# WEB INTELLIGENCE

## Lecture 1: Introduction

### Russa Biswas

September 02, 2025

# Organisational Details

**Lecturer**: Russa Biswas

**Schedule**:

- Every Tuesday at 12:30 – 14:00 in Room 2.2.120

- Exercises: Every Monday at 14:15 – 16:15 in Room 2.2.120/ Group Room

**Exercises**:

Classroom Exercises: smaller problems solved on pen-paper (or white board)

Extended Exercises: Exercises similar to classroom exercises + Practical programming tasks with Python and Jupyter notebooks

# Introduction

# WHAT IS THE WEB?

# What is the Web?

"…The dream behind the Web is of a common information space in which we communicate by sharing information… The web of human-readable document is being merged with a web of machine-understandable data…" – Tim Berners-Lee, https://w3.org/People/Berners-Lee/ShortHistory.html

- The Web is a
  - global **decentralized hypertext-based** information system.
  - **Hyperlinks** are used to navigate from one document to another.
  - This information space build on a set of technical standards for the identification, retrieval and representation of content.

Document 1
Document 2
Document 3
Document 4
Document 5

# World Wide Web (WWW)

- The **World Wide Web** (referred to as WWW or W3 or simply Web), developed by **Tim Berners-Lee** in 1989 at CERN, Switzerland
  - The project document described a "**hypertext project**" called "**WorldWideWeb**" in which a "**web**" of "**hypertext documents**" could be viewed by "**browsers**".[1]

- First Website in 1991: **info.cern.ch**

- First Webpage address:
  http://info.cern.ch/hypertext/WWW/TheProject.html



Tim Berners-Lee, pictured at CERN (Image: CERN)

**Web Intelligence, Lecture 1: Introduction, Russa Biswas**

# World Wide Web (WWW)

- First Browser

- In 1993: Graphical Web Browsers such as Mosiac and Netscape Navigator were made accessible outside of academia.

- In 1994: W3C (World Wide Web Consortium) was founded by Tim Berners Lee. W3C publishes recommendations, that are considered web standards.
  - Web standards are blueprints –or building blocks– of a consistent and harmonious digitally connected world. They are implemented in browsers, blogs, search engines, and other software that power our experience on the web.



A screenshot showing the NeXT world wide web browser created by Tim Berners-Lee (Image: CERN)

# World Wide Web – Data Model

W3 data model enables:

- Information need only be represented once, as a reference may be made instead of making a copy.

- Links allow the topology of the information to evolve, so modelling the state of human knowledge at any time is without constraint.

- The web stretches seamlessly from small personal notes on the local workstation to large databases on other continents.

- Indexes such as phone books are presented as documents, and so may themselves be found by searches and/or following links.

- The documents in the web do not have to exist as files; they can be "virtual" documents generated by a server in response to a query or document name. They can therefore represent views of databases, or snapshots of changing data (such as the weather forecasts, financial information, etc.).

**Advantages:**

- Information access doesn't require expert knowledge

- Information Retrieval via search engines

# World Wide Web - Components

The World Wide Web (WWW) comprise of different components:

- **Identification:** Universal Resource Identifiers (URIs)- Address system; globally unique identification of the web resources.

For e.g., the URI of the main page for the first WWW project is http://info.cern.ch/hypertext/WWW/TheProject.html

- **Interaction:** Hypertext Transfer Protocol (HTTP) - network protocol used for transferring information/interacting between the web resources. The data transferred can be plain text, hypertext, images, etc.

- **Content Format:** Hypertext Markup Language (HTML) - a markup language, used to define the structure and the content of the webpage. HTML supports various content types, including text, images, video, audio, scripts, and hyperlinks for easy web resource access.



