Introduction
00000

Wikipedia dataset
00

WikipediaSearch implementation
00

WikipediaSearch evaluation
00

Conclusion
00

# WikipediaSearch

## An application of Personalized PageRank on a Wikipedia subset

Roberto Corti

Information Retrieval exam

December 8, 2020

UNIVERSITÀ
DEGLI STUDI DI TRIESTE

DATA SCIENCE &
SCIENTIFIC COMPUTING

## Outline

Introduction

Wikipedia dataset

WikipediaSearch implementation

WikipediaSearch evaluation

Conclusion

# Outline

### Introduction

Wikipedia dataset

WikipediaSearch implementation

WikipediaSearch evaluation

Conclusion

# WikipediaSearch
*A brief introduction*

**WikipediaSearch** is a user-interactive tool that computes a Personalized (or Topic Specific) PageRank over a Wikipedia corpus.

# WikipediaSearch
*A brief introduction*



**Input interface**: user specifies the topics in which he/she has more interest

# WikipediaSearch

*A brief introduction*

**WikipediaSearch**

**Result for Roberto Corti**

- Mean (statistics)
- Statistics
- Computer
- Mathematics
- Science
- Computer science
- Message (computer science)
- Mean
- Sampling (statistics)
- Physics
- Algorithm
- Error
- Object-oriented programming
- Football (soccer)
- Finance
- United States 09d4
- Sampling
- Number
- Computer program
- People
- Programmer

**Output page**: a rank of Wikipedia articles in which the user would be interested to read

# The problem
*"Bringing Order to the Web"*

Algorithm developed by L. Page and S. Brin (Google co-founders) used to determine the order of web pages

**Problem:** Is there a rating system that could measure the human interest and attention devoted to web pages?

Roberto Corti

# The problem
*"Bringing Order to the Web"*

Algorithm developed by L. Page and S. Brin (Google co-founders) used to determine the order of web pages
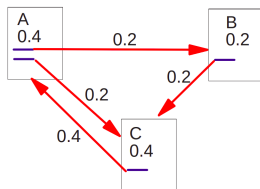
**Problem:** Is there a rating system that could measure the human interest and attention devoted to web pages?

**PageRank idea:** modeling the behavior of a "random surfer" in the Web graph.

# A recap of PageRank theory
*From L.Page and S.Brin (1998)*



$\vec{x}$ : probability distribution vector of the nodes

$M$ : square, stochastic matrix corresponding to the directed grap where $M_{ij} = 1/N_j$
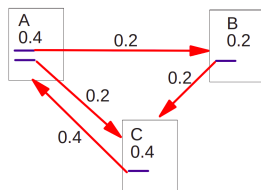
$\Rightarrow$ Find the stationary probability distribution $\Leftrightarrow$ find the unique stochastic eigenvector of to the eigenvalue 1:

$$M\vec{\pi} = \vec{\pi}$$

$\vec{\pi} =$ PageRank vector

# A recap of PageRank theory
*From L.Page and S.Brin (1998)*



Power method solution:

▶ Start with $\vec{x}_0 = \texttt{random}()$

▶ $\vec{x}_i = M\vec{x}_{i-1}$
   until $|\vec{x}_i - \vec{x}_{i-1}| < \epsilon$

# A recap of PageRank theory
*From L.Page and S.Brin (1998)*

Two problems:

▶ **Dangling nodes**: pages without outgoing edges. How to assign probability to them?

▶ **Pages without incoming or outgoing links** : node without incoming edges or group of nodes without outgoing edges. For the first ones we have probability 0 of returning to it once we leave it, while for the others we can never leave them once entered

# A recap of PageRank theory
*From L.Page and S.Brin (1998)*

Allow the random surfer to move to a random page of the graph with probability $\alpha$.

The stochastic matrix will be:

$$
\begin{aligned}
M' &= (1-\alpha)M + \alpha \left[\frac{1}{N}\right]_{N \times N} \\
&= (1-\alpha)M + \alpha \vec{1}^T \cdot \vec{J}
\end{aligned}
$$

where $\vec{J} = (1/N, ..., 1/N)$ is the *jump vector*

# A recap of Topic-Sensitive PageRank
*From Taher H. Haveliwala (2003)*

In addition to the PageRank calculation we can *specialize* the scores of the pages by limiting them to a single topic.

For a given topic-specific set of pages $S$, we allow the random surfer to teleport only to pages that are inside $S$

# A recap of Topic-Sensitive PageRank
*From Taher H. Haveliwala (2003)*

In addition to the PageRank calculation we can *specialize* the scores of the pages by limiting them to a single topic.

For a given topic-specific set of pages $S$, we allow the random surfer to teleport only to pages that are inside $S$

$$M' = (1 - \alpha)M + \alpha \vec{1}^T \cdot \vec{J_S}$$

where the *topic-specific jump vector* is defined as

$$J_{S_i} = \begin{cases} 1/|S|, & \text{if page } i \in S. \\ 0, & \text{otherwise.} \end{cases}$$

Roberto Corti

Introduction
00000

**Wikipedia dataset**
●○

WikipediaSearch implementation
00

WikipediaSearch evaluation
00

Conclusion
00

# Outline

Roberto Corti

Introduction
○○○○○

Wikipedia dataset
○●

WikipediaSearch implementation
○○

WikipediaSearch evaluation
○○

Conclusion
○○

# Wikipedia ...

content...

# Outline

Introduction
00000

Wikipedia dataset
00

WikipediaSearch implementation
0●

WikipediaSearch evaluation
00

Conclusion
00

# PageRank implementation

# Outline

# Standard PageRank

# Outline

Introduction

Wikipedia dataset

WikipediaSearch implementation

WikipediaSearch evaluation

Conclusion

# Conclusion