

The Polarised State of the Swedish Political Twitterverse: Lessons from Ideological Forests of Hate in the 2016 US Presidential Election

Raazesh Sainudiin

Associate Professor, Department of Mathematics, Uppsala University, Uppsala, Sweden, and
Director, Technical Strategy & Research, Combient Mix AB, Stockholm

lamastex.org

PART - I

“Lessons from Ideological Forests ...”

**Scalably Vertex-Programmable Ideological Forests
from Certain Political Twitterverses:
US 2016, UK 2017 & SE 2018 Elections**

1 Questions and Experimental Design

2 Data and Statistics

- Experimental design of twitter streams

3 Models and Methods

4 Results

Disclaimer! This is a highly empirical presentation.

- +-* / Game: *Statistical Hypothesis Testing and Estimation* while limiting oneself to scalable fault-tolerant distributed programs (sort, join and pregel on distributed graphs)

Three Questions

- (Q1) Is Trump preferentially retweeted by various types of hate groups or their leadership relative to other politicians (i.e., Clinton, Sanders, Cruz, or Ryan) against the null random network model of *apathetic retweeting*?

Three Questions

- (Q1) Is Trump preferentially retweeted by various types of hate groups or their leadership relative to other politicians (i.e., Clinton, Sanders, Cruz, or Ryan) against the null random network model of *apathetic retweeting*?
- (Q2) What frequency of unique users retweeted both a politician and a hate group or its leadership more than expected under the null model?

Three Questions

- (Q1) Is Trump preferentially retweeted by various types of hate groups or their leadership relative to other politicians (i.e., Clinton, Sanders, Cruz, or Ryan) against the null random network model of *apathetic retweeting*?
- (Q2) What frequency of unique users retweeted both a politician and a hate group or its leadership more than expected under the null model?
- (Q3) What is the joint distribution of the *degrees of separation* to each user from each of the five politicians and the eight most prolific hateful ideologies on Twitter, measured through the lengths of the most retweeted directed paths in the observed network?

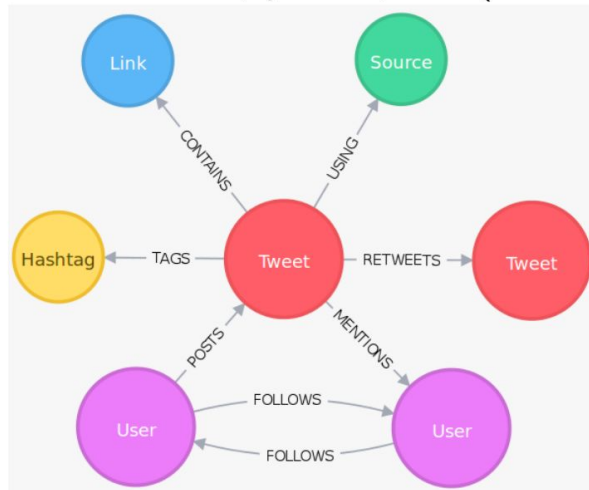
Three Questions

- (Q1) Is Trump preferentially retweeted by various types of hate groups or their leadership relative to other politicians (i.e., Clinton, Sanders, Cruz, or Ryan) against the null random network model of *apathetic retweeting*?
- (Q2) What frequency of unique users retweeted both a politician and a hate group or its leadership more than expected under the null model?
- (Q3) What is the joint distribution of the *degrees of separation* to each user from each of the five politicians and the eight most prolific hateful ideologies on Twitter, measured through the lengths of the most retweeted directed paths in the observed network?
- (Q4) Did the US Hate Networks get help from “Russian Trolls”? — [back to basics in the “Big Data” Age – Scientific Hypothesis Testing](#)

Twitterverse

twitter is a micro-blogging service...

What is a tweet? retweet? reply-tweet, etc. (*status updates*)

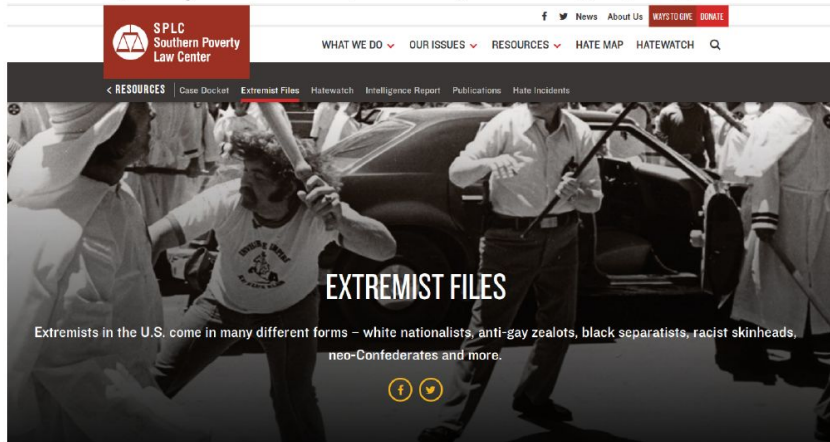


Via public streams and REST APIs we collected ~22M status updates related to 5 politicians and 52 hate groups (retrospective REST-based network augmentations).

