will be there
  - Questions about lectures/labs/CW are there
  - Feel free to answer each other questions
  - Lab support will be mainly there
  - Please share your lab answers there
  - Join NOW: [link](#)
- Microsoft Teams:
  - Live lab support will be there
  - Join NOW: [link](#)

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

37

## FAQ

- How the project would be managed? What if one member does not work?
- I am not that solid in programming, should I take this course?
- Can I audit the course?
- Anything else?

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

38

## Next Lecture

- Definitions of IR main concepts  
(more introduction)

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Definitions

Instructor:  
**Walid Magdy**

22-Sep-2021

1

## Lecture Objectives

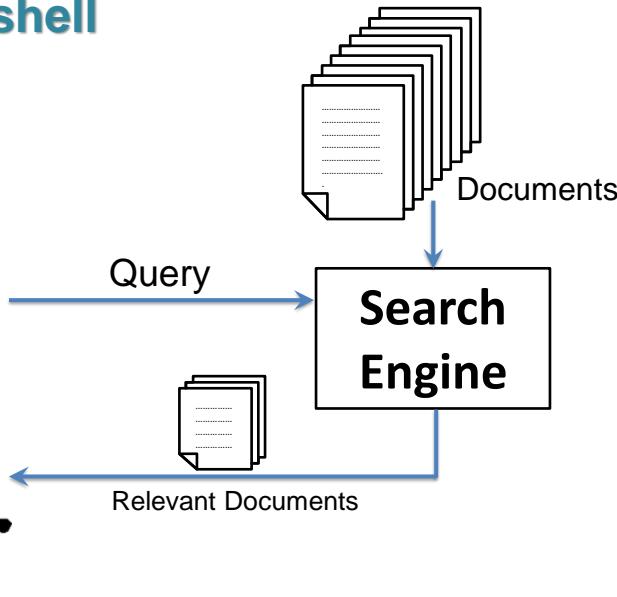
- Learn about main concepts in IR
  - Document
  - Information need
  - Query
  - Index
  - BOW



THE UNIVERSITY  
of EDINBURGH

2

## IR in a nutshell



Walid Magdy, TTDS 2021/2022

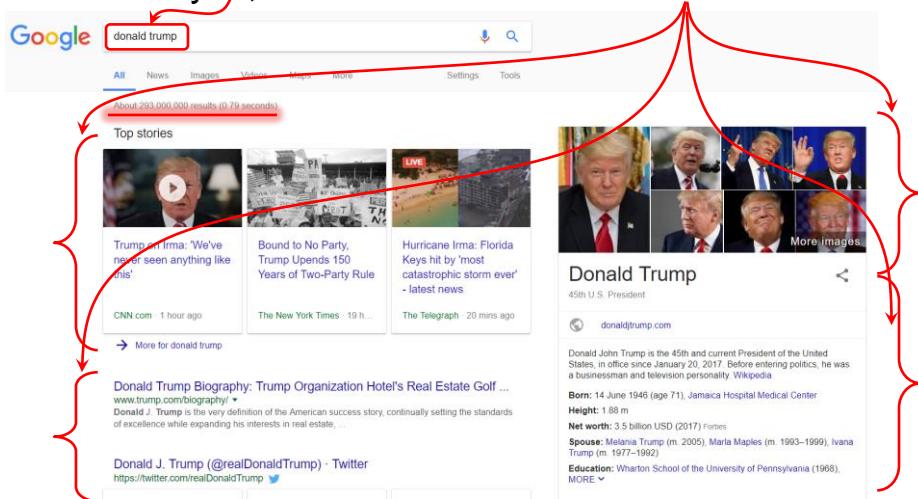


THE UNIVERSITY  
of EDINBURGH

3

## IR, basic form

- Given Query **Q**, find relevant documents **D**



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

4

## Two main Issues in IR

About 293,000,000 results (0.79 seconds)

- Effectiveness
  - need to find **relevant** documents
  - needle in a haystack
  - very different from relational DBs (SQL)
- Efficiency
  - need to find them quickly
  - vast quantities of data (100's billions pages)
  - thousands queries per second (Google, 63,000)
  - data constantly changes, need to keep up
  - compared with other NLP areas, IR is **very fast**

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

5

## IR main components

- Documents
- Queries
- Relevant documents

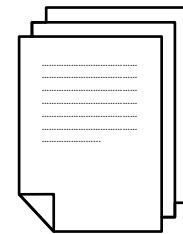
Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

6

## Documents



- The element to be retrieved
  - Unstructured nature
  - Unique ID
  - $N$  documents → Collection
- web-pages, emails, book, page, sentence, tweets
- photos, videos, musical pieces, code
- answers to questions
- product descriptions, advertisements
- may be in a different language
- may not have words at all (e.g. DNA)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Queries

- Free text to express user's information need
- Same information need can be described by multiple queries
  - Latest news on the hurricane in the US
  - North Carolina storm
  - Florence
- Same query can represent multiple information needs
  - Apple
  - Jaguar



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## Queries – different forms

- Web search → keywords, narrative ...
- Image search → keywords, sample image
- QA → question
- Music search → humming a tune
- Filtering/recommendation → user's interest/history
- Scholar search → structured (author, title ..)
  
- Advanced search  
 $\#wsyn(0.9 \#field (title, \#phrase (homer,simpson)) 0.7 \#and (\#> (pagerank,3), \#ow3 (homer,simpson)) 0.4 \#passage (homer, simpson, dan, castellaneta))$

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## Relevance

- At an abstract level, IR is about:
  - does item  $D$  **match** item  $Q$ ? ...or...
  - is item  $D$  **relevant** to item  $Q$ ?
- Relevance a tricky notion
  - will the user like it / click on it?
  - will it help the user achieve a task?  
(satisfy information need)
  - is it novel (not redundant)?
- *Relevance = what is the topic about?*
  - i.e.  $D, Q$  share similar “meaning”
  - about the same topic / subject / issue

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## What is the challenge in relevance?

- No clear semantics, contrast:
  - “*William Shakespeare*”
  - Author history’s? list of plays? a play by him?
- Inherent ambiguity of language:
  - synonymy: “North Carolina storm” = “Florence hurricane”
  - polysemy: “Apple”, “Jaguar”
- Relevance highly subjective
  - Relevance: yes/no
  - Relevance: perfect/excellent/good/fair/bad
- On the web: counter SEOs / spam

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Relevant Items are Similar

- Key idea:
  - Use similar vocabulary → similar meaning
  - Similar documents relevant to same queries
- Similarity
  - String match
  - Word overlap
  - $P(D|Q)$  → retrieval model

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

12

## IR vs. DB

	Databases	IR
What we're retrieving	Structured data. Clear semantics based on a formal model.	Mostly unstructured. Free text with some metadata.
Queries we're posing	Formally-defined (relational algebra, SQL). Unambiguous.	Free text ("natural language"), Boolean
Results we get	Exact (always "correct")	Imprecise (need to measure relevance)
Interaction with system	One-shot queries.	Interaction is important.

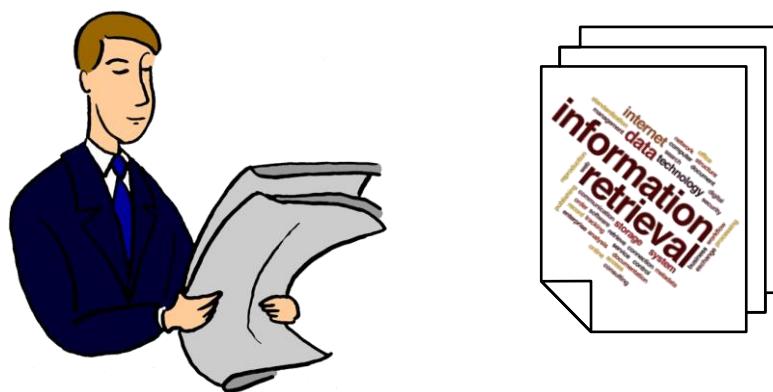
Tamer Elsayed, QU

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

13

## How IR sees documents?



Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

14

## Bag-of-words trick

- Can you guess what this is about:
  - per is salary hour €4,200 Neymar's
  - obesity French is of full cause and fat fries
- Re-ordering doesn't destroy the topic
  - individual words – "building blocks"
  - "bag" of words: a "composition" of "meanings"

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Bag-of-words trick

- Most search engines use BOW
  - treat documents, queries as bags of words
- A "bag" is a set with repetitions
  - match = "degree of overlap" between D,Q
- Retrieval models
  - statistical models (function) that use words as features
  - decide which documents most likely to be relevant
- What should be the top results for Q?
  - BOW makes these models tractable

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

16

## Bag-of-words: Criticism

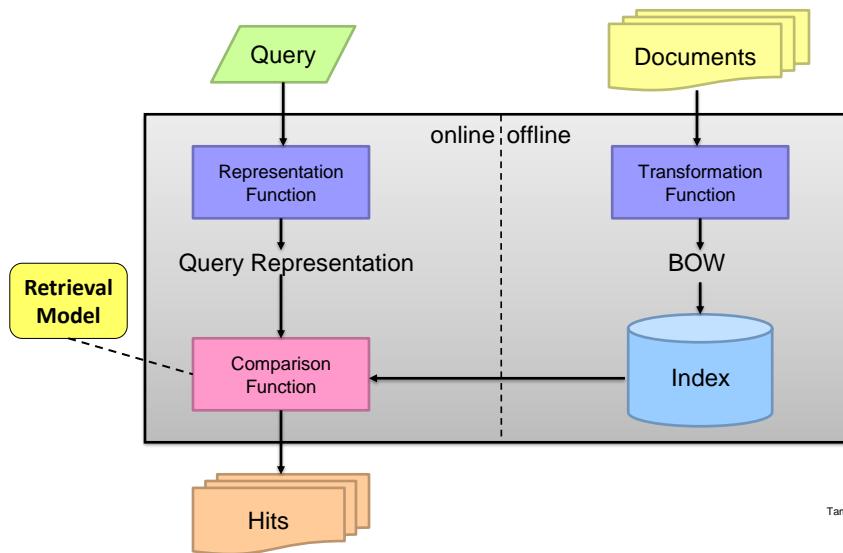
- word meaning lost without context
  - True, but BOW doesn't really discard context
- what about negations, etc.?
  - {no, climate change is real} vs. {climate change is no real}
- does not work for all languages
  - No natural "word" unit for Chinese, images, music
  - Solve by "segmentation" or "feature induction"

Walid Magdy, TTDS 2021/2022



17

## IR Black Box



Walid Magdy, TTDS 2021/2022



18

## Systems perspective on IR

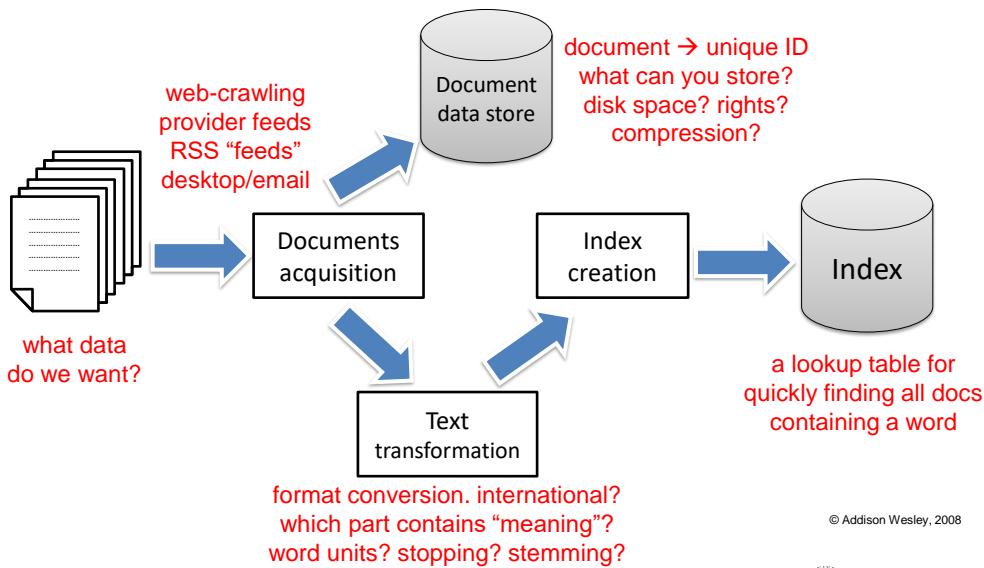
- Indexing Process: (*offline*)
  - get the data into the system
    - acquire the data from crawling, feeds, etc.
    - store the originals (if needed)
    - transform to BOW and “index”
- Search (retrieval) Process: (*online*)
  - satisfy users’ requests
    - assist user in formulating query
    - retrieve a set of results
    - help user browse / re-formulate
    - log user’s actions, adjust retrieval model

Walid Magdy, TTDS 2021/2022



19

## Indexing Process



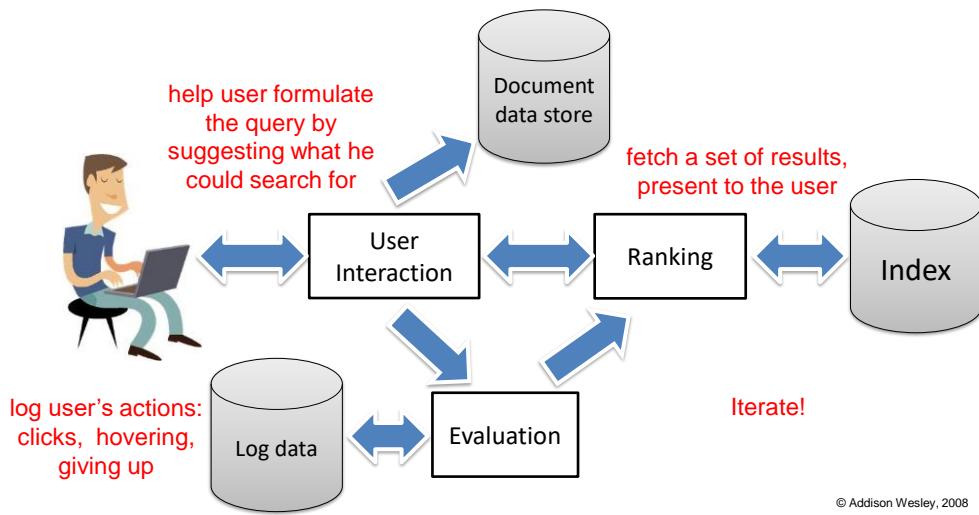
Walid Magdy, TTDS 2021/2022

© Addison Wesley, 2008



20

## Search Process



© Addison Wesley, 2008

Walid Magdy, TTDS 2021/2022



21

## Summary

- Information Retrieval (IR): core technology
  - selling point: IR is very fast, provides context
- Main issues: effectiveness and efficiency
- Documents, queries, relevance
- Bag-of-words trick
- Search system architecture:
  - indexing: get data into the system
  - searching: help users find relevant data

Walid Magdy, TTDS 2021/2022



22

## Resources

- Search Engines: Information Retrieval in Practice, chapter 1 & 2
- Lab 0:
  - You have to be confident doing it!
  - If you have trouble finishing it, think twice before committing to the course

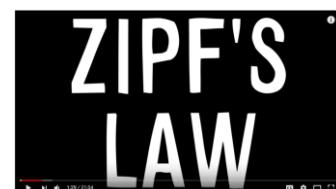
Walid Magdy, TTDS 2021/2022



23

## Questions

- Next time:
  - Laws of text (Zipf ....)
  - Vector space models
- Skill to learn by next time:
  - Read text file from disk
  - Read word by word
- Videos:
  - The Zipf Mystery, Vsauce
- Tools:
  - (Perl) regular expressions: <https://perldoc.perl.org/perlre.html>



Walid Magdy, TTDS 2021/2022



24



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Laws of Text

Instructor:  
**Walid Magdy**

29-Sep-2021

1

## Lecture Objectives

- Learn about some text laws
  - Zipf's law
  - Benford's law
  - Heap's law
  - Clumping/contagion
- This lecture is practical



THE UNIVERSITY  
of EDINBURGH

2

## You can try with me ...

- Shell commands: cat, sort, uniq, grep
- Perl (or alternative)
- Excel (or alternative)
- Download the following:
  - Bible: <http://www.gutenberg.org/cache/epub/10/pg10.txt>

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

3

## Words' nature

- Word → basic unit to represent text
- Certain characteristics are observed for the words we use!
- These characteristics are very consistent, that we can apply laws for them
- These laws apply for:
  - Different languages
  - Different domains of text

*Walid Magdy, TTDS 2021/2022*

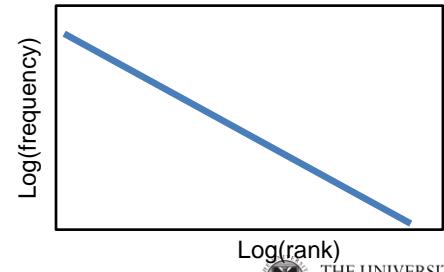


THE UNIVERSITY  
of EDINBURGH

4

## Frequency of words

- Some words are very frequent  
e.g. “the”, “of”, “to”
- Many words are less frequent  
e.g. “schizophrenia”, “baazinga”
- ~50% terms appears once
- Frequency of words has hard exponential decay



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

5

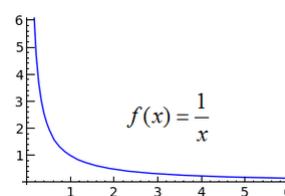
## Zipf's Law:

- For a given collection of text, ranking unique terms according to their frequency, then:

$$r \times P_r \cong \text{const}$$

- $r$ , rank of term according to frequency
- $P_r$ , probability of appearance of term

$$\bullet P_r \cong \frac{\text{const}}{r} \rightarrow f(x) \cong \frac{1}{x}$$



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

6

## Zipf's Law:

Wikipedia abstracts  
→ 3.5M En abstracts

$$r \times P_r \cong \text{const} \rightarrow \\ r \times freq_r \cong \text{const}$$

Term	Rank	Frequency	r × freq
the	1	5,134,790	5,134,790
of	2	3,102,474	6,204,948
in	3	2,607,875	7,823,625
a	4	2,492,328	9,969,312
is	5	2,181,502	10,907,510
and	6	1,962,326	11,773,956
was	7	1,159,088	8,113,616
to	8	1,088,396	8,707,168
by	9	766,656	6,899,904
an	10	566,970	5,669,700
it	11	557,492	6,132,412
for	13	493,374	5,970,456
as	14	480,277	6,413,862
on	15	471,544	6,723,878
from	16	412,785	7,073,160

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Practical

Collection	# words	File size
Bible	824,054	4.24 MB
Wiki abstracts	80,460,749	472 MB

Walid Magdy, TTDS 2021/2022

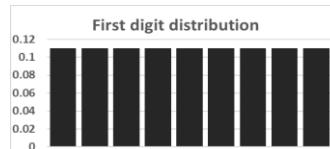


THE UNIVERSITY  
of EDINBURGH

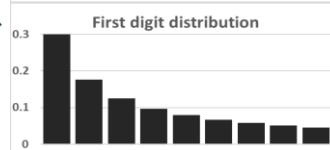
8

## Distribution of first digit in frequencies?

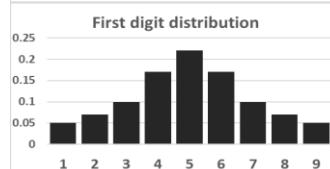
1) Uniform →



2) Exp decay →



3) Normal →



Term	Rank	Frequency
the	1	5134,790
of	2	3102,474
in	3	2607,875
a	4	2492,328
is	5	2181,502
and	6	1962,326
was	7	11159,088
to	8	1088,396
by	9	766,656
an	10	566,970
it	11	557,492
for	13	493,374
as	14	480,277
on	15	471,544
from	16	412,785

Walid Magdy, TTDS 2021/2022



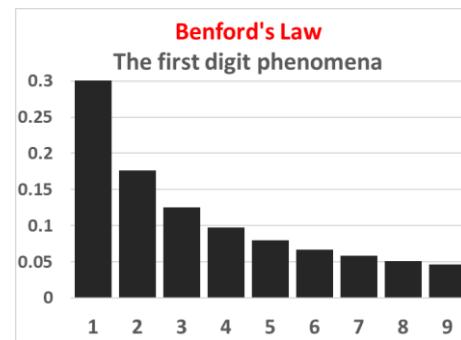
THE UNIVERSITY  
of EDINBURGH

9

## Benford's Law:

- First digit of a number follows a Zipf's like law!
  - Terms frequencies
  - Physical constants
  - Energy bills
  - Population numbers
- Benford's law:

$$P(d) = \log\left(1 + \frac{1}{d}\right)$$



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## Practical

Walid Magdy, TTDS 2021/2022



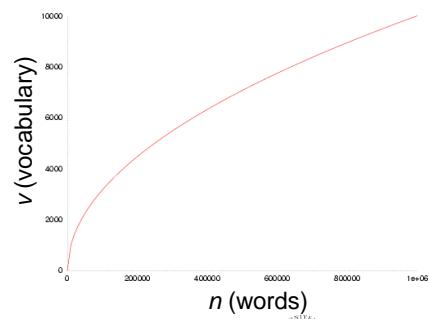
11

## Heap's Law:

- While going through documents, the number of new terms noticed will reduce over time
- For a book/collection, while reading through, record:
  - $n$ : number of words read
  - $v$ : number of news words (unique words)
- Vocabulary growth:  

$$v(n) = k \times n^b$$

where,  $b < 1$   
 typically,  $0.4 < b < 0.7$

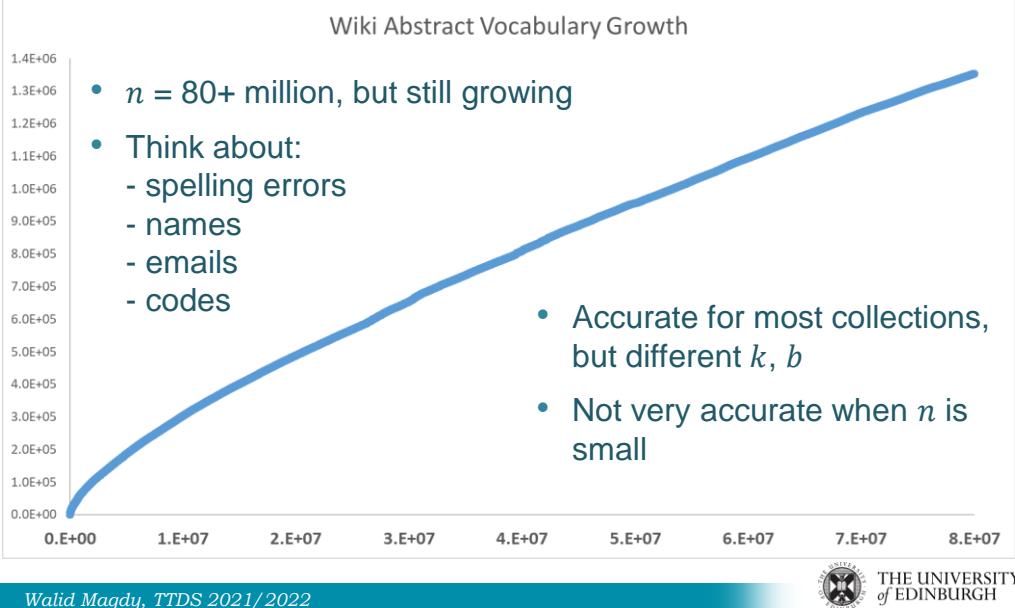


Walid Magdy, TTDS 2021/2022



12

## Heap's Law: shouldn't it saturate?



13

## Practical

Walid Magdy, TTDS 2021/2022



14

## Clumping/Contagion in text

- From Zipf's law, we notice:
  - Most words do not appear that much!
  - Once you see a word once → expect to see again!
  - Words are like:
    - Rare contagious disease
    - Not, rare independent lightning
- Words are rare events, but they are contagious

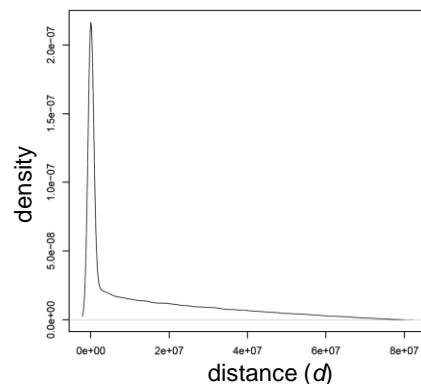
Walid Magdy, TTDS 2021/2022



15

## Clumping/Contagion in text

- Wiki abstract collection
  - Identify terms appeared only twice
  - Measure distance between the two occurrences of the terms:  
 $d = n_{occurrence2} - n_{occurrence1}$
  - Plot density function of  $d$
- Majority of terms appearing only twice appear close to each other.



Walid Magdy, TTDS 2021/2022



16

## Applying the laws

- Given a collection of 20 billion terms,
- What is the number of unique terms?
  
- What is the number of terms appearing once?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

17

## Summary

- Text follows well-known phenomena
- Text Laws:
  - Zipf
  - Heap
  - Contagion in text
  
- Try it on another language ...

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

18

## Recourses

- Text book:
  - Search engines: IR in practice → chapter 4
- Videos:
  - Zipf's law, Vsouce:  
<https://www.youtube.com/watch?v=fCn8zs912OE>
  - Benford's law, Numberphile:  
<https://www.youtube.com/watch?v=XXjIR2OK1kM>
- Tools:
  - Unix commands for windows  
<https://sourceforge.net/projects/unxutils>

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

19

## Next Lecture

- Getting ready for indexing?
- Pre-processing steps before the indexing process
- Reminder: 5-10 mins break after L1
  - Have a break, stretch, get food ... etc.
  - Ask questions on chat
  - Questions on L1 are allowed before starting L2
  - Mind teaser math problem (for fun)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

20



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Preprocessing

Instructor:  
**Walid Magdy**

29-Sep-2021

1

## Lecture Objectives

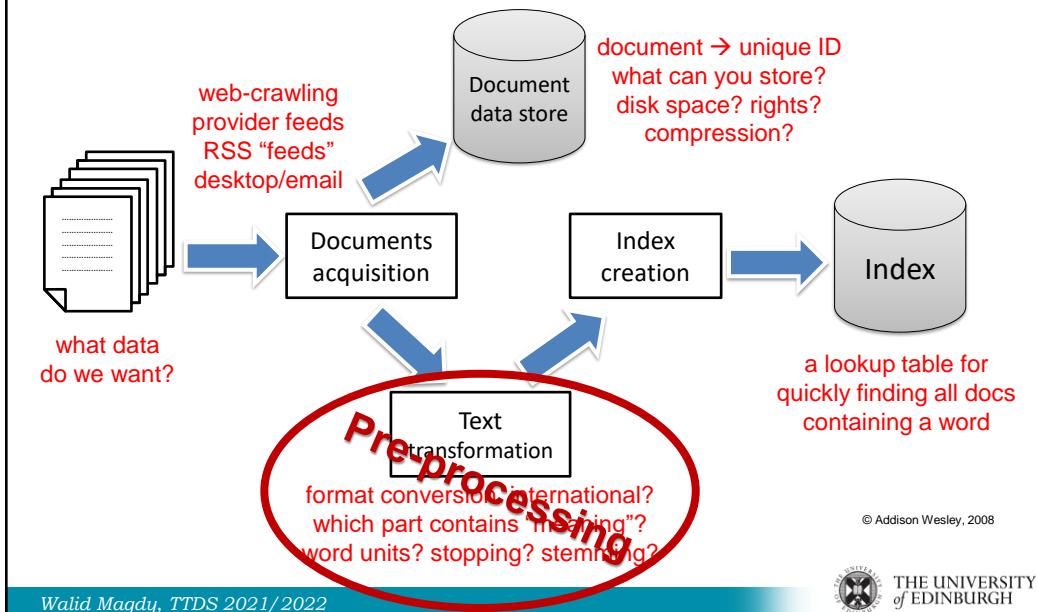
- Learn about and implement
- Standard text pre-processing steps:
  - Tokenisation
  - Stopping
  - Normalisation
  - Stemming



THE UNIVERSITY  
of EDINBURGH

2

## Indexing Process



3

## Preprocessing

Find the best text transformation technique (preprocessing) that will lead to better match between different forms of words in document and query

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY of EDINBURGH

4

## Getting ready for indexing?

- BOW, what is a word?
- In IR, we refer to word-elements as “terms”
  - word “*preprocessing*”
  - part of a word “*pre*”
  - number / code “*INFR11145*”
- Pre-processing steps before indexing:
  - Tokenisation
  - Stopping
  - Stemming
- **Objective** → identify the optimal form of the term to be indexed to achieve the best retrieval performance

Walid Magdy, TTDS 2021/2022



5

## Tokenisation

- Input: “*Friends, Romans; and Countrymen!*”
- Output: Tokens
  - *Friends*
  - *Romans*
  - *and*
  - *Countrymen*
- Sentence → tokenization (splitting) → tokens
- A **token** is an instance of a sequence of characters
- **Typical technique**: split at non-letter characters
- Each such token is now a candidate for an index entry (**term**), after further processing

Walid Magdy, TTDS 2021/2022



6

## Issues in Tokenisation

- “*Finland’s*” capital → *Finland?* *Finlands?* *Finland’s?*
- Hewlett-Packard → one token or two?
  - **state-of-the-art:** break up hyphenated sequence.
  - *co-education*
  - *lowercase, lower-case, lower case ?*
  - It can be effective to get the user to put in possible hyphens
- **Numbers?**
  - 3/20/91 vs. Mar. 20, 1991 vs. 20/3/91
  - This course code is INFR11145
  - (800) 234-2333

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Issues in Tokenisation

- **URLs:**
  - <http://www.bbc.co.uk>
  - <http://www.bbc.co.uk/news/world-europe-41376577>
- **Social Media**
  - Black lives matter
  - #Black\_lives\_matter
  - #BlackLivesMatter
  - #blacklivesmatter
  - @blacklivesmatter
- **San Francisco:** one token or two?
  - How do you decide it is one token?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## Tokenisation for different languages

- French → *L'ensemble* → one token or two?
  - *L*? *L'*? *Le*?
  - Want *L'ensemble* to match with *un ensemble*
  - Until at least 2003, it didn't on Google
- German → compounds
  - *Lebensversicherungsgesellschaftsangestellter*  
'life insurance company employee'
  - German retrieval systems benefit greatly from a **compound splitter** module → Can give a 15% performance boost for German
- Chinese and Japanese → no spaces between words:
  - 莎拉波娃现在居住在美国东南部的佛罗里达
  - Tokenisation → Segmentation

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## Tokenisation: common practice

- Just split at non-letter characters
- Add special cases if required
- Some applications have special setup
  - Social media: hashtags/mentions handled differently
  - URLs: no split, split at domain only, remove entirely!
  - Medical: protein & diseases names

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## Stopping (stop words removal)

- ~~This is a very exciting lecture on the technologies of text~~
- **Stop words:** the most common words in collection  
→ the, a, is, he, she, I, him, for, on, to, very, ...
- There are a lot of them ≈ 30-40% of text
- New stop words appear in specific domains
  - Tweets: RT → “*RT @realDonaldTrump Mexico will ...*”
  - Patents: said, claim → “*a said method that extracts ....*”
- Stop words
  - influence on sentence structure
  - less influence on topic (aboutness)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Stopping: always apply?

- Sometimes very important:
  - Phrase queries: “Let it be”, “To be or not to be”
  - Relational queries:
    - flights to London from Edinburgh
    - flights from London to Edinburgh
- In Web search, trend is to keep them:
  - Good compression techniques means the space for including stop words in a system is very small
  - Good query optimization techniques mean you pay little at query time for including stop words.
  - Probabilistic retrieval models give them low weight.

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

12

## Stopping: stop words

- Common practice in many applications  
→ remove stop words
- There are common stop words list for each language
  - NLTK (python)
  - <http://members.unine.ch/jacques.savoy/clef/index.html>
- There are special stop words list for some applications
- How to create your list:
  - Sort all terms in a collection by frequency
  - Manually select the possible stop words from top  $N$  terms

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Normalisation

- **Objective** → make words with different surface forms look the same
- Document: “this is my CAR!!”  
Query: “car”  
should “car” match “CAR”?
- Sentence → tokenisation → **tokens** → normalisation  
→ **terms** to be indexed
- Same tokenisation/normalisation steps should be applied to documents & queries

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

14

## Case folding and equivalents

- “A” & “a” are different strings for computers
- Case folding: convert all letters to lower case
  - CAR, Car, caR → car
  - Windows → windows, should we do that?
- Diacritics/Accents removal
  - French: Château → chateau
  - German: Tüebingen → tuebingen
  - Arabic: كتب → كتب

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Equivalence Classes

- U.S.A. → USA
- Ph.D. → PhD
- 92.3 → 923? 92 3?
- multi-disciplinary → multidisciplinary ← multi disciplinary
- The most important criteria:
  - Be consistent between documents & queries
  - Try to follow users' most common behaviour

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

16

## Stemming

- Search for: “play”  
should it match: “played”, “playing”, “player”?
- Many morphological variations of words
  - *inflectional* (plurals, tenses)
  - *derivational* (making verbs nouns etc.)
- In most cases, aboutness does not change
- Stemmers attempt to reduce morphological variations of words to a common stem
  - usually involves removing suffixes (in English)
- Can be done at indexing time or as part of query processing (like stopwords)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

17

## Stemming

- Usually, it achieves 5-10% improvement in retrieval effectiveness, e.g. English
- For highly inflected languages, it is more critical:
  - 30% improvement in Finnish IR
  - 50% improvement in Arabic IR

They are Peter's **children**

هؤلاء **أبناء** بيتر

The **children** behaved well

ال**الأبناء** تصرفوا جيدا

Her **children** are cute

**أبناءها** لطاف

My **children** are funny

**أبنائي** ظرفاء

We have to save **our children**

عليينا أن نحمي **أبناءنا**

Patents and **children** are happy

الآباء **والأبناء** سعداء

He loves **his children**

هو يحب **أبناءه**

**His children** loves him

**أبناءه** يحبونه

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

18

## Stemming

- Two basic types
  - Dictionary-based: uses lists of related words
  - Algorithmic: uses program to determine related words
- Algorithmic stemmers
  - suffix-s: remove ‘s’ endings assuming plural
  - e.g., **cats** → **cat**, **lakes** → **lake**, **windows** → **window**
  - Many false negatives: **supplies** → **supplie**
  - Some false positives: **James** → **Jame**



## Porter Stemmer

- Most common algorithm for stemming English
- Conventions + 5 phases of reductions
  - phases applied sequentially
  - each phase consists of a set of commands
  - sample convention:  
of the rules in a compound command, select the one that applies to the longest suffix.
- Example rules in Porter stemmer
 

• <i>sses</i> → <i>ss</i>	(processes → process)
• <i>y</i> → <i>i</i>	(reply → repli)
• <i>ies</i> → <i>i</i>	(replies → repli)
• <i>ement</i> → null	(replacement → replac)



## Stemmed words are misspelled!!

- repli, replac, suppli, inform retriev, anim
- These are not words anymore, these are terms
- These terms are not seen by the user, but just used by the IR system (search engine)
- These represent the optimal form for a better match between different surface forms of a term
  - e.g. `replac` → `replace, replaces, replaced, replacing, replacer, replacers, replacement, replacements.`

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

21

## Pre-processing: Common practice

- Tokenisation: split at non-letter characters
  - One line of code in Perl  
→ process \w and neglect anything else
  - For tweets, you might want to keep "#" and "@"
- Remove stop words
  - find a common list, and filter these words out
- Apply case folding
  - One command in Perl or Python: `lc($string)`
- Apply Porter stemmer
  - Other stemmers are available, but Porter is the most famous with many implementations available in different programming languages

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

22

## Limitations

- Irregular verbs:
  - saw → see
  - went → go
- Different spellings
  - colour vs. color
  - tokenisation vs. tokenization
  - Television vs. TV
- Synonyms
  - car vs. vehicle
  - UK vs. Britain
- Solution → Query expansion ...

