

Interesting Content and User Discovery without Upvotes

An application of centrality measures to Reddit

Rollen S. D'Souza

rollen.dsouza@uwaterloo.ca

University of Waterloo

2018

Outline

- 1 Introduction
 - About Reddit
 - Mathematical Preliminaries
- 2 Technical Notes
- 3 Case Study: /r/uwaterloo
- 4 General Observations
- 5 Summary

What is Reddit?

Purposes

- Social news aggregation

What is Reddit?

Purposes

- Social news aggregation
- Discussion platform

What is Reddit?

Purposes

- Social news aggregation
- Discussion platform
- Community for help or relaxation

What is Reddit?

Basic Structure

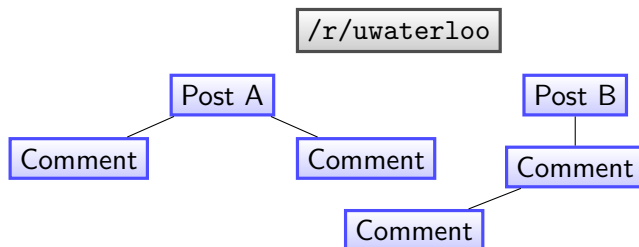


Figure: Graph-like structure of the Reddit platform.

What is Reddit?

↑ 94 ↓ | 🔗 **Infinity** Highschooler proves infinity is odd

↑ **r/badmathematics** · Posted by u/[deleted] 1 year ago 📄

94 **Infinity** Highschooler proves infinity is odd

reddit.com/r/AskR... 🔗

💬 85 Comments Share ... 99% Upvoted

This thread is archived
New comments cannot be posted and votes cannot be cast

SORT BY BEST ▾

👤 jacob8015 68 points · 1 year ago (More than 35 children)

👤 **IcedRoren** **Cursed by Dimensionality** 54 points · 1 year ago

👤 Idk what's worse. The argument, or the fact that the geometry teacher confirmed the proof.

Share Save Edit ...

👤 **Snowayne2** 37 points · 1 year ago

👤 The worst is how arrogant he is about the whole thing.

Is the light bulb on yet?

Share Report Save Give gold

Mathematical Preliminaries

What I assume:

- Centrality Measures
- Properties of Laplacians

The former is how I sort through content and users.

Katz Centrality Measure

Definition

The *Katz Centrality Measure* for a node i in graph G with adjacency matrix A is defined as,

$$c_i = \sum_{k=1}^{\infty} \sum_{j=1}^n \alpha^k (A_{j,i})^k,$$

where $0 < \alpha \leq \rho(A)$.

Katz Centrality Measure

Why this measure?

- Has larger reach in the graph (looks at all paths to a node.)
- Easy to compute.
- Loose requirements on A .
- Intuitive meaning.

On Laplacians

Laplacians can be used for clustering. Recall,

$$L = D_{out} - A.$$

Use eigenvalues near 0 for clustering (eigenvectors associated with nearly disconnected components).

On Laplacians

Laplacians can be used for clustering. Recall,

$$L = D_{out} - A.$$

Use eigenvalues near 0 for clustering (eigenvectors associated with nearly disconnected components).

So...finding eigenvalues is easy right?

On Laplacians

Different Formulation

Problem: The graph can have as many as 10^4 nodes!

On Laplacians

Different Formulation

Problem: The graph can have as many as 10^4 nodes!

Definition

The *normalized Laplacian* of a graph is defined as,

$$L_{rw} = I_{n \times n} - D_{out}^{-1}A$$

On Laplacians

Different Formulation

Problem: The graph can have as many as 10^4 nodes!

Definition

The *normalized Laplacian* of a graph is defined as,

$$L_{rw} = I_{n \times n} - D_{out}^{-1}A$$

Need to ensure the graph has no sinks!

About the Code

Code was implemented using Python 3 and MATLAB. Code is available on my Github. Other than that,

- redis for storing data in-memory.
- numpy, scipy for computation outside of MATLAB.

About the Data

Sourced thanks to Jason Baumgartner who uploaded data dumps of the JSON objects gathered from the Reddit API. Structures contain,

- Submissions with author (id and flair), tags, score, text or link,
- Comments with author (id and flair), text
- and much much more!

About the Data

Sourced thanks to Jason Baumgartner who uploaded data dumps of the JSON objects gathered from the Reddit API. Structures contain,

- Submissions with author (id and flair), tags, score, text or link,
- Comments with author (id and flair), text
- and much much more!

Data for: November 2017 — February 2018

About the Assumptions

In order to do any analysis, we need to make some assumptions.

- 1 Users are interested (positively or negatively) by content they read.

Data for: November 2017 — February 2018

About the Assumptions

In order to do any analysis, we need to make some assumptions.