Spark D3 Demo

Hateful Networks

US Hate Groups by SPLC <https://www.splcenter.org/fighting-hate/extremist-files>



The Extremist Files database contains profiles of various prominent extremists and extremist organizations. It also examines the histories and core beliefs – or ideologies – of the most common types of extremist movements. In addition, it illustrates connections between individuals, groups and extremist ideologies.

- <https://www.splcenter.org/fighting-hate/extremist-files/ideology>
- <https://www.splcenter.org/fighting-hate/extremist-files/individual>
- <https://www.splcenter.org/fighting-hate/extremist-files/groups>

Definition (Hate Group):

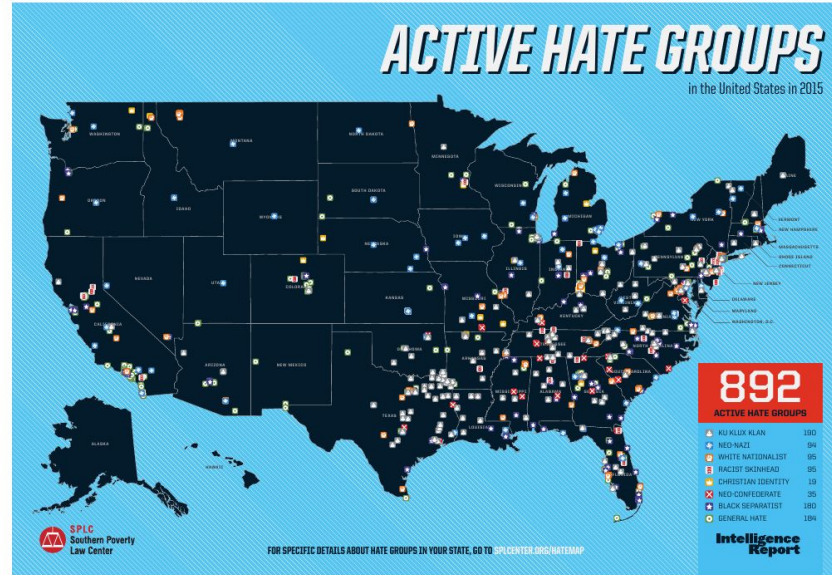
The SPLC does not necessarily consider all groups or individuals on its “Extremist Files” as violent or engaged in criminal activities, but rather identifies any group or individual “**whose beliefs or practices attack or malign an entire class of people, typically for their immutable characteristics**”.

The database does not include foreign hate groups or extremist groups such as ISIS, Al Qaeda, or Boko Haram, as its focus is on American hate groups.

Southern Poverty Law Center (2016) Hate map. SPLC. October 11, 2013 4:00 AM, Available from <https://www.splcenter.org/hate-map>. Accessed on May 28, 2017.

Hateful Networks

US Hate Groups by SPLC <https://www.splcenter.org/fighting-hate/extremist-files>






<https://www.splcenter.org/hate-map>

Hateful Networks

US Hate Groups by SPLC <https://www.splcenter.org/fighting-hate/extremist-files>

ABOUT THE HATE MAP

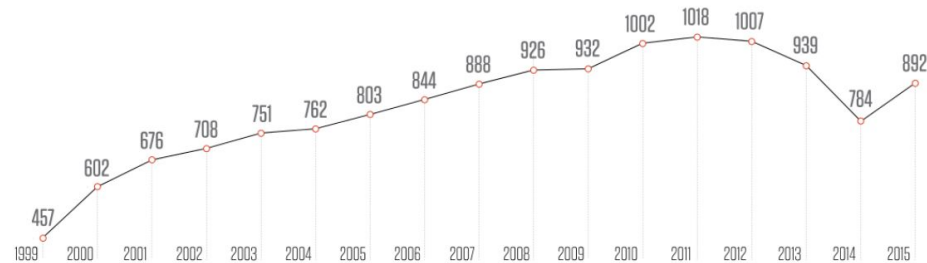
-  All hate groups have beliefs or practices that attack or malign an entire class of people, typically for their immutable characteristics.
-  This list was compiled using hate group publications and websites, citizen and law enforcement reports, field sources and news reports. Groups that appear in the center of states represent statewide groups.
-  Hate group activities can include criminal acts, marches, rallies, speeches, meetings, leafleting or publishing.