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

23

## Asymmetric Expansion

- Maintains relations between unnormalized tokens
- An alternative to equivalence classing
- An example of where this may be useful
  - query: *window*      search: *window, windows*
  - query: *windows*      search: *windows, Windows*
  - query: *Windows*      search: *Windows*
- Potentially more powerful, but less efficient
  - More vocabulary, longer query
- Can be less effective:
  - Inaccurate stats on terms (“car” ≠ “Car”)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

24

## Summary

- Text pre-processing before IR:
  - Tokenisation → Stopping → Stemming

This is an **example sentence** of how the **pre-processing** is applied to **text in information retrieval**. It **includes**: **Tokenization**, **Stop Words Removal**, and **Stemming**



exampl sentenc pre process appli text inform retriev includ token stop word remov stem

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

25

## Practical

Collection	Original		After Pre-processing	
	# words	File size	# words	File size
Bible	824,054	4.24 MB	358,112	2.05 MB
Wiki abstracts	78,137,597	472 MB	47,741,065	309 MB

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

26

## Resources

- Text book 1: Intro to IR, Chapter 2 → 2.2.4
- Text book 2: IR in Practice, chapter 4
- Lab 1 → Implement what learnt in these two lectures
- Optional reading:  
*if you think English pre-processing is hard!*  
- Arabic Information Retrieval. *Darwish & Magdy*

Walid Magdy, TTDS 2021/2022



27

## Next lecture

- Indexing:  
How to build an index!
- Assignment 1 announcement:
  - Build indexing components
  - Today: build your pre-processing module!
  - Next time: build the index

Walid Magdy, TTDS 2021/2022



28



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Indexing

Instructor:  
**Walid Magdy**

07-Oct-2020

1

## Lecture Objectives

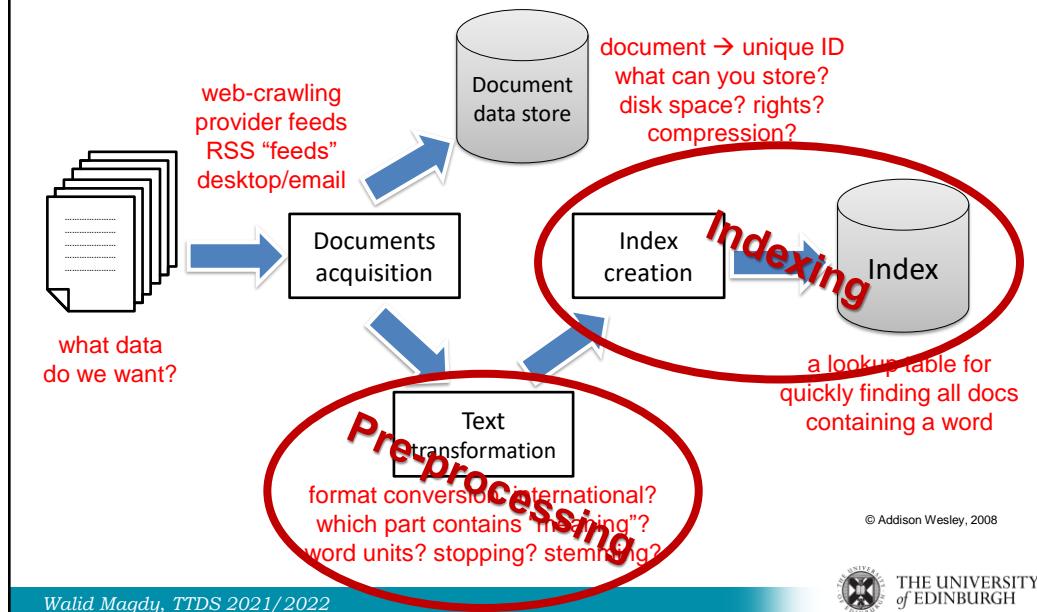
- Learn about and implement
- Boolean search
- Inverted index
- Positional index



THE UNIVERSITY  
of EDINBURGH

2

## Indexing Process



3

## Pre-processing output

This is an example sentence of how the pre-processing is applied to text in information retrieval. It includes: Tokenization, Stop Words Removal, and Stemming

example sentence pre process apply text inform retrieval include token stop word removal stemming

- Add processed terms to index
  - What is “index”?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

4

# Index

- How to match your term in non-linear time?
- Find/Grep:  
Sequential search for term
- Index:  
Find term locations immediately

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

5

# Book Index

## Index

absolute error, 437  
accuracy, 359  
ad hoc search, 3, 280, 423  
adaptive filtering, 425  
adversarial information retrieval, 294  
advertising, 218, 371  
    classifying, 371  
    contextual, 218–221  
agglomerative clustering, 375  
anchor text, 21, 56, 105, 280  
API, 439, 461  
architecture, 13–28  
authority, 21, 111  
automatic indexing, 400  
background probability, *see* collection probability  
bag of words, 345, 451  
Bayes classifier, 245  
Bayes Decision Rule, 245  
Bayes' Rule, 246, 343  
Bayes' rule, 342  
Bayesian network, 268  
bibliometrics, 120  
bidding, 218  
bigram, 100, 253  
BigTable, 57

## 512 Index

binary independence model, 246  
blog, 111  
BM25, 250–252  
BM25F, 294  
Boolean query, 235  
Boolean query language, 24  
Boolean retrieval, 235–237  
boosting, 448  
BPREF, 322  
brute force, 331  
burstiness, 254  
caching, 26, 181  
card catalog, 400  
case folding, 87  
case normalization, 87  
categorization, *see* classification  
CBIR, *see* content-based image retrieval  
character encoding, 50, 119  
checksum, 60  
Chi-squared measure, 202  
CJK (Chinese-Japanese-Korean), 50, 119  
classification, 3, 339–373  
    faceted, 224  
    monothetic, 223, 374  
    polythetic, 223, 374  
classifier, 21

crawler, 17, 32  
cross-language information retrieval, 226  
cross-lingual search, *see* cross-language information retrieval  
cross-validation, 331

Damerau-Levenshtein distance, 194  
dangling link, 107  
data mining, 113  
database system, 459  
DCG, *see* discounted cumulative gain  
deep Web, 41, 448  
delta encoding, 144  
dendrogram, 375  
desktop search, 3, 46  
Dice's coefficient, 192  
digital reference, 447  
Dirichlet smoothing, 258  
discounted cumulative gain, 319  
discriminative model, 284, 360  
distance measure, 374  
distributed hash table, 445  
distributed information retrieval, 438  
distribution, 23  
divisive clustering, 375  
document, 2  
document conversion, 18  
document crawler, 17  
document data store, 19  
document distribution, 180  
document slope curve, 64  
document statistics, 22  
document structure, 101, 269, 459–466  
document summary, 215  
downcasing, 87  
dumping, 366  
duplicate documents, 60  
dwell time, 27  
dynamic page, 42

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

6

## Indexing

- Search engines vs PDF find or grep?
  - Infeasible to scan large collection of text for every “search”
  - Find section that has: “UK and Scotland and Money”?!
- Book Index
  - For each word, list of “relevant” pages
  - Find topic in sub-linear time
- IR Index:
  - Data structure for fast finding terms
  - Additional optimisations could be applied

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Document Vectors

- Represent documents as vectors
  - Vector → document, cell → term
  - Values: term frequency or binary (0/1)
  - All documents → collection matrix

he	drink	ink	likes	pink	think	wink
2	1	0	2	0	0	1
1	3	0	1	0	0	0
1	1	1	1	0	1	0
1	1	1	1	1	0	0
1	1	1	1	1	0	1

← D1: He likes to wink, he likes to drink

← D2: He likes to drink, and drink, and drink

← D3: The thing he likes to drink is ink

← D4: The ink he likes to drink is pink

← D5: He likes to wink, and drink pink ink

number of occurrence of  
a term in a document

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## Inverted Index

- Represent terms as vectors
  - Vector → term, cell → document
  - Transpose of the collection matrix
  - Vector: inverted list

he	drink	ink	likes	pink	think	wink	
2	1	0	2	0	0	1	← D1: He likes to wink, he likes to drink
1	3	0	1	0	0	0	← D2: He likes to drink, and drink, and drink
1	1	1	1	0	1	0	← D3: The thing he likes to drink is ink
1	1	1	1	1	0	0	← D4: The ink he likes to drink is pink
1	1	1	1	1	0	1	← D5: He likes to wink, and drink pink ink

Walid Magdy, TTDS 2021/2022



9

## Boolean Search

- Boolean: exist / not-exist
- Multiword search: logical operators (AND, OR, NOT)
- Example
  - Collection: search Shakespeare's Collected Works
  - Boolean query: Brutus AND Caesar AND NOT Calpurnia
- Build a **Term-Document Incidence Matrix**
  - Which term appears in which document
  - Rows are terms
  - Columns are documents

Walid Magdy, TTDS 2021/2022



10

## Collection Matrix

	Documents					
	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Terms	1	1	0	0	0	1
Antony	1					
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

1 if **document** contains **term**, 0 otherwise

Query: Brutus AND Caesar AND NOT Calpurnia  
 Apply on rows: **110100 AND 110111 AND !(010000) = 100100**

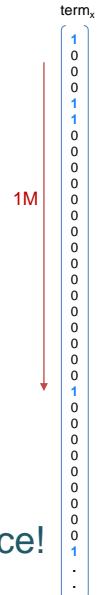
Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

11

## Bigger collections?

- Consider  $N = 1$  million documents, each with about 1000 words.
- $n = 1M \times 1K = 1B$  words  
 $\rightarrow$  Heap's law  $\rightarrow v \approx 500K$
- Matrix size = 500K unique terms x 1M documents = 0.5 trillion 0's and 1's entries!
- If all words appear in many documents  
 $\rightarrow \max\{\text{count(1's)}\} = N * \text{doc. length} = 1B$
- Actually, from Zip's law  $\rightarrow 250k$  terms appears once!
- Collection matrix is extremely sparse. (mostly 0's)



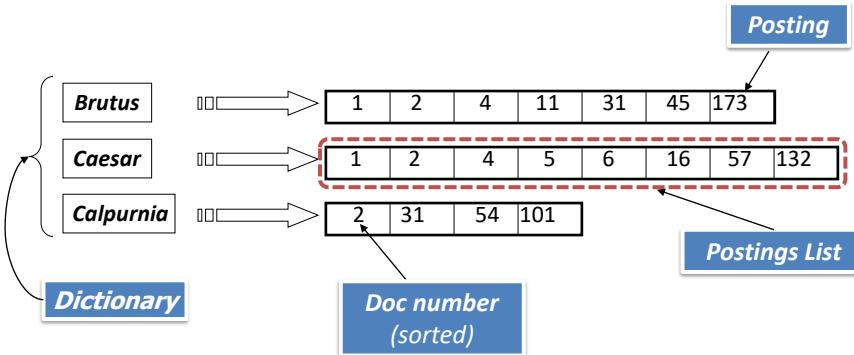
Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

12

## Inverted Index: Sparse representation

- For each term  $t$ , we must store a list of all documents that contain  $t$ .
  - Identify each by a **docID**, a document serial number

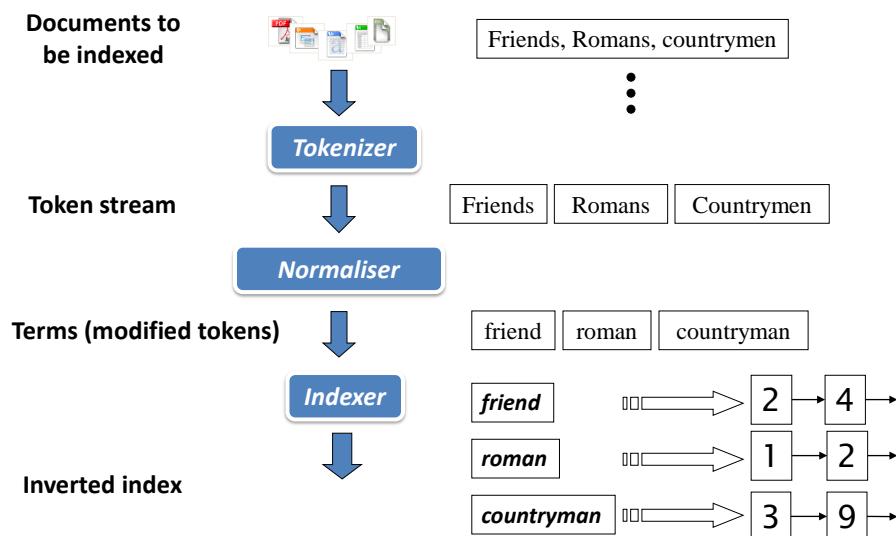


Walid Magdy, TTDS 2021/2022



13

## Inverted Index Construction



Walid Magdy, TTDS 2021/2022



14

## Step 1: Term Sequence

### Doc 1

I did enact Julius Caesar I was killed i' the Capitol; Brutus killed me.

### Doc 2

So let it be with Caesar. The noble Brutus hath told you Caesar was ambitious

Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

## Step 2: Sorting

- Sort by:
  - 1) Term
  - then
  - 2) Doc ID

Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
caesar	2
was	2
was	2
ambitious	2

Sorting

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

## Step 3: Posting

1. Multiple term entries in a single document are merged
2. Split into Dictionary and Postings
3. Doc. Frequency ( $df$ ) information is added

The diagram illustrates the process of creating an inverted index from a term-document matrix. On the left, a table shows terms and their document IDs (docID). An arrow points to the right, where the same terms are shown with their document frequencies (df) and posting lists.

Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
i	1
i	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2

term	doc. freq.	→	postings lists
ambitious	1	→	2
be	1	→	2
brutus	2	→	1 → 2
capitol	1	→	1
caesar	2	→	1 → 2
did	1	→	1
enact	1	→	1
hath	1	→	2
i	1	→	1
i'	1	→	1
it	1	→	2
julius	1	→	1
killed	1	→	1
let	1	→	2
me	1	→	1
noble	1	→	2
so	1	→	2
the	2	→	1 → 2
told	1	→	2
you	1	→	2
was	2	→	1 → 2
with	1	→	2

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

17

## Inverted Index: matrix → postings

he	drink	ink	likes	pink	think	wink
2	1	0	2	0	0	1
1	3	0	1	0	0	0
1	1	1	1	0	1	0
1	1	1	1	1	0	0
1	1	1	1	1	0	1

← D1: He likes to wink, he likes to drink  
 ← D2: He likes to drink, and drink, and drink  
 ← D3: The thing he likes to drink is ink  
 ← D4: The ink he likes to drink is pink  
 ← D5: He likes to wink, and drink pink ink

<i>he</i>	⇒	1	2	3	4	5
<i>drink</i>	⇒	1	2	3	4	5
<i>ink</i>	⇒	3	4	5		
<i>pink</i>	⇒	4	5			
<i>thing</i>	⇒	3				
<i>wink</i>	⇒	1	5			

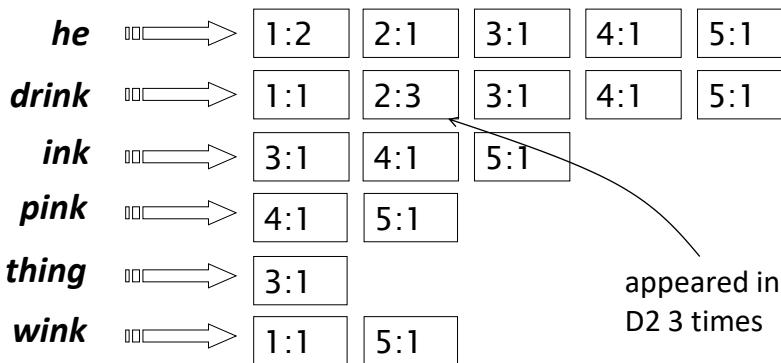
Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

18

## Inverted Index: with frequency

- Boolean: term  $\rightarrow$  DocIDs list
- Frequency: term  $\rightarrow$  tuples (DocID, count(term)) lists



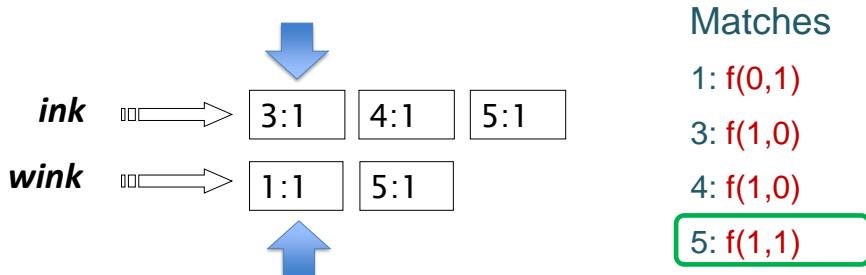
Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

19

## Query Processing

- Find documents matching query {ink **AND** wink}
  - Load inverted lists for each query word
  - Merge two postings lists  $\rightarrow$  **Linear merge**
- Linear merge  $\rightarrow O(n)$   
 $n$ : total number of posts for all query words



Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

20

## Phrase Search

- Find documents matching query “pink ink”

- Find document containing both words
- Both words has to be a phrase

- Bi-gram Index:

He likes to wink, and drink pink ink Convert to bigrams

He\_likes likes\_to to\_wink wink\_and and\_drink drink\_pink pink\_ink

- Bi-gram Index, issues:

- Fast, but index size will explode!
- What about trigram phrases?
- What about proximity? “ink is pink”

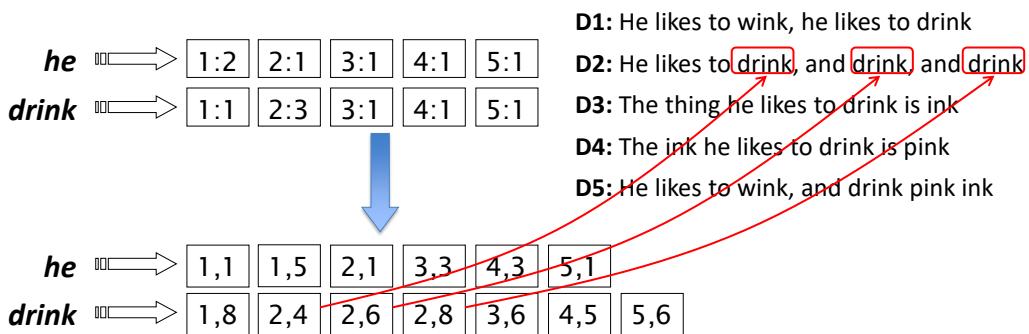
Walid Magdy, TTDS 2021/2022



21

## Proximity Index

- Terms positions is embedded to the inv. Index
  - Called proximity/positional index
  - Enables phrase and proximity search
  - Toubles (DocID, term position)



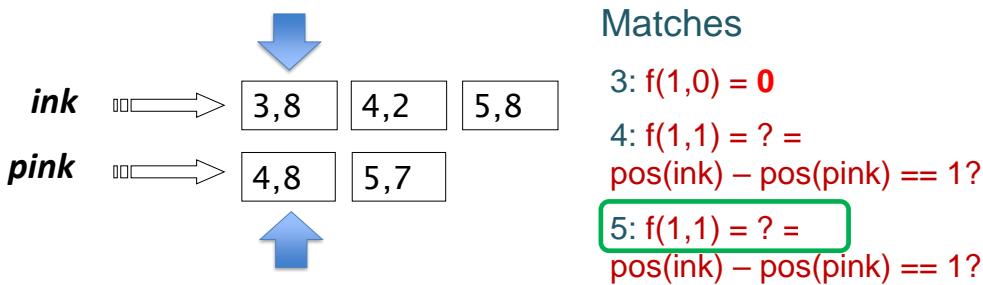
Walid Magdy, TTDS 2021/2022



22

## Query Processing: Proximity

- Find documents matching query “pink ink”
  - Use **Linear merge**
  - Additional step: check terms positions
- Proximity search:**  
 $\text{pos}(\text{term1}) - \text{pos}(\text{term2}) < |w| \rightarrow \#5(\text{pink}, \text{ink})$



Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

23

## Proximity search: data structure

- Possible data structure:  
 $\langle \text{term: df;}$   
 $\text{DocNo: pos1, pos2, pos3}$   
 $\text{DocNo: pos1, pos2, pos3}$   
 $\dots \dots \rangle$
- Example:  
 $\langle \text{be: 993427;}$   
 $1: 7, 18, 33, 72, 86, 231;$   
 $2: 3, 149;$   
 $4: 17, 191, 291, 430, 434;$   
 $5: 363, 367, \dots \rangle$

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

24

## Summary

- Document Vector
- Term Vector
- Inverted Index
- Collection Matrix
- Posting
- Proximity Index
- Query Processing → Linear merge

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

25

## Resources

- Textbook 1: Intro to IR, Chapter 1 & 2.4
- Textbook 2: IR in Practice, Chapter 5
- Lab 2

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

26



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Indexing (2)

Instructor:  
**Walid Magdy**

06-Oct-2021

1

## Lecture Objectives

- Learn more about indexing:
  - Structured documents
  - Extent index
  - Index compression
- Data structure
- Wild-char search and applications

\* You are not asked to implement any of the content in this lecture, but you might think of using some for your course project ☺

## Structured Documents

- Document are not always flat:
    - Meta-data: title, author, time-stamp
    - Structure: headline, section, body
    - Tags: link, hashtag, mention
  - How to deal with it?
    - Neglect!
    - Create separate index for each field
    - Use “extent index”

Walid Maqdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

3

## Extent Index

- Special “term” for each element/field/tag
    - Index all terms in a structured document as plain text
    - Terms in a given field/tag get special additional entry
    - Posting: spans of window related to a given field
    - Allows multiple overlapping spans of different types

*he* [1.1] [1.5] [2.1] [3.3] [4.3] [5.1]

**drink** → 

1.8	2.4	2.6	2.8	3.6	4.5	5.6
-----	-----	-----	-----	-----	-----	-----

*ink*  3,8 4,2 5,8

*pink* → 

4.8	5.7
-----	-----

**Link**  **3,1:2** **4,1:4** **5,7:8**

D1: He likes to work, he likes to drink

D2: He likes to drink, and drink, and drink.

D3: The thing he likes to drink is ink

**D4:** The ink he likes to drink is pink

D5: He likes to wink, and drink pink ink

Walid Maadu, TTDS 2021/2022

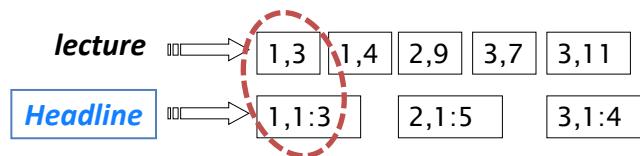


THE UNIVERSITY  
of EDINBURGH

4

## Using Extent

- Doc: 1 → 1 2 3  
Headline: “*Information retrieval lecture*”  
Text: “~~this is lecture 6 of the TTDS course on IR~~”  
4 5 6 7 8
- Query → Headline: lecture



Walid Magdy, TTDS 2021/2022



5

## Index Compression

- Inverted indices are big
  - Large disk space → large I/O operations
- Index compression
  - Reduce space → less I/O
  - Allow more chunks of index to be cached in memory
- Large size goes to:
  - terms? document numbers?
  - Ideas:
    - Compress document numbers, how?

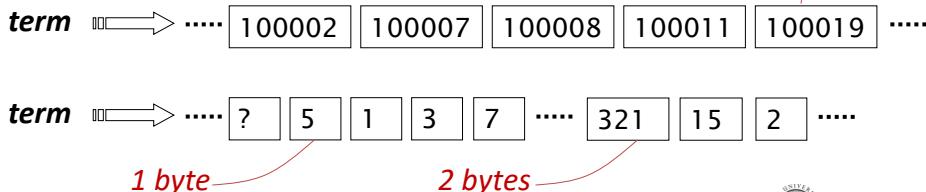
Walid Magdy, TTDS 2021/2022



6

## Delta Encoding

- Large collections → large sequence of doc IDs
  - e.g. Doc IDs: 1, 2, 3, ... 66,032, ..., 5,323,424,235
- Large ID number → more bytes to store
  - 1 byte: 0 → 255
  - 2 bytes: 0 → 65,535
  - 4 bytes: 0 → 4.3 B
- Idea: delta in ID instead of full ID
  - Very useful, especially for frequent terms



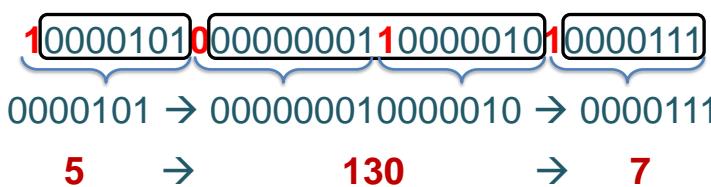
Walid Magdy, TTDS 2021/2022



7

## v-byte Encoding

- Have different byte storage for each delta in index
  - Use fewer bits to encode
  - High bit in a byte → 1/0 = terminate/continue
  - Remaining 7 bits → binary number
  - Examples:
    - "6" → 10000110
    - "127" → 11111111
    - "128" → 0000000100000000
- Real example sequence:



Walid Magdy, TTDS 2021/2022



8

## Index Compression

- There are more sophisticated compression algorithms:
  - Elias gamma code
- The more compression
  - Less storage
  - More processing
- In general
  - Less I/O + more processing > more I/O + no processing  
“>” = faster
  - With new data structures, problem is less severe

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

9

## Dictionary Data Structures

- The dictionary data structure stores the term vocabulary, document frequency, pointers to each postings list ...
- For small collections, load full dictionary in memory.  
In real-life, cannot load all index to memory!
  - Then what to load?
  - How to reach quickly?
  - What data structure to use for inverted index?

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

10

## Hashes

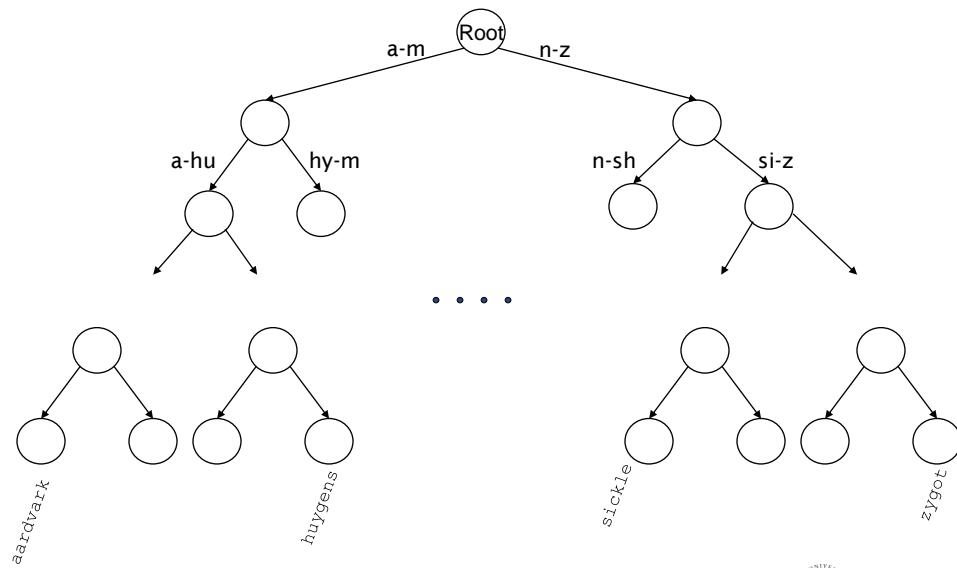
- Each vocabulary term is hashed to an integer
- Pros
  - Lookup is faster than for a tree: O(1)
- Cons
  - No easy way to find minor variants:
    - judgment/judgement
  - No prefix search
  - If vocabulary keeps growing, need to occasionally do the expensive operation of rehashing everything

Walid Magdy, TTDS 2021/2022



11

## Trees: Binary Search Tree

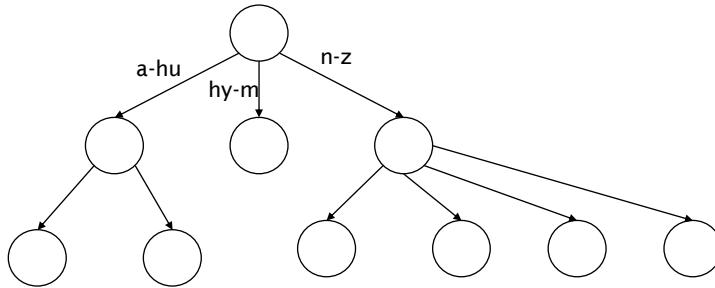


Walid Magdy, TTDS 2021/2022



12

## Trees: B-tree



Every internal node has a number of children in the interval  $[a,b]$  where  $a, b$  are appropriate natural numbers, e.g.,  $[2,4]$ .

## Trees

- Pros?
  - Solves the prefix problem (terms starting with “ab”)
- Cons?
  - Slower:  $O(\log M)$  [and this requires balanced tree]
  - Rebalancing binary trees is expensive
    - But B-trees mitigate the rebalancing problem

## Wild-Card Queries: \*

- $\text{mon}^*$ : find all docs containing any word beginning “mon”.
- Easy with binary tree (or B-tree) lexicon
- $^*\text{mon}$ : find words ending in “mon”: harder
  - Maintain an additional B-tree for terms backwards.
- How can we enumerate all terms meeting the wild-card query  $\text{pro}^*\text{cent}$  ?
- Query processing:  $\text{se}^*\text{ate}$  AND  $\text{fil}^*\text{er}$  ?
  - Expensive

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Permuterm Indexes

- Transform wild-card queries so that the \* occurs at the end
- For term *hello*, index under:
  - *hello\$*, *ello\$h*, *llo\$he*, *lo\$hel*, *o\$hell*, *\$hello*  
where \$ is a special symbol.
- Rotate query wild-card to the end
- Queries:
  - X lookup on X\$
  - X\* lookup on \$X\*
  - \*X lookup on
  - X\*Y lookup on
- Index Size?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

16

## Character n-gram Indexes

- Enumerate all n-grams (sequence of  $n$  chars) occurring in any term
  - e.g., from text “*April is the cruelest month*” we get the 2-grams (bigrams) →  
 $\$a, ap, pr, ri, il, I\$, \$i, is, s\$, \$t, th, he, e\$, \$c, cr, ru, ue, el, le, es, st, t\$,$   
 $\$m, mo, on, nt, h\$$
  - \$ is a special word boundary symbol
- Maintain a second inverted index from bigrams to dictionary terms that match each bigram.
  - Character n-grams → terms
  - terms → documents

Walid Magdy, TTDS 2021/2022

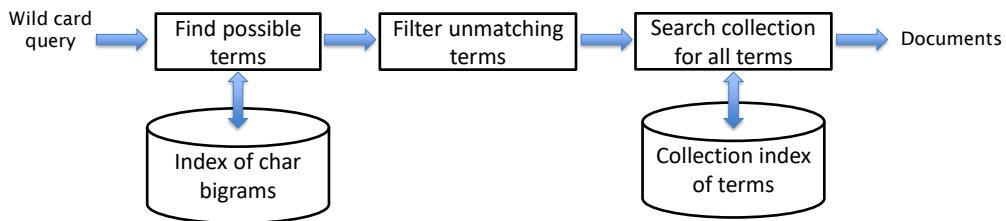


17

## Character n-gram Indexes

- The  $n$ -gram index finds *terms* based on a query consisting of  $n$ -grams (here  $n=2$ ).

$\$m$	→	<b>mace</b>	<b>madden</b>	-----
$mo$	→	<b>among</b>	<b>amortize</b>	-----
$on$	→	<b>almond</b>	<b>among</b>	-----



Walid Magdy, TTDS 2021/2022



18

## Character n-gram Indexes: Query time

- Step 1: Query **mon\*** → \$m AND mo AND on
  - It would still match **moon**.
- Step 2: Must post-filter these terms against query.
  - Phrase match, or post-step1 match
- Step 3: Surviving enumerated terms are then looked up in the term-document inverted index.  
**→ Montreal OR monster OR monkey**
- Wild-cards can result in expensive query execution (very large disjunctions...)

## Character n-gram Indexes: Applications

- Spelling Correction
  - Create n-gram representation for words
  - Build index for words:
    - Dictionary of words → documents (each word is a document)
    - Character n-grams → terms
  - When getting a search term that is misspelled (OOV or not frequent), find possible corrections
    - Possible corrections = most matching results

Query: elepgant → \$e el le ep pg ga an nt t\$

Results:

elegant → \$e el le eg ga an nt t\$

elephant → \$e el le ep ph ha an nt t\$

## Character n-gram Indexes: Applications

- Char n-grams can be used as direct index terms for some applications:
    - Arabic IR, when no stemmer/segmenter is available
    - Documents with spelling mistakes: OCR documents
  - Word char representation can be with multiple n's
    - “elephant” → 2/3-gram →  
“\$e el le ep ph ha an nt t\$ \$el \$ele lep eph pha han ant nt\$”

The **children** behaved well  
Her **children** are cute

**الأبناء** تصرفوا جيدا  
**أبناءها** لطاف

\$ ا ل لا ا ب بن نا اء \$  
\$ ا ب بن نا اء ها \$

Document: Elephant → \$e el le ep pb ba an nt t\$  
Query: Elephant → \$e el le ep ph ha an nt t\$



THE UNIVERSITY  
of EDINBURGH

Walid Magdy, TTDS 2021/2022

21

## Summary

- Index can be multilayer
    - Extent index (multi-terms in one position in document)
  - Index does not have to be formed of words
    - Character n-grams representation of words
  - Two indexes are sometimes used
    - Index of character n-grams to find matching words
    - Index of terms to search for matched words

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

22

## Resources

- Text book 1: Intro to IR, Chapter 3.1 – 3.4
- Text book 2: IR in Practice, Chapter 5

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Ranked IR

Instructor:  
**Walid Magdy**

13-Oct-2021

1

## Lecture Objectives

- Learn about Ranked IR
  - TFIDF
  - VSM
  - SMART notation
- Implement:
  - TFIDF



## Boolean Retrieval

- Thus far, our queries have all been Boolean.
  - Documents either: “match” or “no match”.
- Good for expert users with precise understanding of their needs and the collection.
  - Patent search uses sophisticated sets of Boolean queries and check hundreds of search results  
**(car OR vehicle) AND (motor OR engine) AND NOT (cooler)**
- Not good for the majority of users.
  - Most incapable of writing Boolean queries.
  - Most don’t want to go through 1000s of results.
    - This is particularly true for web search
    - Question: What is the most unused web-search feature?

Walid Magdy, TTDS 2021/2022



3

## Ranked Retrieval

- Typical queries: free text queries
- Results are “ranked” with respect to a query
- Large result sets are not an issue
  - We just show the top k ( $\approx 10$ ) results
  - We don’t overwhelm the user
- Criteria:
  - Top ranked documents are the most likely to satisfy user’s query
  - Score is based on how well documents match a query  
**Score( $d, q$ )**

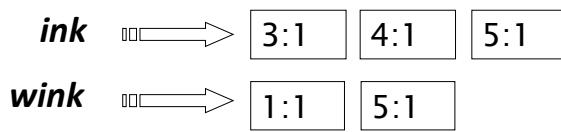
Walid Magdy, TTDS 2021/2022



4

## Old Example

- Find documents matching query {ink wink}
  - Load inverted lists for each query word
  - Merge two postings lists → **Linear merge**
- Apply function for matches
  - Boolean: exist / not exist = 0 or 1
  - Ranked:  $f(tf, df, length, \dots) = 0 \rightarrow 1$



**Matches**  
 1: f(0,1)  
 3: f(1,0)  
 4: f(1,0)  
 5: f(1,1)

## Function example: Jaccard coefficient

- a commonly used measure of overlap of two sets  $A$  and  $B$
- $jaccard(A, B) = \frac{|A \cap B|}{|A \cup B|}$ 
  - D1: He likes to wink, he likes to drink
  - D2: He likes to drink, and drink, and drink
- $jaccard(A, A) = 1$
- $jaccard(A, B) = 0$ , if  $A \cap B = 0$
- Example:
  - $D1 \cup D2 = \{\text{he, likes, to, wink, and, drink}\}$
  - $D1 \cap D2 = \{\text{he, likes, to, drink}\}$
  - $jaccard(D1, D2) = \frac{4}{6} = 0.6667$

## Jaccard coefficient: Issues

- Does not consider **term frequency** (how many times a term occurs in a document)
- It treats all terms equally!
  - How about **rare terms** in a collection?  
more informative than frequent terms.
  - *He likes to drink*, shall “to” == “drink”?
- Needs more sophisticated way of **length normalization**
  - $|D1| = 3, |D2| = 1000!$
  - $D1 \rightarrow Q, D2 \rightarrow D$

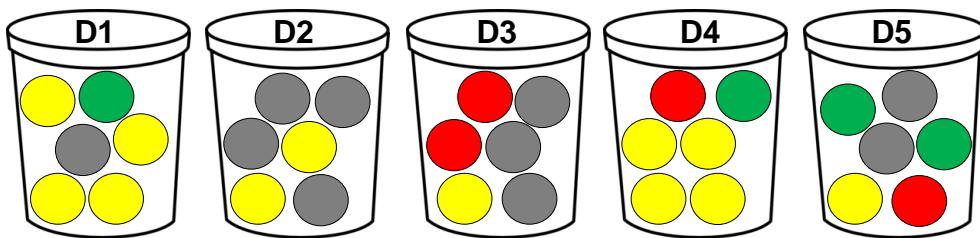
Walid Magdy, TTDS 2021/2022



7

## Should terms be treated the same?

- Collection of 5 documents (balls = terms)
- Query
- Which is the least relevant document?
- Which is the most relevant document?



