



Reddit network analysis

Ema Pajić, Nikola Aleksić, Predrag Obradović, Marko Mišić, Bill Power,
Zoran Obradović

Introduction

- **Reddit** is a social news aggregation, web content rating, and discussion website.
- Posts are organized by subject into user-created boards called **communities** or **subreddits**.
- Some stats:
 - 50M+ daily active users
 - 100K+ active communities
 - 50B+ monthly views



Motivation

- Many possible research directions using Reddit data
- Analyzing online conflicts – [Community Interaction and Conflicts on the Web](#)
- Hate speech analysis - [Dreaddit: A Reddit Dataset for Stress Analysis in Social Media](#)
- Impact of real life events (Covid, elections, finance/crypto)
- Analyzing how people use the website - [Generalists and Specialists: Using Community Embeddings to Quantify Activity Diversity in Online Platforms](#)
- **Prior work either doesn't use user overlap, or uses basic variants**

Pushshift

- **Pushshift Reddit API** collects Reddit data every day and provides easy access to everyone.
- Accessing data via **crawler API** or **large monthly dumps**

Papers Published Using Pushshift Data 2016-2019

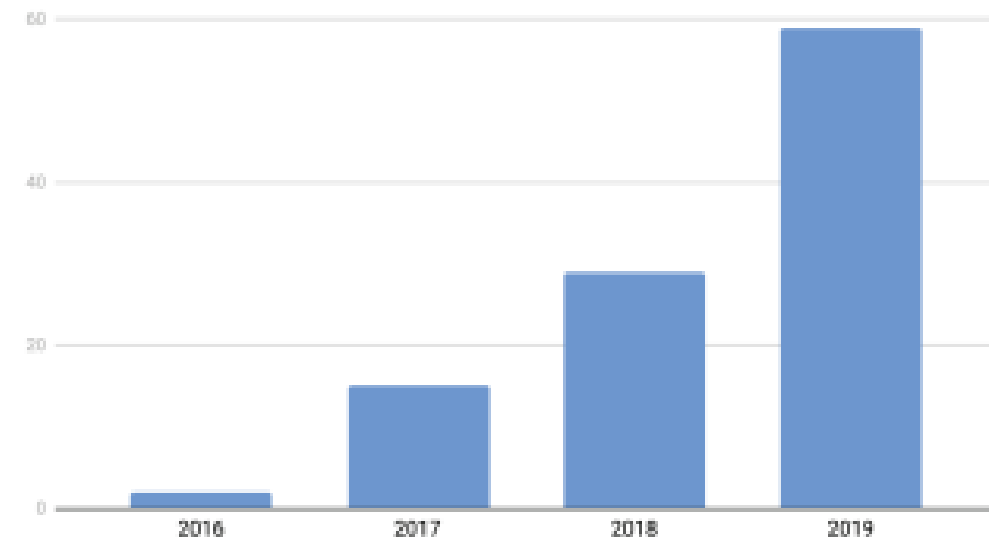
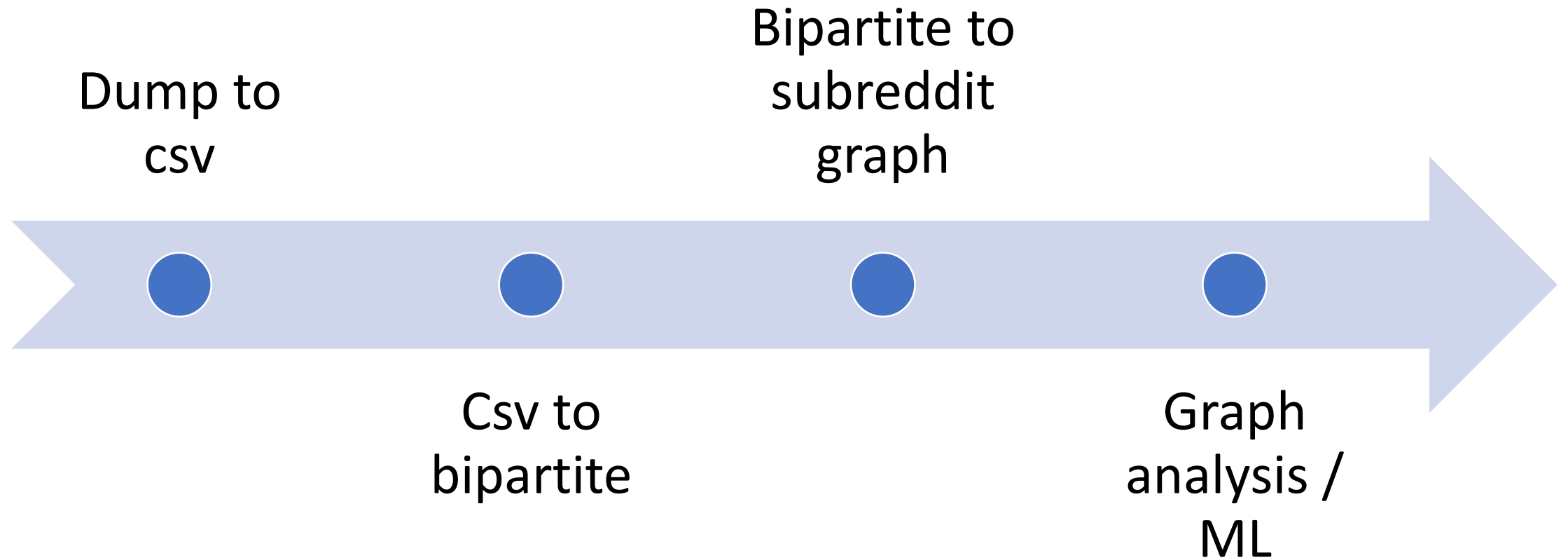


Figure from **The Pushshift Reddit Dataset Paper**

Datasets

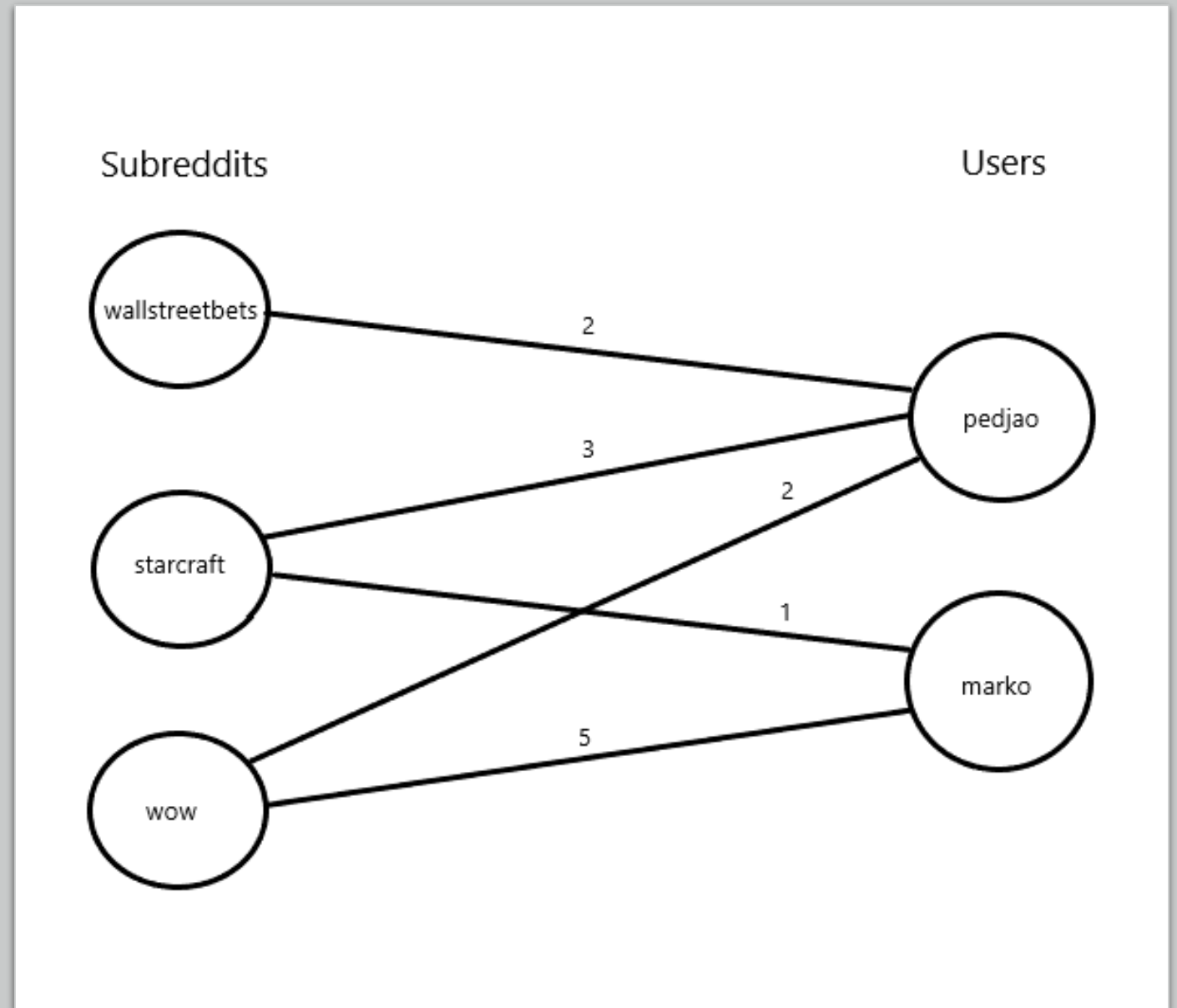
- Crawling data
 - Possible through Pushshift API
 - pmaw, psaw
 - Not reliable
- Pushshift dumps
 - December 2005 – June 2021
 - One dump contains all comments and posts posted in a month
 - 2021 – around 20GB compressed, 200GB uncompressed for each month