<https://www.splcenter.org/hate-map>

Hateful Networks

US Hate Groups by SPLC <https://www.splcenter.org/fighting-hate/extremist-files>

HATE GROUPS 1999-2015



<https://www.splcenter.org/hate-map>

18 US Hateful Ideologies by SPLC

<https://www.splcenter.org/fighting-hate/extremist-files>

Alternative Right

The Alternative Right, commonly known as the Alt-Right, is a set of far-right ideologies, groups and individuals whose core belief is that “white identity” is under attack by multicultural forces using “political correctness” and “social justice” to undermine white people and “their” civilization...



Anti-Immigrant

Anti-immigrant hate groups are the most extreme of the hundreds of nativist and vigilante groups that have proliferated since the late 1990s, when anti-immigration xenophobia began to rise to levels not seen in the United States since the 1920s.



Anti-LGBT

Opposition to equal rights for LGBT people has been a central theme of Christian Right organizing and fundraising for the past three decades – a period that parallels the fundamentalist movement’s rise to political power.



<https://www.splcenter.org/fighting-hate/extremist-files/ideology>

18 US Hateful Ideologies by SPLC

<https://www.splcenter.org/fighting-hate/extremist-files>

Anti-Muslim

Anti-Muslim hate groups are a relatively new phenomenon in the United States, most of them appearing in the aftermath of the World Trade Center terrorist attacks on Sept. 11, 2001. Earlier anti-Muslim groups tended to be religious in orientation and disputed Islam's status as a respectable religion.



Antigovernment Movement

The antigovernment movement has experienced a resurgence, growing quickly since 2008, when President Obama was elected to office. Factors fueling the antigovernment movement in recent years include changing demographics driven by immigration, the struggling economy and the election of the first...



Black Separatist

Black separatists typically oppose integration and racial intermarriage, and they want separate institutions -- or even a separate nation -- for blacks. Most forms of black separatism are strongly anti-white and anti-Semitic, and a number of religious versions assert that blacks are the Biblical "...



<https://www.splcenter.org/fighting-hate/extremist-files/ideology>

18 US Hateful Ideologies by SPLC

<https://www.splcenter.org/fighting-hate/extremist-files>

Christian Identity

Christian Identity is a unique anti-Semitic and racist theology that rose to a position of commanding influence on the racist right in the 1980s. "Christian" in name only, the movement's relationship with evangelicals and fundamentalists has generally been hostile due to the latter's belief that...



General Hate

These groups espouse a variety of rather unique hateful doctrines and beliefs that are not easily categorized. Many of the groups are vendors that sell a miscellany of hate materials from several different sectors of the white supremacist movement.



Holocaust Denial




Deniers of the Holocaust, the systematic murder of around 6 million Jews in World War II, either deny that such a genocide took place or minimize its extent. These groups (and individuals) often cloak themselves in the sober language of serious scholarship, call themselves "historical revisionists..."



<https://www.splcenter.org/fighting-hate/extremist-files/ideology>

18 US Hateful Ideologies by SPLC

<https://www.splcenter.org/fighting-hate/extremist-files>

<u>Ku Klux Klan</u>	
The Ku Klux Klan, with its long history of violence, is the most infamous – and oldest – of American hate groups. Although black Americans have typically been the Klan's primary target, it also has attacked Jews, immigrants, gays and lesbians and, until recently, Catholics.	
<u>Neo-Confederate</u>	
The term neo-Confederacy is used to describe twentieth and twenty-first century revivals of pro-Confederate sentiment in the United States. Strongly nativist, neo-Confederacy claims to pursue Christianity and heritage and other supposedly fundamental values that modern Americans are seen to have...	
<u>Neo-Nazi</u>	
Neo-Nazi groups share a hatred for Jews and a love for Adolf Hitler and Nazi Germany. While they also hate other minorities, gays and lesbians and even sometimes Christians, they perceive "the Jew" as their cardinal enemy.	

<https://www.splcenter.org/fighting-hate/extremist-files/ideology>

Hateful Networks

18 US Hateful Ideologies by SPLC

<https://www.splcenter.org/fighting-hate/extremist-files>

Phineas Priesthood

The Phineas Priesthood is not an actual organization; it has no leaders, meetings, or any other institutional apparatus.