- ① Users are interested (positively or negatively) by content they read.
- ② Users comment only on content they have read.

Data for: November 2017 — February 2018

About the Assumptions

In order to do any analysis, we need to make some assumptions.

- ① Users are interested (positively or negatively) by content they read.
- ② Users comment only on content they have read.
- ③ Deleted content and users are not interesting nor generate interesting content and do not affect other users in a meaningful way.

Data for: November 2017 — February 2018

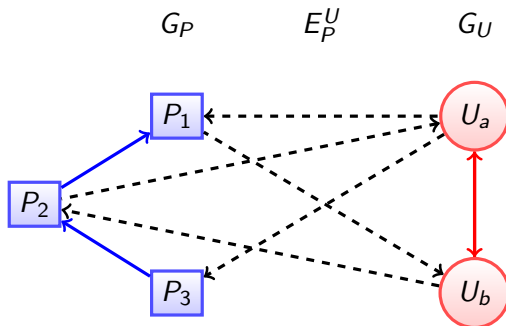
About the Assumptions

In order to do any analysis, we need to make some assumptions.

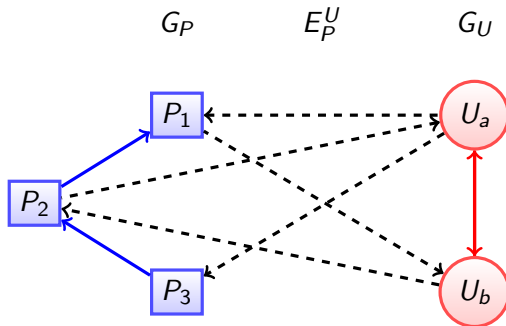
- ① Users are interested (positively or negatively) by content they read.
- ② Users comment only on content they have read.
- ③ Deleted content and users are not interesting nor generate interesting content and do not affect other users in a meaningful way.
- ④ Active engagement of users (through comments or submissions) is a good measure of interest.

Data for: November 2017 — February 2018

Structuring the Graph



Structuring the Graph



Note that there are atleast 3 interesting graphs here. Where are they?

Applied centrality analysis to reveal the following top nodes,

/r/uwaterloo	/r/science	/r/CanadianInvestor
supersonic63	t3_7e1jo1	johnnychi
kw2002	AutoModerator	EyesOnTsx
t3_7f2c1d	mvea	jnf_goonie
t3_7w0dgv	t3_7vxito	langlois44
t3_7i045q	t3_7s6a9z	t3_7ii8af
microwavemasterrace	t3_7my58a	t3_7fapx7
uwsmile	t3_7okc7u	t3_7pawu6
t3_7ld2jb	t3_7hs6ro	helkish
mywaterlooaccount	t3_7sv8vb	Ginhisf
honhonhonFRFR	t3_7qs5xz	MarketStorm

Of Interest to /r/uwaterloo

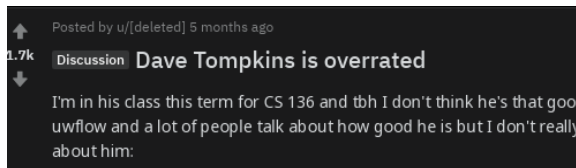
The more generic submissions...

Link	Description	Score	Comment #
f2c1d	Admissions Megathread	311	12k
7i045q	Mr. Goose for good exams...	142	389
7ld2jb	Mr. Goose for good grades...	196	371

A mass thread on WaterlooWorks also shows up in the top 20.

Of Interest to /r/uwaterloo

t3_7w0dgv: A much more interesting post...



Of Interest to /r/uwaterloo

t3_7w0dgv: A much more interesting post...

↑ Posted by u/[deleted] 5 months ago

1.7k Discussion **Dave Tompkins is overrated**

↓

I'm in his class this term for CS 136 and tbh I don't think he's that good uwflow and a lot of people talk about how good he is but I don't really about him:

↑ -dtompkins- CS Lecturer/Advisor 47.7k points · 5 months ago 🏆 12

↓ Dude.... **I ALSO think I'm overrated.**

- -- pauses to think about his response while he has a sip of Coke

I'll be honest, when I started to become infamous for having good thought that student expectations would be way too high, and the actually in my class.