Walid Magdy, TTDS 2021/2022



8

## TFIDF

- **TFIDF:**  
Term Frequency, Inverse DFrequency
- **$tf(t, d)$ :**  
number of times term  $t$  appeared in document  $d$ 
  - As  $tf(t, d) \uparrow\uparrow \rightarrow$  importance of  $t$  in  $d \uparrow\uparrow$
  - Document about IR, contains “retrieval” more than others
- **$df(t)$ :**  
number of documents term  $t$  appeared in
  - As  $df(d) \uparrow\uparrow \rightarrow$  importance if  $t$  in a collection  $\downarrow\downarrow$ 
    - “the” appears in many document  $\rightarrow$  not important
    - “FT” is not important word in financial times articles

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## DF, CF, & IDF

- **DF  $\neq$  CF** (collection frequency)
  - $cf(t)$  = total number of occurrences of term  $t$  in a collection
  - $df(t) \leq N$  ( $N$ : number of documents in a collection)
  - $cf(t)$  can be  $\geq N$
- **DF** is more commonly used in IR than **CF**
  - **CF** is still used
- **$idf(t)$** : inverse of  **$df(t)$** 
  - As  $idf(t) \uparrow\uparrow \rightarrow$  rare term  $\rightarrow$  importance  $\uparrow\uparrow$
  - **$idf(t)$**   $\rightarrow$  measure of the informativeness of  $t$

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## DF vs CF

he	drink	ink	likes	pink	think	wink
2	1	0	2	0	0	1
1	3	0	1	0	0	0
1	1	1	1	0	1	0
1	1	1	1	1	0	0
1	1	1	1	1	0	1

- ← D1: He likes to wink, he likes to drink
- ← D2: He likes to drink, and drink, and drink
- ← D3: The thing he likes to drink is ink
- ← D4: The ink he likes to drink is pink
- ← D5: He likes to wink, and drink pink ink

5 5 3 5 2 1 2 DF

6 7 3 6 2 1 2 CF

Walid Magdy, TTDS 2021/2022



11

## IDF: formula

$$idf(t) = \log_{10}\left(\frac{N}{df(t)}\right)$$

- Log scale used to dampen the effect of IDF

- Suppose  $N = 1$  million →

term	df(t)	idf(t)
calpurnia	1	6
animal	100	4
sky	1,000	3
fly	10,000	2
under	100,000	1
the	1,000,000	0

Walid Magdy, TTDS 2021/2022



12

## TFIDF term weighting

- One the best known term weights schemes in IR
  - Increases with the number of occurrences within a document
  - Increases with the rarity of the term in the collection
- Combines TF and IDF to find the weight of terms

$$w_{t.d} = \left(1 + \log_{10} tf(t, d)\right) \times \log_{10}\left(\frac{N}{df(t)}\right)$$

- For a query  $q$  and document  $d$ , retrieval score  $f(q, d)$ :

$$Score(q, d) = \sum_{t \in q \cap d} w_{t.d}$$

Walid Magdy, TTDS 2021/2022



13

## Document/Term vectors with tfidf

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	6.1	0	1	0	0
Caesar	8.59	2.54	0	1.51	0.25	0
Calpurnia	0	1.54	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	0	1.9	0.12	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

→ Vector Space Model

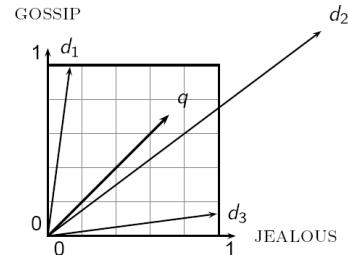
Walid Magdy, TTDS 2021/2022



14

## Vector Space Model

- Documents and Queries are presented as vectors
- Match ( $Q, D$ ) = Distance between vectors
- Example:  $Q = \text{Gossip Jealous}$
- Euclidean Distance?  
*Distance between the endpoints of the two vectors*
- Large for vectors of diff. lengths
- Take a document  $d$  and append it to itself. Call this document  $d'$ .
  - “Semantically”  $d$  and  $d'$  have the same content
  - Euclidean distance can be quite large



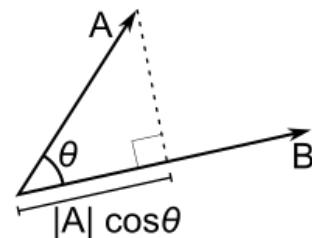
Walid Magdy, TTDS 2021/2022



15

## Angle Instead of Distance

- The angle between the two documents is 0, corresponding to maximal similarity.
- Key idea: Rank documents according to angle with query.
  - Rank documents in increasing order of the angle with query
  - Rank documents in decreasing order of cosine (query, document)
- Cosine of angle = projection of one vector on the other



Walid Magdy, TTDS 2021/2022



16

## Length Normalization

- A vector can be normalized by dividing each of its components by its length – for this we use the L<sub>2</sub> norm:

$$\|\vec{x}\|_2 = \sqrt{\sum_i x_i^2}$$

- Dividing a vector by its L<sub>2</sub> norm makes it a unit (length) vector (on surface of unit hypersphere)
- 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

## Example

- $D1 = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix} \rightarrow \|\vec{D1}\|_2 = \sqrt{1 + 9 + 4} = 3.74$
- $D1_{normalized} = \begin{bmatrix} 0.267 \\ 0.802 \\ 0.535 \end{bmatrix}$
- $D2 = \begin{bmatrix} 3 \\ 9 \\ 6 \end{bmatrix} \rightarrow \|\vec{D2}\|_2 = \sqrt{9 + 81 + 36} = 11.25$
- $D2_{normalized} = \begin{bmatrix} 0.267 \\ 0.802 \\ 0.535 \end{bmatrix}$

## Cosine “Similarity” (Query, Document)

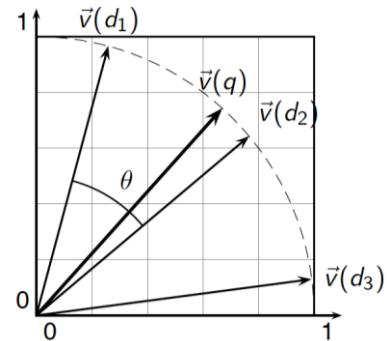
- $\vec{q}_i$  is the tf-idf weight of term  $i$  in the query
- $\vec{d}_i$  is the tf-idf weight of term  $i$  in the document

- For normalized vectors:

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d} = \sum_{i=1}^{|V|} q_i d_i$$

- For non-normalized vectors:

$$\cos(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{\|\vec{q}\| \|\vec{d}\|} = \frac{\vec{q}}{\|\vec{q}\|} \cdot \frac{\vec{d}}{\|\vec{d}\|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$



Walid Magdy, TTDS 2021/2022



19

## Algorithm

```

COSINESCORE( $q$ )
1 float Scores[N] = 0
2 float Length[N]
3 for each query term  $t$ 
4 do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5   for each pair( $d$ ,  $tf_{t,d}$ ) in postings list
6     do  $Scores[d] += w_{t,d} \times w_{t,q}$ 
7 Read the array  $Length$ 
8 for each  $d$ 
9 do  $Scores[d] = Scores[d] / Length[d]$ 
10 return Top  $K$  components of  $Scores[]$ 
```

Walid Magdy, TTDS 2021/2022



20

## TFIDF Variants

Term frequency		Document frequency	Normalization
n (natural)	$tf_{t,d}$	n (no) 1	n (none) 1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf) $\log \frac{N}{df_t}$	c (cosine) $\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$	p (prob idf) $\max\{0, \log \frac{N-df_t}{df_t}\}$	u (pivoted unique) $1/u$
b (boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$		b (byte size) $1/CharLength^\alpha, \alpha < 1$
L (log ave)	$\frac{1+\log(tf_{t,d})}{1+\log(\text{ave}_{t \in d}(tf_{t,d}))}$		

- Many search engines allow for different weightings for queries vs. documents
- SMART Notation:** use notation  $ddd.qqq$ , using the acronyms from the table
- A very standard weighting scheme is: *lnc.ltc*

Walid Magady, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

21

## For Lab and CW

Term frequency		Document frequency	Normalization
n (natural)	$tf_{t,d}$	n (no) 1	n (none) 1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf) $\log \frac{N}{df_t}$	c (cosine) $\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$	p (prob idf) $\max\{0, \log \frac{N-df_t}{df_t}\}$	u (pivoted unique) $1/u$
b (boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$		b (byte size) $1/CharLength^\alpha, \alpha < 1$
L (log ave)	$\frac{1+\log(tf_{t,d})}{1+\log(\text{ave}_{t \in d}(tf_{t,d}))}$		

“OR” operator, then:

$$Score(q, d) = \sum_{t \in q \cap d} \left( 1 + \log_{10} tf(t, d) \right) \times \log_{10} \left( \frac{N}{df(t)} \right)$$

Walid Magady, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

22

## Summary of Steps:

- Represent the query as a weighted *tf-idf* vector
- Represent each document as a weighted *tf-idf* vector
- Compute the cosine similarity score for the query vector and each document vector
- Rank documents with respect to the query by score
- Return the top K (e.g.,  $K = 10$ ) to the user

Walid Magdy, TTDS 2021/2022



23

## Retrieval Output

- For a query  $q_1$ , the output would be a list of documents ranked according to the  $\text{score}(q_1, d)$
- Possible output format:

1,	710,	0.9234
1,	213,	0.7678
1,	103,	0.6761
1,	13,	0.6556
1,	501,	0.4301

Query id                    document id                    score

Walid Magdy, TTDS 2021/2022



24

## Resources

- Text book 1: Intro to IR, Chapter 6.2 → 6.4
- Text book 2: IR in Practice, Chapter 7
- Lab 3

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Ranked Retrieval (2)

Instructor:  
**Walid Magdy**

13-Oct-2021

1

## Lecture Objectives

- Learn about Probabilistic models
  - BM25
- Learn about LM for IR



THE UNIVERSITY  
of EDINBURGH

2

## Recall: VSM & TFIDF term weighting

- Combines TF and IDF to find the weight of terms

$$w_{t,d} = \left(1 + \log_{10} tf(t, d)\right) \times \log_{10}\left(\frac{N}{df(t)}\right)$$

- For a query  $q$  and document  $d$ , retrieval score  $f(q, d)$ :

$$Score(q, d) = \sum_{t \in q \cap d} w_{t,d}$$

- TFIDF observations ***Can we do better?***

- Term appearing more in a doc gets higher weight (TF)
- First occurrence is more important (log)
- Rare terms are more important (IDF)
- Bias towards longer documents

Walid Magdy, TTDS 2021/2022



3

## IR Model

- VSM is very heuristic in nature
  - No notion of relevance is there (still works well)
  - Any weighting scheme, similarity measure can be used
    - Components not interpretable → no guide for what to try next
    - More engineering rather than theory → tweak, run, observe, tweak ...
  - Very popular, hard to beat, strong baseline
    - Easy to adapt good ideas from other models
- Probabilistic Model** of retrieval
  - Mathematical formulisation for relevant / irrelevant sets
    - Explicitly defines random variables (R,Q,D)
    - Specific about what their values are
    - State the assumptions behind each step
    - Watch out for contradictions

Walid Magdy, TTDS 2021/2022



4

## Probabilistic Models

- Concept: Uncertainty is inherent part of IR process
- Probability theory is strong foundation for representing and manipulating uncertainty
- Probability Ranking Principle (1977)



**Stephan Robertson**

Walid Magdy, TTDS 2021/2022



5

## Probability Ranking Principle

- “If a reference retrieval system’s response to each request is a ranking of the documents in the collection in order of decreasing probability of relevance to the user who submitted the request,
- where the probabilities are estimated as accurately as possible on the basis of whatever data have been made available to the system for this purpose,
- the overall effectiveness of the system to its user will be the best that is obtainable on the basis of those data.”
- Basis for most probabilistic approaches for IR

Walid Magdy, TTDS 2021/2022



6

## Formulation of PRP

- Rank docs by probability of relevance
  - $P(R|D_r_1) > P(R|D_r_2) > P(R|D_r_3) > P(R|D_r_4) > \dots$
- Estimate probability as accurate as possible
  - $P_{\text{est}}(R|D) \approx P_{\text{true}}(R|D)$
- Estimate with all possibly available data
  - $P_{\text{est}}(R | \text{doc, session, context, user profile, ...})$
- Best possible accuracy can be achieved with that data
  - → the perfect IR system
  - Is it really doable?
- **How to estimate the probability of relevance?**

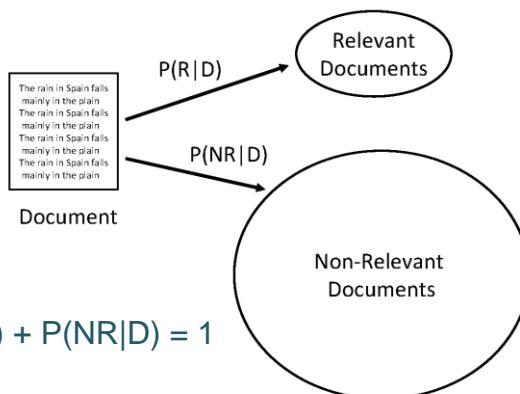
Walid Magdy, TTDS 2021/2022



7

## PRP Concept

- Imagine IR as a classification problem



- Document  $D$  is **relevant** if  $P(R|D) > P(NR|D)$

Walid Magdy, TTDS 2021/2022



8

## Probability of Relevance

- What is  $P_{\text{true}}(\text{rel} | \text{doc, query, session, user, ...})$ ?
  - Isn't relevance just the user's opinion?
  - User decides relevant or not, what is the "probability" thing?
- Search algorithm cannot look into your head (yet!)
  - Relevance depends on factors that algorithm cannot observe
    - SIGIR 2016 best paper award: *Understanding Information Need: an fMRI Study*
- Different users may disagree on relevance of the same doc
  - Even similar users, doing the same task, in the same context
- $P_{\text{true}}(\text{rel} | Q, D)$ :
  - Proportion of all unseen users / context / tasks for which D would have judged relevant to Q
- Similar to:  $P(\text{die}=6 | \text{even and not square})$

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## Okapi BM25 Model

- Based on the probabilistic model
  - A document D is relevant if  $P(R=1|D) > P(R=0|D)$
- Extension to the "binary independence model"
  - **Binary features:** Document represented by a vector of binary features indicating term occurrence
  - Assume **term independence** (Naïve Bayes assumption)  
→ BOW trick
- In 1995, Stephan Robertson with his group came up with the **BM25** Formula as part of the **Okapi** project.
- It outperformed all other systems in TREC
- Popular and effective ranking algorithm

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## Okapi BM25 Ranking Function

- Let  $L_d$  be the number of terms in document  $d$
- Let  $\bar{L}$  be the average number of terms in a document

$$w_{t,d} = \frac{tf_{t,d}}{k \cdot \frac{L_d}{\bar{L}} + tf_{t,d} + 0.5} \times \log_{10} \left( \frac{N - df_t + 0.5}{df_t + 0.5} \right)$$

- Best practices:  $k=1.5$

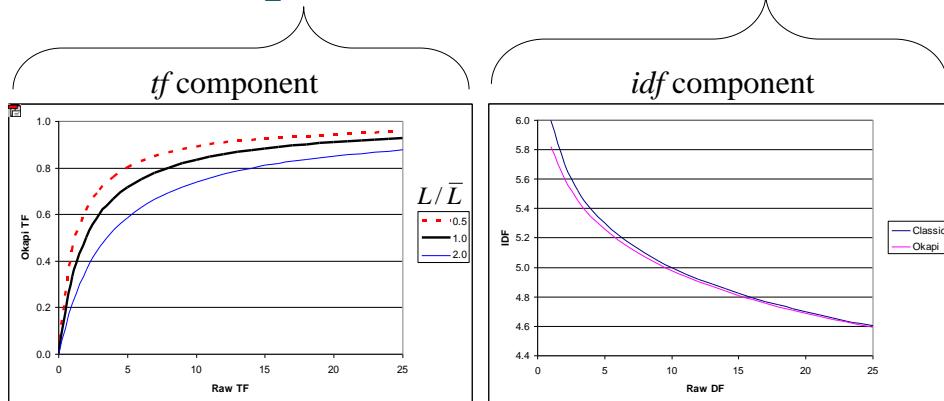
Walid Magdy, TTDS 2021/2022



11

## Okapi BM25 Ranking Function

$$w_{t,d} = \frac{tf_{t,d}}{1.5 \frac{L_d}{\bar{L}} + tf_{t,d} + 0.5} \times \log_{10} \left( \frac{N - df_t + 0.5}{df_t + 0.5} \right)$$



Walid Magdy, TTDS 2021/2022



12

## Probabilistic Model in IR

- Focuses on the probability of relevance of docs
- Could be mathematically proved
- Different ways to apply it
- BM25 is the most common formula for it
- What other models could be still used in IR?

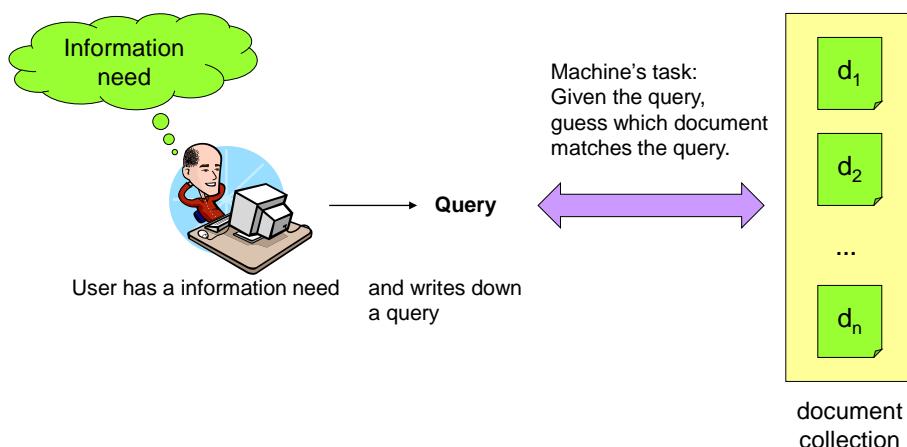
Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## “Noisy-Channel” Model of IR



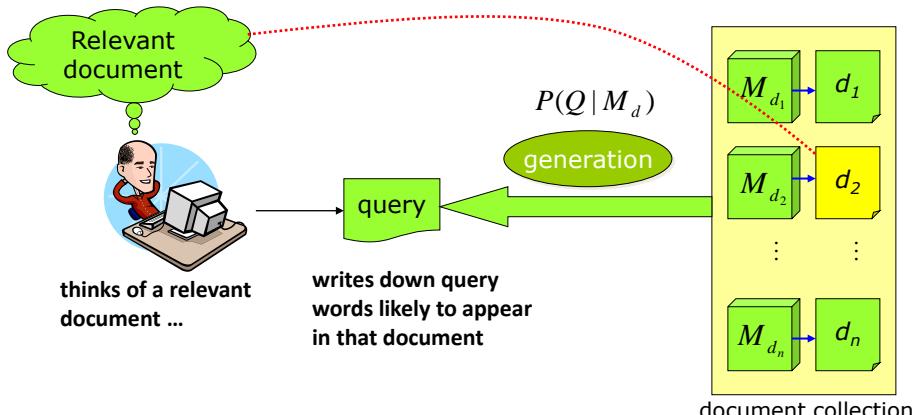
Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

14

## IR based on Language Model (LM)



- **The LM approach directly exploits that idea!**
- a document is a good match to a query if the document model is likely to generate the query

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Concept

- Coming up with good queries?
  - Think of words that would likely appear in a relevant doc
  - Use those words as the query
- The language modeling approach to IR directly models that idea
  - a document is a good match to a query if the document model is likely to generate the query
    - happens if the document contains the query words often.
- Build a probabilistic language model  $M_d$  from each document  $d$
- Rank documents based on the probability of the model generating the query:  $P(q|M_d)$ .

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

16

## Language Model (LM)

- A language model is a probability distribution over strings drawn from some vocabulary
- A topic in a document or query can be represented as a language model
  - i.e., words that tend to occur often when discussing a topic will have high probabilities in the corresponding language model

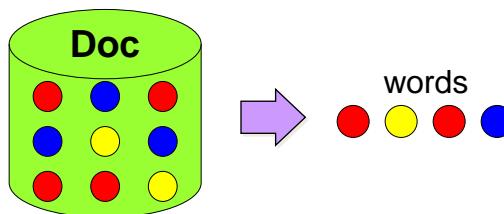
Walid Magdy, TTDS 2021/2022



17

## Unigram LM

- Terms are randomly drawn from a document (with replacement)



$$\begin{aligned} P(\bullet \bullet \bullet \bullet) &= P(\bullet) \times P(\bullet) \times P(\bullet) \times P(\bullet) \\ &= (4/9) \times (2/9) \times (4/9) \times (3/9) \end{aligned}$$

Walid Magdy, TTDS 2021/2022



18

## Example

$w$	$P(w q_1)$	$w$	$P(w q_1)$
STOP	0.2	toad	0.01
the	0.2	said	0.03
a	0.1	likes	0.02
frog	0.01	that	0.04
	...		...

- This is a one-state probabilistic finite-state automaton – a unigram language model.
  - $S = \text{"frog said that toad likes frog STOP"}$
- $$P(S) = 0.01 \times 0.03 \times 0.04 \times 0.01 \times 0.02 \times 0.01 \times 0.02 \\ = 0.00000000000048$$

Walid Magdy, TTDS 2021/2022



## Comparing LMs

- $M_{d1}$   
LM generated from Doc 1
- $M_{d2}$   
LM generated from Doc 2
- Try to generate sentence  $S$  from  $M_{d1}$  &  $M_{d2}$

Model $M_{d1}$	
$P(w)$	$w$
0.2	the
0.0001	yon
0.01	class
0.0005	maiden
0.0003	sayst
0.0001	pleaseth
...	

Model $M_{d2}$	
$P(w)$	$w$
0.2	the
0.1	yon
0.001	class
0.01	maiden
0.03	sayst
0.02	pleaseth
...	

text:	the	class	pleaseth	yon	maiden	$P(S)$
$M_{d1}:$	0.2	0.01	0.0001	0.0001	0.0005	0.0000000000000001
$M_{d2}:$	0.2	0.001	0.02	0.1	0.01	0.000000004

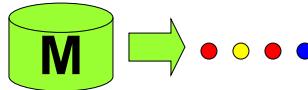
$$P(\text{text}|M_{d2}) > P(\text{text}|M_{d1})$$

Walid Magdy, TTDS 2021/2022



## Stochastic Language Models

- A statistical model for generating text
  - Probability distribution over strings in a given language



$$\begin{aligned}
 P(\bullet \bullet \bullet \bullet | M) &= P(\bullet | M) \\
 &\quad P(\bullet | M, \bullet) \\
 &\quad P(\bullet | M, \bullet \bullet) \\
 &\quad P(\bullet | M, \bullet \bullet \bullet)
 \end{aligned}$$

Walid Magdy, TTDS 2021/2022



21

## Unigram and Higher-order LM

$$\begin{aligned}
 P(\bullet \bullet \bullet \bullet) \\
 &= P(\bullet) P(\bullet | \bullet) P(\bullet | \bullet \bullet) P(\bullet | \bullet \bullet \bullet)
 \end{aligned}$$

- **Unigram Language Models**

$$P(\bullet) P(\bullet) P(\bullet) P(\bullet)$$

- **Bigram** (generally,  $n$ -gram) Language Models

$$P(\bullet) P(\bullet | \bullet) P(\bullet | \bullet) P(\bullet | \bullet)$$

Walid Magdy, TTDS 2021/2022



22

## LM in IR

- Each document is treated as basis for a LM.
- Given a query  $q$ , rank documents based on  $P(d|q)$

$$P(d|q) = \frac{P(q|d)P(d)}{P(q)}$$

- $P(q)$  is the same for all documents → ignore
- $P(d)$  is the prior – often treated as the same for all  $d$ 
  - But we can give a prior to “high-quality” documents, e.g., those with high PageRank (later to be discussed).
  - $P(q|d)$  is the probability of  $q$  given  $d$ .
- So to rank documents according to relevance to  $q$ , ranking according to  $P(q|d)$  and  $P(d|q)$  is equivalent

## LM in IR: Basic idea

- We attempt to model the query generation process.
- Then we rank documents by the probability that a query would be observed as a random sample from the respective document model.
- That is, we rank according to  $P(q|d)$ .

## $P(q|d)$

### Query Likelihood Model

- We will make the conditional independence assumption.

$$P(q|M_d) = P(\langle t_1, \dots, t_{|q|} \rangle | M_d) = \prod_{1 \leq k \leq |q|} P(t_k | M_d)$$

$|q|$ : length of  $q$ ;  $t_k$ : token occurring at position  $k$  in  $q$

- This is equivalent to:

$$P(q|M_d) = \prod_{\text{each term } t \text{ in } q} P(t | M_d)^{tf_{t,q}}$$

$tf_{t,q}$ : term frequency (# occurrences) of  $t$  in  $q$

- Multinomial model (omitting constant factor)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

25

## Parameter estimation

- Probability of a term  $t$  in a LM  $M_d$  using Maximum Likelihood Estimation (MLE)

$$P(t|M_d) = \frac{tf_{t,d}}{|d|}$$

$|d|$ : length of  $d$ ;

$tf_{t,d}$ : # occurrences of  $t$  in  $d$

- Probability of a query  $q$  to be noticed in a LM  $M_d$ :

$$P(q|M_d) = \prod_{\forall t \in q} \left( \frac{tf_{t,d}}{|d|} \right)^{tf_{t,q}}$$

Walid Magdy, TTDS 2021/2022

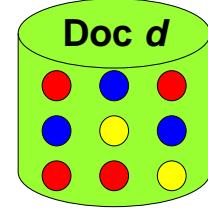


THE UNIVERSITY  
of EDINBURGH

26

## Example

$$\begin{aligned} P(\text{red}, \text{yellow}, \text{red}, \text{blue}) &= P(\text{red})^2 \times P(\text{yellow}) \times P(\text{blue}) \\ &= (4/9)^2 \times (2/9) \times (3/9) = 0.0146 \\ P(\text{red}, \text{yellow}, \text{green}, \text{blue}) & \end{aligned}$$



- Is that fair?
  - In VSM,  $S(Q,D)$  was summation, works more like OR in Boolean search. Missing one term reduces score only
  - In language model,  $S(Q,D)$  is  $P(Q|D)$  → Multiplication of probabilities → missing one term makes score = 0
  - Is there a better way to handle unseen terms?

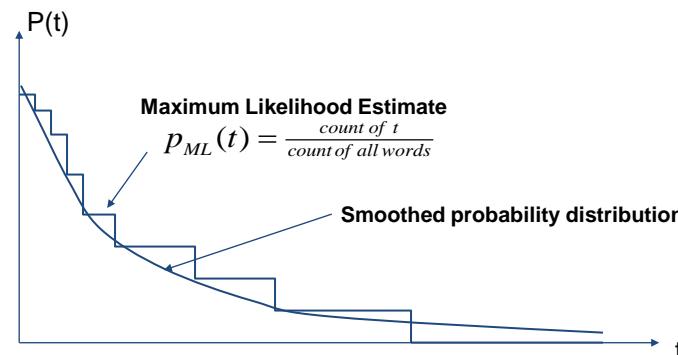
Walid Magdy, TTDS 2021/2022



27

## Smoothing

- Problem: Zero frequency
- Solution: “Smooth” terms probability



Walid Magdy, TTDS 2021/2022



28

## Smoothing

- Document texts are a sample from the language model
- Missing words should not have zero probability of occurring
- A missing term is possible (even though it didn't occur)
  - but no more likely than would be expected by chance in the collection.
- A technique for estimating probabilities for missing (or unseen) words
  - Overcomes data-sparsity problem
  - lower (or discount) the probability estimates for words that are seen in the document text
  - assign that “left-over” probability to the estimates for the words that are not seen in the text (and also on the seen ones)

Walid Magdy, TTDS 2021/2022



29

## Mixture Model

$$P(t|d) = \lambda P(t|M_d) + (1 - \lambda) P(t|M_c)$$

- Mixes the probability from the document with the general collection frequency of the word.
- Estimate for unseen words is  $(1-\lambda) P(t|M_c)$ 
  - Based on collection language model (background LM)
  - $P(t|M_c)$  is the probability for query word  $t$  in the collection language model for collection C (background probability)
  - $\lambda$  is a parameter controlling probability for unseen words
- Estimate for observed words is

$$\lambda P(t|M_d) + (1-\lambda) P(t|M_c)$$

CF

Walid Magdy, TTDS 2021/2022



30

## Jelinek-Mercer Smoothing

$$P(t|d) = \lambda P(t|M_d) + (1 - \lambda)P(t|M_c)$$

- **High value of  $\lambda$ :** “conjunctive-like” search – tends to retrieve documents containing all query words.
- **Low value of  $\lambda$ :** more disjunctive, suitable for long queries
- Correctly setting  $\lambda$  is important for good performance.
- Final Ranking function:

$$P(q|M_d) \propto \prod_{1 \leq k \leq |q|} (\lambda \cdot P(t_k|M_d) + (1 - \lambda) \cdot P(t_k|M_c))$$

Walid Magdy, TTDS 2021/2022



## Example

- **Collection:**  $d_1$  and  $d_2$
- $d_1$ : “Jackson was one of the most talented entertainers of all time”
- $d_2$ : “Michael Jackson anointed himself King of Pop”
- **Query  $q$ :** Michael Jackson
- Use mixture model with  $\lambda = 1/2$
- $P(q|d_1) = \overbrace{[(0/11 + 1/18)/2]} \cdot \overbrace{[(1/11 + 2/18)/2]} \approx 0.003$
- $P(q|d_2) = \overbrace{[(1/7 + 1/18)/2]} \cdot \overbrace{[(1/7 + 2/18)/2]} \approx 0.013$
- Ranking:  $d_2 > d_1$

Walid Magdy, TTDS 2021/2022



## Notes on Query Likelihood Model

- It has similar effectiveness to BM25
- With more sophisticated techniques, it outperforms BM25
  - Topic models
- There are several alternative smoothing techniques
  - That was just an example

Walid Magdy, TTDS 2021/2022



33

## n-grams LMs

- Unigram language model
  - probability distribution over the words in a language
    - associates a probability of occurrence with every word
  - generation of text consists of pulling words out of a “bucket” according to the probability distribution and replacing them
- N-gram language model
  - some applications use bigram and trigram language models where probabilities depend on previous words
  - predicts a word based on the previous n-1 words

Walid Magdy, TTDS 2021/2022



34

## LMs for IR: 3 possibilities

- Probability of generating the query text from a document language model
- Probability of generating the document text from a query language model
- Comparing the language models representing the query and document topics

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

35

## Summary

- **Three ways to model IR**
- VSM  
How query vector aligns with document vector?
- Probabilistic Model  
What is the relevance probability of document D given query Q?
- LM  
How likely is it possible to observe/generate sequence of terms Q in a language model of document D?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

36

## Resources

- Text book 1: Intro to IR, Chapter 12
- Text book 2: IR in Practice, Chapter 7.2, 7.3
- Readings:
  - Robertson, Stephen E., et al.  
*"Okapi at TREC-3."*  
Nist Special Publication Sp 109 (1995): 109.
  - J. Ponte and W. B. Croft.  
*A language modeling approach to information retrieval.*  
In Proceedings on the 21st annual international ACM SIGIR conference, pages 275–281, 1998

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# IR Evaluation

Instructor:  
**Walid Magdy**

20-Oct-2021

1

## Pre-lecture

- Sorry for having last week's lectures recorded
- Discussion session about lectures 7 & 8 available on Learn
- Test collection for CW1 to be released next week
- New CW1 deadline: 31 Oct 2021
- No new lab this week (support to continue for previous labs)



## Lecture Objectives

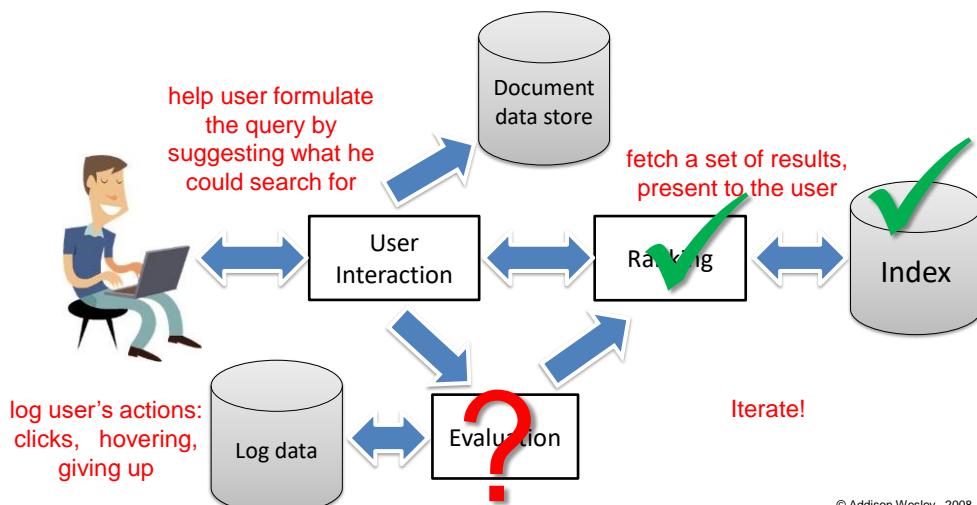
- Learn about how to evaluate IR
  - Evaluation measures
  - P, R, F
  - MAP
  - nDCG
- Implement: (as part of CW2)
  - P, R
  - MAP
  - nDCG

Walid Magdy, TTDS 2021/2022



3

## Search Process



Walid Magdy, TTDS 2021/2022



4

## IR as an Experimental Science!

- Formulate a research question: the hypothesis
- Design an experiment to answer the question
- Perform the experiment
  - Compare with a baseline “control”
- Does the experiment answer the question?
  - Are the results significant? Or is it just luck?
- Report the results!
- Iterate ...
- e.g. stemming improves results? (university → univers)

Walid Magdy, TTDS 2021/2022



5

## Lab 3 output      Is that a good performance?

1, 65, 4.8040	2, 3549, 7.0396	3, 3354, 4.6113
1, 3533, 4.7264	2, 305, 6.8394	3, 3345, 4.5087
1, 3562, 3.5454	2, 288, 6.6742	3, 268, 3.6606
1, 3608, 3.4910	2, 223, 6.1252	3, 328, 3.4825
1, 141, 3.3262	2, 219, 4.8626	3, 21, 3.3984
1, 361, 3.3262	2, 3762, 4.8626	3, 304, 3.3722
1, 92, 3.2311	2, 3663, 4.5415	3, 313, 3.3436
1, 3829, 3.1818	2, 3766, 3.9924	3, 3790, 3.1796
1, 3420, 3.1273	2, 188, 3.8844	3, 55, 3.0462
1, 3734, 3.0561	2, 3360, 3.0988	3, 217, 2.8492
1, 3387, 2.9626	2, 3408, 3.0315	3, 361, 2.8348
1, 3599, 2.9626	2, 3390, 2.8498	3, 3789, 2.7158

Walid Magdy, TTDS 2021/2022



6

## Configure your system

- **About the system:**
  - Stopping? Tokenise? Stemming? n-gram char?
  - Use synonyms improve retrieval performance?
- **Corresponding experiment?**
  - Run your search for a set of queries with each setup and find which one will achieve the best performance
- **About the user:**
  - Is letting users weight search terms a good idea?
- **Corresponding experiment?**
  - Build two different interfaces, one with term weighting functionality, and one without; run a user study

## Types of Evaluation Strategies

- **System-centered studies:**
  - Given documents, queries, and relevance judgments
  - Try several variations of the system
  - Measure which system returns the “best” hit list
  - Laboratory experiment
- **User-centered studies**
  - Given several users, and at least two retrieval systems
  - Have each user try the same task on both systems
  - Measure which system works the “best”

## Importance of Evaluation

- The ability to measure differences underlies experimental science
  - How well do our systems work?
  - Is A better than B?
  - Is it really?
  - Under what conditions?
- Evaluation drives what to research
  - Identify techniques that work and don't work

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

9

## The 3-dimensions of Evaluation

- **Effectiveness**
  - How “good” are the documents that are returned?
  - System only, human + system
- **Efficiency**
  - Retrieval time, indexing time, index size
- **Usability**
  - Learnability, flexibility
  - Novice vs. expert users

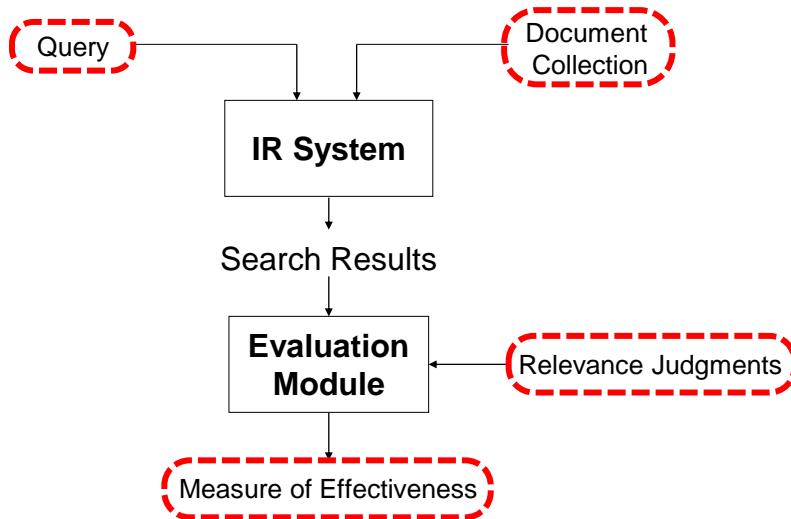
*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

10

## Cranfield Paradigm (Lab setting)



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Reusable IR Test Collection

- **Collection of Documents**
  - Should be “representative” to a given IR task
  - Things to consider: size, sources, genre, topics, ...
- **Sample of information need**
  - Should be “randomized” and “representative”
  - Usually formalized topic statements (query + description)
- **Known relevance judgments**
  - Assessed by humans, for each topic-document pair
  - Binary/Graded
- **Evaluation measure**

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

12

## Good Effectiveness Measures

- Should capture some aspect of what the user wants
  - IR → Do the results satisfy user's information need?
- Should be easily replicated by other researchers
- Should be easily comparable
  - Optimally, expressed as a single number
  - Curves and multiple numbers are still accepted, but single numbers are much easier for comparison
- Should have predictive value for other situations
  - What happens with different queries on a different document collection?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Set Based Measures

- Assuming IR system returns sets of retrieved results without ranking
- Suitable with Boolean Search
- No certain number of results per query

Walid Magdy, TTDS 2021/2022

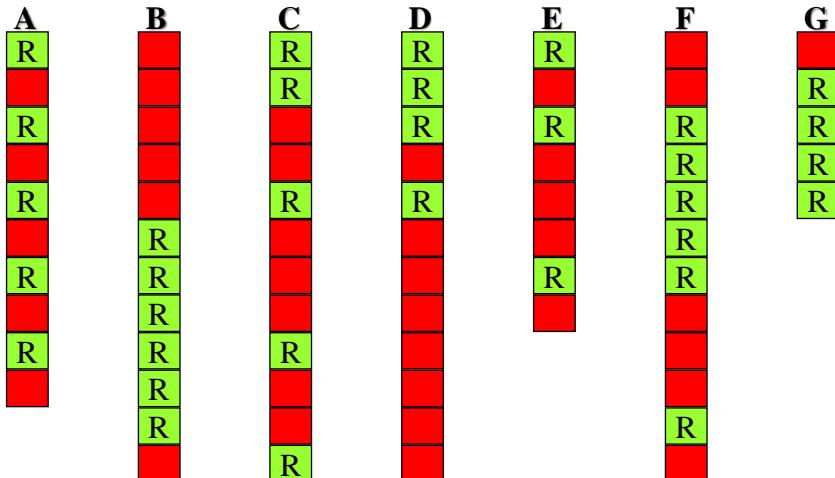


THE UNIVERSITY  
of EDINBURGH

14

## Which looks the best IR system?

- For query Q, collection has **8** relevant documents:



Walid Magdy, TTDS 2021/2022



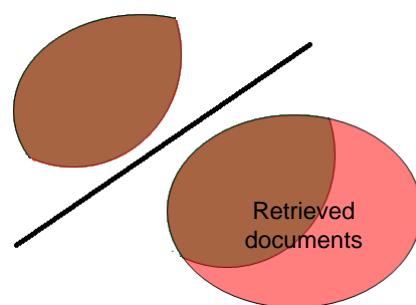
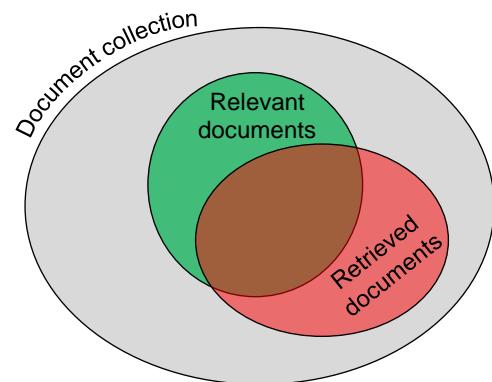
15

## Precision and Recall

- Precision:**

What fraction of these retrieved docs are relevant?

$$P = \frac{rel \cap ret}{retrieved} = \frac{TP}{TP + FP}$$



	relevant		irrelevant	
	retrieved	not retrieved	FP	TN
	retrieved	not retrieved	TP	FN