Racist Music

Racist music groups are typically white power music labels that record, publish and distribute racist music in a variety of genres.



Racist Skinhead

Racist Skinheads form a particularly violent element of the white supremacist movement, and have often been referred to as the "shock troops" of the hoped-for revolution. The classic Skinhead look is a shaved head, black Doc Martens boots, jeans with suspenders and an array of typically racist...



<https://www.splcenter.org/fighting-hate/extremist-files/ideology>

18 US Hateful Ideologies by SPLC

<https://www.splcenter.org/fighting-hate/extremist-files>

Radical Traditional Catholicism

"Radical traditionalist" Catholics, who may make up the largest single group of serious anti-Semites in America, subscribe to an ideology that is rejected by the Vatican and some 70 million mainstream American Catholics.



Sovereign Citizens Movement

The strange subculture of the sovereign citizens movement, whose adherents hold truly bizarre, complex antigovernment beliefs, has been growing at a fast pace since the late 2000s. Sovereigns believe that they get to decide which laws to obey and which to ignore, and they don't think they should...



White Nationalist

White nationalist groups espouse white supremacist or white separatist ideologies, often focusing on the alleged inferiority of nonwhites. Groups listed in a variety of other categories - Ku Klux Klan, neo-Confederate, neo-Nazi, racist skinhead, and Christian Identity - could also be fairly...



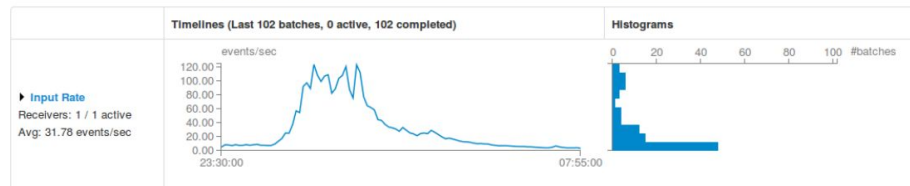
<https://www.splcenter.org/fighting-hate/extremist-files/ideology>

US Presidential Election 2016 - Twitter Streams

Twitter Data — 3rd US Presidential Debate

Streaming Statistics

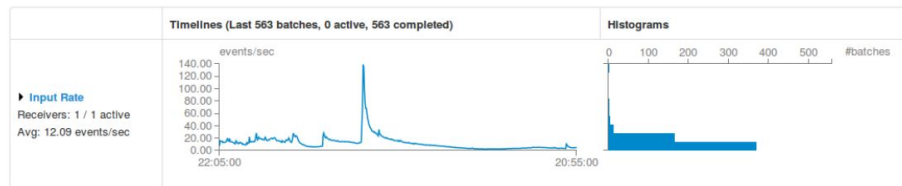
Running batches of 5 minutes for 8 hours 32 minutes 20 seconds since 2016/10/19 23:26:43 (102 completed batches, 972342 records)



Twitter Data — Last 2 Days Around the End of Election

Streaming Statistics

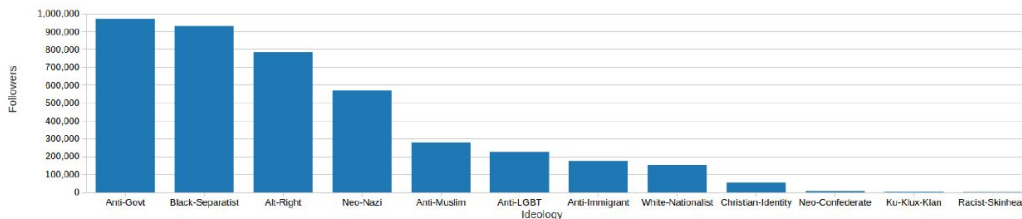
Running batches of 5 minutes for 1 day 22 hours 56 minutes since 2016/11/08 22:02:36 (563 completed batches, 2041501 records)



- public streams of @realDonaldTrump, @HillaryClinton, @BernieSanders, @tedcruz, SpeakerRyan and 52 splc-defined hategroups of their leadership
- collected data includes all mentions, replies, retweets, etc of these twitter accounts of interest for about 9 weeks around the 2016 US Presidential Election

12 SPLC-defined hateful ideologies

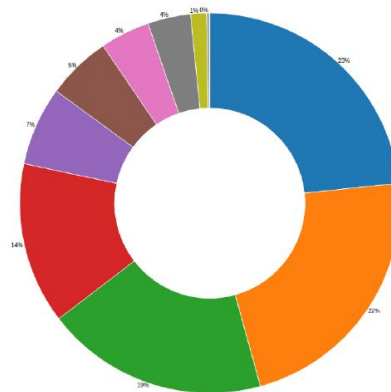
– only 78% of hategroups identified by SPLC were active in Twitter



