# An analysis of scholarly citations in Wikipedia

Alessio Bogon

University of Trento
Department of Information Engineering and Computer Science

Master Degree in Computer Science

March 24, 2016

# Overview

# Introduction

# Wikipedia

- Most used encyclopedia
- Started in 2001
- Studied by many
  - Quality of citations [1]
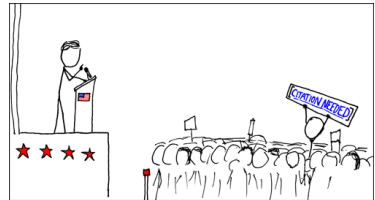  - Illness prediction [2]
  - Stock market moves strategy [3]



Figure: Wikipedian protester
(xkcd.com/285)

# Open questions

- Quality of papers in Wikipedia, in term of:
    - Incoming citations
    - Journals rank
    - Lifetime of citations
- Prediction systems
    - If a paper appears in Wikipedia, will it become more popular in the scientific community?
    - In another way, do researcher use Wikipedia as their primary data source?
    - Predict whether a publication is going to stay on a page

## Available datasets

- Microsoft Academic Graph: a dataset containing data about papers, authors, references, journals, conferences, etc.
- Wikipedia dumps: text of all page revisions since the beginning
- Wikimedia hourly page view statistics

# Microsoft Academic Graph

- Dataset powering the Microsoft Academic Search engine
- Size: 96 GB
- Contains over 120M papers (1800 – 2016)
- Information about authors, references, journals, conferences, keywords, etc.
- Problems
  - Only *computer science* conferences
  - Some of papers' publication dates are incomplete
  - Not all the papers have a DOI (32% of them)

# The missing pieces

- History of papers appearing in Wikipedia (where and when)
- Usable/searchable page views dataset

# Data manipulation

# Extracting citations from Wikipedia

## Problems

- Citations can be structured: *wikimarkup* templates
    - Many different variants
    - Anybody can use custom macros
    - Different templates for each language
- or unstructured: plain text
    - Recognize substrings that appear to be citations
- Entity disambiguation
- Dataset size: 13,3 TB as of September 1st, 2015

## Solution

- Focus on publication identifiers (*DOI*, *PMID*, *arXiv*, *ISBN*)
- The **wikidump** framework

# Wikidump

- Facility framework to extract features from Wikipedia dumps
- Based on libraries by Aaron Halfaker
- Low memory consumption
- Highly parallelizable
- Written in Python

- Processed 445M page revisions (13,3 TB) in 21 hours

| Type | Count |
|------|-------|
| ISBN | 1 153 330 |
| DOI | 651 199 |
| PMID | 372 939 |
| PMC | 79 841 |
| arXiv | 18 832 |

Table: Number of identifiers extracted

# Wikimedia page views

### Problems

- Dataset size: 23 TB (4,7 TB for the 2014)
- Aggregated and ordered by hour (8670 files per year)
- They need to be cleaned
- They need to be **reordered**
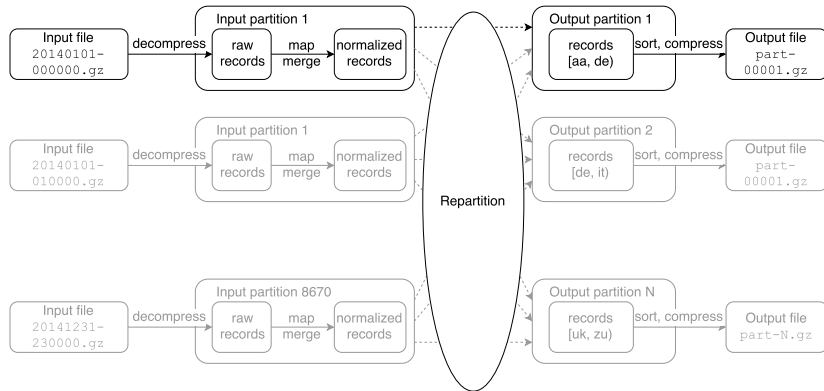- Unfeasible on a single machine

### Solution

- Exploit the UniTN Cisca Cluster

# The Spark job

- UniTN Cisca Cluster
    - 125 workstations
    - 500 CPU cores in total
    - Available only at night and in the weekend
- The job
    - Normalize the content
    - Sample the first file
    - Repartition the keyspace
    - Sort each partition locally

- Took one night for the 2014 dataset
- Took many nights to get it to work

# The Spark job — Workflow



Figure: Workflow showing the processing of Wikimedia page views
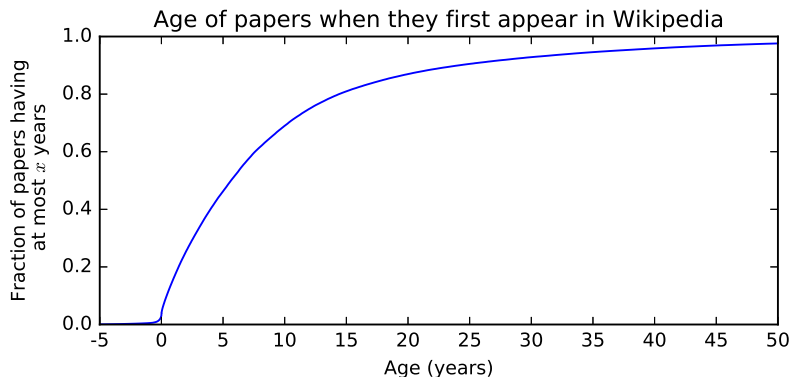
# Results