URI
http://weather.example.com/oaxaca

Identifies

Resource
Oaxaca Weather Report

Represents

Representation

Metadata:
Content-type:
application/xhtml+xml

Data:
<!DOCTYPE html PUBLIC "...
    "http://www.w3.org/...
<html xmlns="http://www...
<head>
<title>5 Day Forecaste for
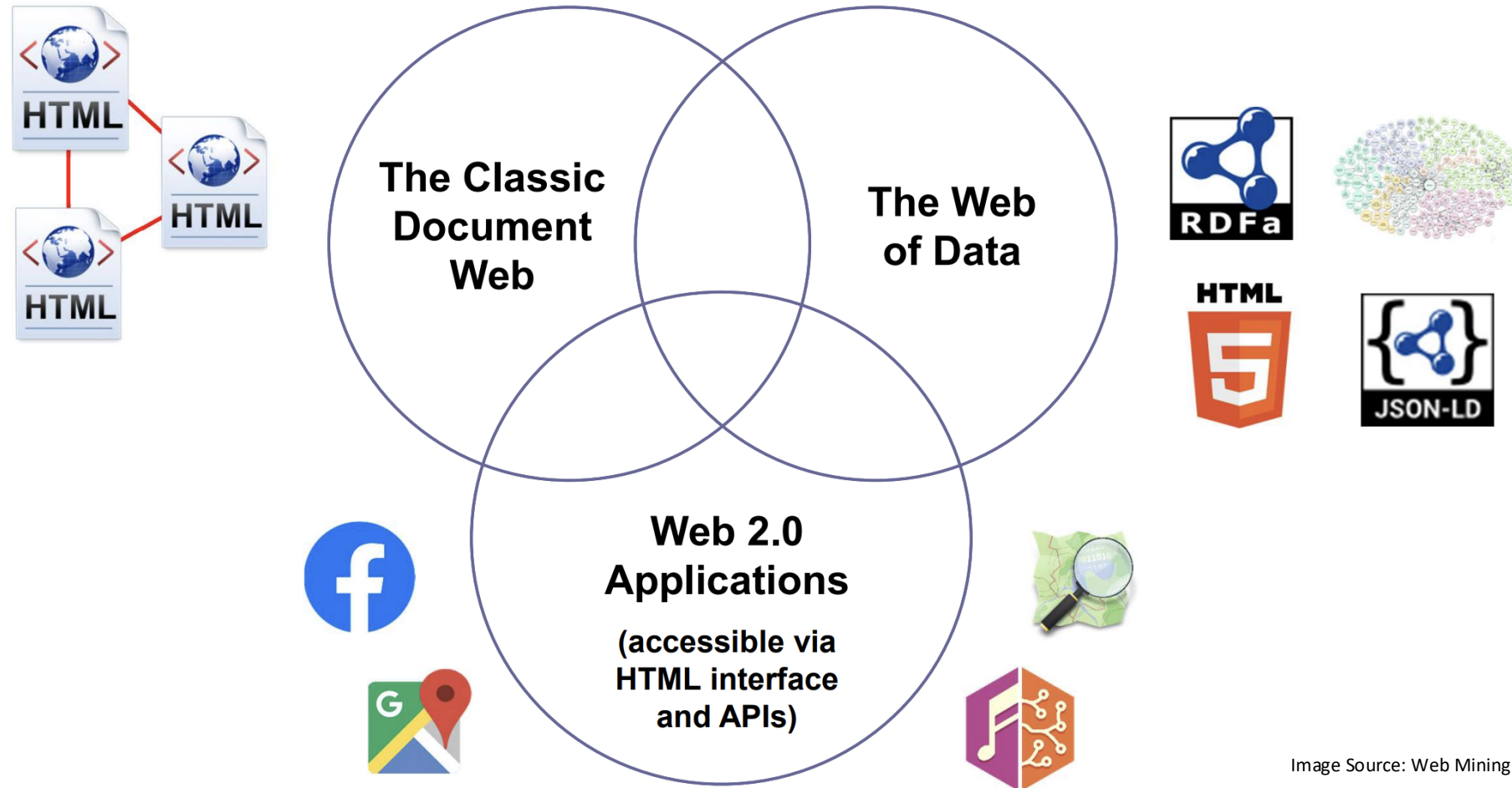Oaxaca</title>
...
</html>

Image source: https://www.w3.org/TR/webarch/

# Topology of the Web



The Classic Document Web

The Web of Data

Web 2.0 Applications

(accessible via HTML interface and APIs)

Image Source: Web Mining FSS2024, University of Mannheim

Web Intelligence, Lecture 1: Introduction, Russa Biswas

10

# Size of the web



The size of the World Wide Web (The Internet)

The Indexed Web contains **at least 4 billion pages** (Sunday, 08 September, 2024).
The Dutch Indexed Web contains **at least 2270.89 million pages** (Sunday, 08 September, 2024).