Ideology	Followers
Anti-Govt	970769
Black-Separatist	931736
Alt-Right	786327
Neo-Nazi	571772
Anti-Muslim	279122
Anti-LGBT	227636
Anti-Immigrant	175441
White-Nationalist	151711
Christian-Identity	56191
Neo-Confederate	6628
Ku-Klux-Klan	3070
Racist-Skinhead	1826

Ideology

- Anti-Govt
- Black-Separatist
- Alt-Right
- Neo-Nazi
- Anti-Muslim
- Anti-LGBT
- Anti-Immigrant
- White-Nationalist
- Christian-Identity
- Others



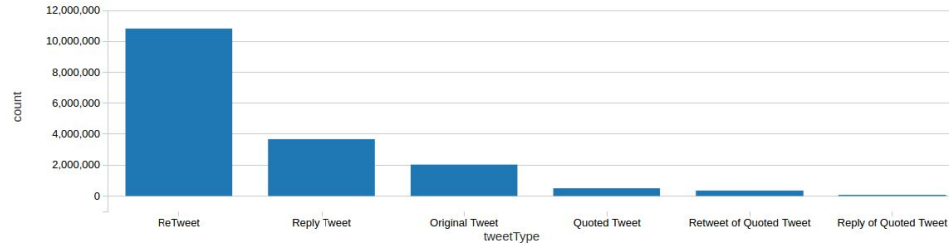
5 prominent Politicians in the USA

Retweet Network statistics of the five political accounts

Politician	in-degree	in-nbhd	out-degree	out-nbhd
Donald Trump	40	12	5,952,257	958,262
Hillary Clinton	225	121	2,774,111	943,995
Bernie Sanders	107	62	762,209	356,718
Paul Ryan	769	158	68,973	28,902
Ted Cruz	322	189	49,479	27,663

Dataset overview

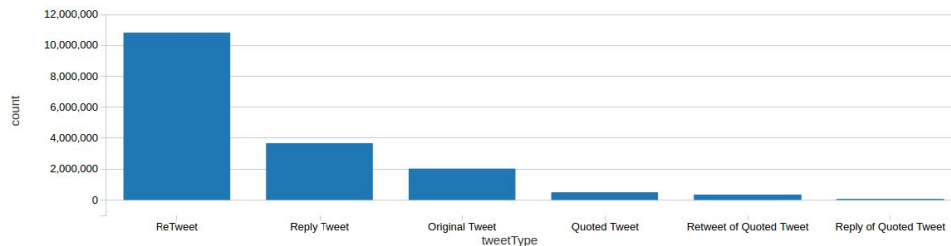
Data collected around the 2016 US Presidential Election



- data = 2.7M tweets, 13.7M retweets, 22M status updates

Dataset overview

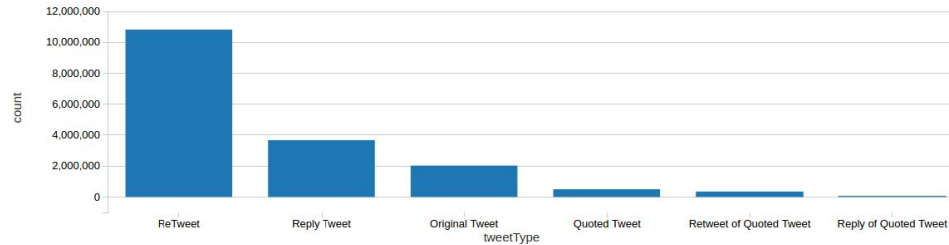
Data collected around the 2016 US Presidential Election



- data = 2.7M tweets, 13.7M retweets, 22M status updates
- 4.4M distinct retweet-pairs: (original-Tweeter, Retweeter)

Dataset overview

Data collected around the 2016 US Presidential Election



- data = 2.7M tweets, 13.7M retweets, 22M status updates
- 4.4M distinct retweet-pairs: (original-Tweeter, Retweeter)
- 2.5M unique users

Dataset overview

ELT Designs by Tweet Anatomy & Transmission Tree

2016, Akinwande Atanda and Raazesh Sainudiin

This notebook describes the structure and key components of a tweet created by Twitter users. The components are used to generate unique categorizations of tweet types and construct the *Tweet Transmission Tree (TTT)*. The purpose of TTT (constructed in a Spark Streaming job) is to encode various types of interactions between twitter users in continuous time by using appropriate attributes based on standard objects returned from [Twitter API for developers](#).

