# Web Information Retrieval

## Academic year 2020/2021

## Instructors

- **Luca Becchetti**
- **Fabrizio Silvestri**

# A quick tour of (tentative) topics

# What is Web IR

- **Information retrieval when the corpus is the Web**

- **Information Retrieval (IR)**

  - Retrieving unstructured material (usually textual documents) meeting an information need from large collections (usually stored on computers)

- **Live example**

# Why is the Web different?

1) **Distributed and larger than traditional information resources**

2) **Linked**

3) **Evolving**

4) **Information is semi-structured → view source of HTML pages**

5) **Multiple-content types (i.e. images, scripts, text etc.) coming in different formats**
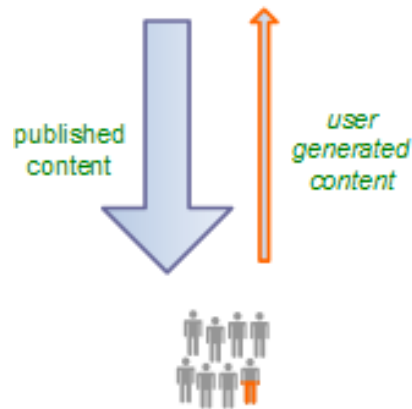
6) **Quality of documents is not homogeneous**

# The Web

- **As it was**
- **10 years later**

## Web 1.0
"the mostly read-only Web"

250,000 sites

published content

user generated content

45 million global users

1996

## Web 2.0
"the wildly read-write Web"

80,000,000 sites

collective intelligence

published content

user generated content

1 billion+ global users

2006

# The Web

- ## As it was
- ## As it is

**Web 1.0** was about…         **Web 2.0** is about…

| reading | → | writing |
| individual users | → | communities |
| taxonomy | → | folksonomy |
| owning | → | sharing |
| consuming | → | producing |
| home pages | → | blogs |
| you find news | → | news finds you |

According to Wikipedia

A Web 2.0 website may allow users to interact and collaborate with each other in a social media dialogue as creators of user-generated content in a virtual community, in contrast to the first generation of Web 1.0-era websites where people were limited to the passive viewing of content. Examples of Web 2.0 features include social networking sites and social media sites (e.g., Facebook), blogs, wikis, folksonomies ("tagging" keywords on websites and links), video sharing sites (e.g., YouTube), hosted services, Web applications ("apps"), collaborative consumption platforms, and mashup applications.

# Course outline

- **Collecting a Web corpus**
- **Pre-processing and organizing a Web corpus**
- **(Web) document retrieval (querying and searching the corpus)**
- **Analyzing documents using NLP**
- **Web page classification and clustering**

# Course outline

- **Collecting a Web corpus**

- **Pre-processing and organizing a Web corpus**

- **(Web) document retrieval (querying and searching the corpus)**

- **Using the Web as a platform to provide services**

# Collecting a Web corpus

# Crawling the Web

- **Exploit link structure**
- **Simplified scheme**
  - Start from an initial page
  - Retrieve all linked pages
  - Iterate on new pages
- **Example**
- **At this point you should have**



Extract

Web Crawling    Database

# Crawling/Caveats and traps

- **Design/algorithmic challenges**
  - E.g.: Multiple Web crawlers
    - How to ensure we are not crawling the same pages?
  - We are not visiting all pages
    - Bias in data
- **Web applications**
  - e.g., social networking platforms
    - Not all pages accessible
    - Specific APIs/restrictions

# Organizing a Web corpus

# We have goals in mind – e.g.

# How Google puts it …

- **https://www.google.com/search/howsearchworks/**

# How Google puts it ...

- **https://www.google.com/search/howsearchworks/**
- **In a nutshell**
    - Crawling
    - Indexing
    - Search algorithms

# Indexing

- **Organize Web corpus so as to efficiently answer (implicit or explicit) queries**

- **Challenging task**
  - Multiple objectives
  - Multiple trade-offs

# Efficient data structures



- **Typically an inverted index**
  - Index construction
  - Search using an inverted index
  - Compression, metadata enrichment …

# Querying the corpus (search)

# Goals and document scoring

- **Return documents that are relevant to the query "web information retrieval"**

- **How to define and measure relevance**

  – Textual analysis

    - Use meta-data when available

  – Link analysis (Web structure)

  – Pages can be "more" or "less" relevant → ranking

- **Relevance vs authority**

# The final picture

# The recent past …

# Providing services over the Web

- **Search engines**

  – Web applications providing search

- **More have emerged over the recent past …**

# Now (social networking only)

# New approaches/challenges

# Personalization



Your profile impacts the outcome

# Recommendations

# ...and much more

# Practical info

# General info

- **Where**
  - Room A3, via Ariosto 25
- **When**
  - Mondays, 2pm – 5pm
  - Thursdays, 2pm – 4pm

- **General info, announcements etc.**
  - https://classroom.google.com/u/0/c/MjcxOTA0NTgyNjEy
    - Please enroll!!

- **Luca Becchetti**
  - becchetti@diag.uniroma1.it
- **Fabrizio Silvestri**
  - fsilvestri@diag.uniroma1.it

# Organization

- **Lectures**
  - New topics
  - Discussions
  - Homeworks
- **Hands on (hopefully)**
  - We try to solve problems together
    - Emphasis on together
    - Bring your laptop if you have one
  - First year – we'll do our best

- **Exam**
  - Possibly: written exam + assignments
  - Details to be decided

# More info

- **Prerequisites**
  - Undergraduate in CS or equivalent

- **Useful things**
  - A laptop
  - Curiosity and independence
  - Presence and participation

- **References**
  - Manning, Christopher D., Prabhakar Raghavan, and Hinrich Schütze. Introduction to information retrieval. Vol. 1. No. 1. Cambridge: Cambridge university press, 2008
    - Thanks to the authors, pdf of chapters is available for free at the book's Web site
  - Scientific papers
  - On-line material, tutorials etc.