Pipeline



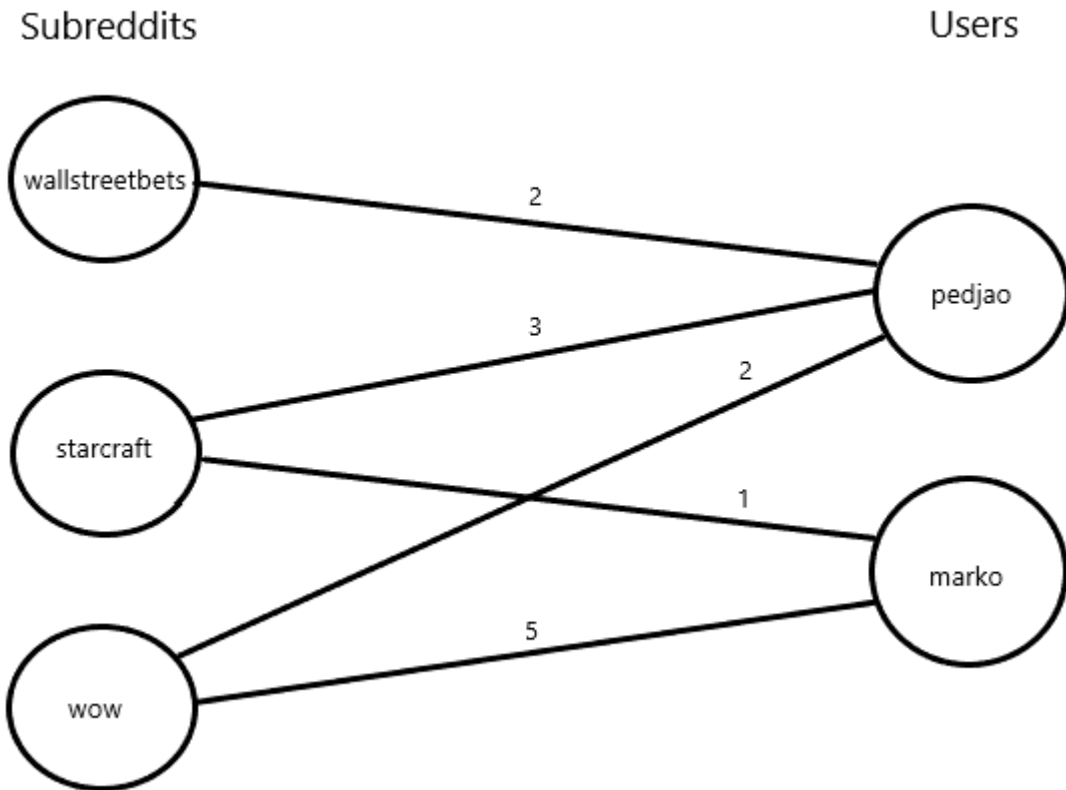
Modelling Reddit

- Bipartite graph – Weight between **Subreddit** and **User** equal to number of comments **User** posted on that **Subreddit**
- Many different ways to transform into Unipartite graph