TTT can be used to:

- define interactions among Twitter users as a tweet status is transmitted in continuous time up to millisecond resolution
- exploit specific types of interactions among users to build networks and detect ideologically aligned communities
- filter appropriate sets of tweets in the context of specific interaction for downstream Natural Language Processing (NLP), including sentiment analysis
- etc.

This is part of [Project MEP: Meme Evolution Programme](#) and supported by databricks academic partners program.

The analysis is available in the following databricks notebook:

- <http://lamastex.org/lmse/mep/src/TweetAnatomyAndTransmissionTree.html>

For details on the mathematical model motivating the anatomy and categorizations of tweet transmission trees in the notebook see:

- [The Transmission Process: A Combinatorial Stochastic Process for the Evolution of Transmission Trees over Networks](#), Raazesh Sainudiin and David Welch, *Journal of Theoretical Biology*, Volume 410, Pages 137–170, 10.1016/j.jtbi.2016.07.038, 2016
 - [preprint of the above paper as PDF 900KB](#).

Other resources that employ transmission trees and networks are summarized here:

- <http://lamastex.org/lmse/mep/>

Copyright 2016 Raazesh Sainudiin and Akinwande Atanda

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at <http://www.apache.org/licenses/LICENSE-2.0>

- SparkSQL: Twitter Experimental Designs via parquet single-column JSON string of each status update as one row
 - future-proofing evolving schema
 - input to generic SparkML pipelines
- GraphX: Pregel-programmed Network Design
- SparkML/lib: various standard algorithms
- Spark Core: distributed sort and join

See [Project MEP: Meme Evolution Programme at <http://lamastex.org/lmse/mep/>](#) and the databricks notebook [<http://lamastex.org/lmse/mep/src/TweetAnatomyAndTransmissionTree.html>](http://lamastex.org/lmse/mep/src/TweetAnatomyAndTransmissionTree.html)

Retweet Network — (3% sample # V = 1205, # E = 29856)

Trump-Clinton Retweet Network — a few samples

CPostUserSN	OPostUserSNinRT	OPostUserSNinOT	favouritesCount	followersCount	friendsCount	IsVerified	IsGeoEnabled	CurrentTweet
georgefayner	realDonaldTrump	null	137811	1466	953	false	true	RT @realDonaldTrump: China is cooking up conspiracy theories that the Olympics are rigged. http://t.co/0ah0hBJt They don't understand why...
KevinCormier10	realDonaldTrump	null	16164	505	367	false	true	RT @realDonaldTrump: EXCLUSIVE: FBI Agents Say Comey 'Stood In The Way' Of Clinton Email Investigation: https://t.co/6n63HvVNo
thuerta	realDonaldTrump	null	13081	128	345	false	true	RT @realDonaldTrump: 'Trump rally disrupter was once on Clinton's payroll' https://t.co/75oLLu4S1
tanladyvolfan	HillaryClinton	null	6316	101	200	false	true	RT @HillaryClinton: Our progress is on the ballot. Tolerance is on the ballot. Democracy is on the ballot. Make a plan to vote:....

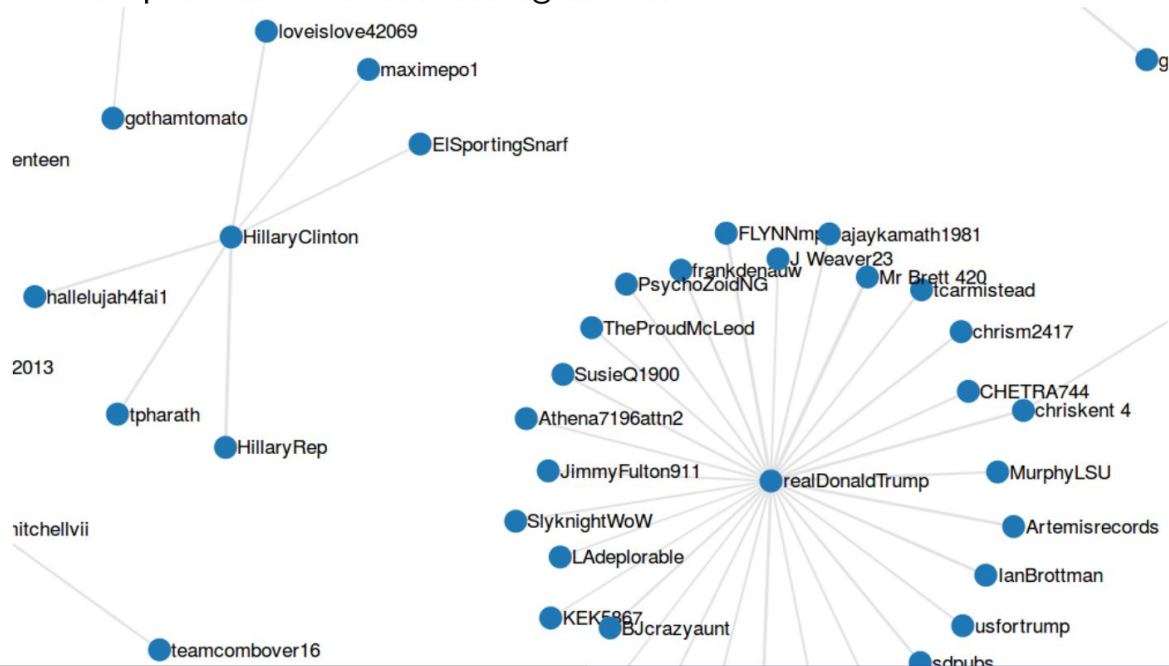
Retweet Network — (3% sample $\#V = 1205$, $\#E = 29856$)