**The Indexed Web** | The Dutch Indexed Web

| Last Month | Last Three Months | Last Year | Last Two Years | Last Five Years | Last Ten Years |

The size of the indexed World Wide Web
(Number of webpages)

GB = Sorted on Google and Bing
BG = Sorted on Bing and Google

The size of the World Wide Web:
Estimated size of Google's index

| Last Month | Last Three Months | Last Year | Last Two Years | Last Five Years | Last Ten Years |

Size Google
(Number of webpages)

# What is Web Intelligence?

- Intelligent ways to extract information and knowledge from the web:

  - finding relevant information available on the web

  - obtaining new knowledge by analyzing web data: the web itself, but also how it evolves, and how users interact on and with the web


- Some applications:

  - Intelligent Search

  - Recommender Systems

  - Business Analytics

  - Crowd Sourcing

  - Not so nice ones: advertising, manipulation, surveillance

# Intelligent Search



Query: Contemporary reaction Puccini Turandot

www.operaphila.org › backstage › opera-blog › watchi... ▾
**Watching Turandot with Modern Eyes - Opera Philadelphia**
Puccini composed Turandot in 1924, at a time when no one had internet and few ... a diff response from that of an audience member from ninety years ago.

People also ask

What is Turandot synopsis?

What are the 3 questions in Turandot?

How much of Turandot did Puccini write?

How do you say Turandot?

www.jstor.org › stable
**Modernism and the Machine Woman in Puccini's 'Turandot'**
Puccini's final opera within the context of contemporary developments in the ... Turando the response to its two heroines in particular, provides a vital key to.
by A Wilson - 2005 - Cited by 16 - Related articles

www.theoperablog.com › what-you-said-turandot-react... ▾
**What you said: Turandot reactions | The Opera Blog**
Jul 31, 2015 - Don't take our word for it – here's what the audience and critics are ... costu fit for a Tarsem Singh film, Puccini earworms & killer vocals.

operawire.com › no-hero-why-calaf-is-the-dramatic-pr... ▾
**No Hero: Why Calaf Is The Dramatic Problem In Puccini's ...**

➡ Close, but not quite the information I was looking for

Query: shortest path algorithms

en.wikipedia.org › wiki › Shortest_path_problem ▾
**Shortest path problem - Wikipedia**
Jump to **Algorithms** - In graph theory, the **shortest path** problem is the problem of finding a path between two vertices (or nodes) in a graph such that the sum ...
Definition · Single-source shortest paths · All-pairs shortest paths · Applications

en.wikipedia.org › wiki › Dijkstra's_algorithm ▾
**Dijkstra's algorithm - Wikipedia**
Dijkstra's **algorithm** (or Dijkstra's **Shortest Path** First **algorithm**, SPF **algorithm**) is an **algorithm** for finding the **shortest paths** between nodes in a graph, which may represent, for example, road networks. It was conceived by computer scientist Edsger W. Dijkstra in 1956 and published three years later.
Description · Pseudocode · Proof of correctness · Running time

www.hackerearth.com › ... › Shortest Path Algorithms ▾
**Shortest Path Algorithms Tutorials & Notes | Algorithms ...**
The **shortest path** problem is about finding a **path** between vertices in a graph such that the total sum of the edges weights is minimum. This problem could be solved easily using (BFS) if all edge weights were ( ), but here weights can take any value.

People also ask

What are the shortest path algorithms?          ⌄

Which is the best shortest path algorithm?          ⌄

How do you solve Dijkstra's shortest path algorithm?          ⌄

Does A * find the shortest path?          ⌄

Feedbac

www.geeksforgeeks.org › dijkstras-shortest-path-algorit... ▾
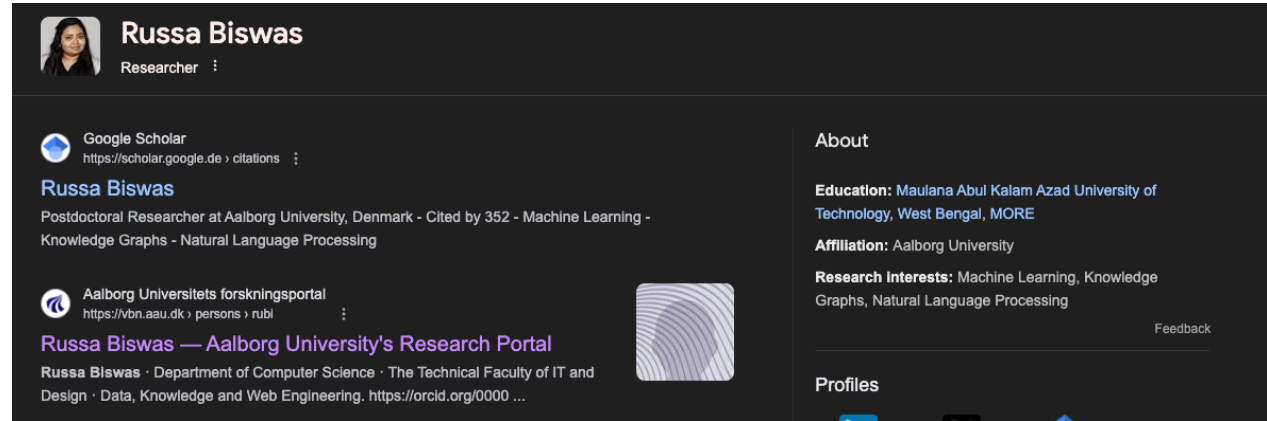**Dijkstra's shortest path algorithm | Greedy Algo-7**
Dijkstra's shortest path algorithm | Greedy Algo-7. Last Updated: 22-04-2020. Given a graph