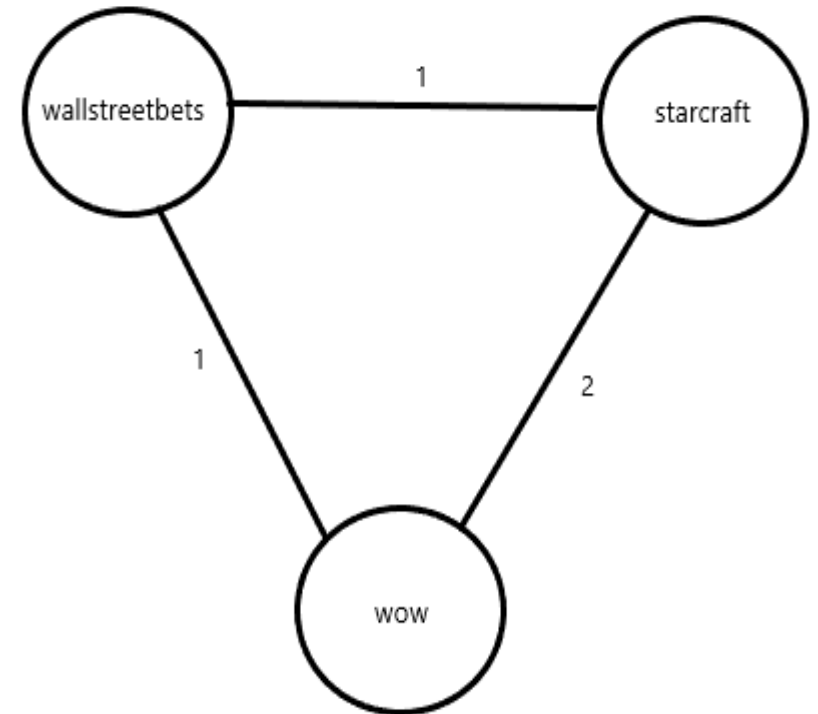


Count overlapping users

Bipartite graph



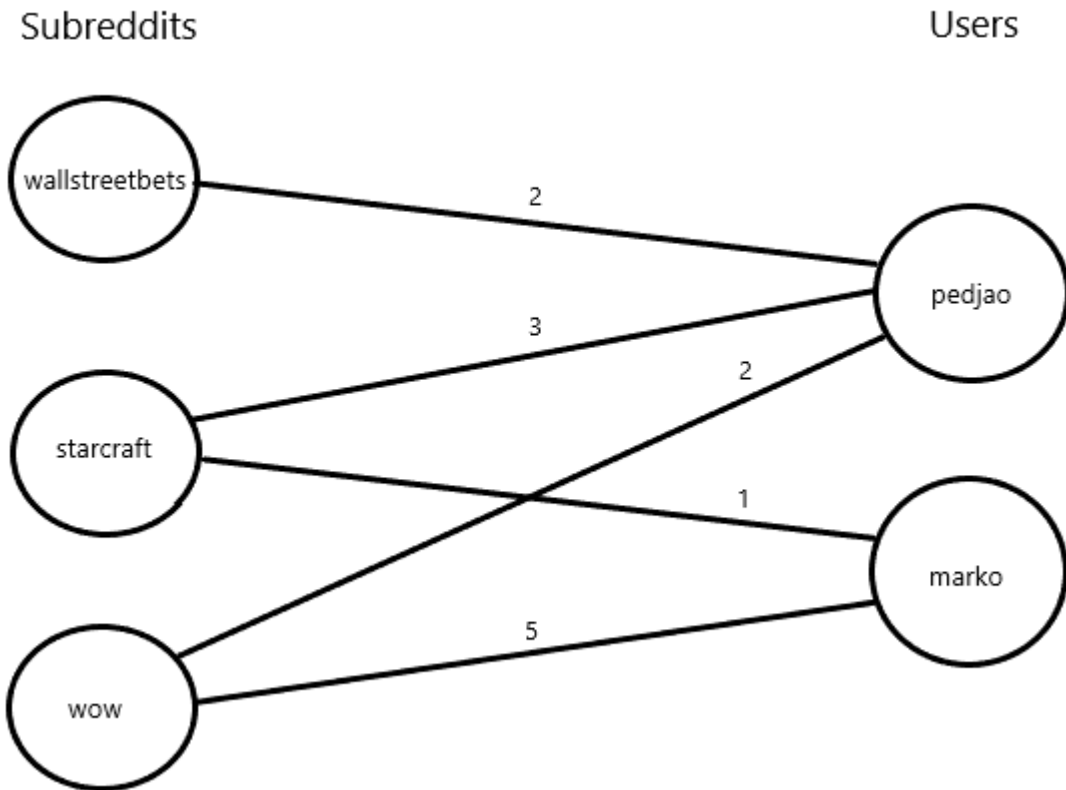
Unipartite graph



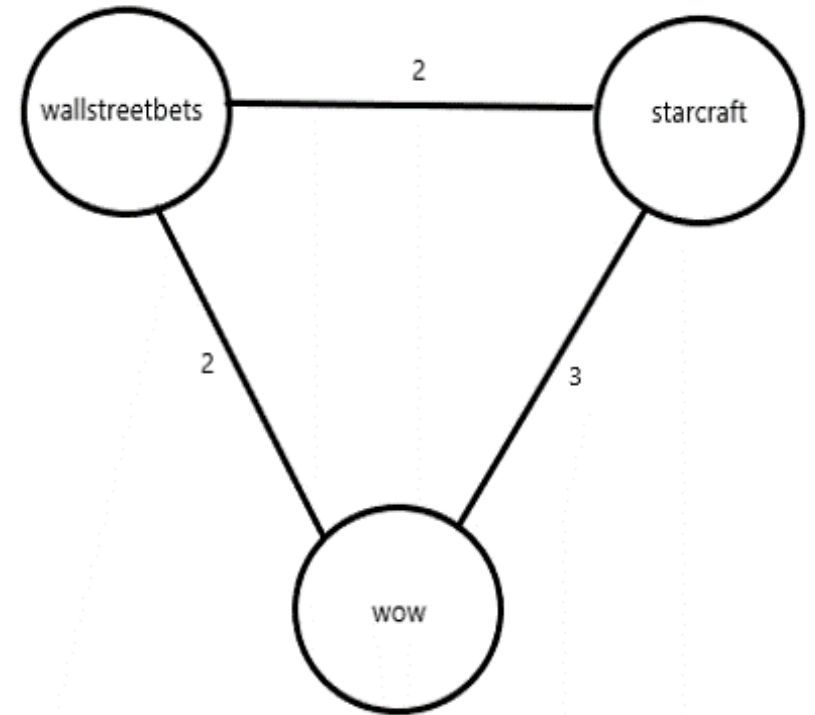
Problem: Two subreddits might have one user posting 1 time on both, while another user posted 50 and 60 times on each, respectively. We don't want to count these the same!!!

Sum minimum between user weights

Bipartite graph



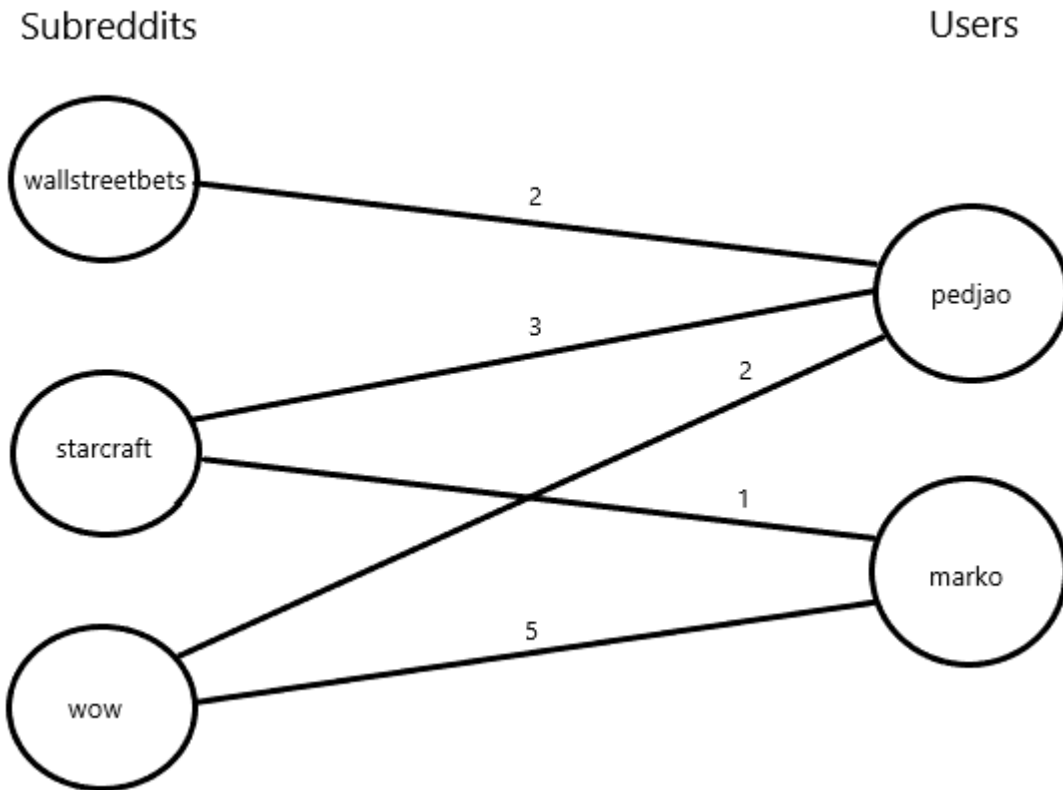
Unipartite graph



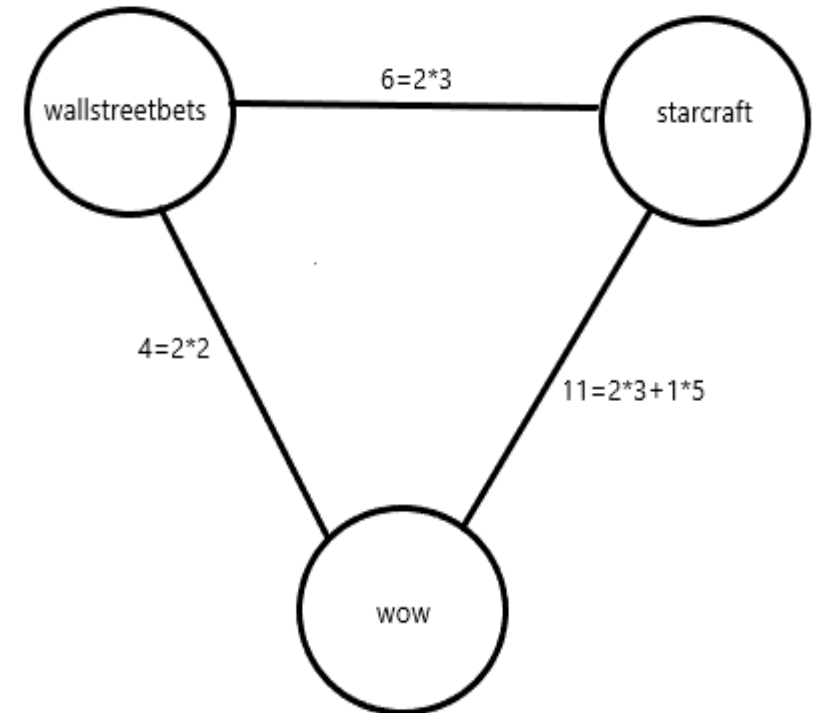
Problem: Strong connection between **marko** and **wow** is disregarded. Also, **pedjao** doesn't visit **starcraft** and **wow** the same but the weight to **wallstreetbets** is the same.

