



THE UNIVERSITY
of EDINBURGH

Text Technologies for Data Science

INFR11145

Introduction

Instructor:
Walid Magdy

18-Sep-2024

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!

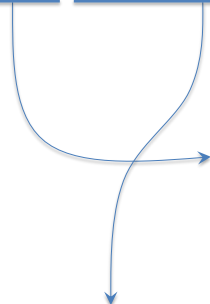
Walid Magdy, TTDS 2024/2025



THE UNIVERSITY
of EDINBURGH

2

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 2024/2025



3

What is Information Retrieval (IR)?

IR is **NOT** just

Google

Google Search

I'm Feeling Lucky

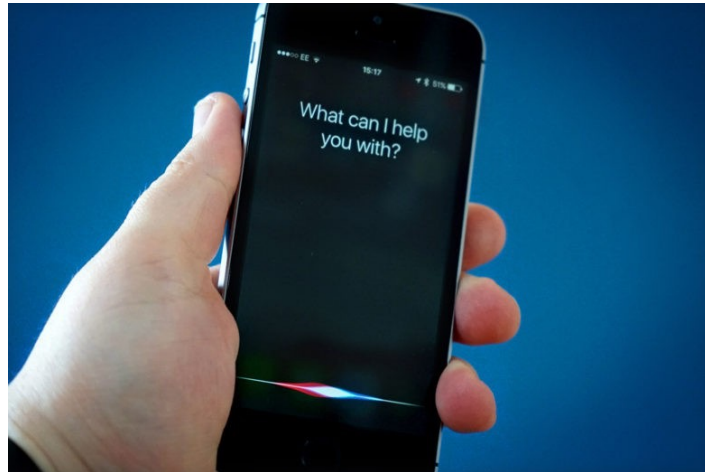
Web search

Walid Magdy, TTDS 2024/2025



4

What is IR?



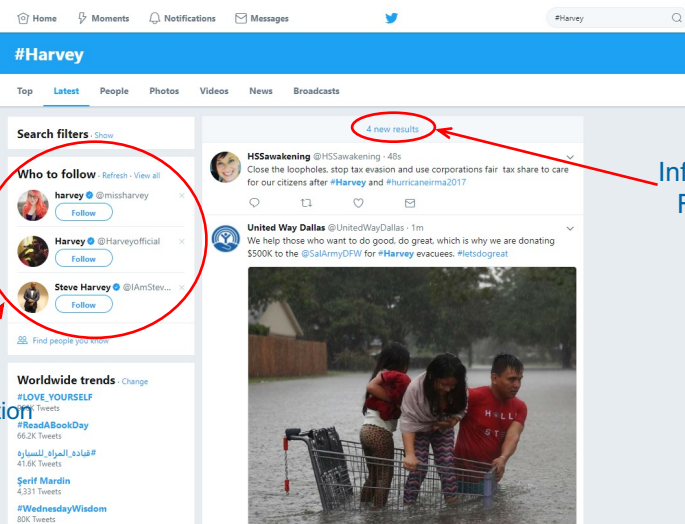
Speech - QA

Walid Magdy, TTDS 2024/2025



5

What is IR?



Recommendation

Information Filtering

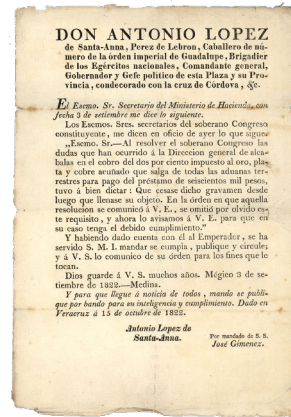
Social search

Walid Magdy, TTDS 2024/2025



6

What is IR?



Library (book) search
1950's

Walid Magdy, TTDS 2024/2025



7

What is IR?



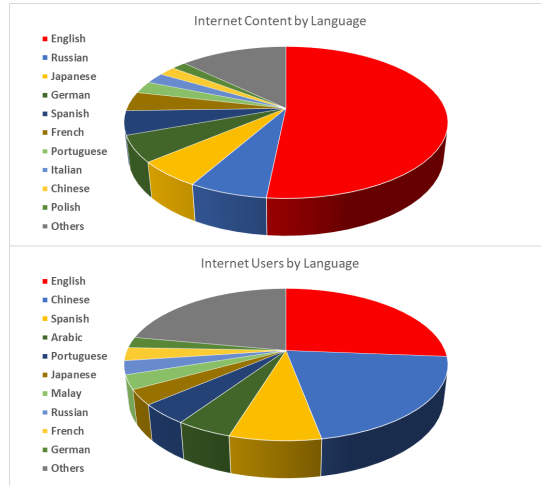
Legal search

Walid Magdy, TTDS 2024/2025



8

What is IR?