➡ Relevant and competent resources

# Intelligent Search

- Keyword Search: the words in the query appear frequently in the document, in any order (*bag of words*).

- Disadvantages:
  - May not retrieve relevant documents that include synonymous terms (e.g., cannot distinguish between "restaurant" and "café")
  - May retrieve irrelevant documents that include ambiguous terms (e.g., cannot distinguish between "bat" mammal and "bat" baseball)

- Beyond Keywords:
  - Considering the meaning of the words used
  - Adapting to user feedback (direct or indirect)
  - Considering the authority of the source

# Recommender Systems

## Looking on Amazon for the Manning IR book:



**Frequently bought together**

Total price: £166.24

**Add all three to Basket**

⏐ Some of these items are dispatched sooner than the others. Show details

☑ **This item:** Introduction to Information Retrieval by Christopher D. Manning  Hardcover  £43.99
☑ Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit by Steven Bird  Paperback  £27.25
☑ Foundations of Statistical Natural Language Processing (The MIT Press) by Christopher Manning  Hardcover  £95.00

**Books you may like**

| Where the Crawdads Sing | Midnight Sun | The Fast 800 Recipe Book: Low-carb,…. | Girl, Woman, Other: WINNER OF THE BOOKER PRIZE 2019 | The Boy, The Mole, The Fox and The Horse | Too Much and Ne Enough: How My Created the Worl |
|---|---|---|---|---|---|
| › Delia Owens | › Stephenie Meyer | Dr Clare Bailey | › Bernardine Evaristo | › Charlie Mackesy | Mary L. Trump Ph. |
| ⭐ 50,230 | ⭐ 3,231 | ⭐ 2,214 | ⭐ 2,122 | ⭐ 14,361 | ⭐ 3,3 |
| Paperback | Hardcover | Paperback | | Hardcover | |
| £5.99 | £11.00 | | | | |

Kinds of recommendations:

- Product Based (collaborative filtering): similar books
- User- Based (content based filtering): based on search history
- Hybrid

# Knowledge Graphs



Search Engines store information about entities in Knowledge Graphs and use them to summarise search results
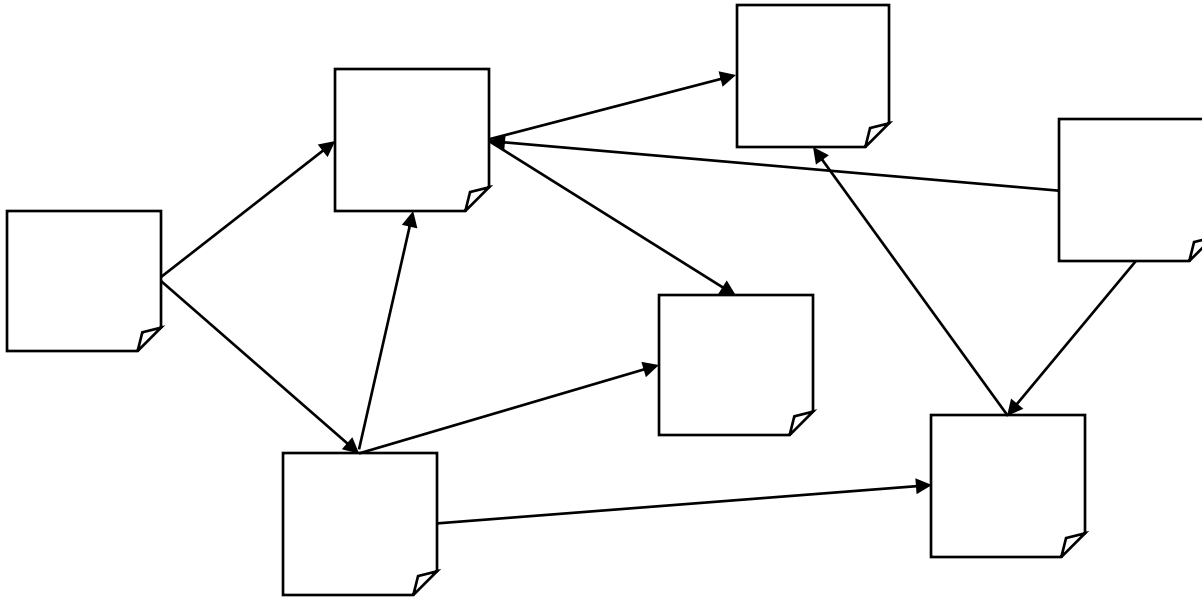
# Question Answering