Sum product between user weights

Bipartite graph



Unipartite graph



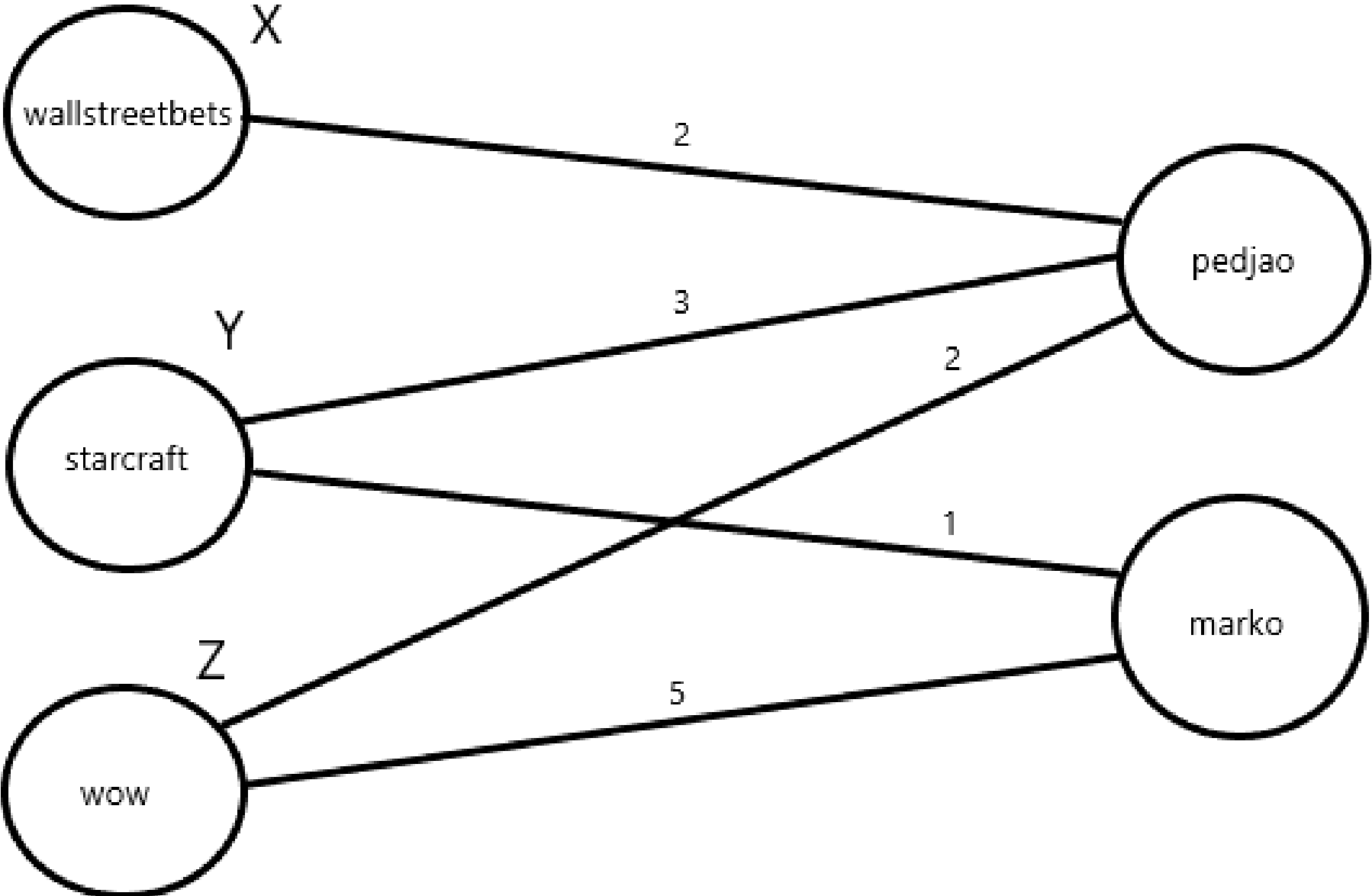
Problem: Weights getting too big. What does the weight of 6 mean, is it big or small, does it correspond to a strong connection?

Two pass aggregation¹

1. Assign 1 unit of **resource** to a subreddit.
2. Redistribute resource to connected users, splitting it proportionally to edge weights.
3. Every user does the same, redistributes the weight to subreddits, proportionally to edge weights.
4. Amount of resource that started in subreddit1 and ended in subreddit2 we will be the edge weight!