Walid Magdy, TTDS 2021/2022



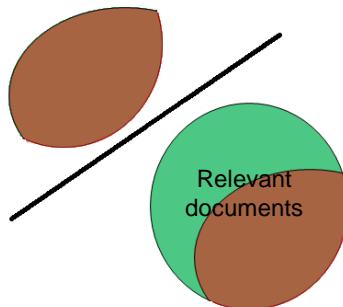
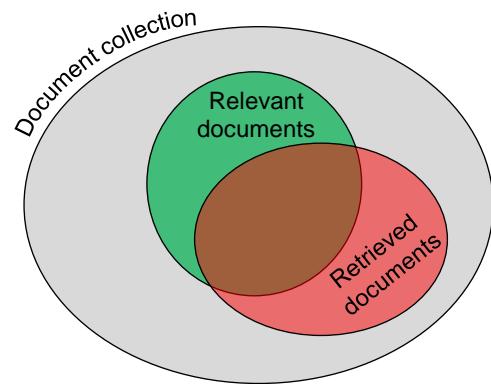
16

## Precision and Recall

- Recall:**

What fraction of the relevant docs were retrieved?

$$R = \frac{rel \cap ret}{relevant} = \frac{TP}{TP + FN}$$



	relevant		not retrieved
	retrieved	not retrieved	
irrelevant	FP	TN	
	TP	FN	

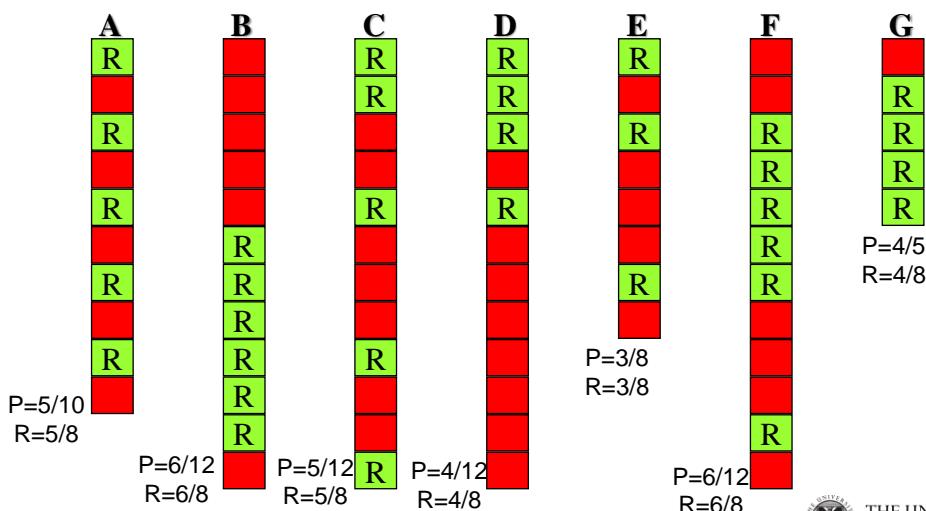
Walid Magdy, TTDS 2021/2022



17

## Which looks the best IR system?

- For query Q, collection has **8 relevant documents**:



Walid Magdy, TTDS 2021/2022



18

## Trade-off between P & R

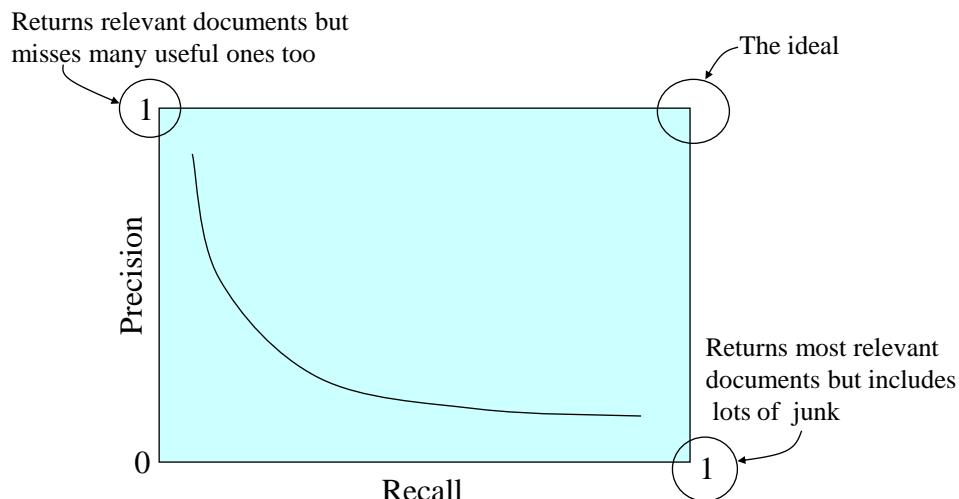
- Precision: The ability to retrieve top-ranked docs that are mostly relevant.
- Recall: The ability of the search to find all of the relevant items in the corpus.
- Retrieve more docs:
  - Higher chance to find all relevant docs →  $R \uparrow\uparrow$
  - Higher chance to find more irrelevant docs →  $P \downarrow\downarrow$

Walid Magdy, TTDS 2021/2022



19

## Trade-off between P & R



Walid Magdy, TTDS 2021/2022



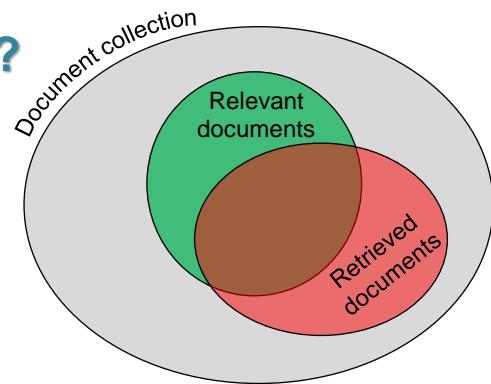
20

## What about Accuracy?

- **Accuracy:**

What fraction of docs was classified correctly?

$$A = \frac{TP + TN}{TP + FP + TN + FN}$$



	relevant	irrelevant
retrieved	FP	TN
not retrieved	TP	FN

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

21

## One Measure? F-measure

$$F1 = \frac{2 \cdot P \cdot R}{P + R}$$

$$F_\beta = \frac{(\beta^2 + 1)P \cdot R}{\beta^2 P + R}$$

- Harmonic mean of recall and precision
  - Emphasizes the importance of small values, whereas the arithmetic mean is affected more by outliers that are unusually large
- Beta ( $\beta$ ) controls relative importance of P and R
  - $\beta = 1$ , precision and recall equally important  $\rightarrow F1$
  - $\beta = 5$ , recall five times more important than precision

Walid Magdy, TTDS 2021/2022

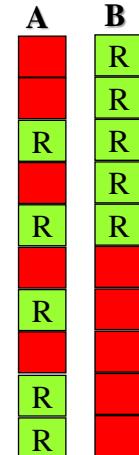


THE UNIVERSITY  
of EDINBURGH

22

## Rank-based IR measures

- Consider systems A & B
  - Both retrieved 10 docs, only 5 are relevant
  - P, R & F are the same for both systems
    - Should their performances considered equal?
- Ranked IR requires taking “ranks” into consideration!
- How to do that?



Walid Magdy, TTDS 2021/2022



23

## Which is the best ranked list?

- For query Q, collection has **8** relevant documents:

	A	B	C	D	E	F	G
1	R	R	R	R	R	R	R
2	R	R	R	R	R	R	R
3	R	R	R	R	R	R	R
4	R	R	R	R	R	R	R
5	R	R	R	R	R	R	R
6	R	R	R	R	R	R	R
7	R	R	R	R	R	R	R
8	R	R	R	R	R	R	R
9	R	R	R	R	R	R	R
10	R	R	R	R	R	R	R
11	R	R	R	R	R	R	R
12	R	R	R	R	R	R	R

Walid Magdy, TTDS 2021/2022



24

## Precision @ K

- $k$  (a fixed number of documents)
- Have a cut-off on the ranked list at rank  $k$ , then calculate precision!
- Perhaps appropriate for most of web search: most people only check the top  $k$  results
- But: averages badly, Why?

Walid Magdy, TTDS 2021/2022



25

## P@5

- For query **Q**, collection has **8 relevant documents**:

A	B	C	D	E	F	G
1 R	1 R	1 R	1 R	1 R	1 R	1 R
2 R	2 R	2 R	2 R	2 R	2 R	2 R
3 R	3 R	3 R	3 R	3 R	3 R	3 R
4 R	4 R	4 R	4 R	4 R	4 R	4 R
5 R	5 R	5 R	5 R	5 R	5 R	5 R
6 R	6 R	6 R	6 R	6 R	6 R	6 R
7 R	7 R	7 R	7 R	7 R	7 R	7 R
8 R	8 R	8 R	8 R	8 R	8 R	8 R
9 R	9 R	9 R	9 R	9 R	9 R	9 R
10 R						
11 R						
12 R						

Walid Magdy, TTDS 2021/2022



26

## R-Precision

- For a query with known  $r$  relevant documents  
→ R-precision is the precision at rank  $r$  ( $P@r$ )
- $r$  is different from one query to another
- Concept:  
It examines the ideal case: getting all relevant documents in the top ranks
- Is it realistic?

Walid Magdy, TTDS 2021/2022



27

## R-Precision

- For query **Q**, collection has **8 relevant documents**:

	A	B	C	D	E	F	G
1	R		R	R	R		
2		R		R		R	
3	R			R			R
4		R			R		R
5	R		R			R	
6		R		R		R	
7	R				R		R
8		R				R	
9		R	R				
10			R				
11				R			
12					R		

Walid Magdy, TTDS 2021/2022



28

## User Satisfaction??

- It is assumed that users needs to find relevant docs at the highest possible ranks  
→ Precision is a good measure
- But, user would cut-off (stop inspecting results) at some point, say rank  $x$   
→  $P@x$
- What is the optimal  $x$ ?  
When you think a user can stop?



## When a user can stop?

- IR objective: “satisfy user information need”
- Assumption: a user will stop once his/her information need is satisfied
- How? user will keep looking for relevant docs in the ranked list, read them, then stop once he/she feels satisfied
- $P@x \rightarrow x$  can be any rank where a relevant document appeared (*assume uniform distribution*)
- **What about calculating the averages over all  $x$ 's?**
  - every time you find relevant doc, calculate  $P@x$ , then take the average at the end



## Average Precision (AP)

**Q<sub>1</sub>**  
(has 4 rel. docs)

1	R	1/1=1.00
2	R	2/2=1.00
3		
4		
5	R	3/5=0.60
6		
7		
8		
9	R	4/9=0.44
10		

$$\text{AP} = \frac{3.04}{4} / 4 \\ = \mathbf{0.76}$$

**Q<sub>2</sub>**  
(has 3 rel. docs)

1		
2		
3	R	1/3=0.33
4		
5		
6		
7	R	2/7=0.29
8		

$$\text{AP} = \frac{0.62}{3} / 3 \\ = \mathbf{0.207}$$

**Q<sub>3</sub>**  
(has 7 rel. docs)

1		
2	R	1/2=0.50
3		
4		
5	R	2/5=0.40
6		
7		
8	R	3/8=0.375
9		

$$\text{AP} = \frac{1.275}{7} / 7 \\ = \mathbf{0.182}$$

Walid Magdy, TTDS 2021/2022



31

## Mean Average Precision (MAP)

**Q<sub>1</sub>**  
(has 4 rel. docs)

1	R	1/1=1.00
2	R	2/2=1.00
3		
4		
5	R	3/5=0.60
6		
7		
8		
9	R	4/9=0.44
10		

$$\text{AP} = 0.76$$

**Q<sub>2</sub>**  
(has 3 rel. docs)

1		
2		
3	R	1/3=0.33
4		
5		
6		
7	R	2/7=0.29
8		

$$\text{AP} = 0.207$$

**Q<sub>3</sub>**  
(has 7 rel. docs)

1		
2	R	1/2=0.50
3		
4		
5	R	2/5=0.40
6		
7		
8	R	3/8=0.375
9		

$$\text{AP} = 0.182$$

$$\text{MAP} = (0.76 + 0.207 + 0.182) / 3 = \mathbf{0.383}$$

Walid Magdy, TTDS 2021/2022



32

## AP & MAP

$$AP = \frac{1}{r} \sum_{k=1}^n P(k) \times rel(k)$$

where,  $r$ : number of relevant docs for a given query

$n$ : number of documents retrieved

$P(k)$  precision @  $k$

$rel(k)$ : 1 if retrieved doc @  $k$  is relevant, 0 otherwise.

$$MAP = \frac{1}{Q} \sum_{q=1}^Q AP(q)$$

where,  $Q$ : number of queries in the test collection



## AP/MAP

$$AP = \frac{1}{r} \sum_{k=1}^n P(k) \times rel(k)$$

- A mix between precision and recall
- Highly focus on finding relevant document as early as possible
- When  $r=1 \rightarrow MAP = MRR$  (mean reciprocal rank  $\frac{1}{k}$ )
- MAP is the most commonly used evaluation metric for most IR search tasks
- Uses binary relevance:  $rel = 0/1$



## Binary vs. Graded Relevance

- Some docs are more relevant to a query than other relevant ones!
  - We need non-binary relevance
- Binary Relevance:
  - Relevant 1
  - Irrelevant 0
- Graded Relevance:
  - Perfect 4
  - Excellent 3
  - Good 2
  - Fair 1
  - Bad 0

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

35

## Binary vs. Graded Relevance

- Two assumptions:
  - Highly relevant documents are more useful than marginally relevant
  - The lower the ranked position of a relevant document, the less useful it is for the user, since it is less likely to be examined
- Discounted Cumulative Gain (DCG)
  - Uses graded relevance as a measure of the usefulness
  - The most popular for evaluating web search

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

36

## Discounted Cumulative Gain (DCG)

- Gain is accumulated starting at the top of the ranking and may be reduced (discounted) at lower ranks
- Users care more about high-ranked documents, so we discount results by  $1/\log_2(\text{rank})$ 
  - the discount at rank 4 is 1/2, and at rank 8 is 1/3
- DCG<sub>k</sub> is the total gain accumulated at a particular rank  $k$  (sum of DG up to rank  $k$ ):

$$DCG_k = \text{rel}_1 + \sum_{i=2}^k \frac{\text{rel}_i}{\log_2(i)}$$

0, 1, 2, 3, ...  
 (graded)

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

37

## DCG

k	G	DG	DCG@k
1	3	3	3
2	2	2	5
3	3	1.89	6.89
4	0	0	6.89
5	0	0	6.89
6	1	0.39	7.28
7	2	0.71	7.99
8	2	0.67	8.66
9	3	0.95	9.61
10	0	0	9.61

Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

38

## Normalized DCG (nDCG)

- DCG numbers are averaged across a set of queries at specific rank values (DCG@ $k$ )
  - e.g., DCG at rank 5 is 6.89 and at rank 10 is 9.61
  - Can be any positive real number!
- DCG values are often normalized by comparing the DCG at each rank with the DCG value for the perfect ranking
  - makes averaging easier for queries with different numbers of relevant documents
- $nDCG@k = DCG@k / iDCG@k$  (divide actual by ideal)
- $nDCG \leq 1$  at any rank position
- To compare DCGs, normalize values so that a ideal ranking would have a normalized DCG of 1.0

Walid Magdy, TTDS 2021/2022



39

## nDCG

$k$	$G$	DG	DCG@ $k$	iG	iDG	iDCG@ $k$	nDCG@ $k$
1	3	3	3	3	3.00	3	1.00
2	2	2	5	3	3.00	6	0.83
3	3	1.89	6.89	3	1.89	7.89	0.87
4	0	0	6.89	2	1.00	8.89	0.78
5	0	0	6.89	2	0.86	9.75	0.71
6	1	0.39	7.28	2	0.77	10.52	0.69
7	2	0.71	7.99	1	0.36	10.88	0.73
8	2	0.67	8.66	0	0.00	10.88	0.80
9	3	0.95	9.61	0	0.00	10.88	0.88
10	0	0	9.61	0	0.00	10.88	0.88

Walid Magdy, TTDS 2021/2022



40

## Summary:

- IR test collection:
  - Document collection
  - Query set
  - Relevant judgements
  - IR measures
- IR measures:
  - R, P, F → not commonly used
  - P@k, R-precision → used sometimes
  - MAP → the most used IR measure
  - nDGC → the most used measure for web search

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

41

## Resources

- Text book 1: Intro to IR, Chapter 8
- Text book 2: IR in Practice, Chapter 8

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

42



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# IR Evaluation (2)

Instructor:  
**Walid Magdy**

20-Oct-2021

1

## Lecture Objectives

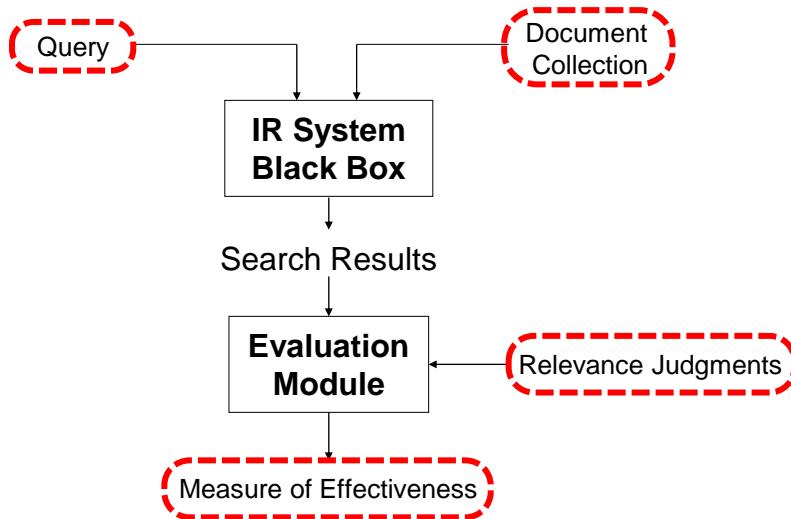
- Learn about how to evaluate IR
  - How to create a test collection?
  - Topic vs. query
  - Relevance judgements
  - Pooling



THE UNIVERSITY  
of EDINBURGH

2

## Cranfield Paradigm



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

3

## Reusable IR Test Collection

- **Collection of Documents**
  - Should be “representative” to a given IR task
  - Things to consider: size, sources, genre, topics, ...
- **Sample of information need**
  - Should be “randomized” and “representative”
  - Usually formalized topic statements (query + description)
- **Known relevance judgments**
  - Assessed by humans, for each topic-document pair
  - Binary/Graded
- **Evaluation measure**

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

4

## Where Do Test Collections Come From?

- For web search, companies apply their own studies to assess the performance of their search engine.
- Web-search performance is monitored by:
  - Traffic
  - User clicks and session logs
  - Labelling results for selected users' queries
- For other search tasks:
  - Someone goes out and builds them (expensive)
  - As the by-product of large scale evaluations
- IR Evaluation Campaigns are created for this reason

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

5

## IR Evaluation Campaigns

- IR test collections are provided for scientific communities to develop best IR methods
- Collections and queries are provided, relevance judgements are built during the campaign
- TREC = Text REtrieval Conferences
  - Main IR eval campaign. Sponsored by NIST (US gov)
  - Series of annual evaluations, started in 1992
  - Organized into “tracks”
- Other evaluation campaigns
  - CLEF: European version (since 2000)
  - NTCIR: Asian version (since 1999)
  - FIRE: Indian version (since 2008)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

6

## TREC Task

- It is a task for search a set of documents of given genre and domain.
- TREC (or other IR eval campaigns) are formed of a set of tracks, each track has a set of search tasks.
- Example
  - TREC Medical track
  - TREC Legal track → CLEF-IP track → NTCIR patent mining track
  - TREC Microblog track
  - Different CLIR tracks in all campaigns

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## TREC Collection

- 100's of collections were released in the different evaluation campaigns covering most of the domains in life
- A set of hundreds of thousands of docs
  - 1B in case of web search (TREC ClueWeb09)
- The typical format:

```
<DOC>
<DOCNO> 1234 </DOCNO>
<TEXT>
Multilines of plain text of the document
</TEXT>
</DOC>
```

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## TREC Topic

- Query sets are provided for each collection. Generated by experts and is associated with additional details. It is called **Topics**, and contains:
  - Query: the query text
  - Description: description of what is meant by the query
  - Narrative: what should be considered relevant

```

<num>189</num>
<title>Health and Computer Terminals</title>
<desc>Is it hazardous to the health of individuals to work with computer terminals on a daily basis?</desc>
<narr>Relevant documents would contain any information that expands on any physical disorder/problems that may be associated with the daily working with computer terminals. Such things as carpel tunnel, cataracts, and fatigue have been said to be associated, but how widespread are these or other problems and what is being done to alleviate any health problems</narr>
```

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## Relevance Judgements

- For each topic, set of relevant docs is required to be known for an effective evaluation!
- **Exhaustive assessment** is usually impractical
  - TREC usually has 50 topics
  - Collection usually has >1 million documents
- **Random sampling** won't work
  - If relevant docs are rare, none may be found!
- **IR systems** can help focus the sample (**Pooling**)
  - Each system finds some relevant documents
  - Different systems find different relevant documents
  - Together, enough systems will find most of them
  - Leverages cooperative evaluations

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## Pooling

1. Systems submit top **1000** documents per topic
2. Top **100** documents from each are judged
  - Single pool, duplicates removed, random ranking
  - Judged by the person who developed the topic
3. Treat unevaluated documents as irrelevant
4. Compute MAP (or others) down to **1000** documents
  - To make pooling work:
    - Large number of reasonable systems participating
    - Systems must not all “do the same thing”

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Pooling, does it work?

- Judgments can't possibly be exhaustive!  
**It doesn't matter: relative rankings of different systems remain the same!**  
Chris Buckley and Ellen M. Voorhees. (2004) Retrieval Evaluation with Incomplete Information. SIGIR 2004.
- This is only one person's opinion about relevance  
**It doesn't matter: relative rankings remain the same!**  
Ellen Voorhees. (1998) Variations in Relevance Judgments and the Measurement of Retrieval Effectiveness. SIGIR 1998.
- What about hits 101 to 1000?  
**It doesn't matter: relative rankings remain the same!**
- We can't possibly use judgments to evaluate a system that didn't participate in the evaluation!  
**Actually, we can!**

Justin Zobel. (1998) How Reliable Are the Results of Large-Scale Information Retrieval Experiments? SIGIR 1998.

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

12

## Who decides a doc is relevant or not?

- The same doc can be seen relevant by me, but not you
- Sometimes, it would be useful to have multiple judgements on relevance on the same document
- How to measure agreement among different assessors?
- Cohen's *kappa*

$$\kappa = \frac{P(A) - P(E)}{1 - P(E)}$$

$P(A)$  – proportion of time judges agree (*inter-annotator agreement*)  
 $P(E)$  – what agreement would be by chance

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Cohen's *kappa*

- Two judges ( $J_1$  &  $J_2$ ) annotating 50 docs for relevance
- $P(A) = \frac{20+15}{50} = 0.7$
- $P(E) = P(J_1, J_2 | rel) + P(J_1, J_2 | irrel)$ 
  - $P(rel) = P(J_1 | rel) \cdot P(J_2 | rel) = \frac{20+10}{50} \cdot \frac{20+5}{50} = 0.6 \times 0.5 = 0.3$
  - $P(irrel) = P(J_1 | irrel) \cdot P(J_2 | irrel) = \frac{20}{50} \cdot \frac{25}{50} = 0.4 \times 0.5 = 0.2$
- $\kappa = \frac{P(A) - P(E)}{1 - P(E)}$
- $= \frac{0.7 - 0.5}{1 - 0.5} = \frac{0.2}{0.5} = 0.4$

		$J_1$	
		Relevant	Irrelevant
$J_2$	Relevant	20	5
	Irrelevant	10	15

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

14

## Cohen's kappa - meaning

- Kappa = 0, for chance agreement,  
= 1, for total agreement.  
< 0, for worse than random!
- Kappa > 0.8 →  
good agreement
- 0.67 < Kappa < 0.8 →  
“fair” agreement
- Kappa < 0.67 →  
seen as data providing a suspicious basis for an evaluation

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Web Search Engines Evaluation

- Search engines have test collections of queries and hand-ranked results
- Recall is difficult to measure on the web – why?
- Search engines often use
  - precision at top k, e.g., k = 10
  - measures that reward you more for getting rank 1 right than for getting rank 10 right (nDCG)
  - non-relevance-based measures:
    - Clickthrough on first result  
not very reliable if you look at a single clickthrough ... but pretty reliable in the aggregate.
    - Studies of user behaviour in the lab
    - A/B testing

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

16

## Web Search Engines: A/B testing

- Purpose: Test a single innovation
- Prerequisite: You have a large search engine up & running.
- Have most users use old system
- Divert a small proportion of traffic (e.g., 1%) to the new system that includes the innovation
- Evaluate with an “automatic” measure like clickthrough on first result
- Now we can directly see if the innovation does improve user satisfaction.
- Probably the evaluation methodology that large search engines trust most

Walid Magady, TTDS 2021/2022



17

## Is system B really better than A?

- Given the results from a number of queries, B achieved better score than A. How can we conclude that ranking algorithm B is really better than algorithm A?

Experiment 1			Experiment 2		
Query	System A	System B	Query	System A	System B
1	0.20	0.40	1	0.02	0.76
2	0.21	0.41	2	0.39	0.07
3	0.22	0.42	3	0.16	0.37
4	0.19	0.39	4	0.58	0.21
5	0.17	0.37	5	0.04	0.02
6	0.20	0.40	6	0.09	0.91
7	0.21	0.41	7	0.12	0.46
Average	0.20	0.40	Average	0.20	0.40

Walid Magady, TTDS 2021/2022



18

## Significance Test

- **Null Hypothesis:**  
No relationship between two observed phenomena
  - Rejecting null hypothesis: observation has a meaning
- A **significance test** enables the rejection of *null hypothesis* (no difference) in favor of the *alternative hypothesis* (B is better than A).
- The power of a test is the probability that the test will reject the *null hypothesis* correctly.
  - increasing the number of queries in the experiment increases the power of test.



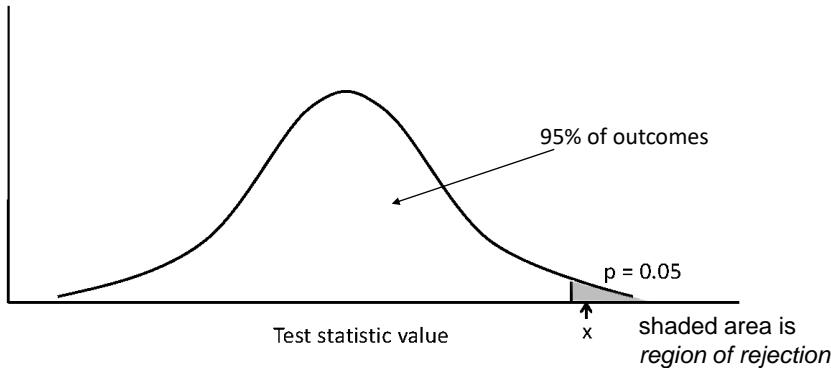
## Significance Test: Steps

- Compute the effectiveness measure for every query for both retrieval systems (note: AP not MAP).
- Compute a **test statistic** based on a comparison of the effectiveness measures for each query.
  - depends on the significance test
- Test statistic is used to compute a **p-value**: reflects the probability that the *null hypothesis* is true.
  - Small p-values suggest that the null hypothesis may be false.
- The null hypothesis (no difference) is rejected in favor of the alternate hypothesis (B is more effective than A) if **p-value  $\leq \alpha$** , where  $\alpha$  is the significance level.
  - Values for  $\alpha$  are small, typically **0.05** or less, to reduce the chance of incorrect rejection.



## One-sided Test Static

- Distribution for the possible values of a test statistic assuming the null hypothesis



Walid Magdy, TTDS 2021/2022



21

## t-test

- Assumption is that the difference between the effectiveness values is a sample from a normal distribution
- Null hypothesis is that the mean of the distribution of differences is zero
- Test statistic

$$t = \frac{\bar{B} - \bar{A}}{\sigma_{B-A}} \cdot \sqrt{N}$$

- t-value to p-value

<http://www.socscistatistics.com/pvalues/distribution.aspx>

Query	A	B	B-A
1	25	35	10
2	43	84	41
3	39	15	-24
4	75	75	0
5	43	68	25
6	15	85	70
7	20	80	60
8	52	50	-2
9	49	58	9
10	50	75	25

$$\bar{B} - \bar{A} = 21.4, \sigma_{B-A} = 29.1, t = 2.33, \text{p-value}=.02$$

Walid Magdy, TTDS 2021/2022



22

## Significance Test

- It is not enough to show that system B achieves better score than system A
  - Significance test is essential
- Two-tailed t-test is highly accepted, with  $\alpha=0.05$ 
  - Sometimes it is required to use others  
Wilcoxon test: does not assume normal distribution
- Meaning of significance test for IR system
  - When a user uses system B that is significantly better than system A, he/she will feel the difference in performance
  - If system B is better than A but not significantly, the user won't notice a difference between the two systems

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

23

## Now, is system B better than A?

Experiment 1			Experiment 2		
Query	System A	System B	Query	System A	System B
1	0.20	0.40	1	0.02	0.76
2	0.21	0.41	2	0.39	0.07
3	0.22	0.42	3	0.16	0.37
4	0.19	0.39	4	0.58	0.21
5	0.17	0.37	5	0.04	0.02
6	0.20	0.40	6	0.09	0.91
7	0.21	0.41	7	0.12	0.46
Average	0.20	0.40	Average	0.20	0.40

t-test p-value = 0

B is statistically significantly better than A

t-test p-value = 0.306

B and A are statistically indistinguishable

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

24

## Summary

- IR test-collection for automatic evaluation
  - Collection of documents
  - Set of topics
    - Topic = query + details on what is meant and what is relevant
    - Recommended minimum number of **25** topics
  - Relevance judgements
    - Pooling is the most common approach for creating judgements
    - Large number of diverse systems are required
  - Evaluation measure
    - Select the proper measure according to the IR task
    - Significance test is essential to confirm that improvement has real meaning
- Web-search uses different evaluation methods that relies on user experience and click-through data

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

25

## Resources

- Text book 1: Intro to IR, Chapter 8
- Text book 2: IR in Practice, Chapter 8
- Pooling:  
Chris Buckley and Ellen M. Voorhees. (2004) Retrieval Evaluation with Incomplete Information. SIGIR 2004

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

26



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Query Expansion

Instructor:  
**Walid Magdy**

27-Oct-2021

1

## Lecture Objectives

- Learn about Query Expansion
  - Query expansion methods
  - Relevance feedback in IR
  - Rocchio's algorithm
  - PRF
- Implement:
  - PRF



THE UNIVERSITY  
of EDINBURGH

2

## Query Expansion

- Query: representation of user's information need
  - Many times it can be suboptimal
- Different words can have the same meaning
  - replacement, replace, replacing, replaced → Stemming
  - go, gone, went → Lemmatisation (NLP)
  - car, vehicle, automobile → ??
  - US, USA, the states, united states of America → ??
- Stemming/Lemmatisation → could be applied to normalise document and queries
  - Research shows that no significant difference between both
- Query Expansion (QE) → add more words of the same meaning to your query for better retrieval

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

3

## Query Expansion: Methods

- Thesaurus
  - Group words into sets of synonyms (synsets)
  - Typically grouping is on the word level (neglects context)
  - Manually built: e.g. WordNet
    - NLTK wordnet: <http://www.nltk.org/howto/wordnet.html>
  - Automatically built:
    - Words co-occurrence
    - Parallel corpus of translations
- Retrieved documents-based expansion
  - Relevance feedback
  - Pseudo (Blind) relevance feedback
- Query logs

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

4

## Automatic Thesaurus: co-occurrence

- Words co-occurring in a document/paragraph are likely to be (*in some sense*) similar or related in meaning
- Built using collection matrix (term-document matrix)
- For a collection matrix  $\mathbf{A}$ , where  $A_{t,d}$  is the normalised weight of term  $t$  in document  $d$ , similarity matrix could be calculated as follows:

$$\mathbf{C} = \mathbf{A} \cdot \mathbf{A}^T$$

where,  $C_{u,v}$  is the similarity score between terms  $u$  and  $v$ . The higher the score, the more similar the terms

- Advantage: unsupervised  
Disadvantage: related words more than real synonyms

Walid Magdy, TTDS 2021/2022



5

## Automatic Thesaurus: co-occurrence

- Example

Word	Nearest neighbors
absolutely	absurd, whatsoever, totally, exactly, nothing
bottomed	dip, copper, drops, topped, slide, trimmed
captivating	shimmer, stunningly, superbly, plucky, witty
doghouse	dog, porch, crawling, beside, downstairs
makeup	repellent, lotion, glossy, sunscreen, skin, gel
mediating	reconciliation, negotiate, case, conciliation
keeping	hoping, bring, wiping, could, some, would
lithographs	drawings, Picasso, Dali, sculptures, Gauguin
pathogens	toxins, bacteria, organisms, bacterial, parasite
senses	grasp, psyche, truly, clumsy, naive, innate

► Figure 9.4 An example of an automatically generated thesaurus. This example is based on the work in Schütze (1998), which employs latent semantic indexing (see Chapter 18 ).

<https://nlp.stanford.edu/IR-book/html/htmledition/automatic-thesaurus-generation-1.html#fig:autothesaurus>

Walid Magdy, TTDS 2021/2022



6

## Automatic Thesaurus: parallel corpus

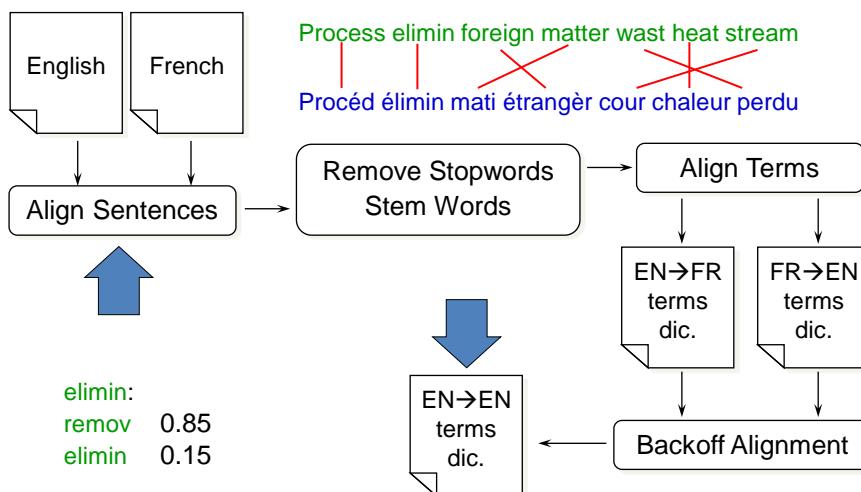
- Parallel corpus are the main training resource for machine translation systems
- Nature: sets of two parallel sentences in two different languages (source and target language)
- Idea:
  - More than one word in language X can be translated into the same word in language Y  
→ these words in language X could be considered synsets
- Requirement: the presence of parallel corpus (training data) → supervised method

Walid Magdy, TTDS 2021/2022



7

## Automatic Thesaurus: parallel corpus



Walid Magdy, TTDS 2021/2022



8

## Automatic Thesaurus: parallel corpus

- Example

motor		weight		travel		color		link	
motor	0.63	weight	0.86	travel	0.67	color	0.56	link	0.4
engin	0.36	wt	0.14	move	0.19	colour	0.25	connect	0.18
				displac		dye	0.19	bond	0.17
								crosslink	0.13
								bind	0.12

cloth		tube		area		game		play	
fabric	0.36	tube	0.88	area	0.4	set	0.6	set	0.3
cloth	0.3	pipe	0.12	zone	0.23	game	0.4	play	0.24
garment	0.2			region	0.2			read	0.17
tissu	0.14			surfac	0.17			game	0.16
								reproduc	0.1

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## Thesaurus-based QE

- Works for very specific applications (e.g. medical domain)
- Many times fails to improve retrieval
  - Sometimes reduces both precision and recall
  - How?
- When it works, it is hard to get a consistent performance over all queries:
  - Improves some, and reduces others. Significant?
- Why it fails?
  - Lack of context
- Current research: word embeddings / BERT
  - No consistent improvement still

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## Relevance Feedback

- Idea: let user give feedback to the IR system about samples of what is relevant and what is not.
- User feedback on relevance of docs in initial results
  - User issues a (short, simple) query
  - The user marks some results as relevant or non-relevant.
  - The system computes a better representation of the information need based on feedback.
  - Relevance feedback can go through one or more iterations
- From user perspective: it may be difficult to formulate a good query when you don't know the collection well, BUT easier to judge particular documents

Walid Magdy, TTDS 2021/2022



11

## Example 1: Image Search

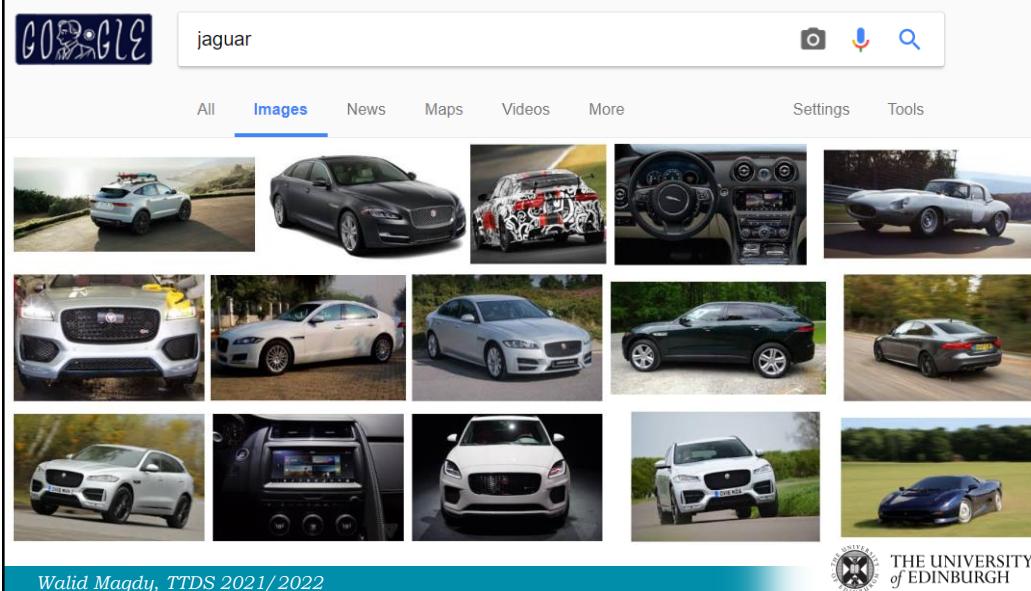
The screenshot shows a Google Images search results page for the query "jaguar". The search bar contains "jaguar". Below the search bar, there are navigation links: All, **Images**, News, Maps, Videos, More, Settings, and Tools. The "Images" link is highlighted in blue. The main content area displays a grid of 15 images. In the first two columns of the first row, three images are highlighted with red boxes: a jaguar walking, a dark-colored car (Jaguar XJ), and a blue SUV (Jaguar F-Pace). The other images in the grid show various jaguars in different poses and environments.