The same knowledge graphs can be used to find answers to user's questions (instead of finding the documents related to the search keywords)

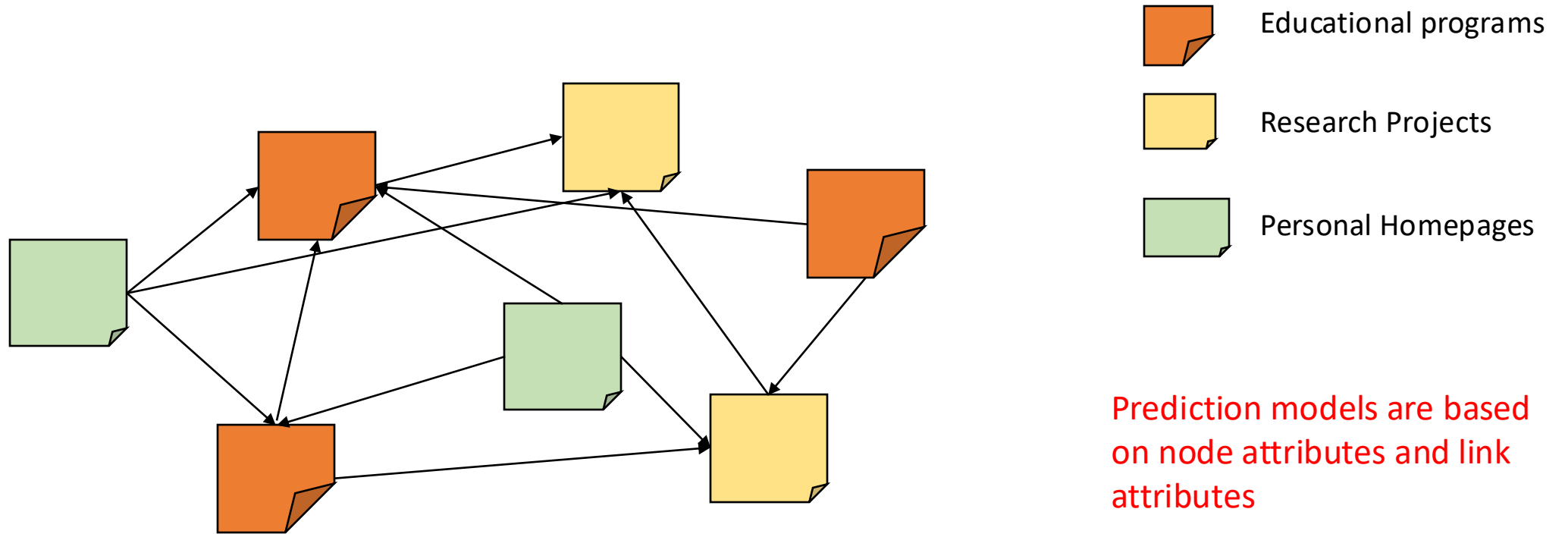[Google - https://www.google.com/search?q=who+is+the+wife+of+the+us+president]

# Node Classification

- Suppose this is aau.dk domain

-  Identify the personal homepages, educational programs, research projects etc.

# Node Classification

- Suppose this is aau.dk domain

- Identify the personal homepages, educational programs, research projects etc.