General Observations


On /r/science

/r/uwaterloo	/r/science	/r/CanadianInvestor
supersonic63	t3_7e1jo1	johnnychi
kw2002	AutoModerator	EyesOnTsx
t3_7f2c1d	mvea	jnf_goonie
t3_7w0dgv	t3_7vxito	langlois44
t3_7i045q	t3_7s6a9z	t3_7ii8af
microwavemasterrace	t3_7my58a	t3_7fapx7
uwsmile	t3_7okc7u	t3_7pawu6
t3_7ld2jb	t3_7hs6ro	helkish
mywaterlooaccount	t3_7sv8vb	Ginhisf
honhonhonFRFR	t3_7qs5xz	MarketStorm

General Observations

On /r/science

 Posted by u/nate **PhD | Chemistry | Synthetic Organic** 8 months ago  7  

124k  **Subreddit Discussion** **Raising the taxes of graduate students by as much as 300% disaster for the USA**

Science and technology development has been the story of the past 100 years. The discoveries and innovations progressing at a dazzling rate, much of this led by researchers at universities in the USA. At these universities a substantial amount of the work is done by graduate students, who work long hours (80 hours weeks are common) for little pay. These graduate students go on to work in good paying jobs, where their innovations make money. Start-ups develop to bring new innovations based on the skills graduate students learn (Google was the work of a couple of Stanford grad students, even Reddit benefited from the skills of a physics grad student/PhD, its current CTO.) Grad school has been for decades a path to prosperity for those who come from humble beginnings, willing to work hard, and make sacrifices, a system that has greatly benefited all of us.

This is why we scientists are shocked and appalled by the recently passed tax bill in congress which will increase the bills of already poor grad students **going up by as much as 300%**, which would see their take-home pay cut in half. As a former grad student myself, I can tell you that I would not have been able to continue if my pay had been cut by \$7,000, and many students would make the same conclusion. Instead, some will not go into science or engineering in the USA to be a grad student in Europe or Asia, most of these students will never return to the USA.

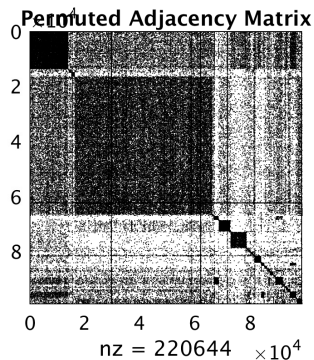
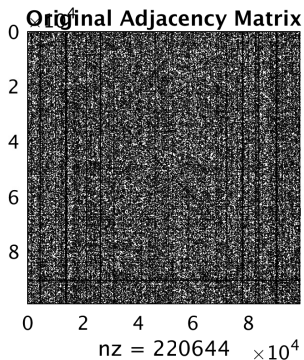
This is why every major science organization has **voiced opposition to the current tax plan**, make no mistake, it will **undermine research and eventually the economy of the USA**.

In comic form from **PhD Comics**.



General Observations

On /r/science



General Observations

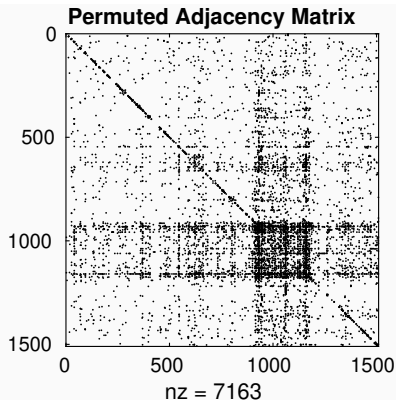
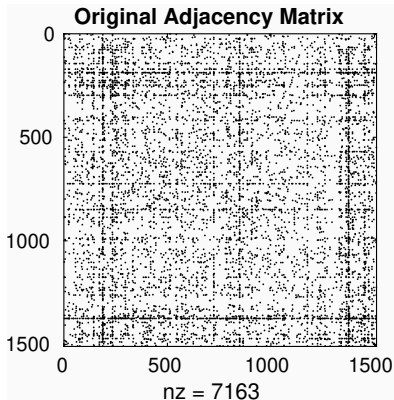
On /r/CanadianInvestor

/r/uwaterloo	/r/science	/r/CanadianInvestor
supersonic63	t3_7e1jo1	johnnychi
kw2002	AutoModerator	EyesOnTsx
t3_7f2c1d	mvea	jnf_goonie
t3_7w0dgv	t3_7vxito	langlois44
t3_7i045q	t3_7s6a9z	t3_7ii8af
microwavemasterrace	t3_7my58a	t3_7fapx7
uwsmile	t3_7okc7u	t3_7pawu6
t3_7ld2jb	t3_7hs6ro	helkish
mywaterlooaccount	t3_7sv8vb	Ginhisf
honhonhonFRFR	t3_7qs5xz	MarketStorm



General Observations

On /r/CanadianInvestor



Summary

- Reddit is cool!

Summary

- Reddit is cool!
- Users totally read what they react to.

Summary

- Reddit is cool!
- Users totally read what they react to.
- Talked about how Dave Tompkins is awesome.

Summary

- Reddit is cool!
- Users totally read what they react to.
- Talked about how Dave Tompkins is awesome.
- Obligatory reference to Trump.

Summary

- Reddit is cool!
- Users totally read what they react to.
- Talked about how Dave Tompkins is awesome.
- Obligatory reference to Trump.
- Importance of investing expertise.