Walid Magdy, TTDS 2021/2022



12

## Example 1: Image Search



13

## Example 2: Text Search

- Initial query: **New space satellite applications**
- Initial Results**
  1. [NASA Hasn't Scrapped Imaging Spectrometer](#)
  2. [NASA Scratches Environment Gear From Satellite Plan](#)
  3. [Science Panel Backs NASA Satellite Plan, But Urges Launches of Smaller Probes](#)
  4. [A NASA Satellite Project Accomplishes Incredible Feat: Staying Within Budget](#)
  5. [Scientist Who Exposed Global Warming Proposes Satellites for Climate Research](#)
  6. [Report Provides Support for the Critics Of Using Big Satellites to Study Climate](#)
  7. [Arianespace Receives Satellite Launch Pact From Telesat Canada](#)
  8. [Telecommunications Tale of Two Companies](#)
- User then marks relevant documents with “+”
- System learns new terms

Walid Magady, TTDS 2021/2022



14

## New terms common in selected docs

2.074 new	15.10 space
30.81 satellite	5.660 application
5.991 nasa	5.196 eos
4.196 launch	3.972 aster
3.516 instrument	3.446 rianespace
3.004 bundespost	2.806 ss
2.790 rocket	2.053 scientist
2.003 broadcast	1.172 earth
0.836 oil	0.646 measure

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Adding new terms to the query

1. NASA Scratches Environment Gear From Satellite Plan
2. NASA Hasn't Scrapped Imaging Spectrometer
3. When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work of Their Own
4. NASA Uses 'Warm' Superconductors For Fast Circuit
5. Telecommunications Tale of Two Companies
6. Soviets May Adapt Parts of SS-20 Missile For Commercial Use
7. Gaping Gap: Pentagon Lags in Race To Match the Soviets In Rocket Launchers
8. Rescue of Satellite By Space Agency To Cost \$90 Million

**Hopefully better results!**

Walid Magdy, TTDS 2021/2022

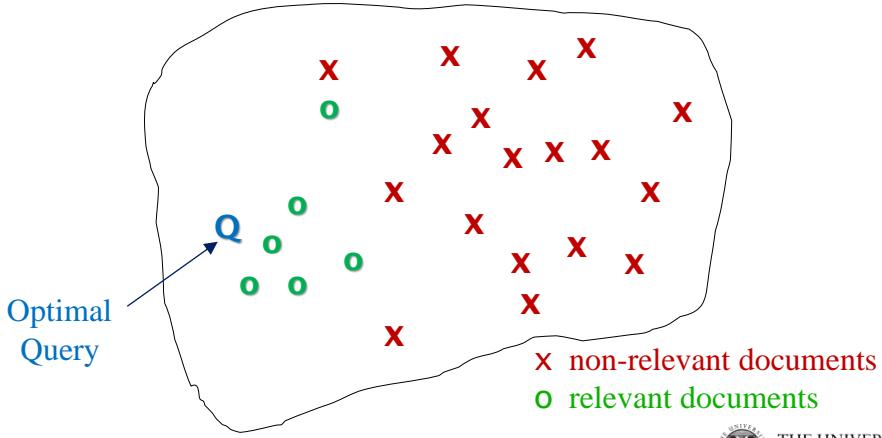


THE UNIVERSITY  
of EDINBURGH

16

## Theoretical Optimal Query

- Found closer to *rel* docs and away from *irrel* ones.
- Challenge: we don't know the truly relevant docs



Walid Magdy, TTDS 2021/2022



17

## Rocchio's Algorithm

- Key Concept: Vector Centroid
- Recall that, in VSM, we represent documents as points in a high-dimensional space
- The centroid is the centre mass of a set of points

$$\vec{\mu}(C) = \frac{1}{|C|} \sum_{\vec{d} \in C} \vec{d}$$

where C is a set of documents.

- Introduced 1963

Walid Magdy, TTDS 2021/2022



18

## Rocchio Algorithm: theory

- Rocchio seeks the query  $\vec{q}_{opt}$  that maximizes

$$\vec{q}_{opt} = \underset{\vec{q}}{\operatorname{argmax}} [sim(\vec{q}, C_{rel}) - sim(\vec{q}, C_{irrel})]$$

- For Cosine similarity

$$\vec{q}_{opt} = \frac{1}{|C_{rel}|} \sum_{\vec{d}_j \in C_{rel}} \vec{d}_j - \frac{1}{|C_{irrel}|} \sum_{\vec{d}_j \notin C_{rel}} \vec{d}_j$$

$$\vec{q}_{opt} = \vec{\mu}(C_{rel}) - \vec{\mu}(C_{irrel})$$

## Rocchio Algorithm: in practice

- Only small set of docs are known to be *rel* or *irrel*

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_{rel}|} \sum_{\vec{d}_j \in D_{rel}} \vec{d}_j - \gamma \frac{1}{|D_{irrel}|} \sum_{\vec{d}_j \in D_{irrel}} \vec{d}_j$$

$\vec{q}_0$  = original query vector

$D_{rel}$  = set of known relevant doc vectors

$D_{irrel}$  = set of known non-relevant doc vectors

$\vec{q}_m$  = modified query vector

$\alpha$  = original query weights (hand-chosen or set empirically)

$\beta$  = positive feedback weight

$\gamma$  = negative feedback weight

- New query moves toward relevant documents and away from non-relevant documents

## Notes about setting weights: $\alpha, \beta, \gamma$

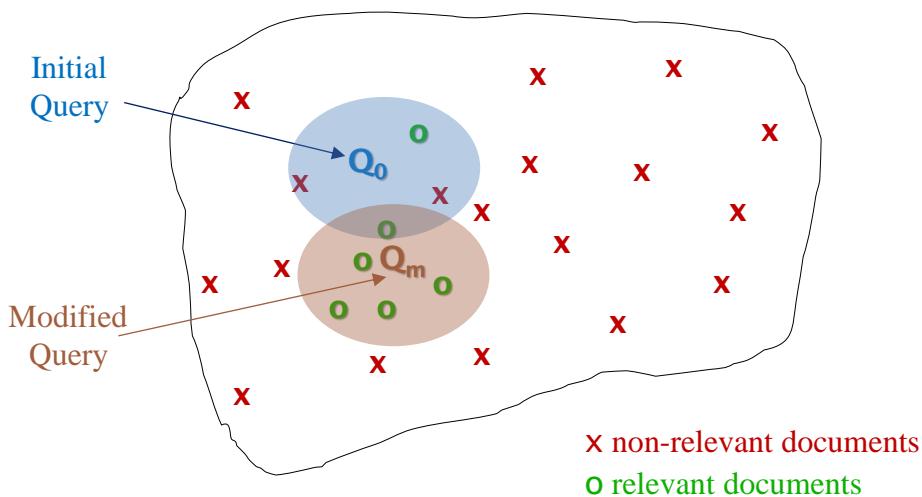
- Values of  $\beta, \gamma$  compared to  $\alpha$  are set high when large judged documents are available.
- In practice, +ve feedback is more valuable than -ve feedback (usually, set  $\beta > \gamma$ )
  - Many systems only allow positive feedback ( $\gamma=0$ ).
  - Or, use only highest-ranked negative document.
- When  $\gamma > 0$ , some weights in query vector can go -ve.
  - “Jaguar”  $\xrightarrow{\text{feedback}}$  jaguar + car + model - animal - jungle
- In practice, top  $n_t$  terms in  $\vec{d}_j \in D_{rel}$  are only selected
  - $n = 5 \rightarrow 50$
  - Top  $n_t$  are identified using e.g. TFIDF

Walid Magdy, TTDS 2021/2022



21

## Effect of Relevance Feedback on Query



Walid Magdy, TTDS 2021/2022



22

## Effect of Relevance Feedback on Retrieval

- Relevance feedback can improve recall and precision
- In practice, relevance feedback is most useful for increasing recall in situations where recall is important.
- Empirically, one round of relevance feedback is often very useful. Two rounds is sometimes marginally useful.

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

23

## Relevance Feedback: Issues

- Long queries are inefficient for typical IR engine.
  - High cost for retrieval system. (why?)
  - Long response times for user.
- It's often harder to understand why a particular document was retrieved after applying relevance feedback
- Users are often reluctant to provide explicit feedback  
→ not practical!

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

24

## Relevance Feedback: Practicality

- User revises and resubmits query
  - Users may prefer revision/resubmission to having to judge relevance of documents.
  - Useful for query suggestion to other users
- Is there a way to apply relevance feedback without user's input?

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

25

## Pseudo (Blind) Relevance Feedback

- Solves the problem of users hate to provide feedback
- Feedback is applied blindly (PRF)
  - Automates the “manual” part of true relevance feedback.
- Algorithm:
  - Retrieve a ranked list of hits for the user’s query
  - Assume that the top  $k$  documents are relevant
  - Do relevance feedback (e.g. Rocchio)
  - Typically applies only positive relevance feedback ( $\gamma=0$ )
- Mostly works
  - Still can go horribly wrong for some queries (when top  $k$  docs are not relevant)
  - Several iterations can lead to query drift

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

26

## PRF (BRF)

- Was proven to be useful for many IR applications
  - News search (learn names and entities)
  - Social media search (learn hashtags)
  - Web search (implicit feedback is used more = clicks)
- Some domains are more challenging
  - Patent search
    - Top documents are usually not relevant
    - Patent text in general is unclear/confusing
- PRF is the most basic QE method for IR
  - Unsupervised
  - Language independent
  - Does not require any kind of language resources

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

27

## PRF (BRF): Evaluation

- In practice, different number of feedback docs ( $n_d$ ) and terms ( $n_t$ ) are usually tested for PRF
  - $n_d$ : 1 → 50
  - $n_t$ : 5 → 50
- Results of PRF are directly compared to baseline (with no PRF)
  - It is not considered cheating.
  - It is essential to show that improvement is significant, and preferred to show the % of queries improved vs degraded.

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

28

## Summary

- QE: automatically add more terms to user's query to better match relevant docs
- QE via thesaurus
  - Manual/automatic thesaurus: useful for specific applications
  - Fail when context is important
- Relevance feedback
  - Get samples of *rel/irrel* docs for extracting QE useful terms
  - Rocchio's is one of the most common algorithms for query modification
- PRF
  - Skips user's input for the feedback process
  - Found to be useful in many applications

Walid Magdy, TTDS 2021/2022



29

## Resources

- Text book 1: Intro to IR, Chapter 9
- Text book 2: IR in Practice, Chapter 6.2, 6.3
- Reading:  
Magdy W. and G. J. F. Jones.  
A Study on Query Expansion Methods for Patent Retrieval.  
*PAIR 2011 - CIKM 2011* ([link](#))
- Lab 5

Walid Magdy, TTDS 2021/2022



30



THE UNIVERSITY  
of EDINBURGH

## Search is not only the Web IR Applications

Walid Magdy

School of Informatics  
University of Edinburgh

27 Oct 2021

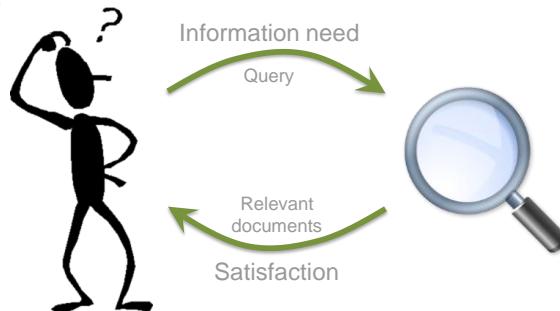
1

## Objectives

- Main objective of IR
- Two search tasks
  - Printed documents search
  - Patent search
- Possible ideas for Group Project ☺

## Information Retrieval Objective

- IR is finding material of an unstructured nature that satisfies an information need from within large collections.
- **Information need**
  - Expected search scenario?
  - Modeling the task?
- **Data nature**
  - Approach?
  - Scalable? Fast?
- **User Satisfaction**
  - More relevant documents?
  - Effective evaluation?



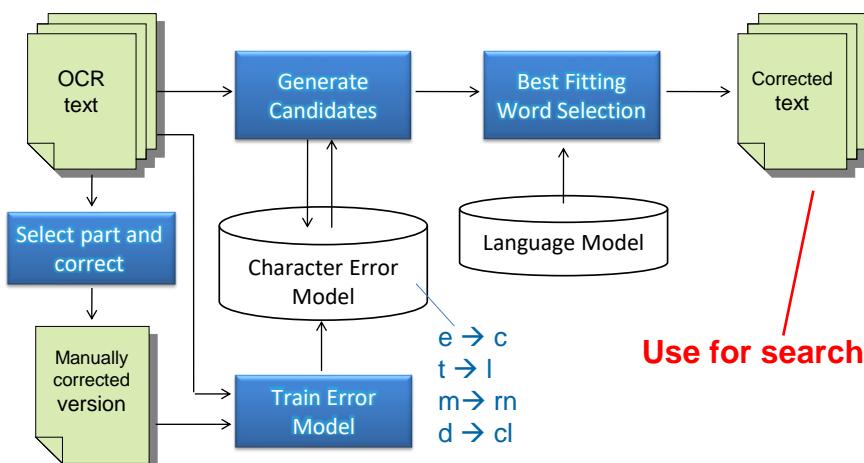
## Printed Documents Retrieval



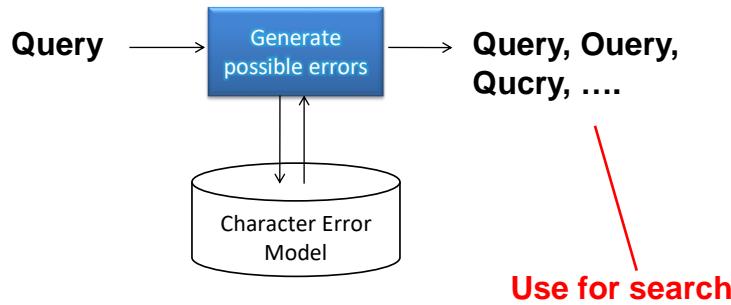
## n-gram Char Representation of OCR

- **Original:**  
example sentence
- **OCR output:**  
example senlcnce
- **3-gram char representation:**  
\$ex exa xar arn rnp npl ple le\$ \$se sen enl nlc lcn cnc nce ce\$
- **Query:**  
example sentence →  
\$ex exa xam amp mpl ple le\$ \$se sen ent nte ten enc nce ce\$
- **Matching:**  
\$ex exa xar arn rnp npl ple le\$ \$se sen enl nlc lcn cnc nce ce\$  
\$ex exa xam amp mpl ple le\$ \$se sen ent nte ten enc nce ce\$

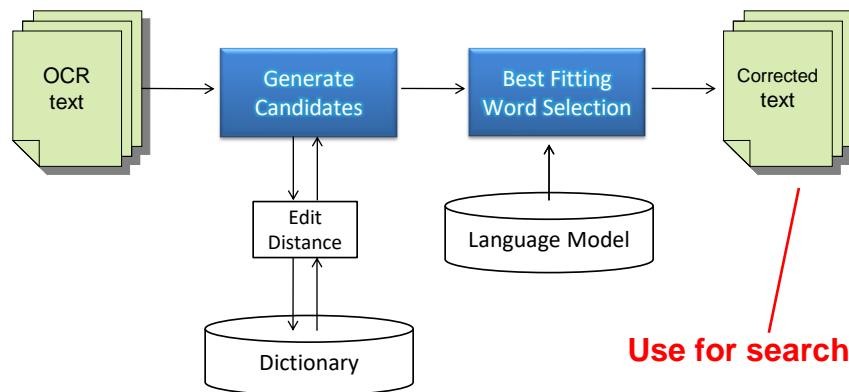
## OCR Correction using Error Model



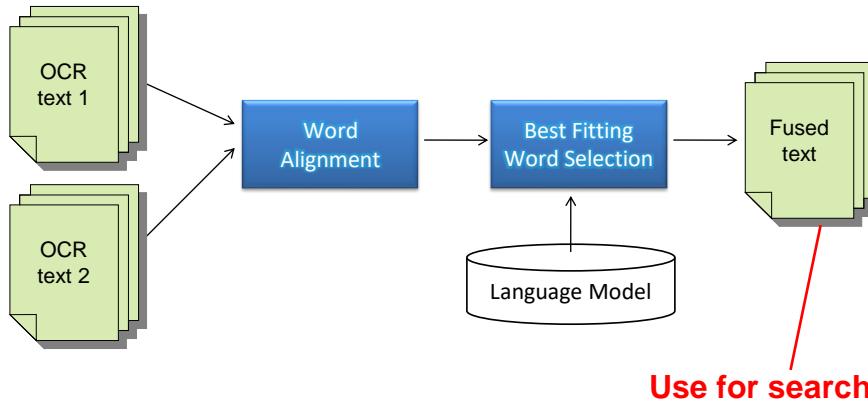
## Query Garbling using Error Model



## OCR Correction using Edit Distance



## Multi-OCR Text Fusion

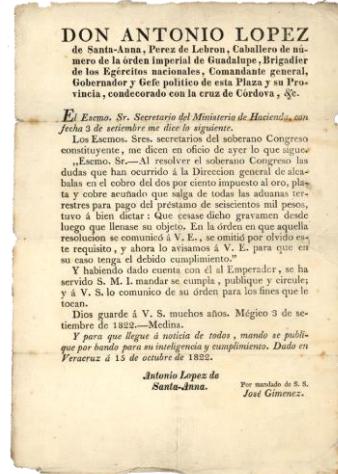


## OCR Search

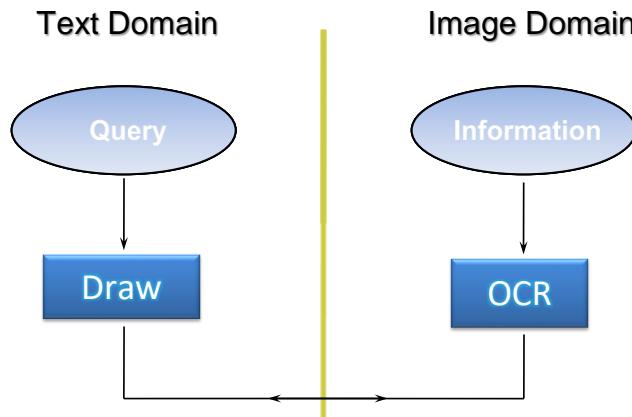
- Recognition errors in OCR text degrades retrieval
- Different methods of text processing can overcome the negative effect on retrieval and improves search
- n-gram character representation improves retrieval, but not that much
- Some training and resources are needed which can be manual correction, trained language model, or both
- Previous methods fail when errors are large ( $WER > 50\%$ )

## Solution – back to Information Need

- **Information need:**  
the printed papers
- **Question:**  
Why convert image to text?
- **Related work:**  
Word Spotting



## Modeling the Problem



## OCRless Search

قال : صدقت ، قال : فعجبنا له يسأله ويفسده ، قال : فأخبرني عن الإيمان ؟ قال :  
 أَنْ تُؤْمِنَ بِاللَّهِ وَبِالْمَلَائِكَةِ وَرَبِّكِهِ ، وَرَسُولِهِ ، وَالْأَئِمَّةِ ، وَتَوَمَّنَ بِالظَّرْفِ خَيْرِهِ  
 وَشَرِّهِ ، قَالَ صدقت ، قال : فأخبرني عن الإيمان ؟ قال : أَنْ تَعْمَدَ اللَّهُ كَافِرَكَ  
 تَرَاهُ ، فَإِنْ لَمْ تَكُنْ تَرَاهُ ، فَإِنَّهُ بِرَبِّكَ ، قَالَ صدقت ، قال : فأخبرني عن السَّاعَةِ ؟  
 قَالَ : مَا أَمْرُ رَبِّكَ عَنْهَا يَعْلَمُ بِهَا ، قَالَ فأخبرني عن أَمْرَاتِها ؟ قَالَ : أَنْ  
 تَلِدَ الْأَمَّةَ رَبِّهَا ، وَأَنْ تَرِي الْخَلَقَ الْمَرْأَةَ الْمَلَأَةَ رَعَاءَ الشَّاءِ يَطَّاولُونَ فِي الْبَيَانِ .

Segment to elements

فَعْجَبَنَا لَهُ وَصَدَقَتْ لَهُ

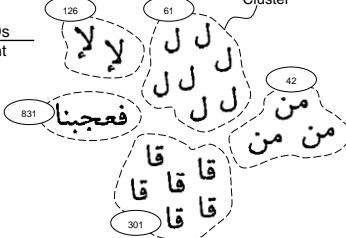
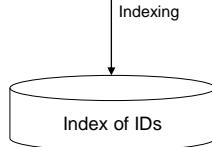
Clustering

Cluster ID

Cluster

213 31 89 32 2 213 31 3341  
1190 23 802 ...

Create IDs document

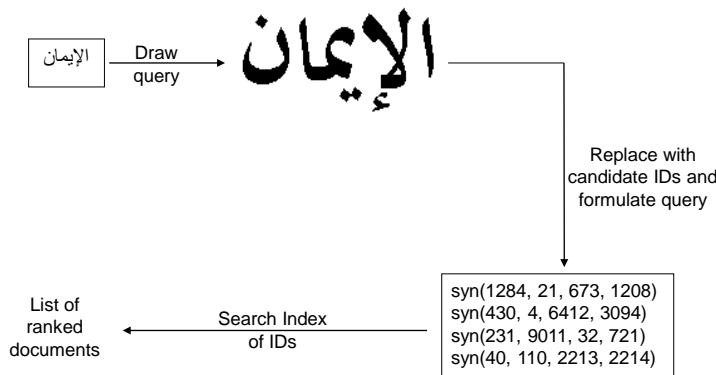


Institute for Language, Cognition and Computation  
ILCC

THE UNIVERSITY  
of EDINBURGH

15

## Solution – OCRless Search



Institute for Language, Cognition and Computation  
ILCC

THE UNIVERSITY  
of EDINBURGH

16

## Solution – OCRless Search

- Effective and fast
- Robust to OCR errors (*v1de0*)
- No training resources required
- Language independent

العربية חִמְצִין

English 熊貓



- Microsoft TechFest Demo

The same engine for searching printed documents in:  
Arabic, English, Chinese, Hebrew, and Hieroglyphic

## Printed Documents Retrieval

- Text-based solutions: correction
- Image-based: clustering

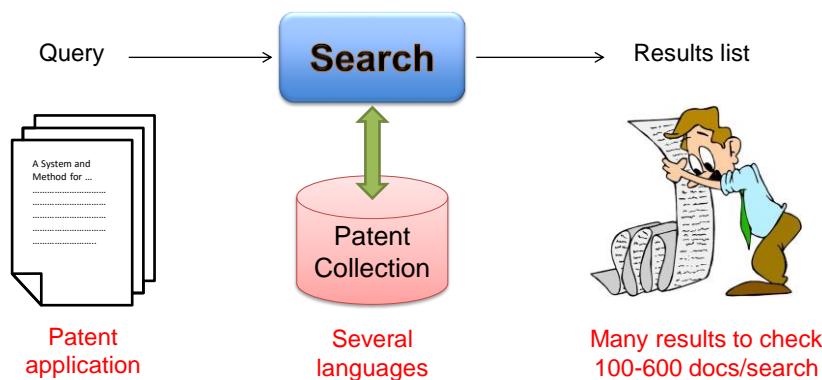
- Current State-of-the-art?

- Information need → Approach

# Patent Search

## Patent Search

- Given a patent application, check if the invention described is novel



## Patent Search – User Satisfaction

- NTCIR, CLEF, TREC
- Recall-oriented → Try not to miss a relevant document
  - Recall is the objective
- Precision is also important
- Huge # documents checked (100-600 documents)
  
- Evaluation: average precision (AP)!!
  - Focuses on finding relevant docs early in ranked list
  - Less focus on recall

## Example

For a topic with 4 relevant docs and 1<sup>st</sup> 100 docs to be examined:

System1: relevant ranks = {1}

System2: relevant ranks = {50, 51, 53, 54}

System3: relevant ranks = {1, 2, 3, 4}

$$AP_{\text{system1}} = 0.25$$

$$AP_{\text{system2}} = 0.0481$$

$$AP_{\text{system3}} = 1$$

$$R_{\text{system1}} = 0.25$$

$$R_{\text{system2}} = 1$$

$$R_{\text{system3}} = 1$$

- We need a metric that reflects recall and ranking quality in one measure

## PRES: Patent Retrieval Evaluation Score

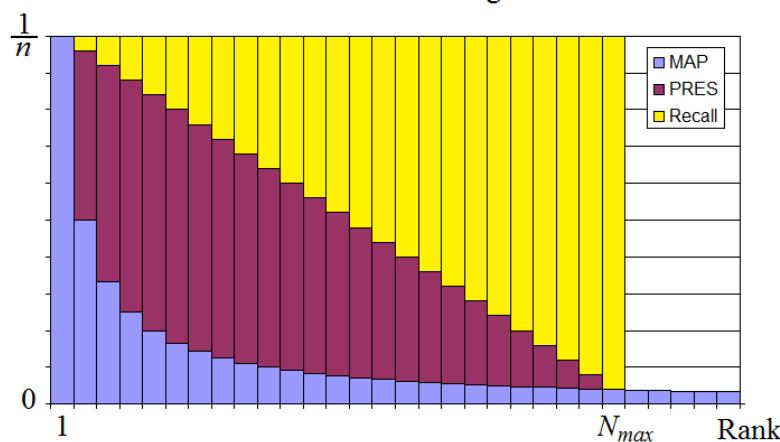
$$PRES = 1 - \frac{\frac{\sum r_i}{n} - \frac{n+1}{2}}{N_{\max}}$$

n: number of relevant docs  
 $r_i$ : rank of the  $i^{\text{th}}$  relevant document  
 $N_{\max}$ : max number of checked docs

- Derived from  $R_{\text{norm}}$  (Rocchio, 1964)
- Gives higher score for systems achieving higher recall and better average relative ranking
- Dependent on user's potential/effort ( $N_{\max}$ )
- Robust to incomplete relevance judgements

## PRES: as a cumulative gain

Value added to score when finding relevant document



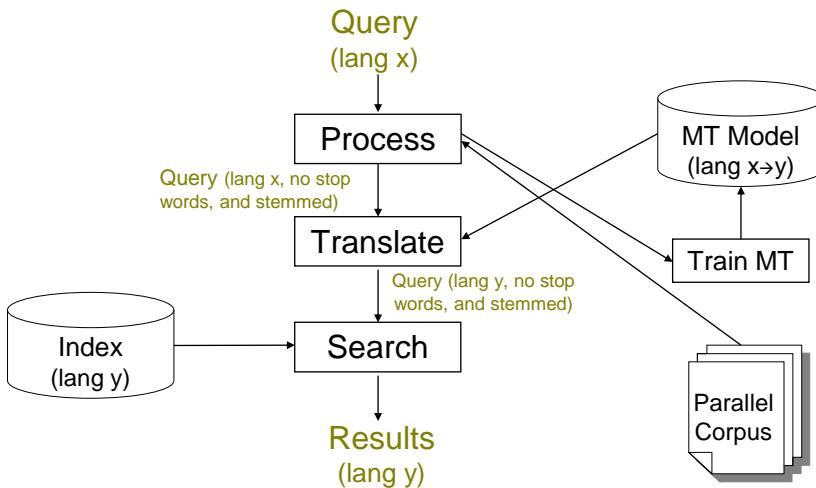
## Patent Search – CLIR

- **Query:** Full patent application
- **Common approach:** MT (the best)
- **Challenge:** training resources + speed!
- **Ideal:** Query + Document translation

## Patent Search – CLIR – Objective?

- **Manual translation**  
It is a great idea to apply stemming in information retrieval
- **MT output**  
he are an great ideas to applied stem by information retrieving
- **MT evaluation: MT sucks**
- **IR evaluation: MT rocks ☺**
- **MT4IR:** An efficient MT that neglects morphological and syntactic features of output

## Ordinary MT vs. MT4IR



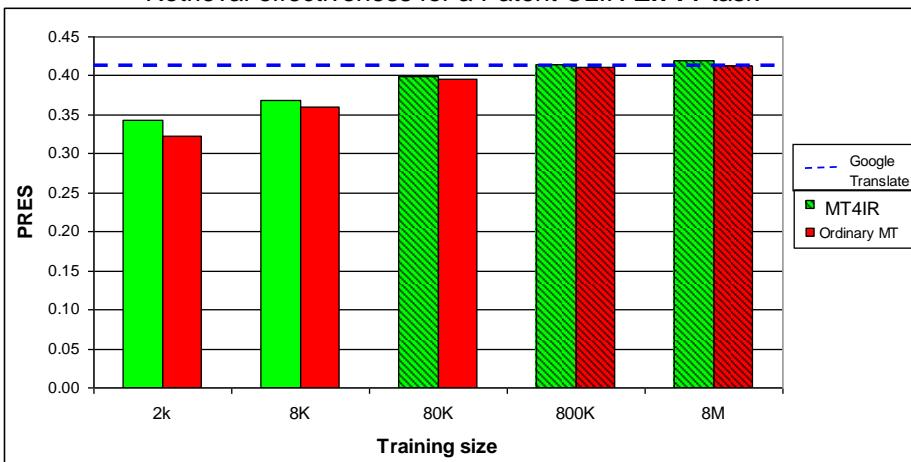
Institute for Language, Cognition and Computation  
ILCC



27

## Patent Search – MT4IR

Retrieval effectiveness for a Patent CLIR En-Fr task

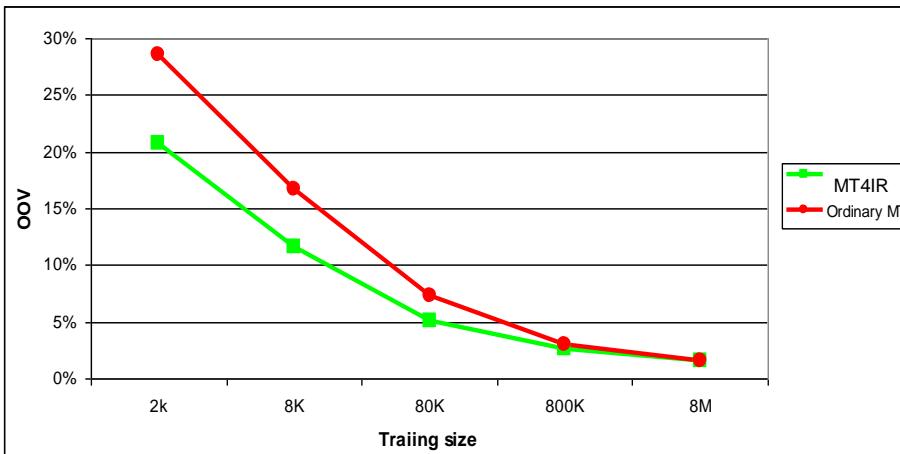


Institute for Language, Cognition and Computation  
ILCC



28

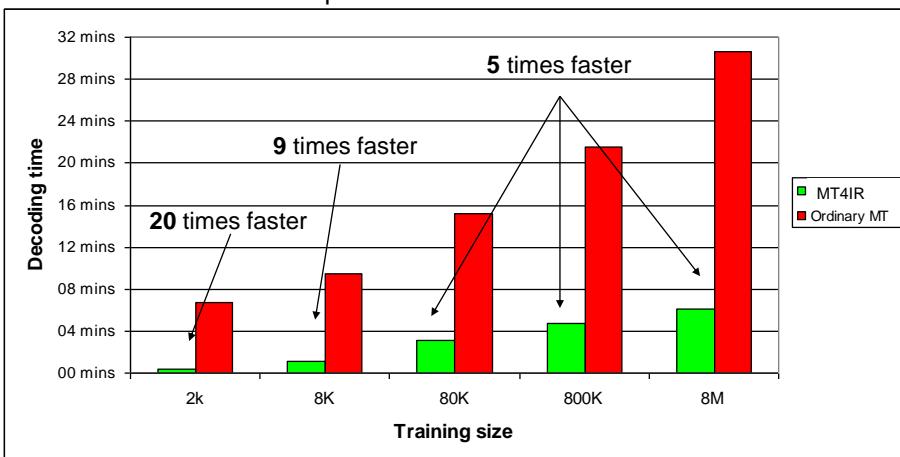
## Patent Search – MT4IR



E.g. play, plays, played, playing

## Patent Search – MT4IR

Translation speed for a Patent CLIR En-Fr task



## Summary

- The objective is IR is “User Satisfaction”
- Understand the user needs well
- Design the IR task carefully
- You do not have to stick to the path in the literature
- Are you sure performance is measured correctly?

## Readings

- Magdy W. and G. J. F. Jones. Studying Machine Translation Technologies for Large-Data CLIR Tasks: A Patent Prior-Art Search Case Study. *Springer, Information Retrieval, 2013*
- Magdy W. and G. J. F. Jones. PRES: A Score Metric for Evaluating Recall-Oriented Information Retrieval Applications. *SIGIR 2010*
- Magdy W. , K. Darwish, and M. El-Saban. Efficient Language-Independent Retrieval of Printed Documents without OCR. *SPIRE 2009*
- Magdy W. and K. Darwish. Effect of OCR Error Correction on Arabic Retrieval. *Springer, Information Retrieval, 2008*



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Web Search

Instructor:  
**Walid Magdy**

3-Nov-2021

1

## Lecture Objectives

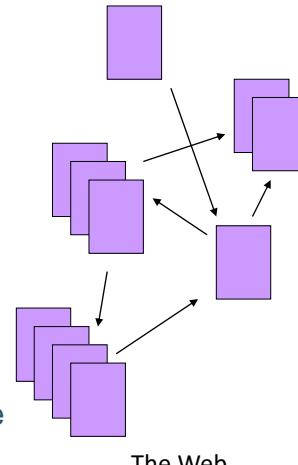
- Learn about:
  - Working with Massive data
  - Link analysis (PageRank)
  - Anchor text



2

## The Web Document Collection

- Huge / Massive
- Graph / Connected
- No design/co-ordination
- Distributed content publishing
- Content includes truth, lies, obsolete information, contradictions ...
- Unstructured (text, html, ...), semi-structured (XML, annotated photos), structured (DB) ...
- Growth – slowed down from initial “volume doubling every few months” but still expanding
- Content can be dynamically generated



Walid Magdy, TTDS 2021/2022



3

## Effect of Massive data

- Web search engines work with huge amount of data
  - 20 PB/day in 2008 → 160 PB/day in 2013 → now??
  - 1 PB = 1,000 TB = 1,000,000 GB
- How this would affect a search engine?
  - Very challenging (storage, processing, networking, ...)
  - Very useful still (makes stuff easier), how?
- Assume two good search engines the collects two sub-sets of the web
  - Search engine A collected N docs → precision@10 = 40%
  - Search engine B collected 4N docs → precision@10??

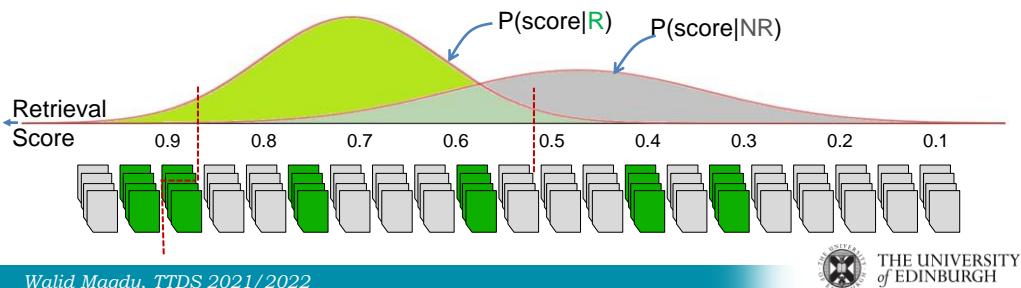
Walid Magdy, TTDS 2021/2022



4

## Effect of Massive data on Precision

- Assume two good search engines that collect two sub-sets of the web
  - Search engine A collected N docs  $\rightarrow$  precision@10 = 40%
  - Search engine B collected 4N docs  $\rightarrow$  precision@10???
    - Distribution of positive/negative scores stays the same
    - Precision/Recall at a given score stays the same
    - In any decent IR system: more relevant docs exist at the top  
 $\rightarrow P@n \uparrow \uparrow \rightarrow$  precision@10 = 60% (increases)



5

## Big Data or Clever Algorithm?

- For Web search, larger index usually would beat a better retrieval algorithm
  - Google Index vs Bing Index
- Similar to other applications
  - Google MT vs IBM MT
    - Statistical methods trained over **10x** training data beat deep NLP methods with **1x** training data
  - In general ML, the more data, the better the results
    - Tweets classification: using **100x** of noisy training data beats **1x** of well prepared training data, even with absence of stemming & stopping
  - Question answering task:
    - IBM Watson vs Microsoft experiment

6

## Big Data or Clever Algorithm?

- Question answering task:
  - **Q:** Who created the character of Scrooge?
  - **A:** Scrooge, introduced by **Charles Dickens** in “A Christmas Carol”
  - Requires heavy linguistic analysis, lots of research in TREC
- 2002, Microsoft
  - Identify (**subj verb obj**), rewrite as queries:
    - Q1: “created the character of Scrooge”
    - Q2: “the character of Scrooge was created by”
  - Search the web for exact phrase, get top 500 results
  - Extract phrase: ■Q1 or Q2■ , get most frequent ■
  - Very naive approach, ignores most answers patterns
  - Who cares!! Web is huge, you will find matches anyway

117	Dickens
78	Christmas Carol
75	Charles Dickens
72	Disney
54	Carl Banks
...	

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Search “Microsoft”

Doc1

Microsoft.com

“Microsoft” mentioned  
5 times

Doc2

Tutorial.com  
*Tutorial on MS word*

“Microsoft” mentioned  
35 times

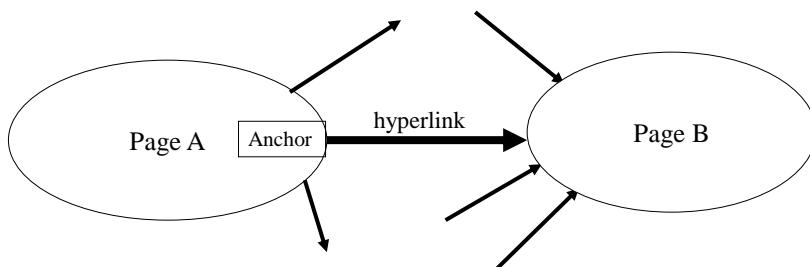
Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## The Web as a Directed Graph



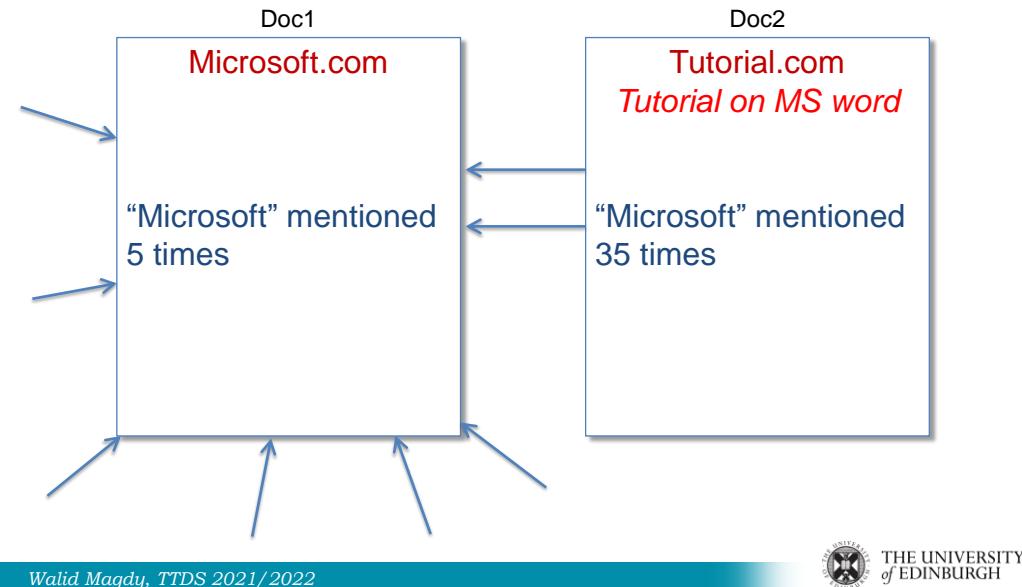
**Assumption 1:** A hyperlink between pages denotes author perceived relevance (quality signal)

**Assumption 2:** The text in the anchor of the hyperlink describes the target page (textual context)

## Links between Pages

- Google Description of PageRank:
  - Relies on the “**uniquely democratic**” nature of the web
  - Interprets a link from page A to page B as “**a vote**”
- A → B: means A thinks B worth something
  - “**wisdom of the crowds**”: many links means B must be good
  - **Content-independent** measure of quality of B
- Use as ranking feature, combined with content
  - But not all pages that link to B are of equal importance!
    - Importance of a link from CNN >> link from blog page
- Google PageRank, 1998
  - How many “good” pages link to B?