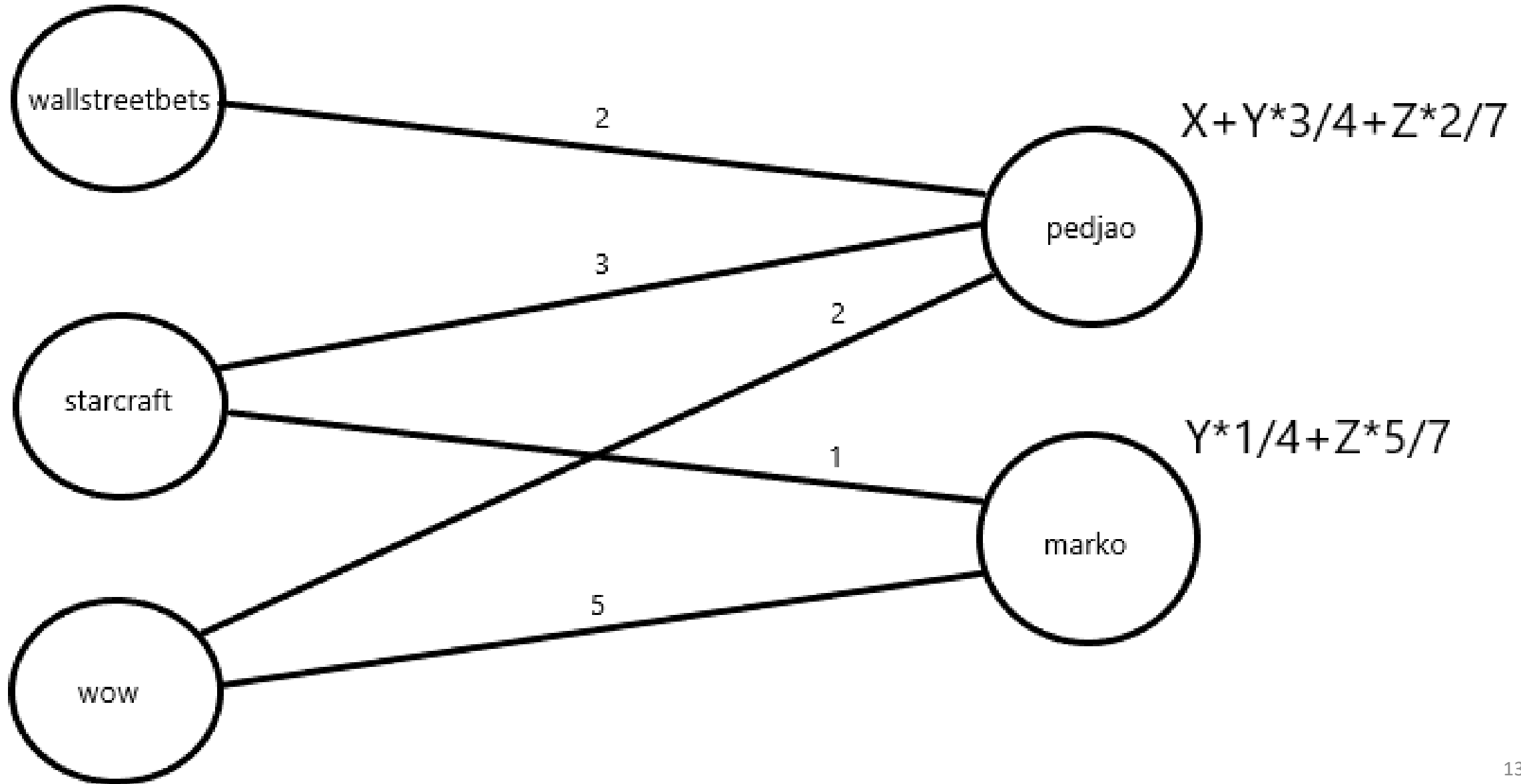
Subreddits

Users



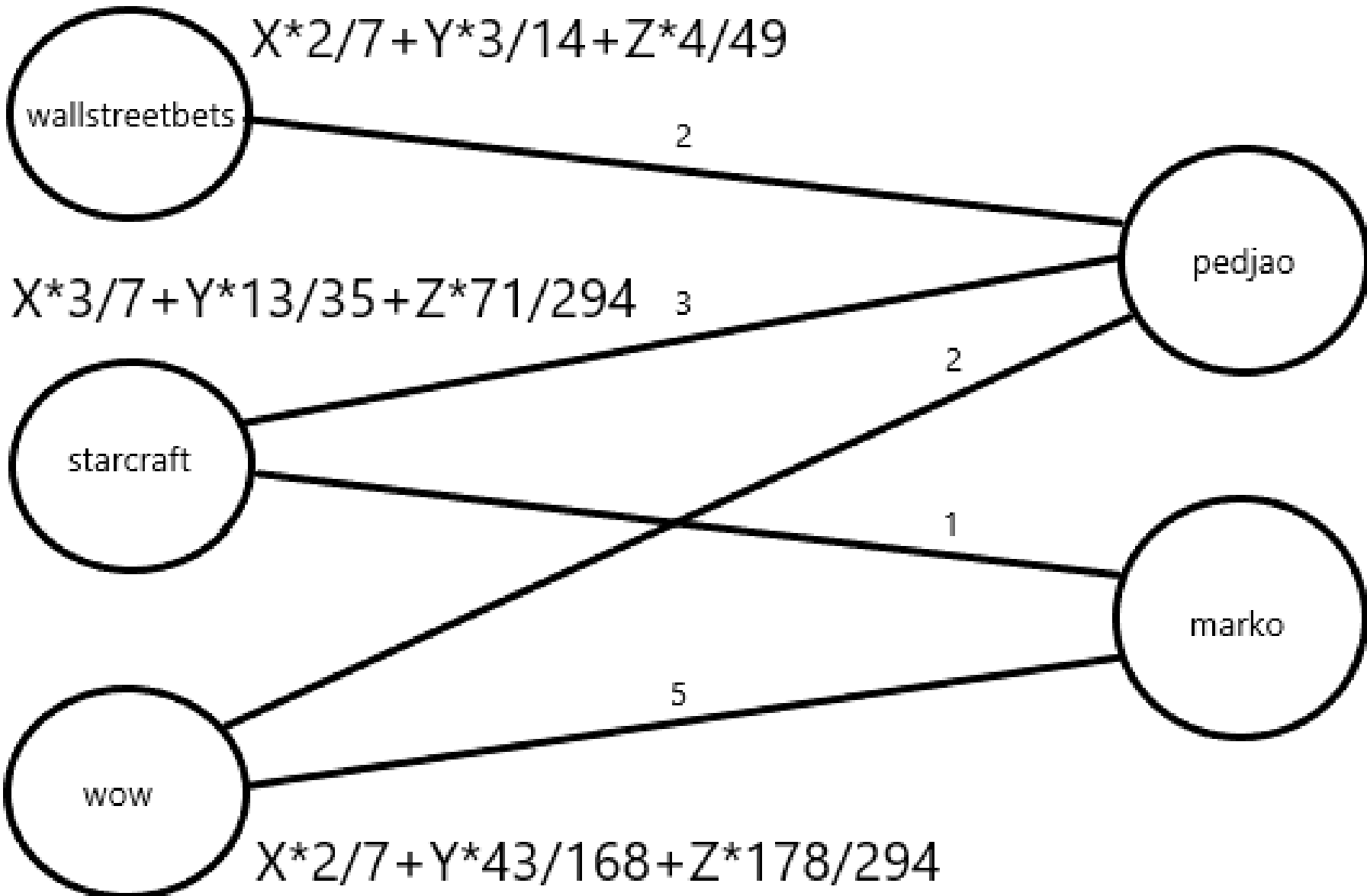
Subreddits

Users



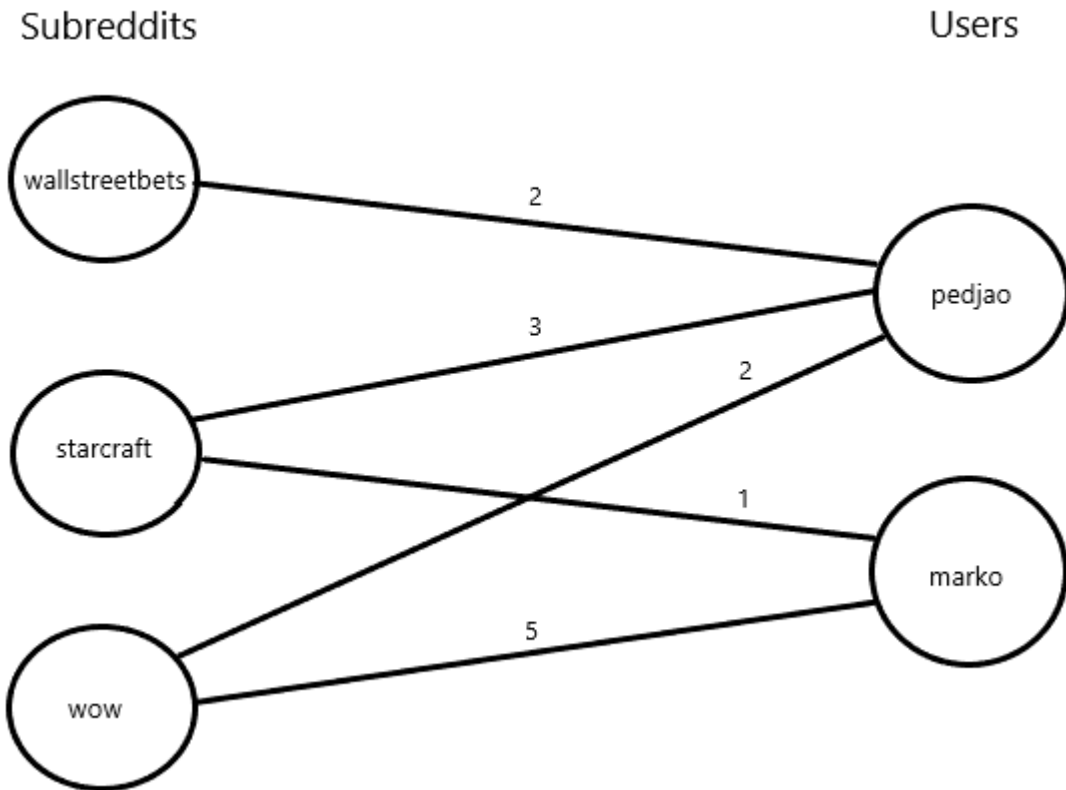
Subreddits

Users

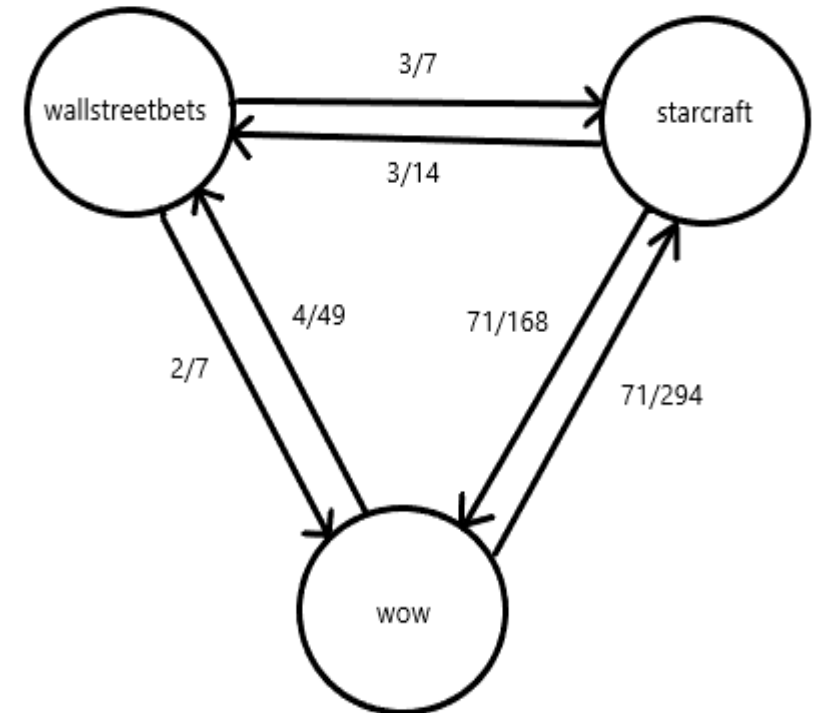


- Two pass aggregation

Bipartite graph



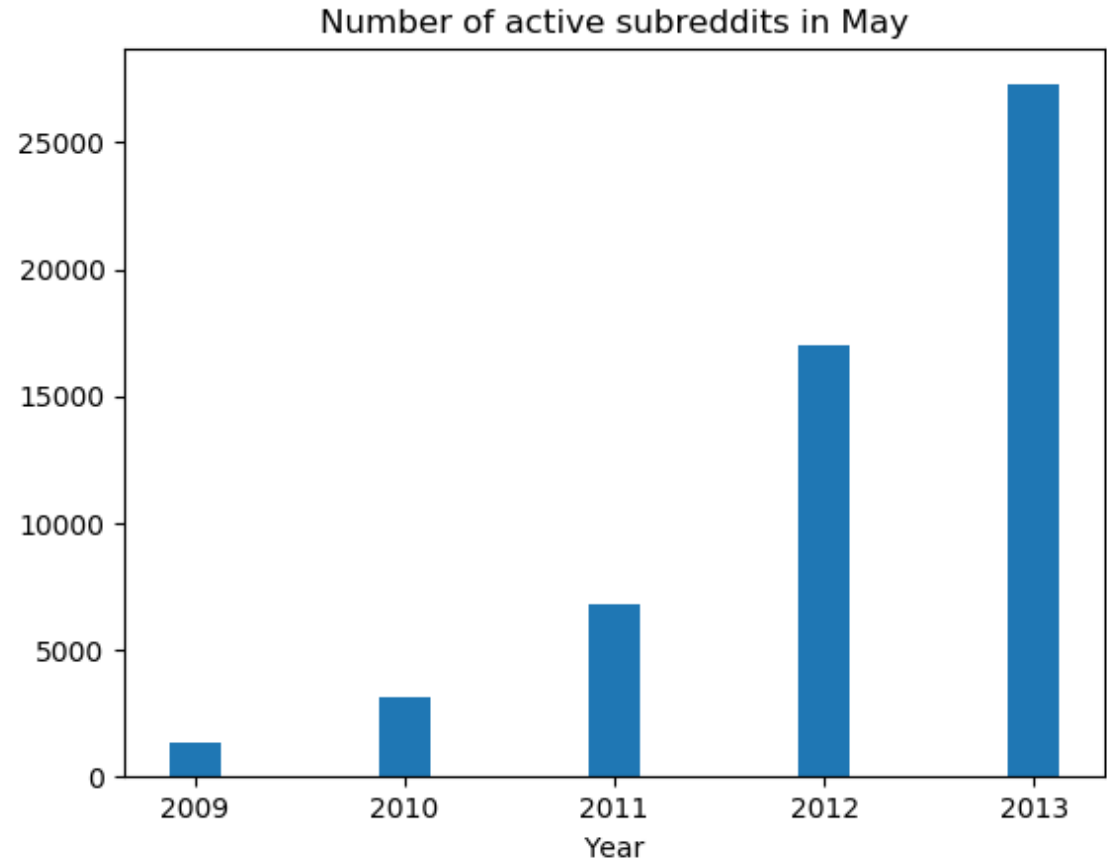
Unipartite graph



Awesome properties: Similar to using product as earlier, but normalizes by both user and subreddit sizes. **Weights are between 0 and 1!!!**

Applying on real data

- Reddit data is huge