Trump-Clinton Retweet Network weighted by Retweet counts

userCreatedAtDate	daysSinceUserCreated	OPostUserSNinRT	CPostUserSN	max(favouritesCount)	max(followersCount)	max(friendsCount)	RetweetCount
2011-12-13T19:10:28.000+0000	1781	realDonaldTrump	Mr_Brett_420	3294	78	194	100
2015-04-30T09:13:34.000+0000	181	HillaryClinton	HillaryRlap	4196	2168	4954	158
2011-03-22T13:09:23.000+0000	2047	realDonaldTrump	FLYNHrpc	1653	48	75	146
2014-08-25T17:02:48.000+0000	795	realDonaldTrump	mikenyc499	17427	183	155	132
2009-04-26T07:07:03.000+0000	2742	yottapoint	gcomrking	5076	797	1826	120
2014-06-20T21:37:33.000+0000	851	BUILDseriesNYC	suzannebuzz	30604	1705	485	112
2009-06-28T18:51:31.000+0000	2710	realDonaldTrump	ohitskent_4	838	284	85	112
2009-05-08T12:59:18.000+0000	2791	realDonaldTrump	Artemisrecords	2000	2777	5000	112
2012-06-25T16:09:37.000+0000	1434	realDonaldTrump	lanBjottfson	1	89	151	107
2011-03-31T09:54:09.000+0000	2038	realDonaldTrump	frankdenauw	43	55	18	102
2015-07-17T21:30:47.000+0000	103	HillaryClinton	lovelslove42069	3510	168	593	90
2015-09-01T18:52:08.000+0000	423	realDonaldTrump	bjorazyaurit	1064	1296	1432	95
2011-12-24T03:52:02.000+0000	1770	HillaryClinton	iprerath	703	36	183	91
2015-03-08T23:47:05.000+0000	600	HillaryClinton	haleelujah4ta1	16765	227	270	88
2014-06-30T18:44:10.000+0000	851	realDonaldTrump	ajaykamat1981	8309	2667	3910	85
2012-04-29T21:49:30.000+0000	1643	realDonaldTrump	MurphyLSU	85	26	47	84
2010-06-05T16:02:11.000+0000	2276	realDonaldTrump	sdpubs	23674	123	34	83
2011-07-34T19:55:57.000+0000	1923	realDonaldTrump	chris-m3417	3012	182	1112	81
2015-02-03T23:58:01.000+0000	258	realDonaldTrump	SusieQ1900	6797	366	415	81

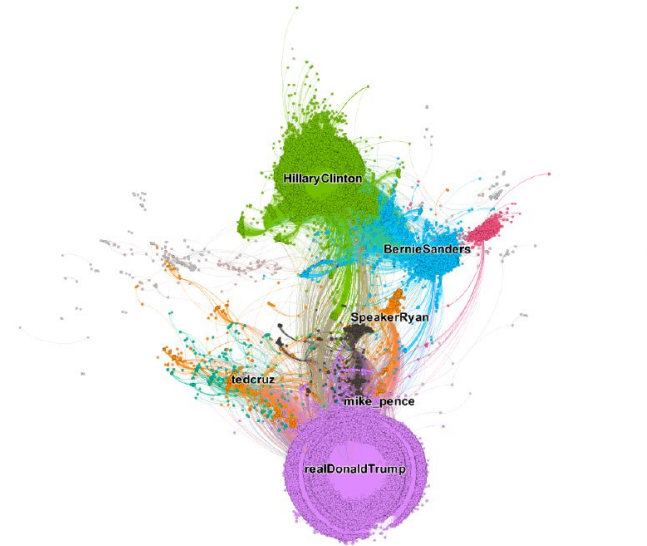
Retweet Network — (3% sample #V = 1205, #E = 29856)