Educational programs

Research Projects

Personal Homepages

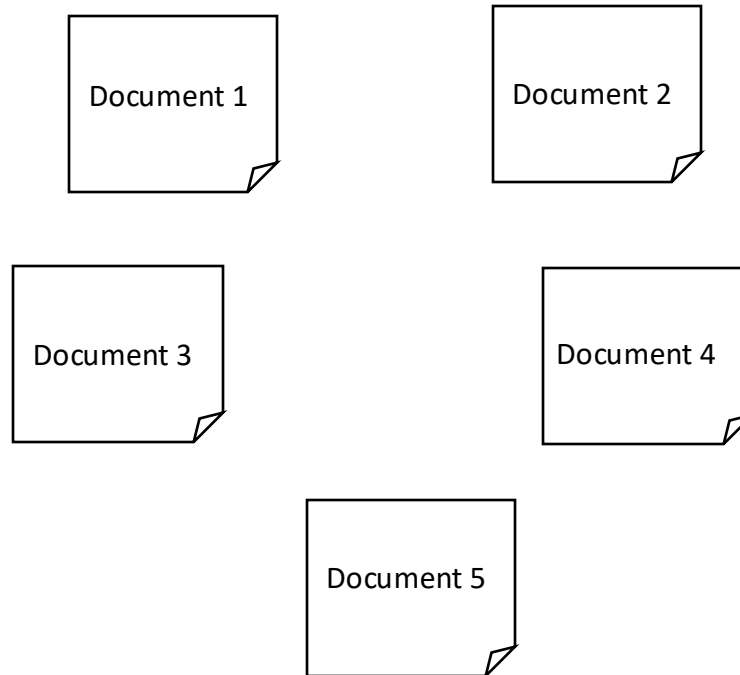Prediction models are based on node attributes and link attributes

# Information diffusion

- Information diffusion: process by which a piece of information or knowledge is spread and reaches individuals through interactions
  - Studied in Sociology, epidemiology (the study of the determinants, occurrence, and distribution of health and disease in a defined population)
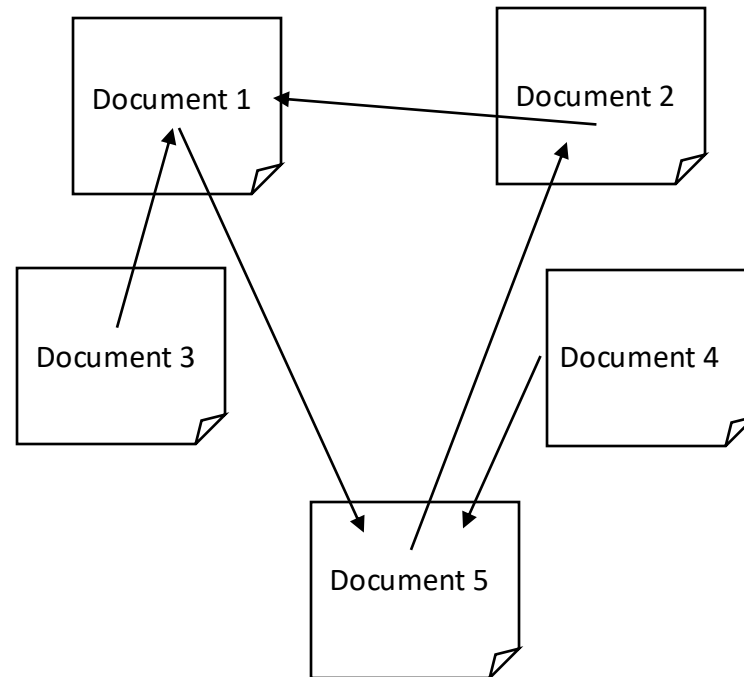
# Different levels of Web

- We can distinguish 3 levels (perspectives) of modeling and analytics

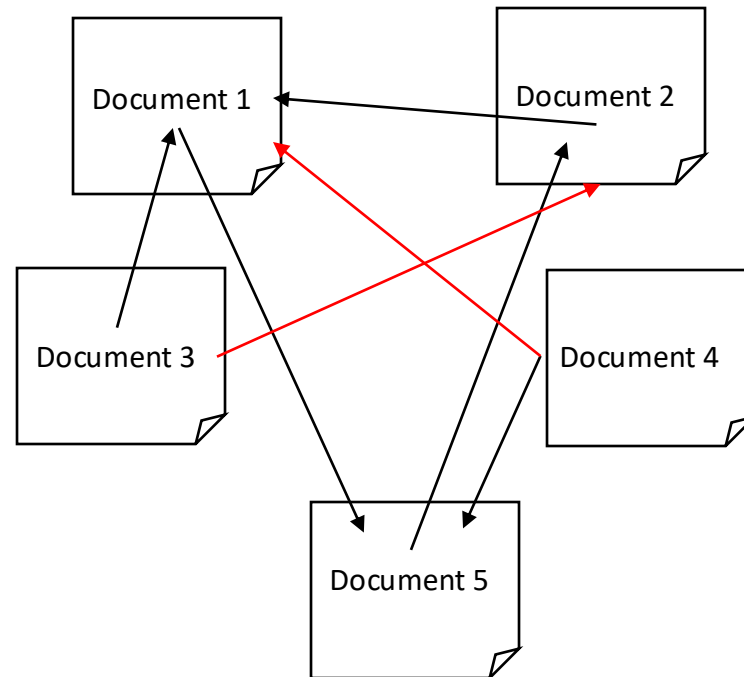- Level 1: Web Content – The web is a collection of documents