- Even for one month in early Reddit years, $10^6 - 10^9$ edges.
- To perform analysis in NetworkX, limit on one topic

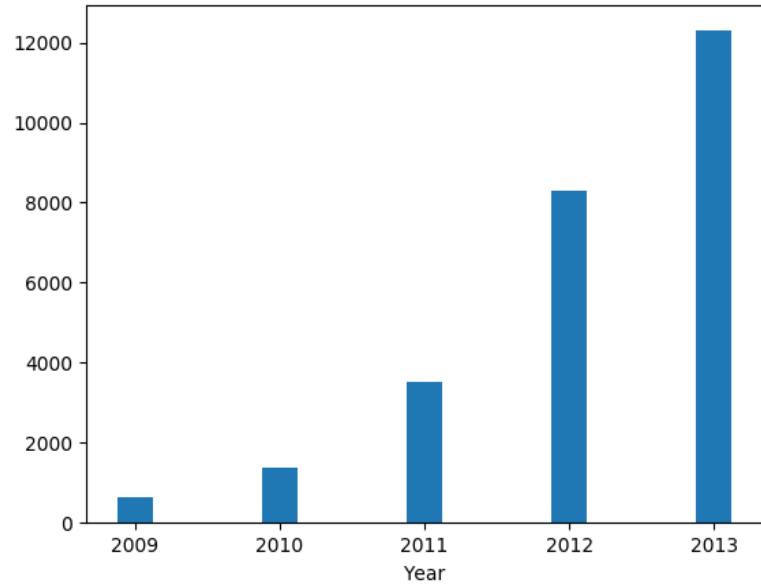




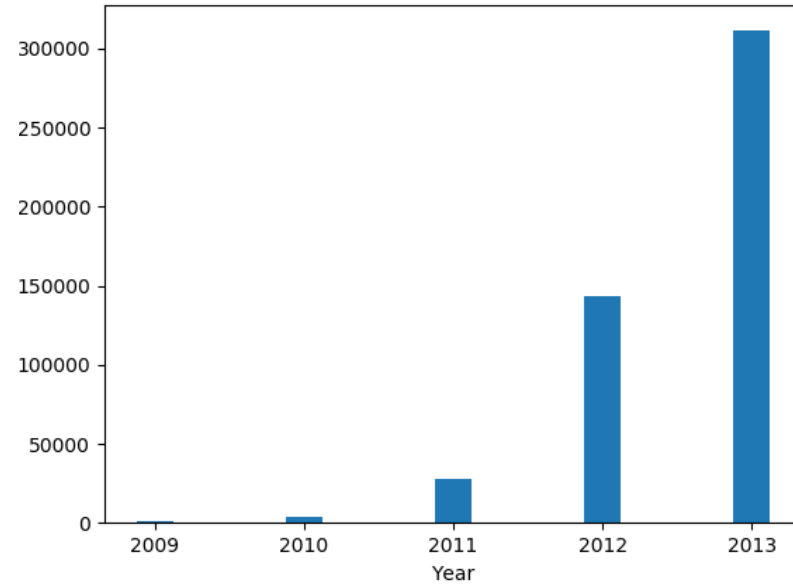
Blizzard subreddit analysis

- **Blizzard** – one of the biggest video game companies
- Manually collected ~100 of Blizzard subreddits
- Computed edges between Blizzard subreddits, and Blizzard subreddits and other subreddits
- Added connected subreddits to the list
- Analyzing only May 2009-2013