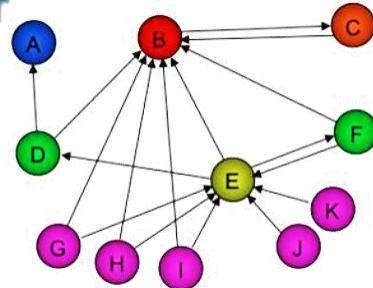
## Search “Microsoft”



11

## PageRank: Random Surfer

- **Analogy:**
  - User starts browsing at a random page
  - Pick a random outgoing link  
→ goes there → repeat forever
  - Example:  
 $G \rightarrow E \rightarrow F \rightarrow E \rightarrow D \rightarrow B \rightarrow C$
  - With probability  $1-\lambda$  jump to a random page
    - Otherwise, can get stuck forever A, or B ↔ C
- **PageRank of page x**
  - Probability of being at page x at a random moment in time



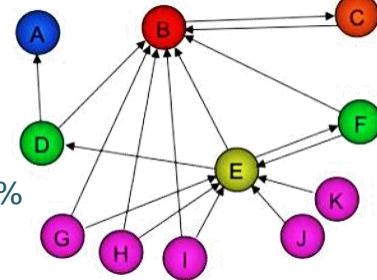
Walid Magdy, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

12

## PageRank: Algorithm

- Initialize  $PR_0(x) = \frac{100\%}{N}$ 
  - $N$ : total number of pages
  - $PR_0(A) = \dots = PR_0(K) = \frac{100\%}{11} = 9.1\%$



- For every page  $x$

$$PR_{t+1}(x) = \frac{1 - \lambda}{N} + \lambda \sum_{y \rightarrow x} \frac{PR_t(y)}{L_{out}(y)}$$

- $y \rightarrow x$  contributes part of its PR to  $x$
- Spread PR equally among out-links
- Iterate till converge  $\rightarrow$  PR scores should sum to 100%

Walid Magdy, TTDS 2021/2022

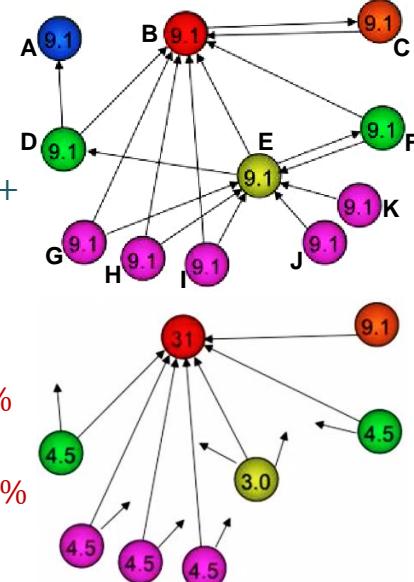


THE UNIVERSITY  
of EDINBURGH

13

## PageRank: Example

- Let  $\lambda = 0.82$
- $PR(B) = \frac{0.18}{11} + 0.82 \times [PR(C) + \frac{1}{2}PR(D) + \frac{1}{3}PR(E) + \frac{1}{2}PR(F) + \frac{1}{2}PR(G) + \frac{1}{2}PR(H) + \frac{1}{2}PR(I)] \approx 0.31 = 31\%$
- $PR(C) = \frac{0.18}{11} + 0.82 \times PR(B) = 0.18 \times 9.1\% + 0.82 \times 31\% = 9.1\%$
- $PR_{t+1}(C) = 0.18 \times 9.1\% + 0.82 \times 31\% \approx 26\%$



Walid Magdy, TTDS 2021/2022

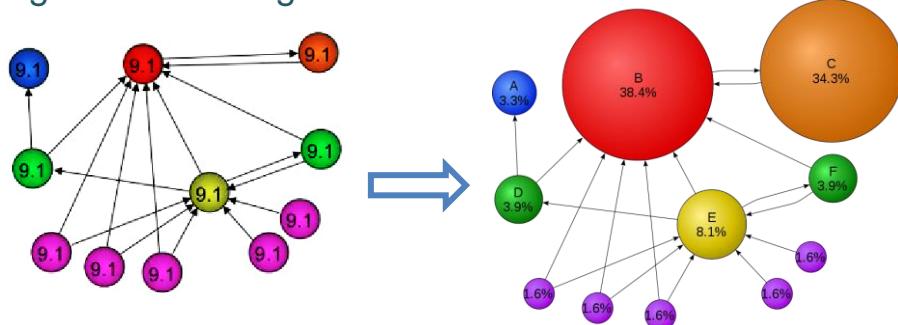


THE UNIVERSITY  
of EDINBURGH

14

## PageRank: Example result

- Algorithm converges after few iterations



- Observations
  - Pages with no inlinks:  $PR = (1 - \lambda)/N = 0.18/11 = 1.6\%$
  - Same (or symmetric) inlinks → same PR (e.g. D and F)
  - One inlink from high PR >> many from low PR (e.g. C vs E)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Anchor Text

- Anchor Text (text of a link):
  - Description of destination page
  - Short, descriptive like a query
  - Re-formulated in different ways
    - Human “query expansion”
- Used when indexing page content
  - Add text of all anchor text linking the page
  - Different weights for different anchor text
    - Weighted according to PR of linking page
- Significantly improves retrieval

[International Business Machines](#) announced today a deal of \$100M ...

We support:  
[Sun](#)  
[HP](#)  
[IBM](#)

[www.ibm.com](#)  
[Big Blue](#) today announced record profits for the quarter

Walid Magdy, TTDS 2021/2022

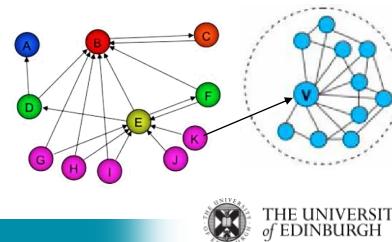


THE UNIVERSITY  
of EDINBURGH

16

## Link Spam

- Trackback links (blogs that link to me)
  - Based on `$HTTP_REFERER`
  - Artificial feedback loops
  - Similar to “follow back” in Twitter
- Links from comments on sites with high PR
  - Links in comments on CNN
  - One solution: insert `rel=nofollow` into links
    - Link ignored when computing PR
- Link farms
  - Fake densely-connected graph
  - Hundreds of web domains / IPs can be hosted on one machine



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

17

## The Reality

- **PageRank** is used in Google, but is hardly the full story of ranking
  - A big hit when initially proposed, but just one feature now
  - Many sophisticated features are used
  - Machine-learned ranking heavily used
    - Learning to Rank (L2R)
    - Many features are used, including PR
  - Still counted as a very useful feature

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

18

## Summary

- Web data is massive
  - Challenging for efficiency, but useful for effectiveness
- PageRank:
  - Probability than random surfer is currently on page x
  - The more powerful pages linking to x, the higher the PR
- Anchor text:
  - Short concise description of target page content
  - Very useful for retrieval
- Link Spam
  - Trackable links, link farms

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

19

## Resources

- Text book 1: Intro to IR, Section 21.1
- Text Book 2: IR in Practice: 4.5, 10.3
- Page Rank Paper:  
 Page, L., Brin, S., Motwani, R., & Winograd, T. (1999).  
*The PageRank citation ranking: Bringing order to the web.*  
 Stanford InfoLab.
- Additional reading:  
 Dumais, S., Banko, M., Brill, E., Lin, J., & Ng, A. (2002)  
 Web question answering: Is more always better?.  
 SIGIR 2002.
- YouTube Video: How Search Works  
<https://www.youtube.com/watch?v=BNHR6IQJGZs>

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

20



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Web Search (2)

Instructor:  
**Walid Magdy**

03-Nov-2021

1

## Lecture Objectives

- Learn about:
  - Basics of Web search
  - Brief History of web search
  - SEOs
  - Web Crawling (intro)



THE UNIVERSITY  
of EDINBURGH

2

## Brief History

- Early keyword-based engines (1995-1997)
  - Altavista, Excite, Infoseek, Lycos, AOL
  - Traditional IR techniques
  - Scalability is an issue
- Paid search ranking: Goto (morphed into Overture.com → Yahoo!)
  - Your search ranking depended on how much you paid
  - Auction for keywords
  - Called “sponsored search”
    - CPC (Cost Per Click)
    - CPM (Cost Per Thousand Impressions)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

3

## CPC / CPM / RPM

- With new services on the web → RPM
- RPM: Revenue per 1000 video views
- Read more:  
Understand ad revenue analytics  
<https://support.google.com/youtube/answer/9314357>

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

4

## Brief (non-technical) History

- 1998+: Link-based ranking pioneered by Google
  - Blew away all early engines
  - Great user experience in search of a business model
  - Meanwhile Goto/Overture's annual revenues: ~ \$1 billion
- Result: Google added paid search “ads” to the side, independent of search results
  - Yahoo followed, acquiring Overture (for paid placement) and Inktomi (for search)
- 2005+: Google gains search share, dominating in Europe and very strong in North America
  - 2009: Yahoo! and Microsoft combined paid search offering

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

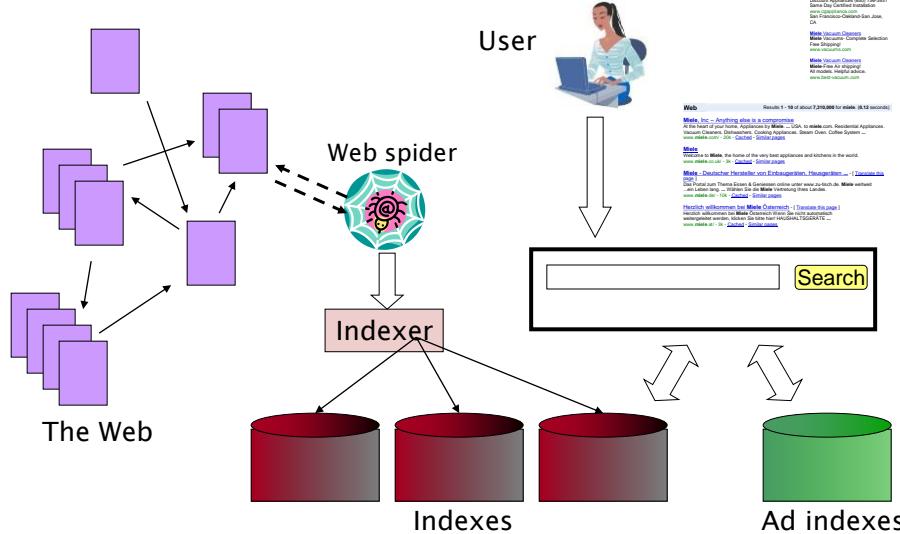
5

The screenshot shows a Google search results page with the following details:

- Search Query:** nigritude ultramarine
- Results Count:** 1 - 10 of about 185,000
- Section Headers:**
  - Algorithmic Search Results
  - Sponsored Search Ads
- Content Preview:**
  - Anil Dash: [Nigritude Ultramarine](#)
  - Nigritude Ultramarine FAQ**
  - SEO contest - Wikipedia, the free encyclopedia**
  - Slashdot | How To Get Googled, By Hook Or By Crook**
  - The Nigritude Ultramarine Search Engine Optimization Contest**
  - Sponsored Links:**
    - [Business Blogging Seminar](#)
    - [Full-Time SEO & SEM Jobs](#)
    - [SEO Contests](#)
    - [The SEO Book](#)
    - [Ultramarine - Companion](#)

6

## Web Search Basics



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## User Need on Web Search

- **Informational** – want to learn about something (~40% / 65%)  
**Information Retrieval**
- **Navigational** – want to go to that page (~25% / 15%)  
**United Airlines**
- **Transactional** – want to do something (web-mediated) (~35% / 20%)
  - Access a service **Seattle weather**
  - Downloads **Mars surface images**
  - Shop **Canon S410**
- **Gray areas**
  - Exploratory search “see what’s there”

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## Search Engine Optimization (SEO)

- The Trouble with Paid Search Ads:  
It costs money. What's the alternative?
- **Search Engine Optimization (SEO):**
  - “Tuning” your web page to rank highly in the algorithmic search results for selected keywords
  - Alternative to paying for placement
  - Thus, intrinsically a marketing function
- Performed by companies, webmasters and consultants (“Search engine optimizers”) for their clients
- Some perfectly legitimate, some very shady

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

## SEO: Simplest Form

- First generation engines relied heavily on *tf/idf*
  - The top-ranked pages for the query **maui resort** were the ones containing the most **maui's** and **resort's**
- SEOs responded with dense repetitions of chosen terms
  - e.g., **maui resort maui resort maui resort**
  - Misleading meta-tags, excessive repetition
  - Often, the repetitions would be in the same color as the background of the web page
    - Repeated terms got indexed by crawlers
    - But not visible to humans on browsers

*Pure word density cannot be trusted  
as an IR signal*



Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## SEO word manipulating examples

- XYZ Hotel in ABC city
  - Accommodation, hotel, room, flat, travel, sights, attractions, vacation, holiday, in ABC ABC ABC
- XYZ for family advices
  - Family, couples, parents, spouse, wife, husband, fights, relationship, cheating, communication, kids, children
- XYZ Umbrellas
  - Raining, rainy, wet, weather, day

Walid Magdy, TTDS 2021/2022

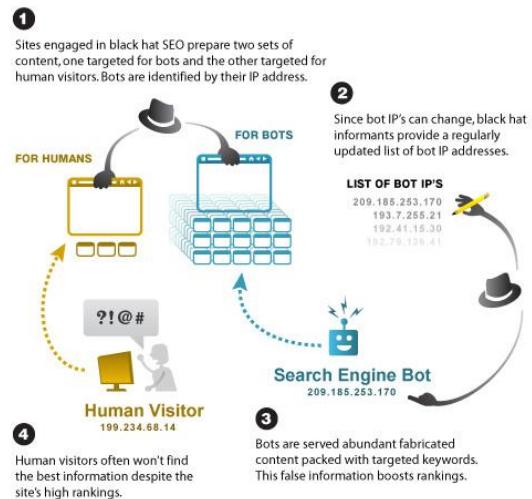


11

## SEO: Cloaking

- Serve fake content to search engine spider
- Famous technique: **Black Hat**
- Kind of a spam!

### Black Hat Cloaking Explained



Walid Magdy, TTDS 2021/2022



12

## Duplicate Detection

- The web is full of duplicated content
- Strict duplicate detection = exact match
  - Not as common
  - can be detected with fingerprints
- But many, many cases of **near duplicates**
  - e.g., last modified date the only difference between two copies of a page
- *Near-Duplication:* Approximate match
  - Use similarity threshold to detect near-duplicates
    - e.g., Similarity > 80% => Documents are “near duplicates”
    - Not transitive though sometimes used transitively
      - $A \approx B \ \& \ B \approx C \rightarrow$  doesn't have to mean  $A \approx C$

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Duplicate Detection: MiniHash

- Features of similarity:
  - Segments of a document (natural or artificial breakpoints)
  - **Shingles** (word n-grams)
  - **a rose is a rose is a rose** →  
 a\_rose\_is\_a  
 rose\_is\_a\_rose  
 is\_a\_rose\_is  
 a\_rose\_is\_a
- Similarity measure between two docs (= sets of shingles)
  - Set intersection
  - Specifically ( $\text{Size\_of\_Intersection} / \text{Size\_of\_Union}$ )

Walid Magdy, TTDS 2021/2022

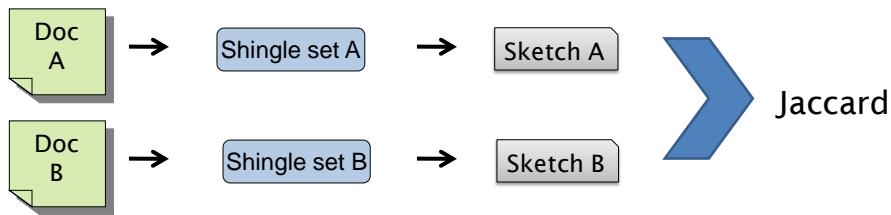


THE UNIVERSITY  
of EDINBURGH

14

## Shingles + Set Intersection

- Computing exact set intersection of shingles between all pairs of documents is expensive/intractable
- Approximate using a cleverly chosen subset of shingles from each (a sketch)
- Estimate  $\frac{\text{size of intersection}}{\text{size of union}}$  based on a short sketch

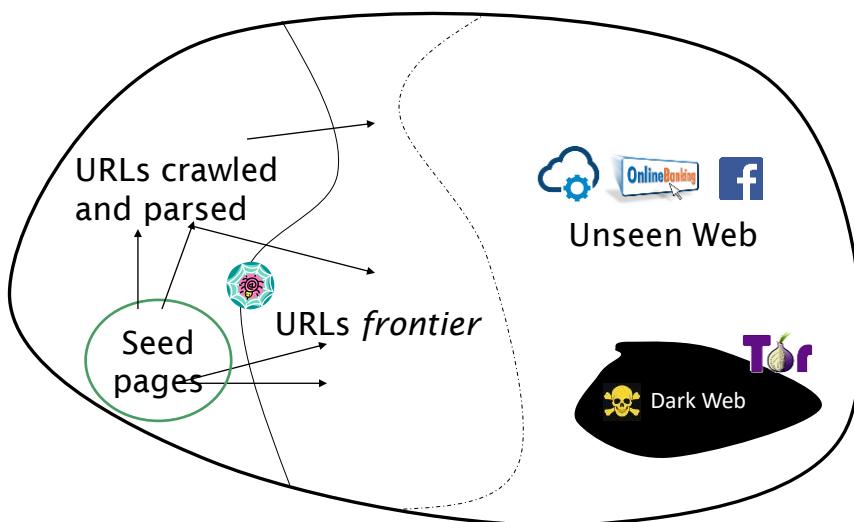


Walid Magdy, TTDS 2021/2022



15

## Web Crawling



Walid Magdy, TTDS 2021/2022



16

## Basic Crawler Operation

- Begin with known “seed” URLs
- Fetch and parse them
  - Extract URLs they point to
  - Place the extracted URLs on a queue
- Fetch one URL from the queue
- Repeat

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

17

## What Any Crawler Must Do

- Be Polite: Respect implicit and explicit politeness considerations
  - Only crawl allowed pages
    - respect `robots.txt`
  - Avoid hitting any site too often
- Be Robust: Be immune to spider traps and other malicious behaviour from web servers
  - Be careful to spams (link farms)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

18

## What Any Crawler Should Do

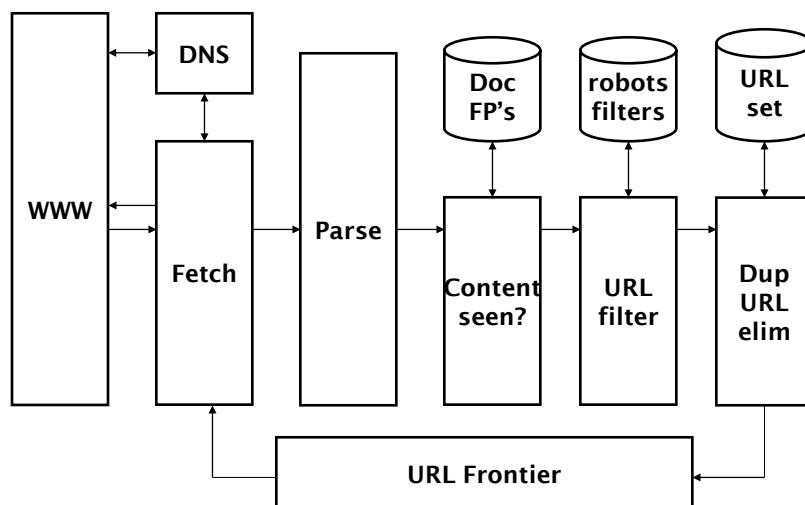
- Be capable of distributed operation
  - designed to run on multiple distributed machines
- Be scalable: designed to increase the crawl rate by adding more machines
- Performance/efficiency: permit full use of available processing and network resources
- Fetch pages of “higher quality” first
- Freshness/Continuous operation: Continue fetching fresh copies of a previously fetched page
- Extensible: Adapt to new data formats, protocols

Walid Magdy, TTDS 2021/2022



19

## Basic Crawler Architecture



Walid Magdy, TTDS 2021/2022



20

## Processing Steps in Crawling

1. Pick a URL from the frontier
2. Fetch the document at the URL
3. Parse the document
  1. Extract links from it to other docs (URLs)
4. Check if document has content already seen
  1. If not, add to indexes
5. For each extracted URL
  1. Ensure it passes certain URL filter tests
  2. Check if it is already in the frontier (duplicate URL elimination)

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

21

## URL Frontier

- Can include multiple pages from the same host
- Must avoid trying to fetch them all at the same time
- Must try to keep all crawling threads busy

Walid Magdy, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

22

## Explicit and Implicit Politeness

- Explicit politeness: specifications from webmasters on what portions of site can be crawled
  - `robots.txt`
- Implicit politeness: even with no specification, avoid hitting any site too often

```
User-agent: *
Disallow: /yoursite/temp/

User-agent: searchengine
Disallow:
```

- No robot should visit any URL starting with "/yoursite/temp/", except the robot called "searchengine"



## URL Frontier: 2 Main Considerations

- Politeness: do not hit a web server too frequently
- Priority/Freshness: crawl some pages more often than others
  - Pages whose content changes often (e.g. News sites)
- These goals may conflict each other.
  - e.g., simple priority queue fails – many links out of a page go to its own site, creating a burst of accesses to that site.
- Even if we restrict only one thread to fetch from a host, can hit it repeatedly
- Common heuristic: insert time gap between successive requests to a host that is  $>>$  time taken in most recent fetch from that host



## Summary

- History of Web search
- Basics of web search
- Usage of web search
- SEO
- Web crawling

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

25

## Resources

- Text book 1: Intro to IR, Chapter 19
- Text Book 2: IR in Practice: Chapter 3
- YouTube Videos (nice to watch)
  - How Search Works. Google  
<https://www.youtube.com/watch?v=BNHR6IQJGZs>
  - The Evolution of Search. Google  
<https://www.youtube.com/watch?v=mTBShtTwCnD4>
  - What Is The Deep Web?. Mashable  
[https://www.youtube.com/watch?v=\\_UOK7aRmUtw](https://www.youtube.com/watch?v=_UOK7aRmUtw)
  - Most popular search engines over time  
<https://www.youtube.com/watch?v=1a3WL1iOvnE>
  - This is How Much YouTube Pays Me  
<https://www.youtube.com/watch?v=l3MeCEwVxB0>

*Walid Magdy, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

26



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Comparing Text Corpora

Instructor:  
**Björn Ross**

10-Nov-2021

1

## Initial Text Analysis

- Scenario: you are given access to a new dataset
  - 2 corpora, each contains thousands of plain text files
  - You want to understand and quantify:
    - What is the *content* of these documents? What are they *about*?
    - How does the content of these corpora *differ*?
- What are some things you might try first?



THE UNIVERSITY  
of EDINBURGH

## Lecture Objectives

- Analyze text corpora
  - Content analysis background
  - Word-level differences
  - Dictionaries and Lexicons
  - Topic modeling
  - Annotation + classification



## Content Analysis

- Goal: given some documents determine
  - What are the types of content present? (themes/topics)
  - Which documents contain which topics?
- Traditionally a manual process
  1. Read a subset of documents, define themes/topics
  2. Determine consistent coding\* methodology
  3. Read all documents and label them according to codes
  4. Check agreement between human coders
  5. Settle disagreements via a third-party
  6. Analyze resulting annotations



## Content Analysis

- Can this process be automated?
  - Yes, to an extent
- Should this process be automated?
  - Humans are better than machines at this task (for now?)
  - Computers are *much, much* faster
    - Avg. human reading speed: 250 wpm
    - Assume 1K words/document, 50K documents...
      - Average person needs > 4 months to read
      - This is a **relatively small** corpus for modern NLP
    - Modern computers can process millions of words/second



## Automated Content Analysis

- Single corpus/class
    - Word frequency analysis
    - Dictionaries & Lexicons
    - Topic modelling
  - Multiple corpora/classes
    - Word-level differences
    - Dominance Scores
    - Topic-level differences
- 
- ```
graph LR; A[Single corpus/class] <--> B[Multiple corpora/classes]; A <--> C[Word frequency analysis]; A <--> D[Dictionaries & Lexicons]; A <--> E[Topic modelling]; B <--> F[Word-level differences]; B <--> G[Dominance Scores]; B <--> H[Topic-level differences]
```



## Word Level Analysis

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

9

# Word frequency analysis

- Very simple starting point
  1. Preprocess as usual (lowercasing? stemming?...)
  2. Count words
  3. Normalize by document length
  4. Average across all documents



*Björn Ross, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

10

## Word-level Differences

- Which words best characterize set of documents (such as a corpus or class)?
  - Need a reference corpus
- Some methods to do this:
  - Mutual information
  - Chi squared
- Can also be used for *feature selection*

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Mutual Information

- $I(X;Y)$ 
  - How much can I learn about Y by observing X?
  - Is the same as *information gain*
  - Is **not** the same as *pointwise mutual information*
- We want to learn about important words in our class
- What should X and Y be?
  - $X = U =$  document contains term t (Boolean)
  - $Y = C =$  class is the target class (Boolean)

$$I(U;C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U = e_t, C = e_c) \log_2 \frac{P(U = e_t, C = e_c)}{P(U = e_t)P(C = e_c)}$$

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

12

## Mutual Information

$$I(U;C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U = e_t, C = e_c) \log_2 \frac{P(U = e_t, C = e_c)}{P(U = e_t)P(C = e_c)}$$

- Given a corpus and a term, how do we estimate the probability of this term appearing in a random document in the corpus?

Source: Manning, Raghavan, and Schütze, 2008

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Mutual Information

$$I(U;C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U = e_t, C = e_c) \log_2 \frac{P(U = e_t, C = e_c)}{P(U = e_t)P(C = e_c)}$$

- Given count data for 2 classes, can be computed as:

$$\begin{aligned} I(U;C) &= \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ &\quad + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}} \end{aligned}$$

Source: Manning, Raghavan, and Schütze, 2008

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

14

## Mutual Information

$$I(U;C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_1.N_1} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_0.N_1} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_1.N_0} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_0.N_0}$$

- Example:
  - What is  $I(U;C)$  given these values?

|                               |                         |                         |
|-------------------------------|-------------------------|-------------------------|
|                               | $e_c = e_{poultry} = 1$ | $e_c = e_{poultry} = 0$ |
| $e_t = e_{\text{export}} = 1$ | $N_{11} = 49$           | $N_{10} = 27,652$       |
| $e_t = e_{\text{export}} = 0$ | $N_{01} = 141$          | $N_{00} = 774,106$      |

Example: Manning, Raghavan, and Schütze, 2008

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Mutual Information for News Data

| UK            |        | China            |        | poultry       |        |
|---------------|--------|------------------|--------|---------------|--------|
| london        | 0.1925 | china            | 0.0997 | poultry       | 0.0013 |
| uk            | 0.0755 | chinese          | 0.0523 | meat          | 0.0008 |
| british       | 0.0596 | beijing          | 0.0444 | chicken       | 0.0006 |
| stg           | 0.0555 | yuan             | 0.0344 | agriculture   | 0.0005 |
| britain       | 0.0469 | shanghai         | 0.0292 | avian         | 0.0004 |
| plc           | 0.0357 | hong             | 0.0198 | broiler       | 0.0003 |
| england       | 0.0238 | kong             | 0.0195 | veterinary    | 0.0003 |
| pence         | 0.0212 | xinhua           | 0.0155 | birds         | 0.0003 |
| pounds        | 0.0149 | province         | 0.0117 | inspection    | 0.0003 |
| english       | 0.0126 | taiwan           | 0.0108 | pathogenic    | 0.0003 |
| <i>coffee</i> |        | <i>elections</i> |        | <i>sports</i> |        |
| coffee        | 0.0111 | election         | 0.0519 | soccer        | 0.0681 |
| bags          | 0.0042 | elections        | 0.0342 | cup           | 0.0515 |
| growers       | 0.0025 | polls            | 0.0339 | match         | 0.0441 |
| kg            | 0.0019 | voters           | 0.0315 | matches       | 0.0408 |
| colombia      | 0.0018 | party            | 0.0303 | played        | 0.0388 |
| brazil        | 0.0016 | vote             | 0.0299 | league        | 0.0386 |
| export        | 0.0014 | poll             | 0.0225 | beat          | 0.0301 |
| exporters     | 0.0013 | candidate        | 0.0202 | game          | 0.0299 |
| exports       | 0.0013 | campaign         | 0.0202 | games         | 0.0284 |
| crop          | 0.0012 | democratic       | 0.0198 | team          | 0.0264 |

Example: Manning, Raghavan, and Schütze, 2008

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

17

## Chi-squared

- Hypothesis testing approach
- $H_0$ : Term appearance is independent from a document's class
  - i.e.,  $P(U = 1, C = 1) = P(U = 1)P(C = 1)$
- Compute:

$$X^2(\mathbb{D}, t, c) = \sum_{e_t \in \{0,1\}} \sum_{e_c \in \{0,1\}} \frac{(N_{e_t e_c} - E_{e_t e_c})^2}{E_{e_t e_c}}$$

- Or to directly plug in values like before:

$$X^2(\mathbb{D}, t, c) = \frac{(N_{11} + N_{10} + N_{01} + N_{00}) \times (N_{11}N_{00} - N_{10}N_{01})^2}{(N_{11} + N_{01}) \times (N_{11} + N_{10}) \times (N_{10} + N_{00}) \times (N_{01} + N_{00})}$$

## Chi-squared

$$X^2(\mathbb{D}, t, c) = \frac{(N_{11} + N_{10} + N_{01} + N_{00}) \times (N_{11}N_{00} - N_{10}N_{01})^2}{(N_{11} + N_{01}) \times (N_{11} + N_{10}) \times (N_{10} + N_{00}) \times (N_{01} + N_{00})}$$

- Example
  - What is the value of  $X^2$  given the example data?

|                                |                                                           |
|--------------------------------|-----------------------------------------------------------|
| $e_c = e_{\text{poultry}} = 1$ | $e_c = e_{\text{poultry}} = 0$                            |
| $e_t = e_{\text{export}} = 1$  | $N_{11} = 49$                                             |
| $e_t = e_{\text{export}} = 0$  | $N_{10} = 27,652$<br>$N_{01} = 141$<br>$N_{00} = 774,106$ |

# Dictionaries and Lexicons

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

21

## Dictionaries and Lexicons

- What if we know what we are looking for?
- Dictionaries (lexicons) are prebuilt mappings
  - Category -> word list
  - E.g., a tiny sentiment lexicon:
    - Positive: good, great, happy, amazing, wonderful, best, incredible
    - Negative: terrible, horrible, bad, awful, nasty, gross, worst, poor
- Domain can be important
  - “**unpredictable** movie plot” ✓
  - “**unpredictable** coffee pot” ✗

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

22

## Dictionaries and Lexicons

- How to get a score per category?

$$\frac{\text{num\_dictionary\_words\_in\_document}}{\text{num\_total\_words\_in\_document}}$$

- That's it!
- Can also be used as machine learning features
- A more advanced approaches to quantifying categories (optional reading)
  - <https://www.ncbi.nlm.nih.gov/pubmed/28364281>

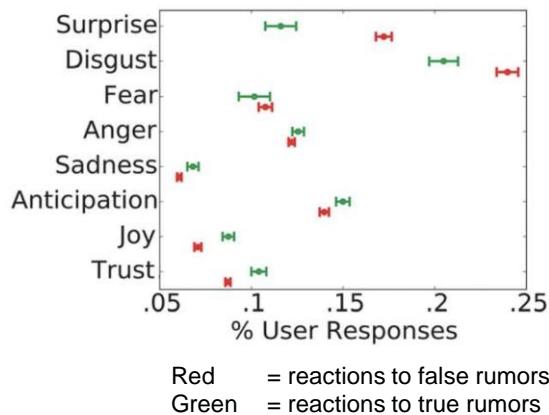


## Some Dictionaries

- LIWC (Pennebaker et al. 2015)
- General Inquirer (Stone 1997)
- Roget's Thesaurus Categories
- VADER (Hutto and Gilbert, 2014)
- Sentiwordnet (Esuli and Sebastiani 2006)
- Wordnet Domains (Magnini and Cavaglia, 2000)
- EmoLex (Mohammad and Turney, 2010)
- Empath (Fast et al., 2016)
- Personal Values Lexicon (Wilson et al., 2018)
- ...



## Reactions to Rumor Tweets with EmoLex



Vosoughi, Roy, and Aral, 2018

Björn Ross, TTDS 2021/2022



25

## Dominance Scores

- The dominance score for a category w.r.t. a corpus:

$$\frac{\text{category\_score\_in\_target\_corpus}}{\text{category\_score\_in\_background\_corpus}}$$

- From Mihalcea and Pulman, 2009

Björn Ross, TTDS 2021/2022



26

## LIWC category dominance scores

| Truthful   |       |         |       | Deceptive  |       |          |       |
|------------|-------|---------|-------|------------|-------|----------|-------|
| Interviews |       | Trials  |       | Interviews |       | Trials   |       |
| Class      | Score | Class   | Score | Class      | Score | Class    | Score |
| Metaphor   | 2.98  | You     | 3.99  | Assent     | 4.81  | Anger    | 2.61  |
| Money      | 2.74  | Family  | 3.07  | Past       | 2.59  | Anxiety  | 2.61  |
| Inhibition | 2.74  | Home    | 2.45  | Sexual     | 2.00  | Certain  | 2.28  |
| Home       | 2.13  | Humans  | 1.87  | Other      | 1.87  | Death    | 1.96  |
| Humans     | 2.02  | Posemo  | 1.81  | Motion     | 1.68  | Physical | 1.77  |
| Family     | 1.96  | Insight | 1.64  | Negemo     | 1.44  | Negemo   | 1.52  |

Pérez-Rosas et al, 2015

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

27

## Topic Level Analysis

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

28

12

## Intro to Topic Modelling

- Goals are similar to traditional content analysis:
  - What are the main themes/topics in this corpus?
  - Which documents contain which topics?

Björn Ross, TTDS 2021/2022



29

## Topic Models

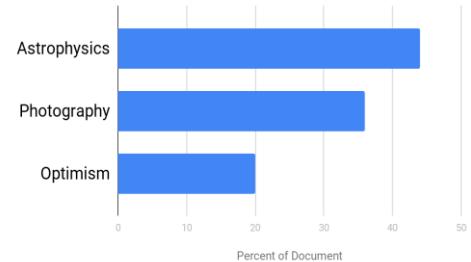
The New York Times

### Expected Soon: First-Ever Photo of a Black Hole

Have astronomers finally recorded an image of a black hole? The world will know on Wednesday.