Document 1

Document 2

Document 3

Document 4

Document 5

# Different levels of Web

Level 2: Web Structure – The web as a network of documents:
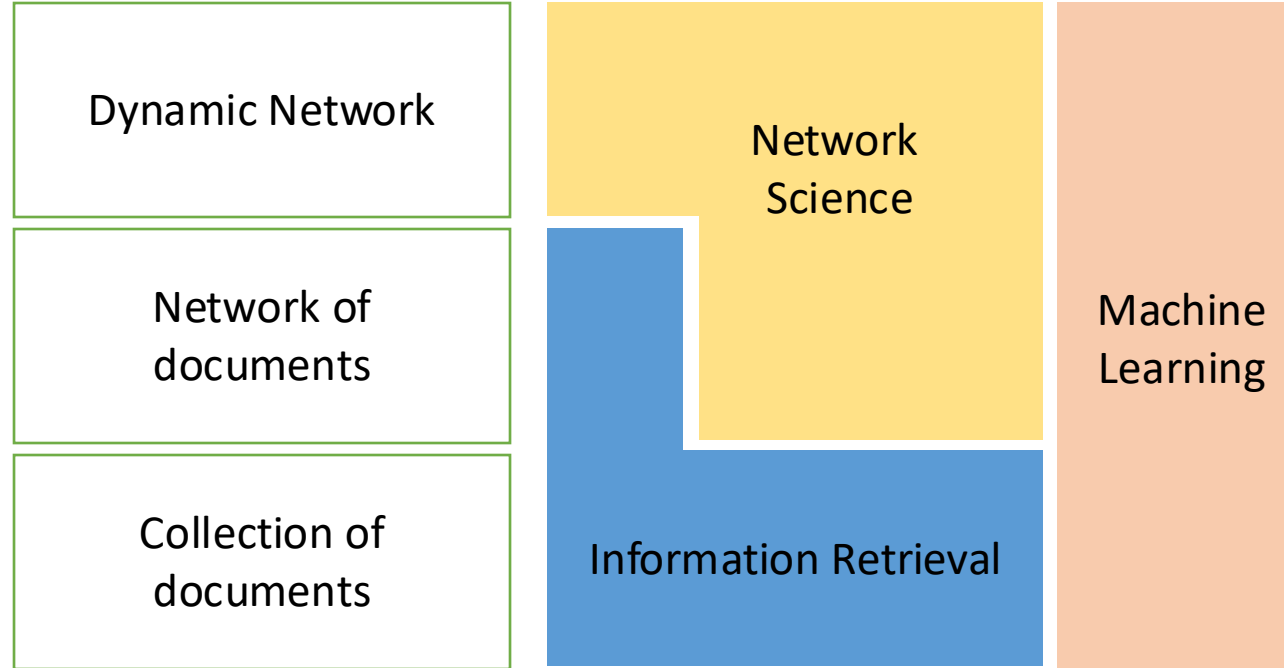
# Different levels of Web

Level 3: Web Dynamics – The web as a dynamic network

- Network evolution
- Dynamic processes on the web
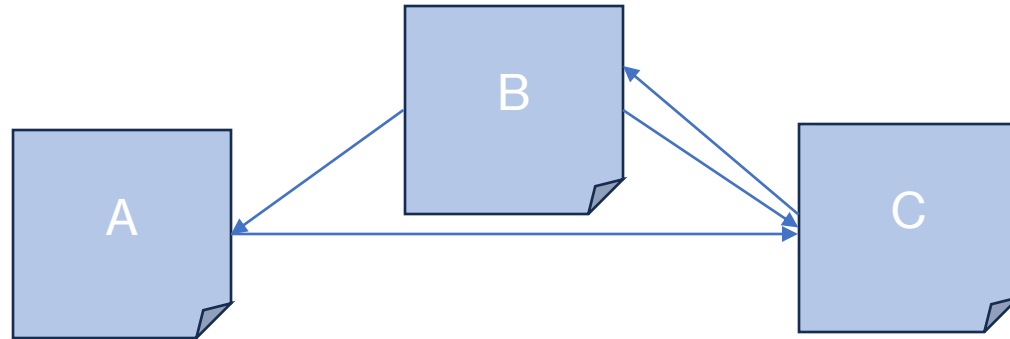
# Relationship to Scientific disciplines