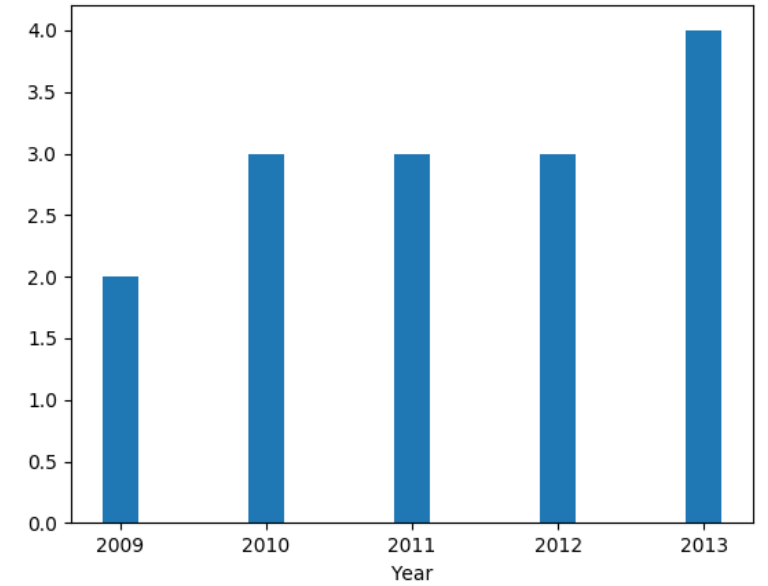
Number of subreddits



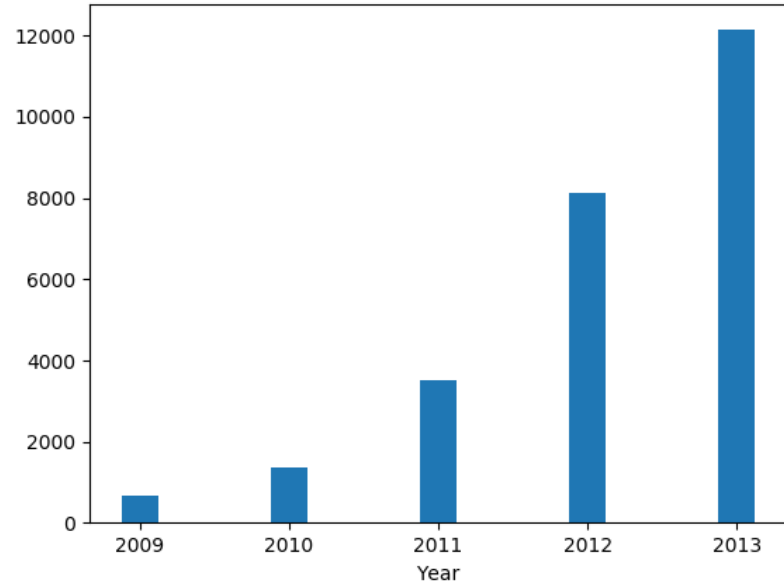
Number of edges



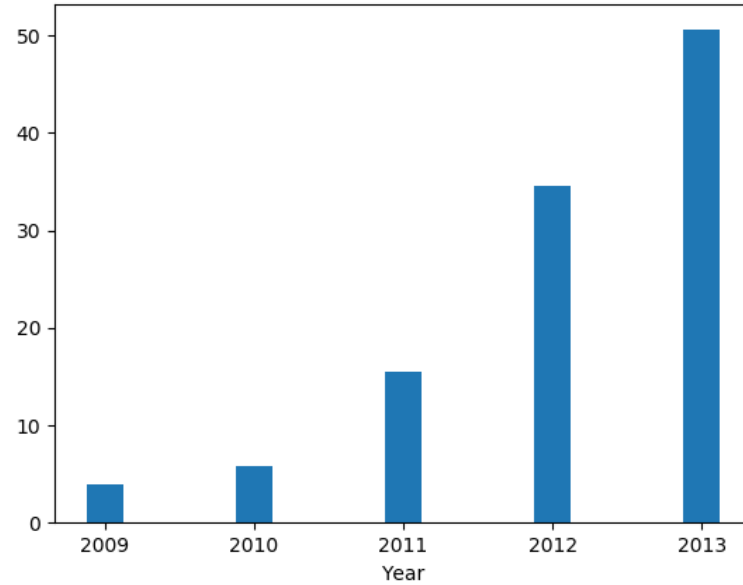
Diameter



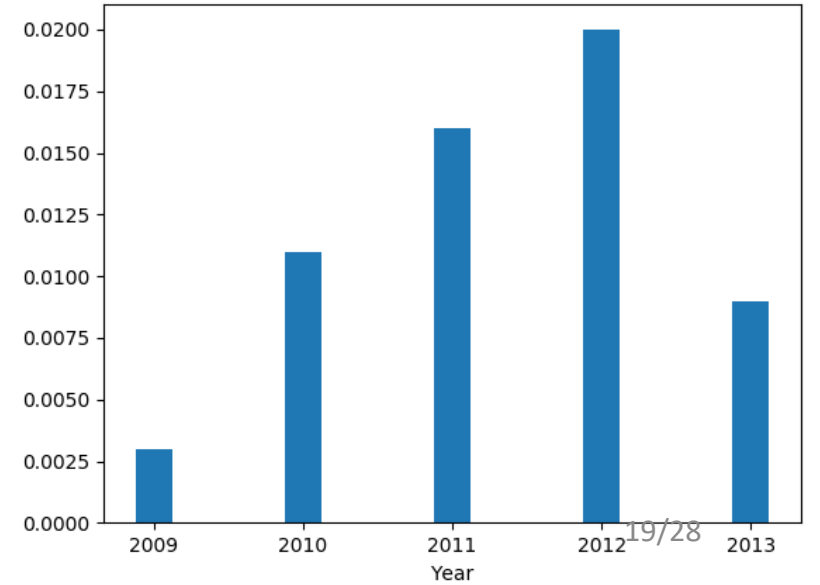
Max degree



Average degree



Average weighted degree



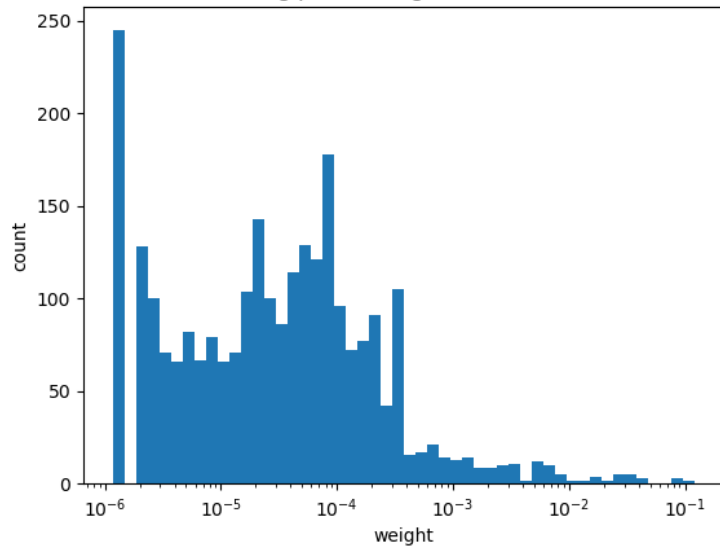
Extracting a backbone

- After projecting a bipartite graph, we are left with many really weak edges
- **Keep strong edges, remove weak ones**
- Threshold approach: $\begin{cases} \textit{weight} > T, \textit{keep the edge} \\ \textit{weight} < T, \textit{remove the edge} \end{cases}$
- More advanced approaches: Look at a distribution of that edge weight for large number of randomly generated graphs, keep the edge if the weight is statistically significant.¹

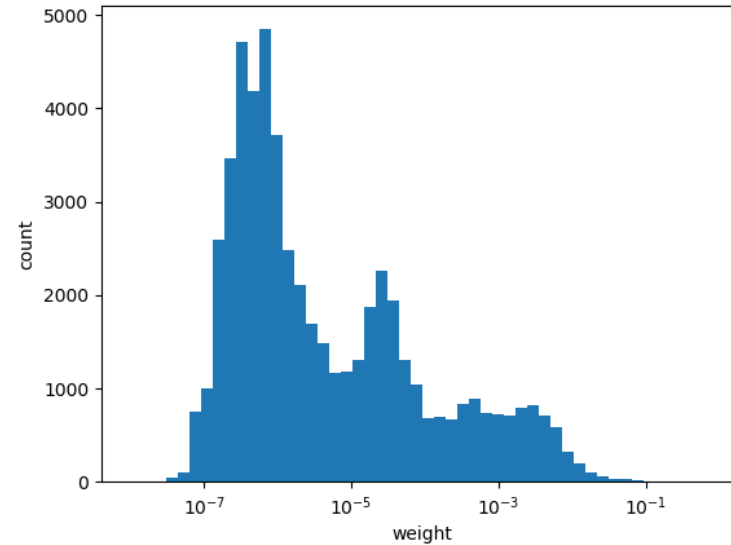
Choosing the threshold

- As we want a small number of important edges, we need to choose a strict threshold that removes >99% edges
- We see that values between 10^{-2} and 10^{-1} satisfy that. We choose 0.04.