# Quality of papers in Wikipedia

- Age of papers when inserted
- Incoming citations distribution
- Journals rank
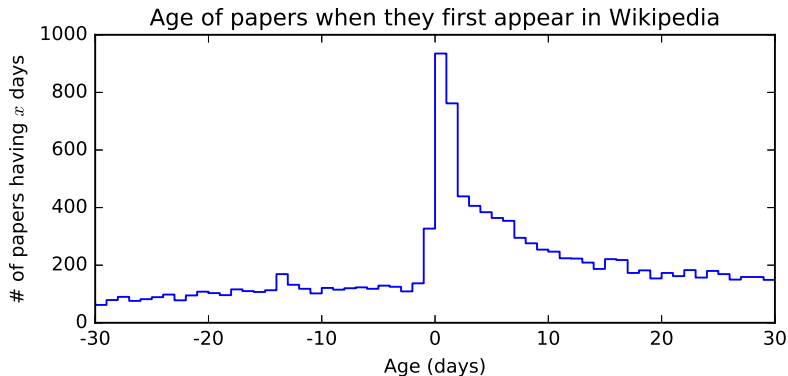- Lifetime of citations

# Age of papers when inserted

- How old is a paper when it is insterted in Wikipedia for the first time?
- Interesting behavior of papers having few days
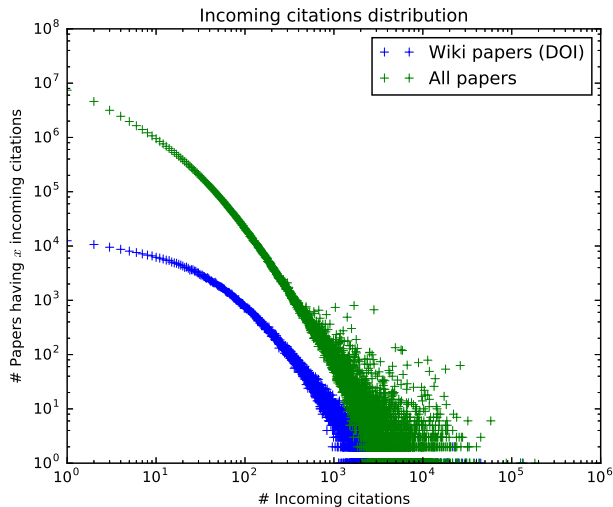
# Age of papers when inserted



Age of papers when they first appear in Wikipedia

# Age of papers when inserted — Detail



Age of papers when they first appear in Wikipedia

1.3% of papers are inserted within 7 days after the publication

# Incoming citations distribution

- Arguably follows a power law [4]: $N(x) \sim x^{-\alpha}$
- How well papers in Wikipedia perform?



Incoming citations distribution

# Incoming citations distribution

- Papers in Wikipedia behave like *Genome Research* and *PNAS*

- They outclass *Nature* and *Science*

- 75% of papers in Wikipedia have more than 10 incoming citations
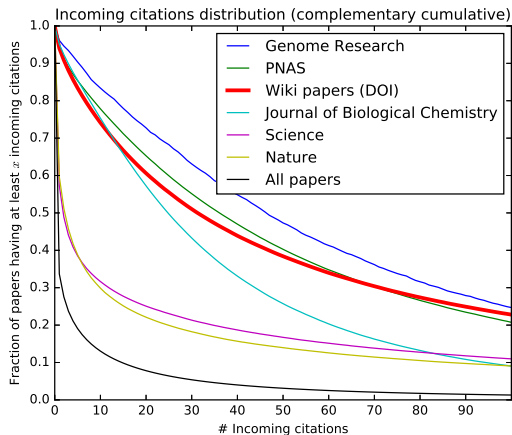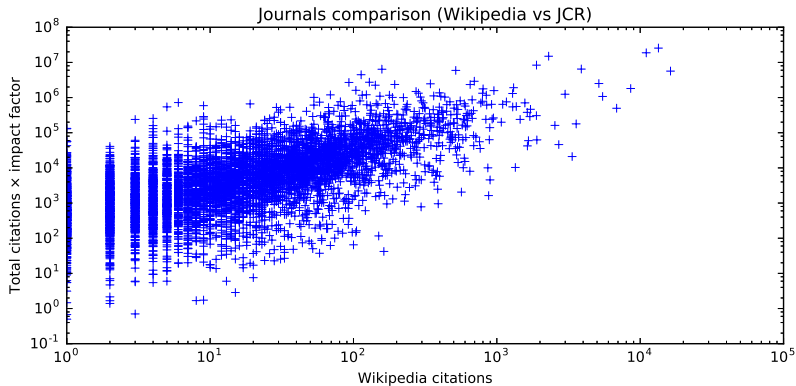