Topic Distribution



30

Björn Ross, TTDS 2021/2022

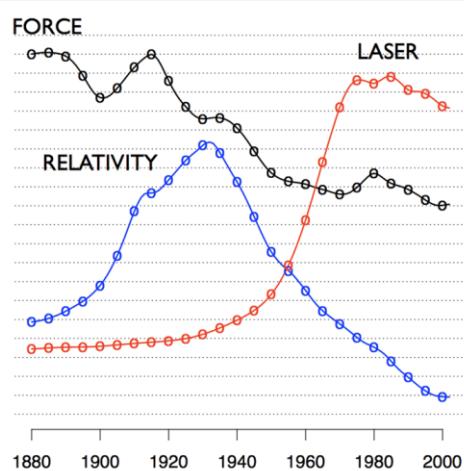
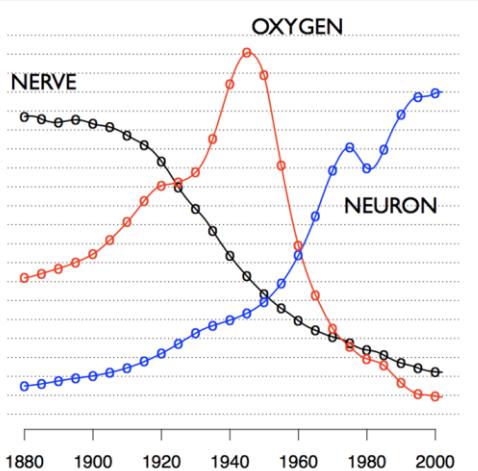


30

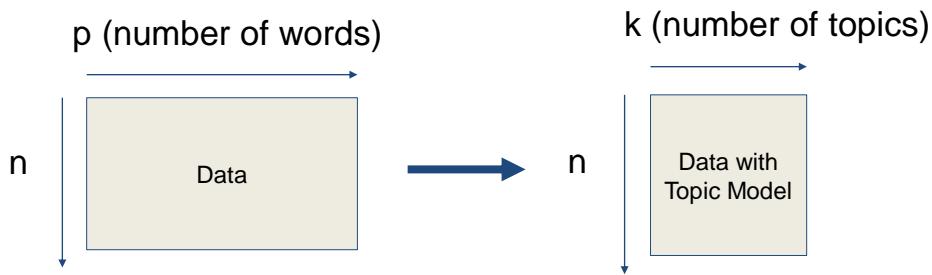
|             |              |              |             |
|-------------|--------------|--------------|-------------|
| human       | evolution    | disease      | computer    |
| genome      | evolutionary | host         | models      |
| dna         | species      | bacteria     | information |
| genetic     | organisms    | diseases     | data        |
| genes       | life         | resistance   | computers   |
| sequence    | origin       | bacterial    | system      |
| gene        | biology      | new          | network     |
| molecular   | groups       | strains      | systems     |
| sequencing  | phylogenetic | control      | model       |
| map         | living       | infectious   | parallel    |
| information | diversity    | malaria      | methods     |
| genetics    | group        | parasite     | networks    |
| mapping     | new          | parasites    | software    |
| project     | two          | united       | new         |
| sequences   | common       | tuberculosis | simulations |

Example from  
David Blei

31

**"Theoretical Physics"****"Neuroscience"**Example from  
David Blei

## Dimensionality Reduction



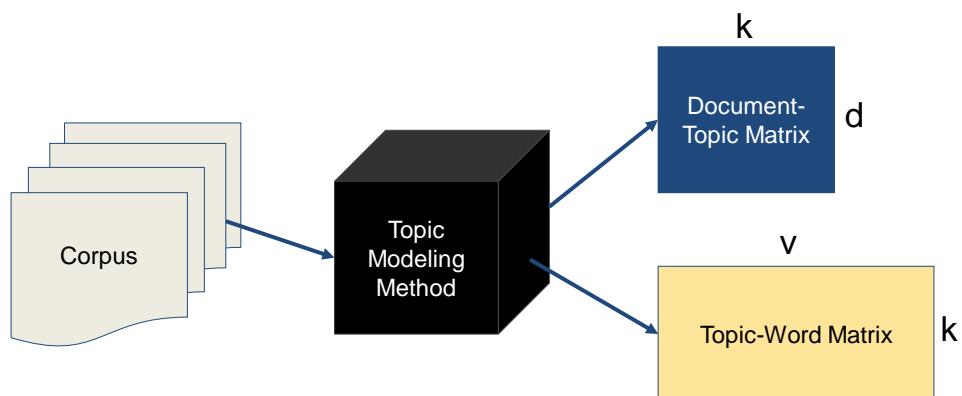
33

Björn Ross, TTDS 2021/2022



33

## Topic Modeling



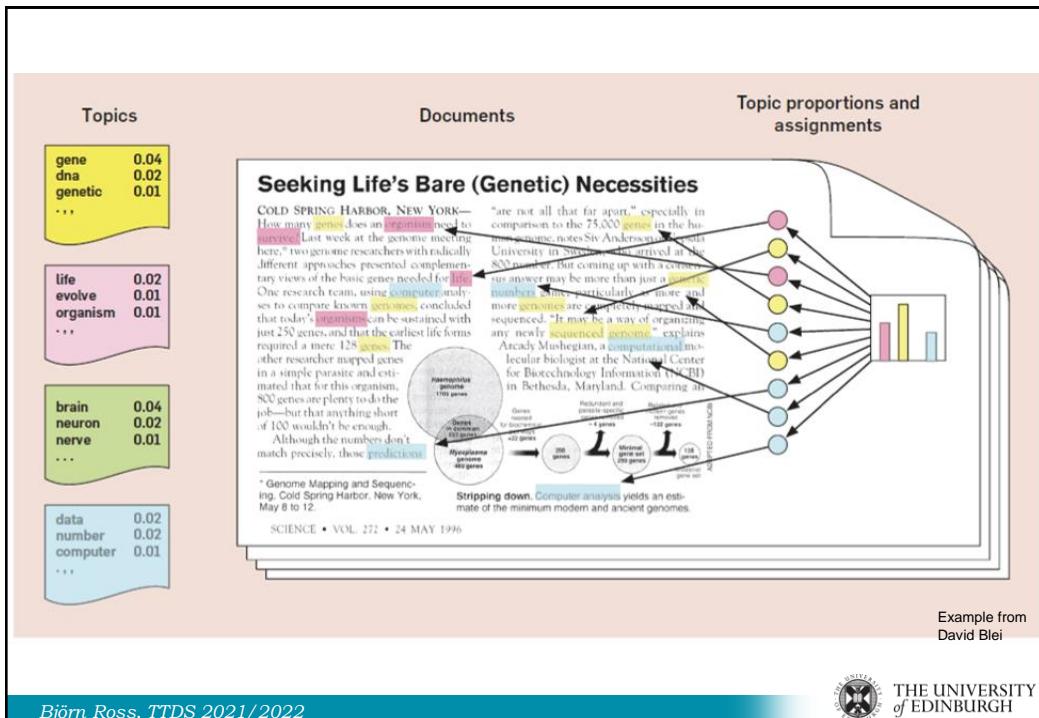
34

Björn Ross, TTDS 2021/2022



34

15



35

## Topic Models

- Most often used for text data, but can also be applied in other settings:
  - Bioinformatics (Liu et al. 2016)
  - Computer code (McBurney et al. 2014)
  - Music (Hu and Saul 2009)
  - Network data (Cha and Cho 2014)

## Topic Modeling Methods

- Most popular: Latent Dirichlet Allocation (LDA)
  - Introduced by David Blei, Andrew Ng, and Michael Jordan (2003)
- Other methods include
  - pLSI
  - PCA-based methods
  - Non-negative matrix factorization
  - Deep learning based topic modeling
  - ...

37

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

37

## Topic Modeling Methods

- Most popular: Latent Dirichlet Allocation (LDA)
  - Introduced by David Blei, Andrew Ng, and Michael Jordan (2003)
- Other methods include
  - pLSI
  - PCA-based methods
  - Non-negative matrix factorization
  - Deep learning based topic modeling
  - ...

38

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

38

17

## Latent Dirichlet Allocation (LDA)

- More details coming up in next lecture...

39

*Björn Ross, TTDS 2021/2022*



THE UNIVERSITY  
of EDINBURGH

39



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Comparing Text Corpora (2)

Instructor:  
**Björn Ross**

10-Nov-2021

1

## LDA Overview



## Background: Plate Notation

Make a  
basket



3

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

3

## Background: Plate Notation

Basketball  
shooting  
accuracy

Make a  
basket



4

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

4

2

## Background: Plate Notation

Basketball shooting accuracy

Make first basket

Make second basket

Make third basket



Björn Ross, TTDS 2021/2022



5

## Background: Plate Notation

Basketball shooting accuracy

Make nth basket

N

Björn Ross, TTDS 2021/2022



6

3

## Latent Dirichlet Allocation

- Let's start with a very simple model
- We will work our way up to the full LDA model

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Unigram Model

w is a word  
N words in a document

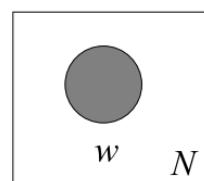


Figure from  
Blei et al 2003

8

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

4

## Unigram Model

w is a word  
N words in a document  
M documents in a corpus

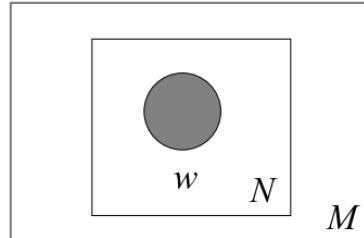


Figure from  
Blei et al 2003

9

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

## Unigram Model

w is a word  
N words in a document  
M documents in a corpus  
**w** is a vector of words (i.e. doc)

$$p(\mathbf{w}) = \prod_{n=1}^N p(w_n)$$

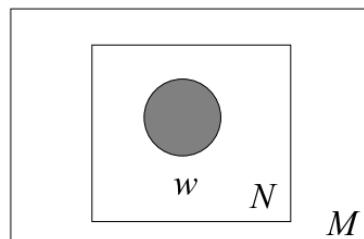


Figure from  
Blei et al 2003

10

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

10

## Probability with a Unigram Model

$$p(\mathbf{w}) = \prod_{n=1}^N p(w_n)$$

What is the probability of the example sentence?

“My dog barked at another dog.”

| word        | my  | at  | dog | another | barked |
|-------------|-----|-----|-----|---------|--------|
| probability | .10 | .10 | .05 | .04     | .03    |

11

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Unigram Model...

- What is the point of making these models more complex?
- Why not just use the basic unigram model for everything?
- Remember:
  - Higher text probability **doesn't imply a better model**
  - We want to **accurately describe** the data
    - → higher probability for *real* documents, lower probability for noise

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Mixture of Unigrams Model

$z$  is the topic of a document

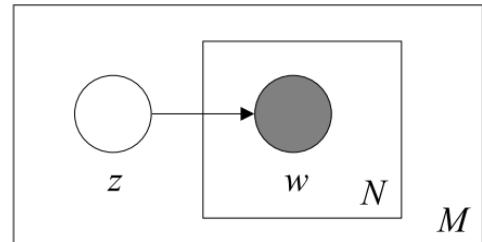


Figure from  
Blei et al 2003

14

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

14

## Mixture of Unigrams Model

$z$  is the topic of a document

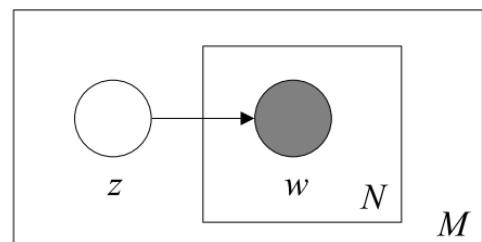


Figure from  
Blei et al 2003

15

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

15

## Probability with Mixture of Unigrams

$$p(\mathbf{w}) = \sum_z p(z) \prod_{n=1}^N p(w_n | z).$$

What is the probability of the sentence?

Ignore stopwords: "my", "after", "the"

"My dog chased after the bus."

| $w_i$                          | cat | dog | chased | car | bus |
|--------------------------------|-----|-----|--------|-----|-----|
| $P(w_i   z = \text{pets})$     | .20 | .30 | .10    | .01 | .01 |
| $P(w_i   z = \text{vehicles})$ | .01 | .01 | .10    | .30 | .20 |

$p(z = \text{pets}) = 0.6,$  16  
 $p(z = \text{vehicles}) = 0.4$

Björn Ross, TTDS 2021/2022



16

## Probabilistic Latent Semantic Indexing

d is a document ID

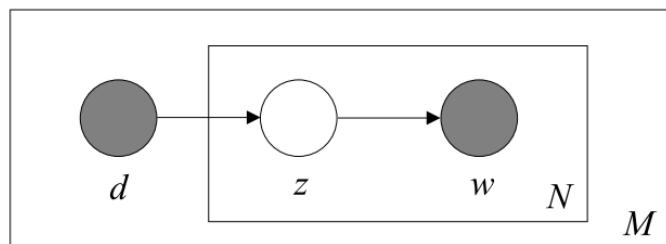


Figure from  
Blei et al 2003

18

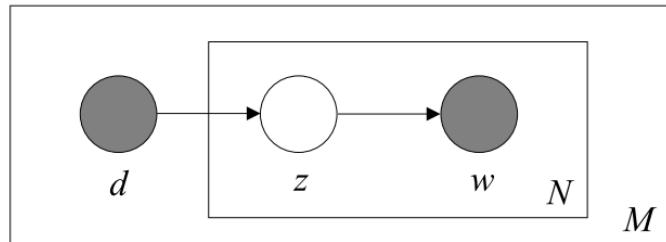
Björn Ross, TTDS 2021/2022



18

# Probabilistic Latent Semantic Indexing

$d$  is a document ID



$$p(d, w_n) = p(d) \sum_z p(w_n | z) p(z | d)$$

Figure from Blei et al 2003

19

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

## Probability with pLSI

| $w_i$              | cat | sat | down | car | broke |
|--------------------|-----|-----|------|-----|-------|
| $p(w_i   z = t_1)$ | .2  | .1  | .05  | .01 | .1    |
| $p(w_i   z = t_2)$ | .01 | .05 | .1   | .3  | .1    |

$d_1$  “The **cat** sat down.”

|                        |     |
|------------------------|-----|
| $p(d = d_1)$           | .01 |
| $p(z = t_1   d = d_1)$ | .6  |
| $p(z = t_2   d = d_1)$ | .4  |

| $w_i$              | cat | sat | down | car | broke |
|--------------------|-----|-----|------|-----|-------|
| $p(w_i   z = t_1)$ | .2  | .1  | .05  | .01 | .1    |
| $p(w_i   z = t_2)$ | .01 | .05 | .1   | .3  | .1    |

What is the joint probability of the document and the word “cat”? <sup>20</sup>

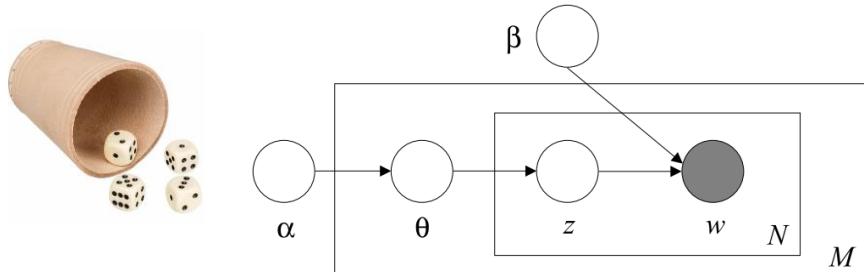
Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

20

## Latent Dirichlet Allocation



$\theta$  is the distribution over topics in a document  
 $\alpha$  is the parameter of a Dirichlet distribution giving possible topic distributions within documents  
 $\beta$  gives word distributions within topics

Figure from Blei et al 2003

22

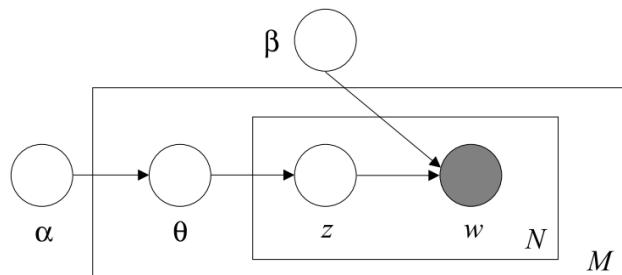
Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

22

## Latent Dirichlet Allocation



$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

Figure from Blei et al 2003

23

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

23

## Latent Dirichlet Allocation

$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

$$p(\mathbf{w} | \alpha, \beta) = \int p(\theta | \alpha) \left( \prod_{n=1}^N \sum_{z_n} p(z_n | \theta) p(w_n | z_n, \beta) \right) d\theta$$

$$p(\mathcal{D} | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) d\theta_d$$

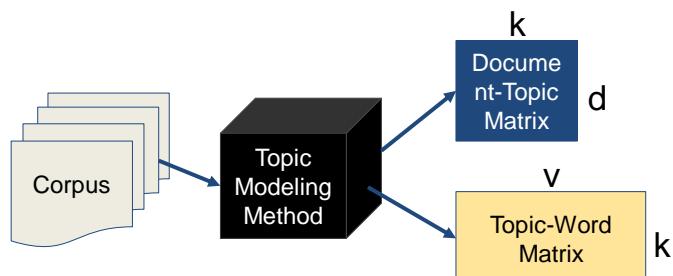
Björn Ross, TTDS 2021/2022



26

## Model Inference

- Want to learn the model parameters
- Exact inference becomes intractable



27

Björn Ross, TTDS 2021/2022



27

## Model Inference

- Instead, use an approximate method such as:
  - Gibbs sampling
  - Variational Inference

28

Björn Ross, TTDS 2021/2022



28

## Gibbs Sampling for LDA

Goal: Learn  $\Phi, \theta$  given a set of documents D

$\Phi$  = topic-word probabilities

$\theta$  = document-topic probabilities

Known:

corpus,  $\alpha, \beta$  and the probability that a word is from a topic conditional on the assignments of all other words to topics

$$P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot) \propto \frac{C_{w_i j}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \cdot \frac{C_{d_i j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i t}^{DT} + T\alpha}$$

29

Björn Ross, TTDS 2021/2022



29

## Gibbs Sampling for LDA

Want to learn  $\Phi, \theta$  given a set of documents D

1. Assign each word a topic randomly
  2. Calculate count matrices
  3. Repeat until convergence:
    - For every document d
      - For every word i
        - Decrement count matrices  $C^{WT}$  and  $C^{DT}$  for current topic assignment
        - Sample a new topic assignment
        - Increment count matrices  $C^{WT}$  and  $C^{DT}$  for new topic assignment
4. Calculate  $\Phi$  and  $\theta$

30

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

30

## Gibbs Sampling for LDA

d1 Green eggs and ham.  
d2 Ham and green peppers.  
d3 Ham and cheese.

Green eggs and ham.  
Ham and green peppers.  
Ham and cheese.

Random  
initialization.

31

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

31

13

## Gibbs Sampling for LDA

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 1     | 1    | 1   | 1   | 1       | 1      |
| t2       | 1     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
 Ham and green peppers.  
 Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 1  | 1  | 1  |
| t2       | 1  | 1  | 1  |

32

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

32

## Gibbs Sampling for LDA

Assume (for the moment)  $\alpha = \beta = 0$

| $\theta$ | green | eggs | and  | ham  | peppers | cheese |
|----------|-------|------|------|------|---------|--------|
| t1       | 0.17  | 0.17 | 0.17 | 0.17 | 0.17    | 0.17   |
| t2       | 0.20  | 0.00 | 0.40 | 0.40 | 0.00    | 0.00   |

Green eggs and ham.  
 Ham and green peppers.  
 Ham and cheese.

| $\Phi$ | d1   | d2   | d3   |
|--------|------|------|------|
| t1     | 0.50 | 0.50 | 0.66 |
| t2     | 0.50 | 0.50 | 0.33 |

33

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

33

## Gibbs Sampling for LDA

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 1     | 1    | 1   | 1   | 1       | 1      |
| t2       | 1     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
 Ham and green peppers.  
 Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 2  | 2  | 2  |
| t2       | 2  | 2  | 1  |

34

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

34

## Gibbs Sampling for LDA

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 1     | 1    | 1   | 1   | 1       | 1      |
| t2       | 1     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
 Ham and green peppers.  
 Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 2  | 2  | 2  |
| t2       | 2  | 2  | 1  |

35

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

35

15

## Gibbs Sampling for LDA

$$\frac{C_{w_j}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_j}^{DT} + \alpha}{\sum_{t=1}^T C_{dj}^{DT} + T\alpha}$$

Assume (for the moment)  $\alpha = \beta = 0$

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 0     | 1    | 1   | 1   | 1       | 1      |
| t2       | 1     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
Ham and green peppers.  
Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 1  | 2  | 2  |
| t2       | 2  | 2  | 1  |

36

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

36

## Gibbs Sampling for LDA

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 0     | 1    | 1   | 1   | 1       | 1      |
| t2       | 2     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
Ham and green peppers.  
Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 1  | 2  | 2  |
| t2       | 3  | 2  | 1  |

37

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

37

## Gibbs Sampling for LDA

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 0     | 1    | 1   | 1   | 1       | 1      |
| t2       | 2     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
Ham and green peppers.  
Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 1  | 2  | 2  |
| t2       | 3  | 2  | 1  |

38

Björn Ross, TTDS 2021/2022



38

## Gibbs Sampling for LDA

$$\frac{C_{w_j}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_j}^{DT} + \alpha}{\sum_{t=1}^T C_{dj}^{DT} + T\alpha}$$

Assume (for the moment)  $\alpha = \beta = 0$

| $C^{WT}$ | green | eggs | and | ham | peppers | cheese |
|----------|-------|------|-----|-----|---------|--------|
| t1       | 0     | 0    | 1   | 1   | 1       | 1      |
| t2       | 2     | 0    | 2   | 2   | 0       | 0      |

Green eggs and ham.  
Ham and green peppers.  
Ham and cheese.

| $C^{DT}$ | d1 | d2 | d3 |
|----------|----|----|----|
| t1       | 0  | 2  | 2  |
| t2       | 3  | 2  | 1  |

39

Björn Ross, TTDS 2021/2022



39

## Gibbs Sampling for LDA

$$\frac{C_{w_j}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_j}^{DT} + \alpha}{\sum_{t=1}^T C_{dj}^{DT} + T\alpha}$$

| $C^{WT} + \alpha$ | green | eggs | and  | ham  | peppers | cheese |
|-------------------|-------|------|------|------|---------|--------|
| t1                | 0.01  | 0.01 | 1.01 | 1.01 | 1.01    | 1.01   |
| t2                | 2.01  | 0.01 | 2.01 | 2.01 | 0.01    | 0.01   |

Green eggs and ham.  
 Ham and green peppers.  
 Ham and cheese.

| $C^{DT} + \beta$ | d1   | d2   | d3   |
|------------------|------|------|------|
| t1               | 0.01 | 2.01 | 2.01 |
| t2               | 3.01 | 2.01 | 1.01 |

40

Björn Ross, TTDS 2021/2022



40

## Gibbs Sampling for LDA

- Repeat until convergence
- Probabilistic algorithm – results depend on random initialisation and random samples!

Björn Ross, TTDS 2021/2022



41

## Topic Modeling Examples

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

42

## What do students look for in a professor?

| Topic                 | Sample words                                                                                                      |
|-----------------------|-------------------------------------------------------------------------------------------------------------------|
| Approachability       | prof, fair, clear, helpful, teaching, approachable, nice, organized, extremely, friendly, super, amazing          |
| Clarity               | understand, hard, homework, office, material, clear, helpful, problems, explains, accent, questions, extremely    |
| Course Logistics      | book, study, boring, extra, nice, credit, lot, hard, attendance, make, fine, attention, pay, mandatory            |
| Enthusiasm            | teaching, passionate, awesome, enthusiastic, professors, loves, cares, wonderful, fantastic, passion              |
| Expectations          | hard, work, time, lot, comments, tough, expects, worst, stuff, avoid, horrible, classes                           |
| Helpfulness           | helpful, nice, recommend, cares, super, understanding, kind, extremely, effort, sweet, friendly, approachable     |
| Humor                 | guy, funny, fun, awesome, cool, entertaining, humor, hilarious, jokes, stories, love, hot, enjoyable              |
| Interestingness       | interesting, material, recommend, lecturer, engaging, classes, knowledgeable, enjoyed, loved, topics              |
| Readings/ Discussions | readings, papers, writing, ta, interesting, discussions, grader, essays, boring, books, participation             |
| Study Material        | exams, notes, questions, material, textbook, hard, slides, study, answer, clear, tricky, attend, long, understand |

Azab, Mihalcea, and Abernathy, 2016

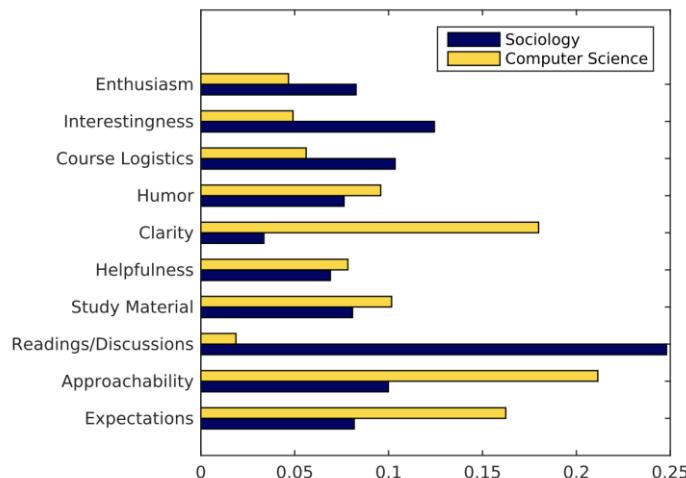
Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

43

## What do students look for in a professor?



Azab, Mihalcea, and Abernathy, 2016

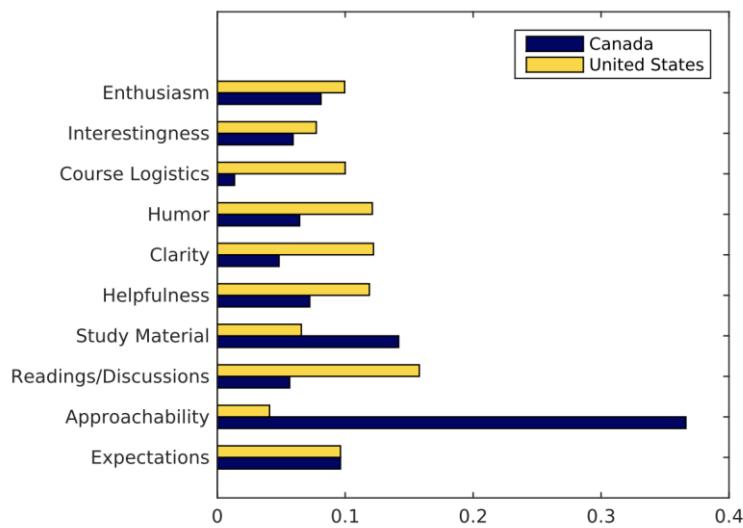
Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

44

## What do students look for in a professor?



Azab, Mihalcea, and Abernathy, 2016

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

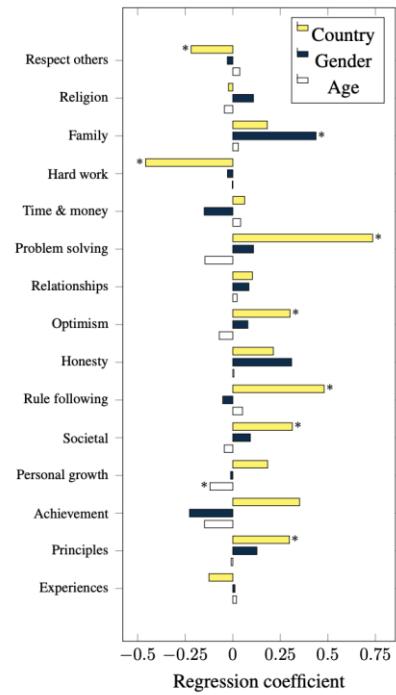
45

# How do personal attributes relate to values?

| Theme           | Example Words                           |
|-----------------|-----------------------------------------|
| Respect others  | people, respect, care, human, treat     |
| Religion        | god, heart, belief, religion, right     |
| Family          | family, parent, child, husband, mother  |
| Hard Work       | hard, work, better, honest, best        |
| Time & Money    | money, work, time, day, year            |
| Problem solving | consider, decision, situation, problem  |
| Relationships   | family, friend, relationship, love      |
| Optimism        | enjoy, happy, positive, future, grow    |
| Honesty         | honest, truth, lie, trust, true         |
| Rule following  | moral, rule, principle, follow          |
| Societal        | society, person, feel, thought, quality |
| Personal Growth | personal, grow, best, decision, mind    |
| Achievement     | heart, achieve, complete, goal          |
| Principles      | important, guide, principle, central    |
| Experiences     | look, see, experience, choose, feel     |

Wilson, Mihalcea, Boyd, and Pennebaker 2016

Björn Ross, TTDS 2021/2022



46

## Annotation + Classification

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

47

## Annotation + Classification

- Method 1: Traditional Supervised Learning
  - Annotate representative samples
  - Train a classifier
  - Apply to rest of data
- Method 2: Transfer Learning
  - Find another large, but similar dataset
  - Train a classifier on that dataset
  - *Optionally: fine-tune classifier to your smaller dataset*
  - Apply to rest of your data

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

48

## After Classification

- Which features are most relevant for each class?
- What are common words/topics for each class?
- How do predicted classes relate to other variables?
- *More about text classification coming up next week!*

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

49

## Wrap-up

- Content analysis background
- Word-level differences
- Dictionaries and Lexica
- Topic modeling
- Annotation + classification

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

50

## Readings

- [Manning: IR book section 13.5](#)
- [“Probabilistic Topic Models” by David Blei](#)
- [“Latent Dirichlet Allocation” by David Blei, Andrew Y. Ng, and Michael I. Jordan](#)
- [“Probabilistic Topic Models” by Mark Steyvers and Tom Griffiths](#)

To watch:

- [Guest lecture \(2017\) by David Blei at University of Edinburgh School of Informatics](#)

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

51



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Text Classification

Instructor:  
**Björn Ross**

17-Nov-2021

1

## Lecture Objectives

- Learn about text basics of text classification
  - Definition
  - Types
  - Methods
  - Evaluation



THE UNIVERSITY  
of EDINBURGH

2

## Text Classification

- **Text classification** is the process of classifying documents into predefined categories based on their content.
- Input: Text (document, article, sentence)
- Task: Classify into predefined one/multiple categories
- Categories:
  - Binary: relevant/irrelevant, spam .. etc.
  - Few: sports/politics/comedy/technology
  - Hierarchical: patents



## Classification is and is not

- **Classification** (a.k.a. “**categorization**”): a common technology in data science; studied within pattern recognition, statistics, and machine learning.
- Definition:  
the activity of **predicting** to which among a **predefined finite set** of groups (“classes”, or “categories”) a data item belongs to
- Formulated as the task of generating a hypothesis (or “**classifier**”, or “**model**”)

$$h : D \rightarrow C$$

where  $D = \{x_1, x_2, \dots\}$  is a domain of data items and  
 $C = \{c_1, \dots, c_n\}$  is a finite set of classes (the **classification scheme**)



## Classification is and is not

- Different from clustering, where the groups (“clusters”) and their number are not known in advance
- The membership of a data item into a class must not be determinable with certainty
  - e.g., predicting whether a natural number belongs to *Prime* or *Non-Prime* is not classification
- In text classification, data items are
  - **Textual:** e.g., news articles, emails, sentences, queries, etc.
  - **Partly textual:** e.g., Web pages

## Types of Classification

- **Binary:**  
item to be classified into one of two classes  
 $h : D \rightarrow C, C = \{c_1, c_2\}$ 
  - e.g., Spam/not spam, offensive/not offensive, rel/irrel
- **Single-Label Multi-Class (SLMC)**  
item to be classified into only one of  $n$  possible classes.  
 $h : D \rightarrow C, C = \{c_1, \dots, c_n\}, \text{ where } n > 2$ 
  - e.g., Sports/politics/entertainment, positive/negative/neutral
- **Multi-Label Multi-Class (MLMC)**  
item to be classified into none, one, two, or more classes  
 $h : D \rightarrow 2^C, C = \{c_1, \dots, c_n\}, \text{ where } n > 1$ 
  - e.g., Assigning CS articles to classes in the ACM Classification System
  - Usually be solved as  $n$  independent binary classification problems

## Dimension of Classification

- Text classification may be performed according to several dimensions (“axes”) orthogonal to each other
- by **topic**; by far the most frequent case, its applications are global
- by **sentiment**; useful in market research, online reputation management, social science and political science
- by **language** (a.k.a. “language identification”); useful, e.g., in query processing within search engines
- by **genre**; e.g., AutomotiveNews vs. AutomotiveBlogs, useful in website classification and others;
- by **author** (a.k.a. “authorship attribution”), by native language (“native language identification”), or by gender; useful in forensics and cybersecurity
- by **usefulness**; e.g., product reviews
- .....

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Rule-based classification

- An old-fashioned way to build text classifiers was via knowledge engineering, i.e., manually building classification rules
  - E.g., (Viagra or Sildenafil or Cialis) → Spam
  - E.g. (#MAGA or America great again) → support Trump
- Common type: dictionary-based classification
- Disadvantages:
  - Expensive to setup and to maintain
  - Depends on few keywords → bad coverage (recall)

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## Supervised-learning classification

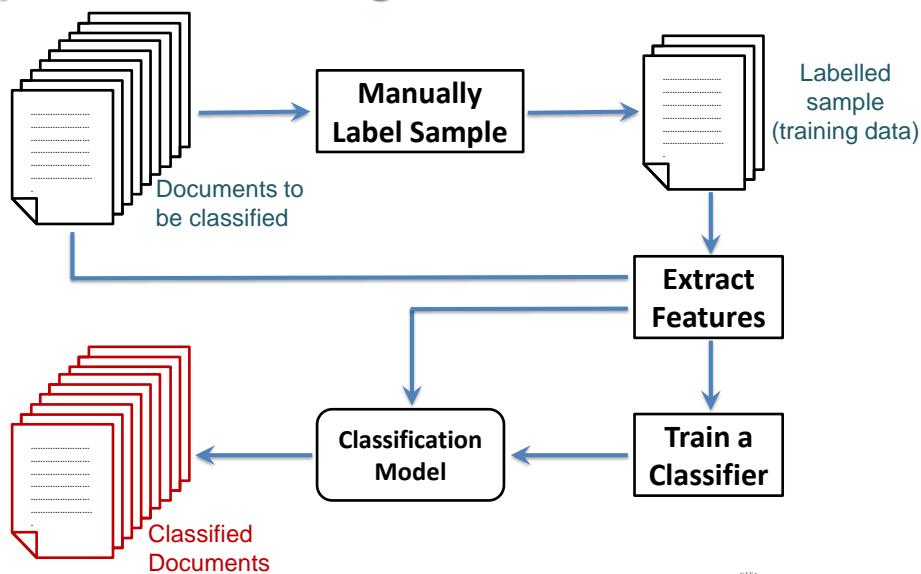
- A generic (task-independent) learning algorithm is used to train a classifier from a set of manually classified examples
- The classifier learns, from these training examples, the characteristics a new text should have in order to be assigned to class c
- Advantages:
  - Generating training examples cheaper than writing classification rules
  - Easy update to changing conditions (e.g., addition of new classes, deletion of existing classes, shifted meaning of existing classes, etc.)

Björn Ross, TTDS 2021/2022



9

## Supervised-learning classification



Björn Ross, TTDS 2021/2022



10

## Extract Features

- In order to be input to a learning algorithm (or a classifier), all training (or unlabeled) documents are converted into **vectors** in a common **vector space**
- The dimensions of the vector space are called **features**
- In order to generate a vector-based representation for a set of documents  $D$ , the following steps need to be taken
  1. Feature Extraction
  2. Feature Selection or Feature Synthesis (optional)
  3. Feature Weighting

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

11

## Step 1: Feature Extraction

- What are the features that should be different from one class to another?
- Simplest form: Bag-of-words (BOW)
  - Each term in a document is a feature
  - Feature space size = vocabulary in all docs
  - Standard IR preprocessing steps are usually applied
    - Tokenisation, stopping, stemming
- Other simple features forms:
  - Word n-grams (bigrams, trigrams, ....)
    - Much larger + more sparse
  - Sometimes char n-grams are used
    - Especially for degraded text (OCR or ASR outputs)

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

12

## Step 1: Feature Extraction

- What other text features could be used?
- Sentence structure:
  - POS (part-of-speech tags)
  - Syntactic tree structure
- Topic-based features:
  - LDA topics
  - NEs (named entities) in text
  - Links / Linked terms
- Non-textual features:
  - Average doc\sentence\word length
  - % of words start with upper-case letter
  - % of links/hashtags/emojis in text

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

13

## Step 1: Feature Extraction

- What preprocessing to apply?
  - Case-folding? `really` vs `Really` vs `REALLY`
  - Punctuation? “?”, “!”, “@”, “#”
  - Stopping? “`he`”, “`she`”, “`what`”, “`but`”
  - Stemming? “`replaced`” vs “`replacement`”
- Other Features:
  - Starts with capital letter, all caps
  - Repeated characters “`congraaaaaats`” “`help!!!!!!`”
  - Scores from dictionaries and lexicons (e.g. LIWC)
- Which to choose?
  - Classification task/application

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

14

## Step 2: Feature Selection

- Number of distinctive features = length of feature vector
- Vector can be of length in the order of  $10^6$ , and might be sparse
  - High computational cost
  - Overfitting
- What are the most important features among those?
  - e.g. Reduce from  $10^6$  to  $10^4$
- For each class, find the top representative  $k$  features for it → get the Union over all classes → reduced feature space

Björn Ross, TTDS 2021/2022



15

## Step 2: Feature Selection Functions

- Document frequency
  - % of docs in class  $c_i$  that contain the term  $t_k$
  - Very basic measure. Will select stop words as features
 
$$\#(t_k, c_i) = P(t_k | c_i)$$
- Mutual Information
  - How much we learn from the presence or absence of term  $t_k$  about whether or not a document is in class  $c_i$
  - Often used in feature selection in text classification
 
$$MI(t_k, c_i) = \sum_{c \in \{c_i, \bar{c}_i\}} \sum_{t \in \{t_k, \bar{t}_k\}} P(t, c) \cdot \log_2 \frac{P(t, c)}{P(t) \cdot P(c)}$$
- Pearson's Chi-squared ( $\chi^2$ )
  - used more in comparisons between classes

Björn Ross, TTDS 2021/2022



16

## Step 2: Feature Selection Functions

| Function               | Denoted by         | Mathematical form                                                                                                                                                                    |
|------------------------|--------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Document frequency     | $\#(t_k, c_i)$     | $P(t_k c_i)$                                                                                                                                                                         |
| DIA association factor | $z(t_k, c_i)$      | $P(c_i t_k)$                                                                                                                                                                         |
| Information gain       | $IG(t_k, c_i)$     | $\sum_{c \in \{c_i, \bar{c}_i\}} \sum_{t \in \{t_k, \bar{t}_k\}} P(t, c) \cdot \log \frac{P(t, c)}{P(t) \cdot P(c)}$                                                                 |
| Mutual information     | $MI(t_k, c_i)$     | $\log \frac{P(t_k, c_i)}{P(t_k) \cdot P(c_i)}$                                                                                                                                       |
| Chi-square             | $\chi^2(t_k, c_i)$ | $\frac{ Tr  \cdot [P(t_k, c_i) \cdot P(\bar{t}_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(\bar{t}_k, c_i)]^2}{P(t_k) \cdot P(\bar{t}_k) \cdot P(c_i) \cdot P(\bar{c}_i)}$             |
| NGL coefficient        | $NGL(t_k, c_i)$    | $\frac{\sqrt{ Tr  \cdot [P(t_k, c_i) \cdot P(\bar{t}_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(\bar{t}_k, c_i)]}}{\sqrt{P(t_k) \cdot P(\bar{t}_k) \cdot P(c_i) \cdot P(\bar{c}_i)}}$ |
| Relevancy score        | $RS(t_k, c_i)$     | $\log \frac{P(t_k c_i) + d}{P(\bar{t}_k \bar{c}_i) + d}$                                                                                                                             |
| Odds Ratio             | $OR(t_k, c_i)$     | $\frac{P(t_k c_i) \cdot (1 - P(t_k \bar{c}_i))}{(1 - P(t_k c_i)) \cdot P(t_k \bar{c}_i)}$                                                                                            |
| GSS coefficient        | $GSS(t_k, c_i)$    | $P(t_k, c_i) \cdot P(\bar{t}_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(\bar{t}_k, c_i)$                                                                                              |

Björn Ross, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

17

## Step 2: Feature Synthesis

- Matrix decomposition techniques (e.g., PCA, SVD, LSI) can be used to synthesize new features that replace the features discussed above
- Principle of distributional semantics: semantics of a word “is” the words it co-occurs with
  - Pros: the synthetic features in the new vector representation do not suffer from problems such as polysemy and synonymy
  - Cons: computationally expensive
- Word embeddings: the new wave of distributional semantics, modern approaches are based on neural networks
  - PCA: Principle component analysis
  - SVD: Singular value decomposition
  - LSA: latent semantic analysis

Björn Ross, TTDS 2021/2022

THE UNIVERSITY  
of EDINBURGH

18

## Step 2: Feature Synthesis

- Deep learning?
- Language modelling “features”
  - Tokenize text and pass to neural network layer
    - E.g., recurrent layer, convolutional layer, self-attention layer
  - Stack on 3+ more layers
  - Train a model to predict the next word (or a missing word) given previous words
  - Penultimate layer of network can be used to generate features for other language-based tasks
  - Basis for many state-of-the-art text classifiers
    - BERT, GPT, Electra, XLNet, etc.