Cross-Language search

What is IR?



Content-based music search

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

What is IR?

Query suggestion / correction

Snippet selection / summarisation

Categorisation (search verticals)

Advertising

Walid Magdy, TTDS 2024/2025

THE UNIVERSITY of EDINBURGH

11

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

What is IR?

AI Generated Result

Retrieval Augmented Generation (RAG)

Walid Magdy, TTDS 2024/2025

THE UNIVERSITY of EDINBURGH

12

What is IR? Find?

ifferent 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

After scoring all terms, the top k expansion terms, are added to the topic vector. To avoid topic drift, the topic vector is reset to the original topic vector before each expansion, as shown in the following equation:

$$\vec{q}^i = \vec{q} + \beta \vec{e}$$

where \vec{e} is the normalized vector of the k expansion terms, and β 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](https://dev.twitter.com/rest/public/search)

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

IR ≠ Find

- Sequential
- Exact match

Walid Magdy, TTDS 2024/2025

THE UNIVERSITY of EDINBURGH

13

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 2024/2025

THE UNIVERSITY of EDINBURGH

14

Text classification

Walid Magdy, TTDS 2024/2025

THE UNIVERSITY of EDINBURGH

15

Text classification

Walid Magdy, TTDS 2024/2025

THE UNIVERSITY of EDINBURGH

16

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 2024/2025



THE UNIVERSITY
of EDINBURGH

17

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 2024/2025



THE UNIVERSITY
of EDINBURGH

18

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, ...)
 - evaluate classification quality
- Apply text analytics
 - Find what makes a set of document different from others
- RAG systems (LLMs with IR)

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

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
- RAG

Walid Magdy, TTDS 2024/2025



21

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 2024/2025



22

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.

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

Walid Magdy, TTDS 2024/2025

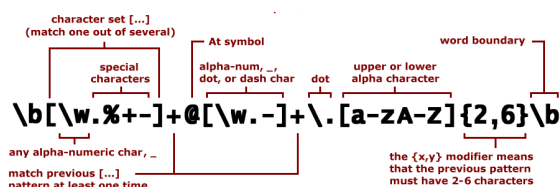


23

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.
- We **DO NOT** teach coding skills in this course!
We assume you can code!



Walid Magdy, TTDS 2024/2025



24

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 2024/2025



25

Skills to be gained !!!

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

Walid Magdy, TTDS 2024/2025



26

Course Structure

- 19 Lectures:
 - 2 lectures → Introduction (today)
 - 13 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 2024/2025



27

Course Instructors



Walid Magdy
Reader
(11 lectures)



Bjorn Ross
Lecturer
(5 lectures)



Tj Elmas
Lecturer
(3 lectures)

+ 1 guest lecture

Walid Magdy, TTDS 2024/2025



28

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 1st 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 2024/2025



29

Labs

- How it works:
 - 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 in-person labs for those still require support
- Optional in-person labs:
 - Location: AT 6.06
 - Times: Tuesday, 11am, 12pm
- Demonstrators:
Wendy Zheng and Zahra Bokae

Walid Magdy, TTDS 2024/2025



30

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
 - Count the number of occurrence of few words
- If Lab 0 challenging →
 - 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

Assessments

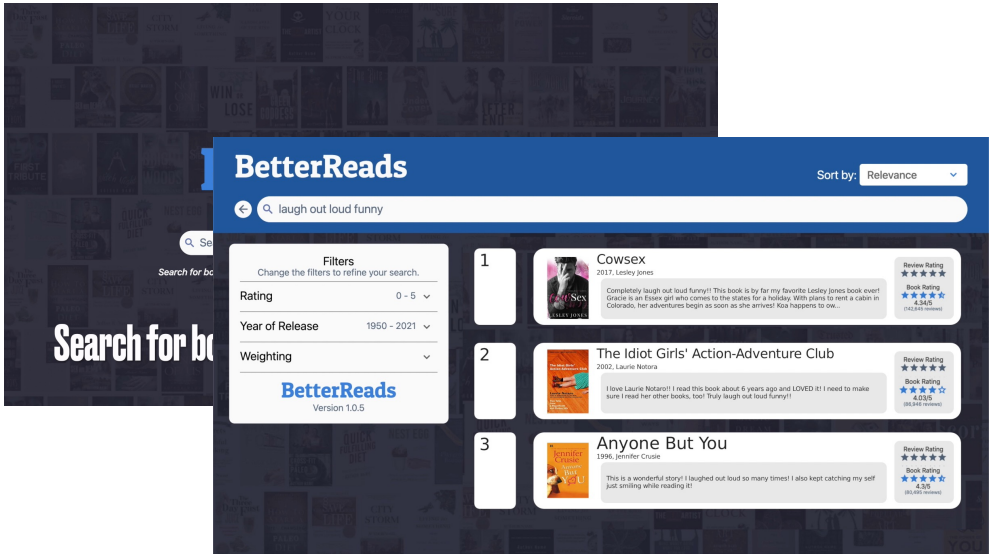
- Coursework 1: **10%**
The same as labs 1-3 → Build your first search engine
- Coursework 2: **20%**
IR Evaluation, Text classification/analytics
- Group project: **40%**
A full running search engine supported by text technologies
- Final Exam: **30%**