Figure: Papers cited in Wikipedia vs papers in top journals
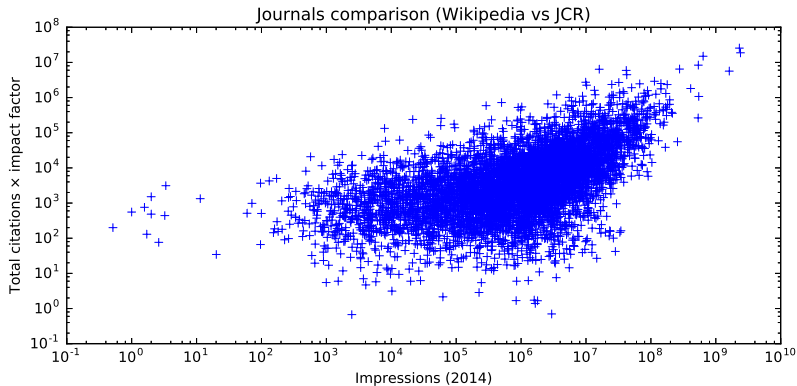
# Journals rank

- First proposed by Nielsen in 2007 [1]
- Most cited journals in Wikipedia are also the most important ones

- Journals rank by impact factor versus:
    - citations in Wikipedia
    - visualizations in Wikipedia (in 2014)
- Measured in term of Kendall rank correlation coefficient ($-1 \leq \tau \leq 1$)

# Journals impact factor vs Wikipedia citations



Journals comparison (Wikipedia vs JCR)

Kendall rank correlation coefficient: 0.464
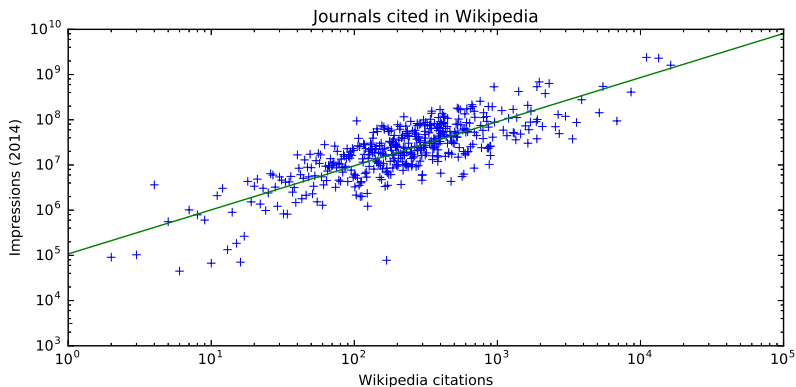
# Journals impact factor vs Wikipedia impressions



Journals comparison (Wikipedia vs JCR)

Kendall rank correlation coefficient: 0.401

# Maybe: Journals citations vs impressions in Wikipedia
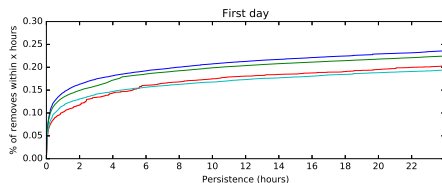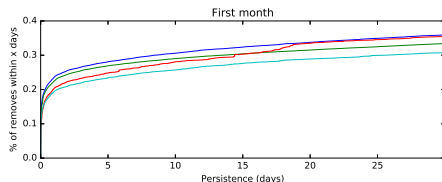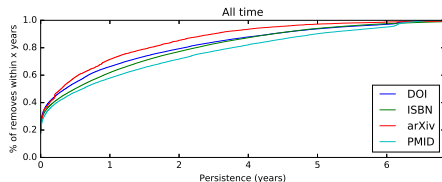


Journals cited in Wikipedia

Linear regression model: $\hat{y}_i = 10^{5.03} \times x^{0.98}$ ($r^2 = 0.65$)

# Lifetime of irrelevant identifiers

- How long does it take for a Wikipedia contributor to discover and remove an "irrelevant" paper from an article?
- An identifier is "irrelevant" for an article if it appeared on that page and was then removed.

- Fraction of identifiers removed in time

- 60% of irrelevant DOI removed within one year
- 35% within one month
- 20% within one day

# Conclusion

## References I

📄 Finn Arup Nielsen.
Scientific citations in Wikipedia.
*First Monday*, 12(8):1–5, August 2007.

📄 David J. McIver and John S. Brownstein.
Wikipedia Usage Estimates Prevalence of Influenza-Like Illness in the
United States in Near Real-Time.
*PLoS Computational Biology*, 10(4), 2014.

📄 Helen Susannah Moat, Chester Curme, Adam Avakian, Dror Y.
Kenett, H. Eugene Stanley, and Tobias Preis.
Quantifying Wikipedia Usage Patterns Before Stock Market Moves.
*Scientific Reports*, 3:1–5, 2013.

# References II

Sidney Redner.
How popular is your paper? An empirical study of the citation distribution.
*The European Physical Journal B - Condensed Matter and Complex Systems*, 4(2):131–134, 1998.