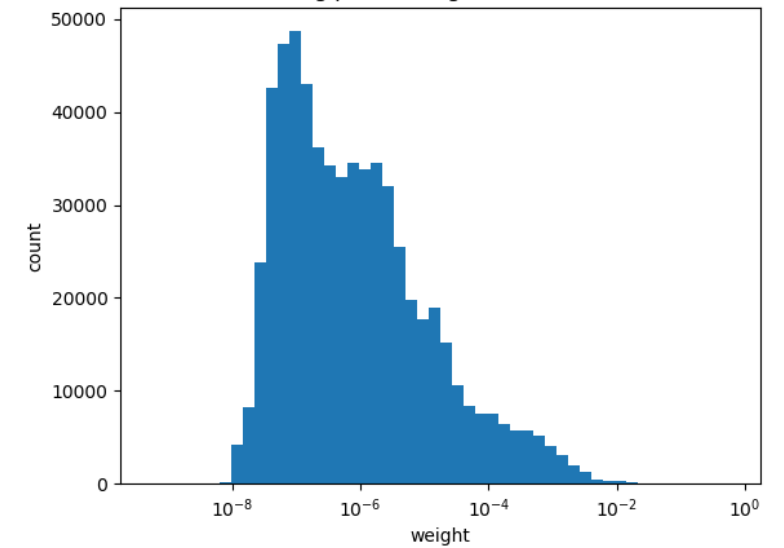
Log plot of weights for 2009



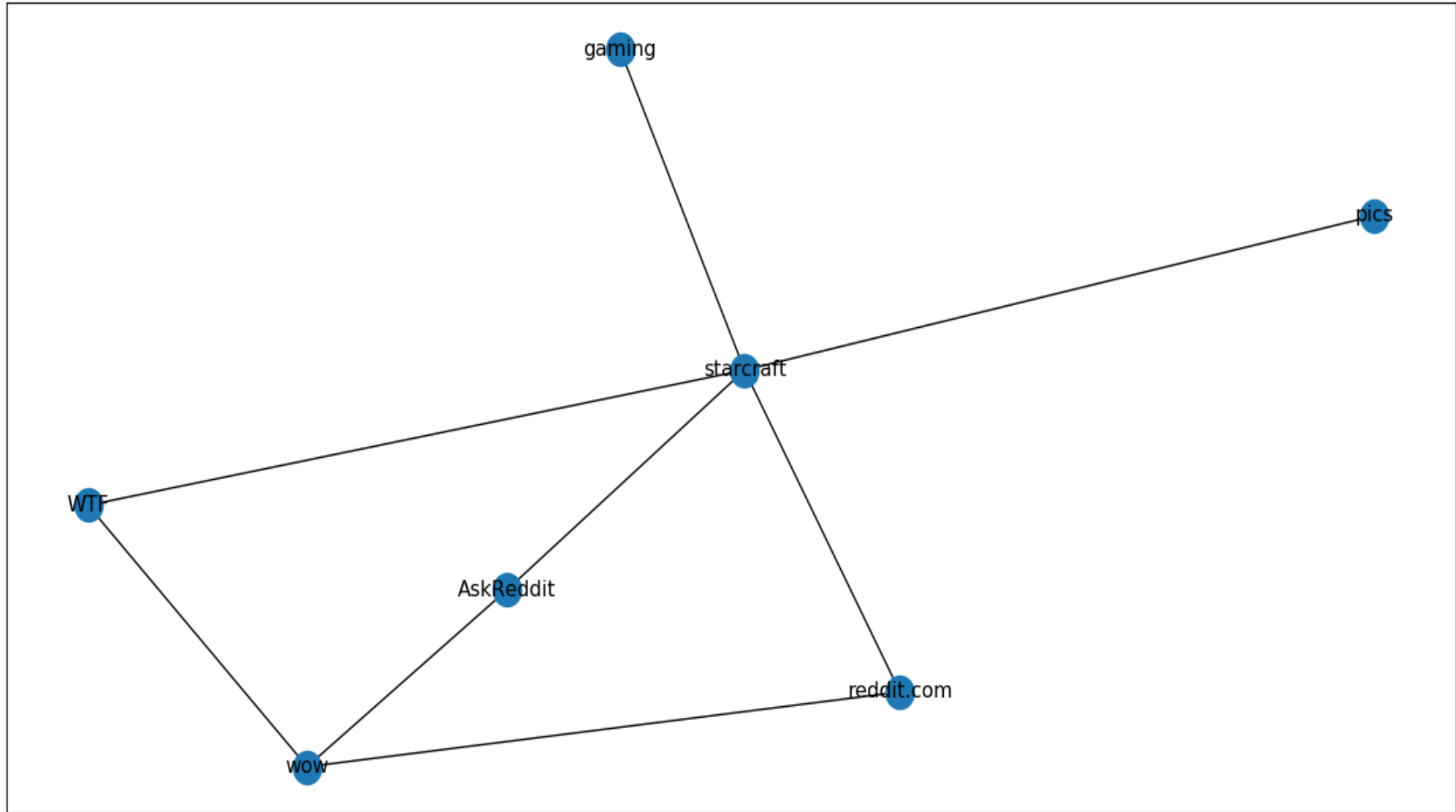
Log plot of weights for 2011



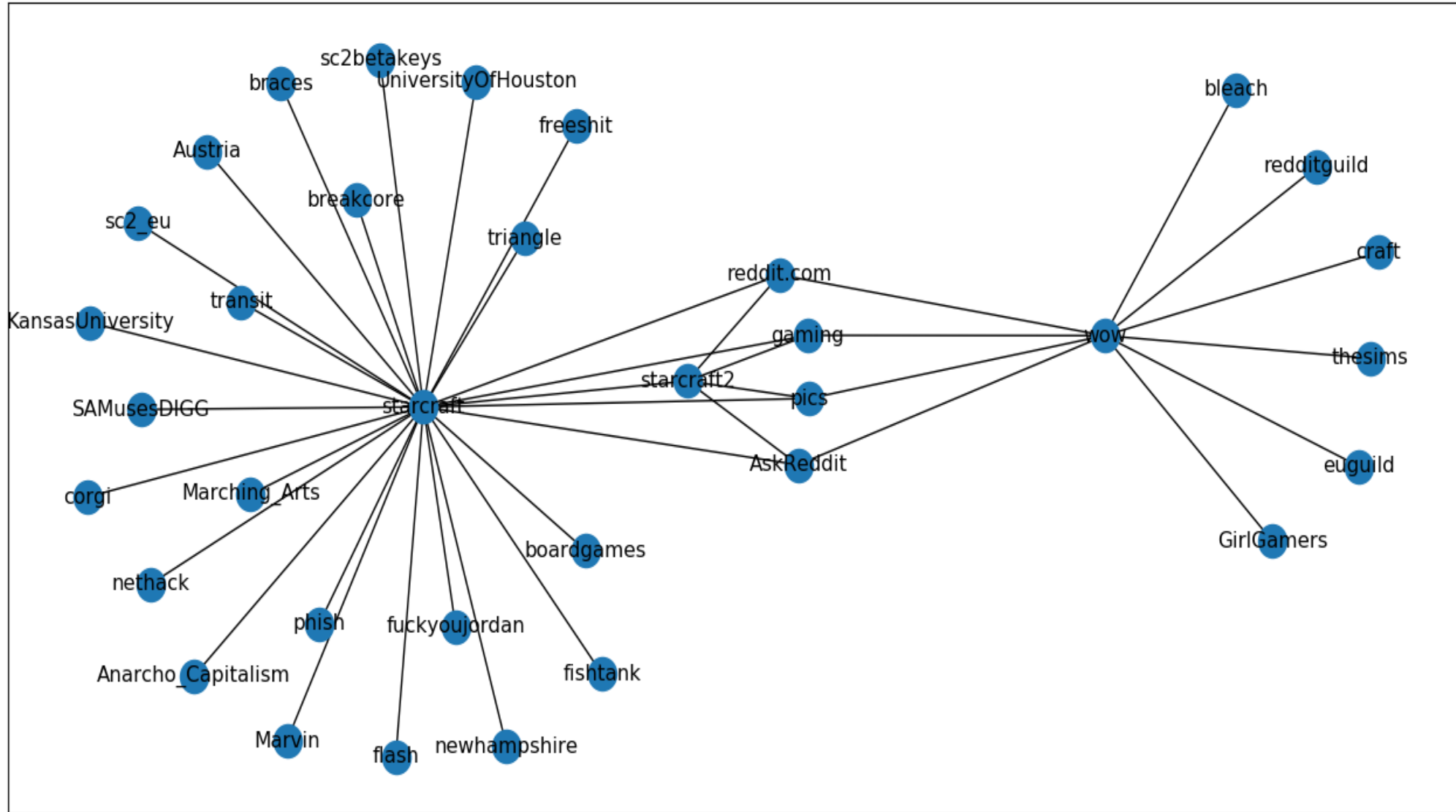
Log plot of weights for 2013



Blizzard graph May 2009



Blizzard graph May 2010



Future work

- GDELT Events dataset - <https://www.gdeltproject.org/>
- Dataset containing all significant news since 1970
- Find correlation between Reddit numbers and real world news
- Predict something happened, evaluate the significance of the news

Q&A

Any questions?

nikolaaleksic44@gmail.com

ema.p25@gmail.com