Group Project

- The largest weight: 40% of the total mark
- Teamwork → Group 5-6 (you select your own 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

33

Example: BetterReads



BetterReads Version 1.0.5

Sort by: Relevance

Search for: laugh out loud funny

Filters
Change the filters to refine your search.

Rating: 0 - 5

Year of Release: 1950 - 2021

Weighting: [dropdown]

- Cowsex**
2013, Lucy Jones

Completely laugh out loud funny!! This book is by far my favorite Lucy Jones book ever!! Grace is an Essex girl who comes to the states for a holiday. With plans to rent a cabin in Colorado, her adventures begin as soon as she arrives! Kee happens to ow...

Review Rating: ★★★★★
Book Rating: ★★★★★ (12,845 reviews)
- The Idiot Girls' Action-Adventure Club**
2002, Laurie R. King

I love Laurie R. King!! I read this book about 6 years ago and LOVED it! I need to make sure I read her other books, too! Truly laugh out loud funny!!

Review Rating: ★★★★★
Book Rating: ★★★★★ (8,428 reviews)
- Anyone But You**
1996, Jennifer Cruise

This is a wonderful story! I laughed out loud so many times! I also kept catching my self just smiling while reading it!

Review Rating: ★★★★★
Book Rating: ★★★★★ (4,375 reviews)

34

Example: BetterReads

- 11.5M Book reviews from Good reads
- Average query time: 1 secs
- New reviews are crawled and indexed automatically every day
- Ranking: Relevance + Sentiment
- Engine hosted on Google cloud compute
- *Note: we will provide credit to Google cloud to host your engine*

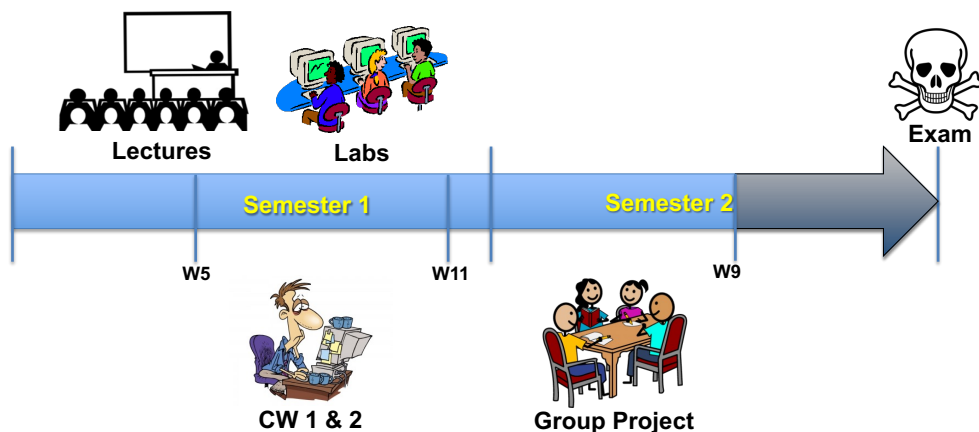
Walid Magdy, TTDS 2024/2025



35

Timeline

- 2 Semesters (or one?)



Walid Magdy, TTDS 2024/2025



36

Logistics

- Lectures:
 - Two lectures on Wednesdays, 15.00-17.00
 - Recording will be available on learn
 - Handouts to be posted on the day of the lecture
- Course webpage:
 - Link: <https://opencourse.inf.ed.ac.uk/ttds/>
 - Handouts, Labs, CW details
- Learn:
 - Lecture recordings
 - Deadlines
- Note: all course materials are made public including recordings. Feel free to share with anyone interested

Walid Magdy, TTDS 2024/2025



37

Piazza

- 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
- Tag each question/post by its relevant topic (lab, CW ... etc)
- Join NOW: [link](#)

Walid Magdy, TTDS 2024/2025



38

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?

Next Lecture

- Definitions of IR main concepts
(more introduction)