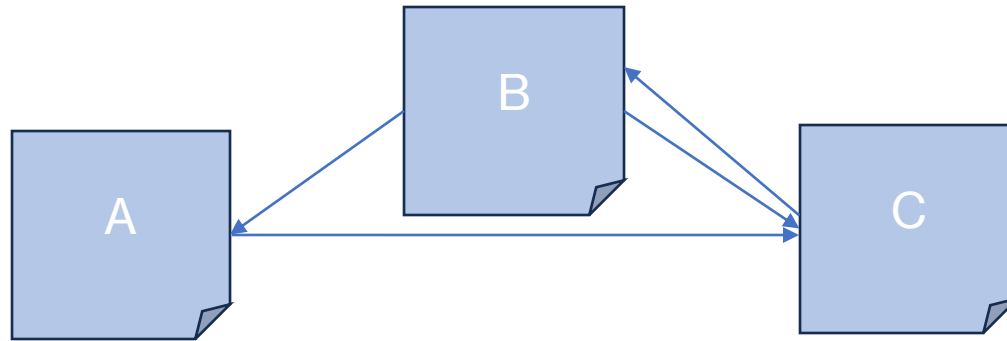Open research Area: Hetereogeneous networks

# Types of Links in Webpages

Types of links:

- Inbound links: The links into the webpage from the outside

- Outbound links: Links from a webpage to other pages in a site or other websites

- Dangling links: Dangling links that point to any page with no outgoing links

# Types of Links in Webpages



Inbound A =

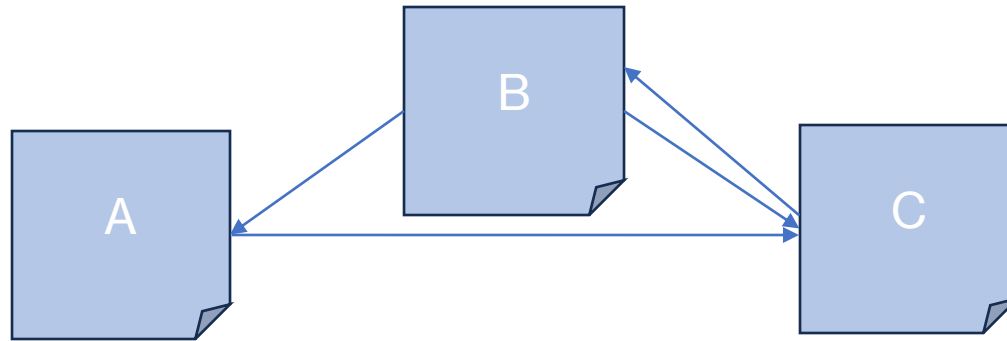Inbound B =

Inbound C =

Outbound A =

Outbound B =

Outbound C =

Dangling =

# Types of Links in Webpages



Inbound A = 1

Inbound B = 1

Inbound C = 2

Outbound A = 1

Outbound B = 2

Outbound C = 1

Dangling = 0

**Web Intelligence, Lecture 1: Introduction, Russa Biswas**

# Course Coverage

- Why do we need this data?

- How do we extract the knowledge available on the Web?

- What are the challenges?

- What are the different forms of data available?

- How do we use the data effectively?

- How do we model the different data formats efficiently?

- How can we use this data to make our life easier?

# Literature

1. https://www.w3.org/TR/webarch/

2. Berners-Lee, T., Cailliau, R., Luotonen, A., Nielsen, H.F. and Secret, A., 1994. The world-wide web. *Communications of the ACM*, *37*(8), pp.76-82.

3. Berners-Lee, T., Cailliau, R., Groff, J.F. and Pollermann, B., 1992. World-Wide Web: the information universe. *Internet Research*, *2*(1), pp.52-58.

4. https://www.w3.org/standards/

5. Berners-Lee, T.J., 1992. The world-wide web. *Computer networks and ISDN systems*, *25*(4-5), pp.454-459

6. https://home.cern/science/computing/birth-web/short-history-web