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

19

## Step 3: Feature Weighting

- Attributing a value to feature  $t_k$  in document  $d_i$   
 This value may be
  - **binary** (representing presence/absence of  $t_k$  in  $d_i$ );
  - **numeric** (representing the importance of  $t_k$  for  $d_i$ );  
 obtained via feature weighting functions in the following two classes:
    - **unsupervised**: e.g., tfidf or BM25,
    - **supervised**: e.g.,  $tf^* MI$ ,  $tf^* x^2$
- Similarity between two vectors may be computed e.g.  
 via **cosine similarity**
- **Scaling** can be important!

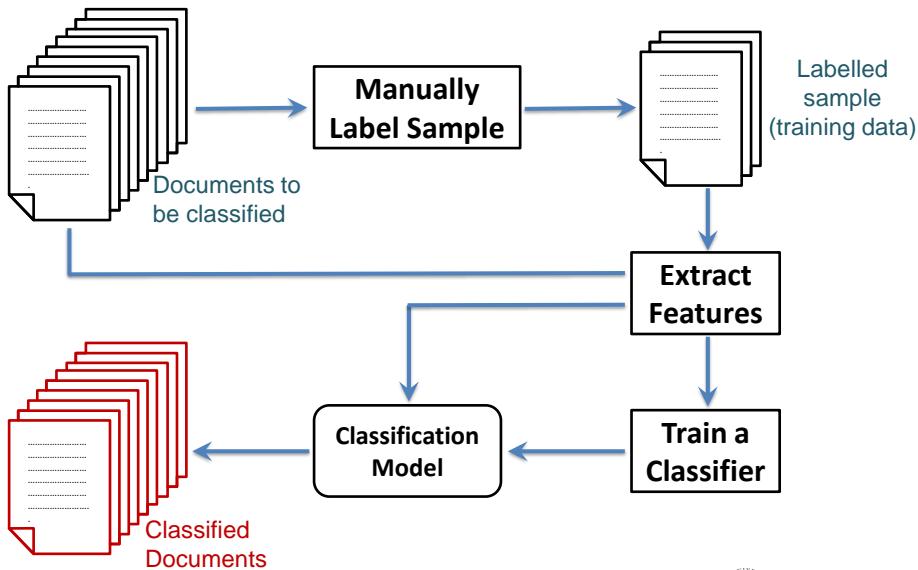
Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

20

## Supervised-learning classification



Björn Ross, TTDS 2021/2022



21

## Training a Classifier

- For **binary** classification, essentially any supervised learning algorithm can be used for training a classifier; classical choices include
  - Support vector machines (SVMs)
  - Random forests
  - Naïve Bayesian methods
  - Lazy learning methods (e.g., k-NN)
  - Logistic Regression
  - ....
- The “**No-free-lunch principle**” (Wolpert, 1996) → *there is no learning algorithm that can outperform all others in all contexts*
- Implementations need to cater for
  - the very high dimensionality
  - the sparse nature of the representations involved

Björn Ross, TTDS 2021/2022



22

## Training a Classifier

- For **Multiclass classification**, some learning algorithms for binary classification are “SLMC-ready”; e.g.
  - Decision trees
  - Random forests
  - Naive Bayesian methods
  - Lazy learning methods (e.g., k-NN)
  - Neural networks
- For other learners (notably: SVMs) to be used for SLMC classification, combinations / cascades of the binary versions need to be used
  - e.g. multi-class classification SVM
  - Could be directly used for MLMC as well

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

23

## Holding out test data

- It's important to avoid overfitting
- Labelled data could be split into **two parts**
  - **Training:** used to train the classifier (e.g. **80%** of the data)
  - **Test:** used to test the performance of the trained classifier on unseen data (e.g. **20%** of the data)

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

24

## Hyperparameter optimisation

- Most classifiers have some hyperparameters to be optimized (we will usually refer to the ones we set manually as “hyperparameters” to distinguish from the “learned” parameters/weights of the model)
  - The C parameter in soft-margin SVMs
  - The  $r$ ,  $d$  parameters of non-linear kernels
  - Decision threshold for binary SVM
- Usually labelled data is split into **three parts**
  - **Training:** used to train the classifier (typically **80%** of the data)
  - **Development:** used to optimise hyperparameters. Apply the classifier on this data with different values of the hyperparameters and report the one that achieves the highest results (usually **10%** of the data)
  - **Test:** used to test the performance of the trained classifier with the optimal hyperparameters on these unseen data (usually **10%** of the data)
- Optimising the hyperparameters on test data is cheating!

Björn Ross, TTDS 2021/2022



## Evaluation

- Effectiveness (e.g. accuracy, precision, recall, F1):
  - Global effectiveness measures
  - Per class effectiveness measures
- Efficiency:
  - Speed in learning
    - SVM with linear kernel is known to be fast
    - DNNs are known to be much slower (specially with large # layers)
  - Speed in classification
    - K-NNs are known to be one of the slowest
  - Speed in feature extraction
    - BOW vs POS vs Link analysis features
- Importance of baselines

Björn Ross, TTDS 2021/2022

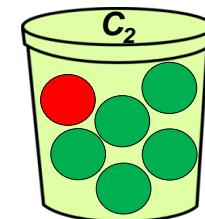
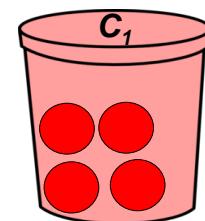
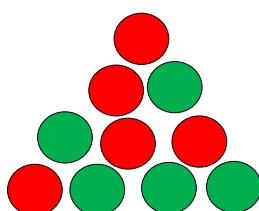


## Evaluation: Baselines

- There are standard methods for creating baselines in text classification to compare your classifier with
- Most popular/simplest baselines
  - Random classification
    - Classes are assigned randomly
    - How much better is the classifier doing than random?
  - Majority class baseline
    - Assign all elements to the class that appears the most
    - How much better you are doing than if you always picked the same thing output regardless of input?
  - Simple algorithm, e.g. BOW
    - Usually used when you introduce new interesting features

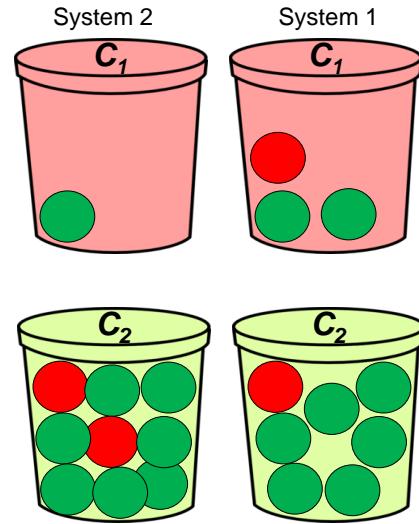
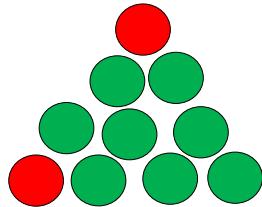
## Evaluation: Binary Classification

- Accuracy:
  - How many of the samples are classified correctly?
- $A = (4+5)/10 = 0.9$



## Evaluation: Binary Classification

- $A = (1+6)/10 = 0.7$  System 1
- $A = (0+7)/10 = 0.7$  System 2
- When classes are highly unbalanced
  - Precision/recall/F1 for the rare class
  - e.g. Spam classification (detection)



Björn Ross, TTDS 2021/2022

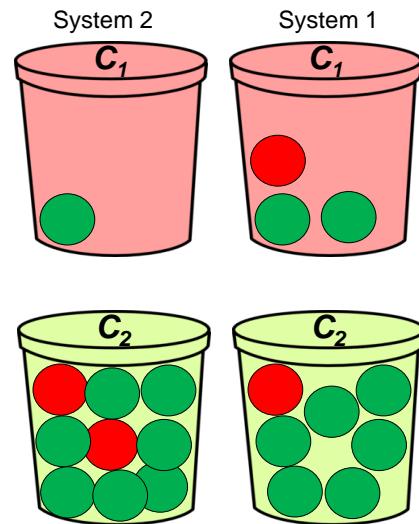
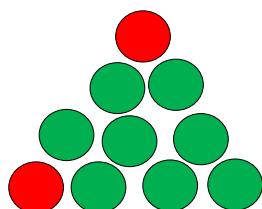


THE UNIVERSITY  
of EDINBURGH

30

## Evaluation: Binary Classification

|           | System 1     | System 2  |
|-----------|--------------|-----------|
| Precision | $1/3 = 0.33$ | $0/1 = 0$ |
| Recall    | $1/2 = 0.5$  | $0/2 = 0$ |
| F1        | 0.4          | 0         |



Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

31

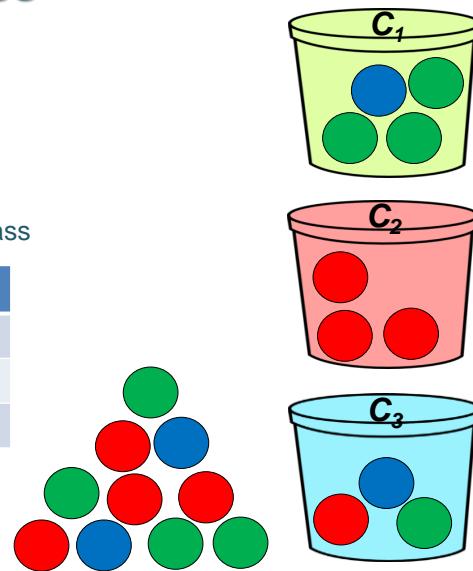
## Evaluation: Multi-class

- Accuracy =  $(3+3+1)/10 = 0.7$
- Good measure when
  - Classes are nearly balanced
- Preferred:
  - Precision/recall/F1 for each class

|    | Green | Red  | Blue  |
|----|-------|------|-------|
| P  | 0.75  | 1    | 0.333 |
| R  | 0.75  | 0.75 | 0.5   |
| F1 | 0.75  | 0.86 | 0.4   |

- Macro-F1
 
$$= (0.75+0.86+0.4)/3$$

$$= 0.67$$



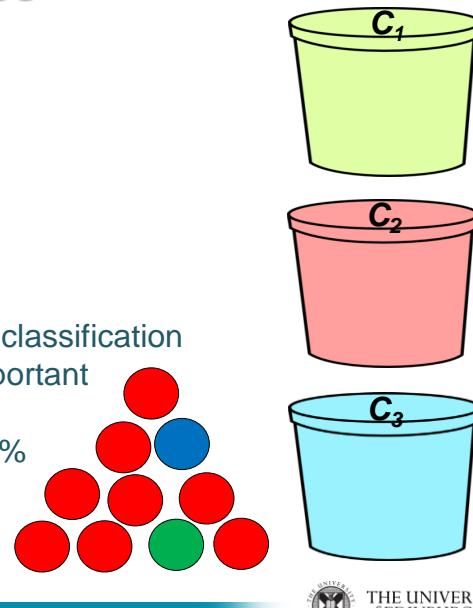
Björn Ross, TTDS 2021/2022



32

## Evaluation: Multi-class

- Majority class baseline
- Accuracy = 0.8
- Macro-F1 = 0.296
- Macro-F1:
  - Should be used in binary classification when two classes are important
  - e.g.: males/females while distribution is 80/20%



Björn Ross, TTDS 2021/2022

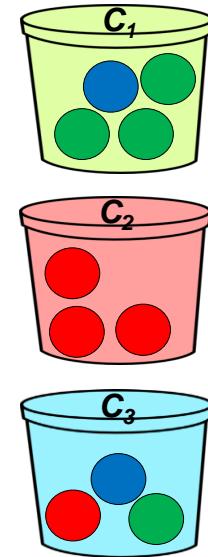


33

## Error Analysis

- Confusion Matrix

|              |       | Predicted class |     |      |
|--------------|-------|-----------------|-----|------|
|              |       | Green           | Red | Blue |
| Actual class | Green | 3               | 0   | 1    |
|              | Red   | 0               | 3   | 1    |
|              | Blue  | 1               | 0   | 1    |



- Useful:
  - Find classes that are confused with others
  - Develop better features to solve the problem

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

34

## Summary

- Text Classification tasks
- Feature extraction/selection/synthesis/weighting
- Learning algorithms
- Cross-validation
- Baselines
- Evaluation measures
  - Accuracy/precision/recall/Macro-F1

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

35

## Resources

- *Fabrizio Sebastiani*

**Machine Learning in Automated Text Categorization**

*ACM Computing Surveys, 2002*

*Link: <https://arxiv.org/pdf/cs/0110053>*

- *Yoav Goldberg*

**A Primer on Neural Network Models for Natural Language Processing**

*Link: <https://arxiv.org/abs/1510.00726>*





THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Text Classification (2)

Instructor:  
**Björn Ross**

17-Nov-2021

1

## Lecture Objectives

- Implement your first text classifier easy steps
- This is practical lecture  
No equations this time ☺



THE UNIVERSITY  
of EDINBURGH

2

## My first text classifier: Ingredients

- Text elements to be classified
  - Document, paragraph, sentence
- Set of predefined classes (classification task)
  - At least two (binary)
  - Topical, spam, relevance, sentiment, ...
- Training set
  - Enough samples of text elements for each class
- Test set (+ possible validation set)
  - Some samples of each class that not used in training
- Features set
  - A set of features extracted from the text to train the classifier
- Classifier
  - The ML module that learns a classification model



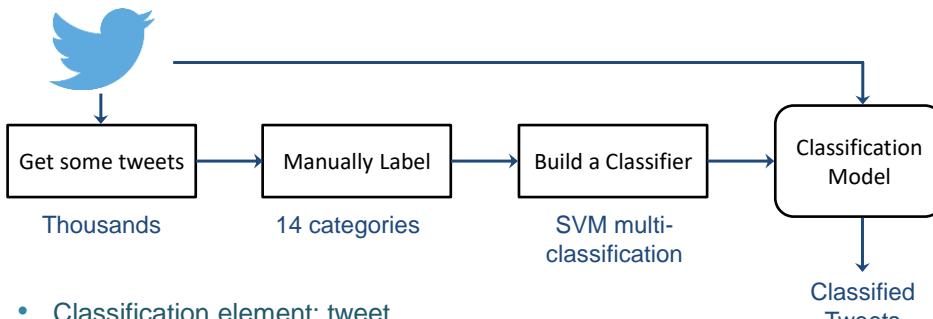
Björn Ross, TTDS 2021/2022



3

## My first text classifier: Application

- Classifying tweets into general-purpose categories



- Classification element: tweet
- Classes: 14 categories: sports, politics, comedy, ...
- Training/test set: 3129 tweets → 80/20% for train/test
- Features: BOW
- Classifier: SVM multiclass classifier

Björn Ross, TTDS 2021/2022



4

# My first text classifier: Steps

1. Prepare training data  
required: piece of text (tweet) + label to class
2. Extract features
  1. Pre-process text: lowercase, tokenise, remove useless strings
  2. Create a list of all unique terms in the training data. Give each term a unique ID.
  3. Convert the text into features, by replacing each term with its corresponding feature ID.  
Add value to the feature (simplest: value "1" if exists, or count of occurrences)
3. Prepare test file  
Convert test file text into features using the same mapping from the training data. For terms that are not in the features list, it could be neglected, or assigned to an ID representing OOV.
4. Run the learning process on the training data features to create a model
5. Run the classification on the features of the test data and get predictions
6. Evaluate performance

Björn Ross, TTDS 2021/2022



5

## Examples

- Tweet + Label
 

|                                                   |
|---------------------------------------------------|
| Kobe passes Wilt for 4th on all-time scoring list |
| Sports                                            |
  - Learned features (BOW) from training data
  - After converting text to feature vectors
- |     | 0 | 2943 | 2944 | 2945 | 2946 | ... | 8330 | 8331 | 10000 | ... |
|-----|---|------|------|------|------|-----|------|------|-------|-----|
| 0   | 0 | 0    | 1    | 0    | 0    | ... | 1    | 0    | 0     | ... |
| 1   | 0 | 1    | 0    | 0    | 1    | ... | 0    | 0    | 0     | ... |
| ... |   |      |      |      |      |     |      |      |       |     |
- Feature ID      Corresponding word

|      |         |
|------|---------|
| 2944 | kobe    |
| 2945 | rapping |
| ..   |         |
| 4525 | 4th     |
| 4526 | trevi   |
| ..   |         |
| 8330 | passes  |
| 8331 | ducks   |
| ..   |         |
| 9929 | 17      |
| 9930 | wilt    |
| ...  |         |
- SVM prediction output
 

|                    |
|--------------------|
| 7                  |
| Predicted Class ID |

Björn Ross, TTDS 2021/2022



6

## Practical

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

7

## Possible Improvements

- Feature Extraction
  - Apply stemming & stopping
  - Duplicate hashtags words (#car → #car car)
  - Expand tweet text that has link with the page title of that link
  - Add new set of features to the terms appearing in the profile description of the author of the tweet
    - E.g. tweets terms features: ID range: 1 → 12000
    - profile terms features: ID range: 12001 → 20000
    - If a term appeared in the tweet and in the profile description, these are two different features with two different IDs
  - Try non-textual features
    - Tweet length, presence of hashtags, links, emojis ...

Björn Ross, TTDS 2021/2022



THE UNIVERSITY  
of EDINBURGH

8

## Possible Improvements

- Feature weighting
  - Using tfidf, BM25 as the feature value instead of binary
- Learning method
  - Test other ML learning methods other than SVM
    - Random forest
    - Decision trees
    - Naive Bayesian
  - Test DNNs with word embeddings
    - Google Collab: you can experiment with using GPUs for free!  
<https://colab.research.google.com/>
- Add more training data
  - Think about a way to create more training data

Björn Ross, TTDS 2021/2022



9

## Resources

- Magdy W., H. Sajjad, T. El-Ganainy and F. Sebastiani. (2015) Bridging Social Media via Distant Supervision.  
*Springer SNAM 2015* [link](#), [arXiv](#)
- Additional reading:  
Nguyen, D. P., Gravel, R., Trieschnigg, R. B., & Meder, T. How old do you think I am? A study of language and age in Twitter.  
*ICWSM 2013*

Björn Ross, TTDS 2021/2022



10



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Learning to Rank

Instructor:  
**Walid Magdy**

01-Dec-2021

1

## Lecture Objectives

- Learn about:
  - IR as a classification task
  - Learning to Rank approaches



THE UNIVERSITY  
of EDINBURGH

2

## Classical Models vs. ML in IR

- Classical Models:
  - Features (factors): only a few, e.g., TF, IDF, |D|,  $P(t|\text{corpus})$  etc.
  - Structure: optimized for the a few particular features
  - Parameter & training
    - Often 1-2; not every factor has a parameter controlling its influence
    - Hand-tuning or data-based; can exhaustive since just 1-2 parameters
  - ***tfidf or BM25 or LMIR? PRF? What  $n_d$ ,  $n_t$ ?***
- ML in IR
  - Features: can include up to hundreds, thousands, or even more
  - Define the basic structure of a model
  - Quite generic: such as a weighted linear combination of all features
  - Parameters & training
    - Many; control the influence of each feature and their combinations
    - Impossible to tune by hand; Must be data-driven
  - **Let the ML decide what is better!**

Walid Magdy, TTDS 2021/2022



3

## Text Classification in IR

- Text Classification:
  - Classify a document into one of two or more classes
  - Different features could be used, e.g. BOW
- Can we model IR as classification?
  - Classify document to C1: R or C2: NR
  - Challenges?
    - Training data?
    - Features? BOW?
- BOW features cannot work
  - Spam? **Viagra, @ed.ac.uk**
  - Sentiment? **happy, sad**
  - Relevant? **Trump, hurricane**
  - Relevance is a query-dependent class

Walid Magdy, TTDS 2021/2022



4

2

## Getting Classification to IR

- Transforming features
  - Text classification: Input (D) → output (yes/no)
  - Information Filtering: Input (D|Q) → output (yes/no)
- Features set:
  - Independent of absolute words
  - More on relation between doc and query
  - Mostly are numbers (formulas, frequencies, ...)
  - Consistent as much as possible among different Q,D pairs
  - e.g.:
    - TFIDF, BM25
    - Query in page title? Heading?
    - Query in anchor text linking pages
    - PageRank of doc
    - Number of times page clicked for the same query

Walid Magdy, TTDS 2021/2022



5

## Popular Features

| Column in Output | Description                             | Column in Output | Description                |
|------------------|-----------------------------------------|------------------|----------------------------|
| 1                | TF(Term frequency) of body              | 24               | LMIR.JM of body            |
| 2                | TF of anchor                            | 25               | BM25 of anchor             |
| 3                | TF of title                             | 26               | LMIR.ABS of anchor         |
| 4                | TF of URL                               | 27               | LMIR.DIR of anchor         |
| 5                | TF of whole document                    | 28               | LMIR.JM of anchor          |
| 6                | IDF(Inverse document frequency) of body | 29               | BM25 of title              |
| 7                | IDF of anchor                           | 30               | LMIR.ABS of title          |
| 8                | IDF of title                            | 31               | LMIR.DIR of title          |
| 9                | IDF of URL                              | 32               | LMIR.JM of title           |
| 10               | IDF of whole document                   | 33               | BM25 of URL                |
| 11               | TF*IDF of body                          | 34               | LMIR.ABS of URL            |
| 12               | TF*IDF of anchor                        | 35               | LMIR.DIR of URL            |
| 13               | TF*IDF of title                         | 36               | LMIR.JM of URL             |
| 14               | TF*IDF of URL                           | 37               | BM25 of whole document     |
| 15               | TF*IDF of whole document                | 38               | LMIR.ABS of whole document |
| 16               | DL(Document length) of body             | 39               | LMIR.DIR of whole document |
| 17               | DL of anchor                            | 40               | LMIR.JM of whole document  |
| 18               | DL of title                             | 41               | PageRank                   |
| 19               | DL of URL                               | 42               | Inlink number              |
| 20               | DL of whole document                    | 43               | Outlink number             |
| 21               | BM25 of body                            | 44               | Number of slash in URL     |
| 22               | LMIR.ABS of body                        | 45               | Length of URL              |
| 23               | LMIR.DIR of body                        | 46               | Number of child page       |

Walid Magdy, TTDS 2021/2022



6

## Training Data

- Training data:  $\{R, X\}$ 
  - $X$ : feature representation of  $(D, Q)$  pairs
  - $R = \{-1, +1\}$  ... is  $D$  relevant to  $Q$  or no
- Samples:
  - Large set of  $(D, Q)$  pairs
  - Wide range of  $Q$ 's (long/short, frequent/rare, ...)
  - Wide range of  $D$ 's for each  $Q$  (top深深 ranked, recent/old pages, ...)
- Labels:
  - Manually labelled: assessors judge relevance of docs to queries (similar to standard IR)
  - Automatically labelled: click-through data

## Classification or Ranking?

- Click-through data
  - User clicks can give indication of relevance
  - What about non-relevance?
  - A list of ranked results:  $D_1 \rightarrow D_2 \rightarrow D_3$   
 user clicked on  $D_3$  and neglected  $D_1$  &  $D_2$   
 what does it mean?
    - $D_3$  is relevant and  $D_1$  &  $D_2$  are not relevant?
    - Relevance:  $D_3 > D_1 \& D_2$ ?
- It might be better to model the problem as ranking
  - Label → Ranking preference (e.g. gain= $\{4, 3, 2, 1, 0\}$ )
  - Learning → to optimize  $Doc_X > Doc_Y$   
 not to classify them to R/NR
  - Input: features for set of docs for a given query  
 Objective: rank them (sort by relevance)

## ML & IR: History

- Considerable interaction between these fields
  - Rocchio algorithm (60s) is a simple learning approach
  - 80s, 90s: learning ranking algorithms based on user feedback
  - 2000s: text categorization
- Limited by amount of training data
- Web query logs have generated new wave of research
  - L2R (LTR): “Learning to Rank”

Walid Magdy, TTDS 2021/2022



9

## What is Learning-to-Rank?

- Purpose
  - Learn a function automatically to rank results effectively
- Point-wise approach
  - Classify document to R / NR
- List-wise
  - The function is based on a ranked list of items
  - given two ranked list of the same items, which is better
- Pair-wise
  - The function is based on a pair of item
  - e.g., given two documents, predict partial ranking

Walid Magdy, TTDS 2021/2022



10

## Point-wise Approaches

- The function is based on features of a single object
  - e.g., regress the rel. score, classify docs into R and NR
- Very similar to classification
  - Examples of (D,Q) pairs with labels 1 or 0
- Classic retrieval models are also point-wise:
  - Calculate score(Q, D)
  - If  $\text{score}(Q,D) > \theta \rightarrow$  relevant  
else, irrelevant
- Referred to as *information filtering*
  - Standing query + new documents coming
  - Decide whether a new document is R or NR

Walid Magdy, TTDS 2021/2022



11

## List-based Approaches

- Given: ranked list A and ranked list B  
Task: decide which is better
- Need a loss function on a list of documents
- Challenge is scale
  - Huge number of potential lists
- Can develop tricks
  - Consider only possible re-rankings of top N retrieved by some fixed method
- Still expensive
  - No clear benefits over pairwise ones (so far)

Walid Magdy, TTDS 2021/2022



12

## Pair-wise Approaches

- Trying to classify
  - Which document of two should be ranked at a higher position?
- Optimize based on:
  - Margin between decision hyperplane and instances
  - Errors
  - Weighted based on some hyper-parameter C
  - Evaluation metric
- Example: SVM-rank
  - A generalization of SVM that supports ranking  
[Herbrich et al. 1999, 2000; Joachims et al. 2002]

Walid Magdy, TTDS 2021/2022



13

## Pair-wise Approaches

- The most popular approach
- Learning method: SVM-rank, RankBoost, GBRank, Ranknet, LambdaRank, LambdaMART
- Several issues of ranking SVM
  - Still, it does not directly optimize an evaluation metric
  - But pairwise ranking error often has better correlations with evaluation metrics than the loss/objective functions in point-wise approaches
    - Why: evaluation measures only care about rankings!
    - e.g., ground-truth:  $\text{rel}(D1) = 2$ ,  $\text{rel}(D2) = 1$ 
      - Regression model 1:  $\text{pred.rel}(D1) = 2$ ,  $\text{pred.rel}(D2) = 3$
      - Regression model 2:  $\text{pred.rel}(D1) = 1$ ,  $\text{pred.rel}(D2) = 0$
      - Model 1 is better than model 2 by criterion of evaluation regression (the prediction error), but model 2 yields a correct ranking of docs

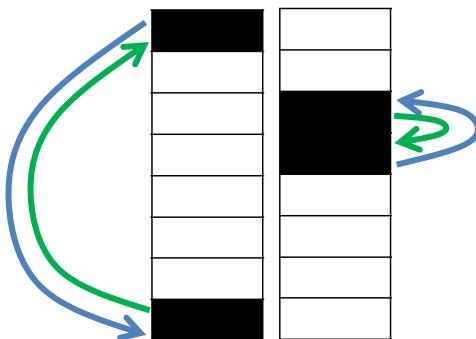
Walid Magdy, TTDS 2021/2022



14

## Pair-wise Approaches

- LambdaMART:
  - Misordered pairs are not equally important
  - Depends on how much they contribute to the changes in the target evaluation measure



Walid Magdy, TTDS 2021/2022



15

## Pair-wise Approaches

- Optimizing for an evaluation metric
  - The general idea is to weight loss/objective function or gradient with pairwise changes in evaluation measure.
  - e.g., in LambdaMART: lambda gradient
- Can we optimize all measures?
  - Not necessarily
  - For some measures, pairwise change do not only relate to the two documents themselves, but also others ...
    - Position-based measures do not have the issues (pairwise change only depends on the two documents)
    - Cascade measures may have issues

Walid Magdy, TTDS 2021/2022



16

## Pair-wise Approaches: Example

- Experiments
  - 1.2k queries, 45.5K documents with 1890 features
  - 800 queries for training, 400 queries for testing

|            | MAP           | P@1           | ERR           | MRR           | NDCG@5        |
|------------|---------------|---------------|---------------|---------------|---------------|
| ListNET    | 0.2863        | 0.2074        | 0.1661        | 0.3714        | 0.2949        |
| LambdaMART | <b>0.4644</b> | <b>0.4630</b> | <b>0.2654</b> | <b>0.6105</b> | <b>0.5236</b> |
| RankNET    | 0.3005        | 0.2222        | 0.1873        | 0.3816        | 0.3386        |
| RankBoost  | 0.4548        | 0.4370        | 0.2463        | 0.5829        | 0.4866        |
| RankingSVM | 0.3507        | 0.2370        | 0.1895        | 0.4154        | 0.3585        |
| AdaRank    | 0.4321        | 0.4111        | 0.2307        | 0.5482        | 0.4421        |
| pLogistic  | 0.4519        | 0.3926        | 0.2489        | 0.5535        | 0.4945        |
| Logistic   | 0.4348        | 0.3778        | 0.2410        | 0.5526        | 0.4762        |

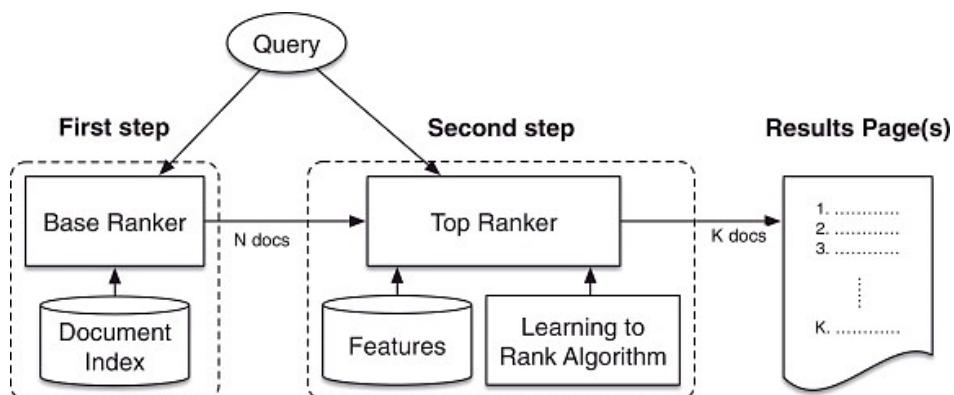
Honglin Wang Slides

Walid Magdy, TTDS 2021/2022



17

## L2R in Practice



Capannini, G., et al.  
 Quality versus efficiency in document scoring with learning-to-rank models.  
 IP&M 2016.

Walid Magdy, TTDS 2021/2022



18

## SVM-rank Example

```

Q1 { 3 qid:1 1:1 2:1 3:0 4:0.2 5:0 # 1A
     2 qid:1 1:0 2:0 3:1 4:0.1 5:1 # 1B
     1 qid:1 1:0 2:1 3:0 4:0.4 5:0 # 1C
     1 qid:1 1:0 2:0 3:1 4:0.3 5:0 # 1D
Q2 { 1 qid:2 1:0 2:0 3:1 4:0.2 5:0 # 2A
     2 qid:2 1:1 2:0 3:1 4:0.4 5:0 # 2B
     1 qid:2 1:0 2:0 3:1 4:0.1 5:0 # 2C
     1 qid:2 1:0 2:0 3:1 4:0.2 5:0 # 2D
Q3 { 2 qid:3 1:0 2:0 3:1 4:0.1 5:1 # 3A
     3 qid:3 1:1 2:1 3:0 4:0.3 5:0 # 3B
     4 qid:3 1:1 2:0 3:0 4:0.4 5:1 # 3C
     1 qid:3 1:0 2:1 3:1 4:0.5 5:0 # 3D
      }
      
```

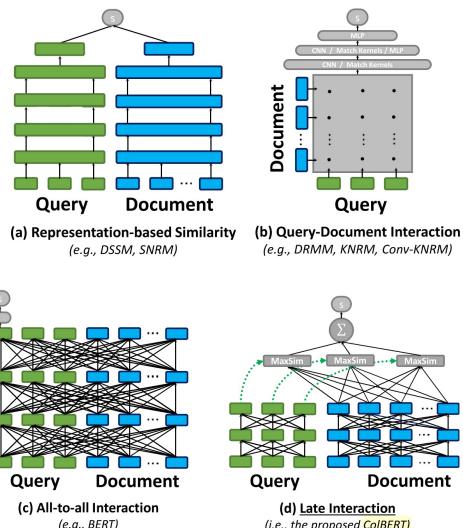
relevance  
(rank importance)      Set of features  
for (D,Q) pair

- Q3: 3C>3A, 3C>3B, 3C>3D, 3B>3A, 3B>3D, 3A>3D

## Current work in L2R

- Deep learning models are mainly used
- No manual feature extraction is applied
- Try to use word-embeddings to represent queries and docs, then learn the features automatically
- Content-independent models: try to learn the pattern of relations between terms in Q and D
- Content dependent: dependent on the terms

## Types of Deep LTR Models



Walid Magdy, TTDS 2021/2022

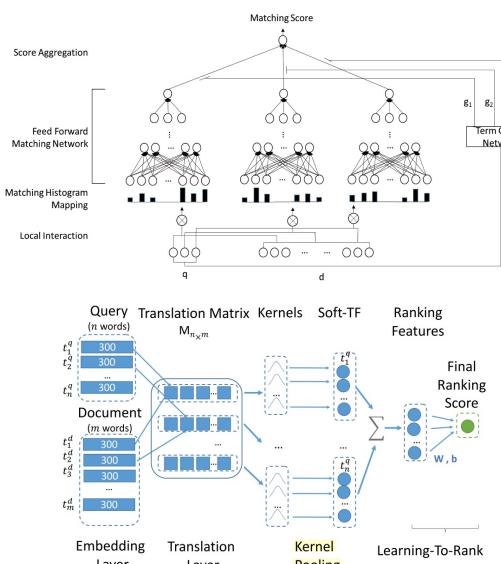
- Early Interaction-based: Learn on the signals from a query-document interaction.
- Late Interaction (Representation) based: Learn independent representations of queries and documents and then consider the interaction between them.
- Early interaction based approaches, e.g. DRMM, are relatively independent of the content (terms themselves) – tend to generalize well.
- Late interaction based approaches, e.g. ColBERT, are usually data hungry approaches – hence likely not to generalize well on standard ad-hoc IR collections.

Ganguly, guest lecture

THE UNIVERSITY  
of EDINBURGH

21

## DRMM & KNRM



- DRMM (up): *Deep Relevance Matching Model*. Uses histograms of word pair similarities (between doc and query) terms as inputs to a feed-forward network.
- The model seeks to utilize inherent patterns in these histograms to distinguish relevance from non-relevance.
- KNRM (down): *Kernel-based Neural Ranking Model*. It does not need to rely on histograms. Instead it applies 1D convolution.

Walid Magdy, TTDS 2021/2022

Ganguly, guest lecture

THE UNIVERSITY  
of EDINBURGH

22

## Summary

- IR as a classification task
- Learning to rank (L2R) approaches
  - Point-wise
    - Information Filtering
  - List-wise
  - Pair-wise
    - Ranking SVM
    - LambdaMART
- Current work in L2R depends on deep learning models and word-embedding representations

Walid Magdy, TTDS 2021/2022



23

## Resources

- Nallapati, Ramesh.  
Discriminative models for information retrieval.  
*SIGIR* 2004.
- Burges, C. J. (2010).  
From ranknet to lambdarank to lambdamart: An overview.  
*Learning*, 11(23-581), 81.
- Xiong, Chenyan, et al. End-to-end neural ad-hoc ranking with kernel pooling. *SIGIR* 2017
- Guest Lecture by Debasis Ganguly
- SVMRank: <http://svmlight.joachims.org/>
- L2R test sets:
  - Microsoft's LETOR project  
<http://research.microsoft.com/en-us/um/beijing/projects/letor//default.aspx>
  - Microsoft L2R datasets  
<http://research.microsoft.com/en-us/projects/mslr/default.aspx>

Walid Magdy, TTDS 2021/2022



24