Trump-Clinton Retweet Ideological Network



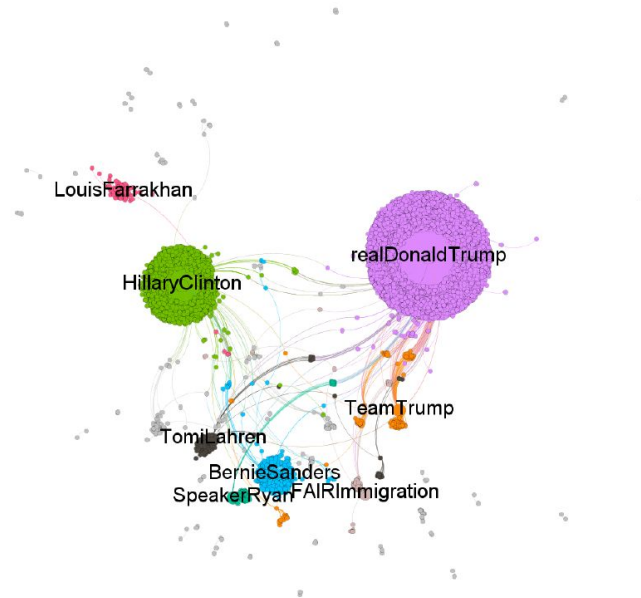
Strong Community Structure – samples of retweet networks

The 3rd US Presidential Debate 22K Retweet Network



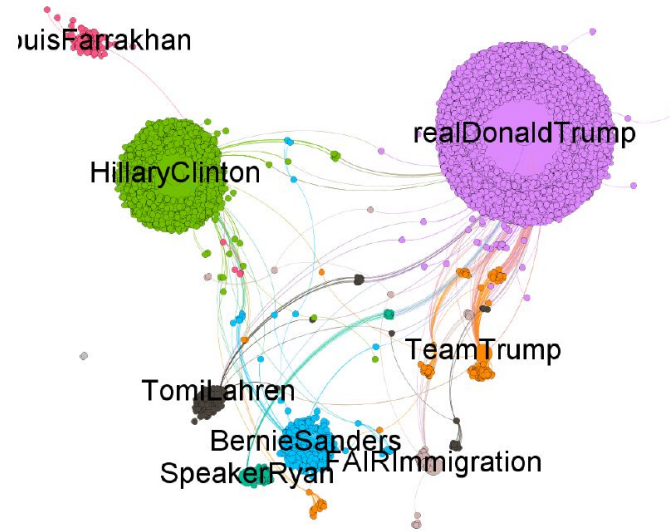
Strong Community Structure – samples of retweet networks

5% random sampled retweet networks for October 19 2016



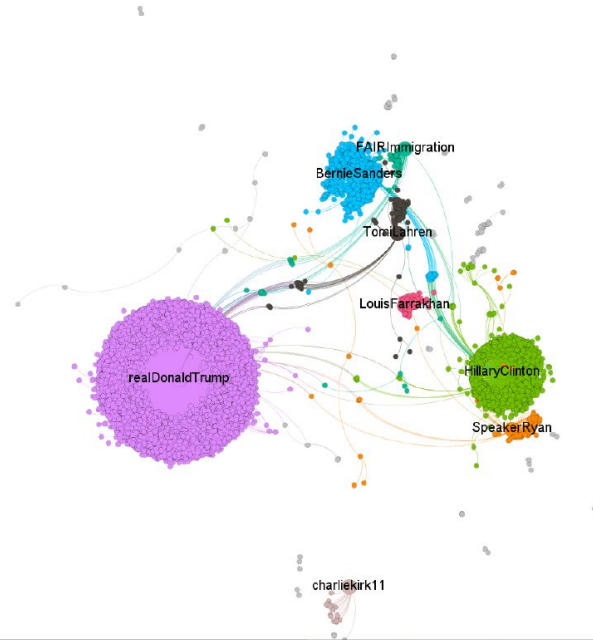
Strong Community Structure – samples of retweet networks

5% random sampled retweet networks for October 19 2016 – top 10



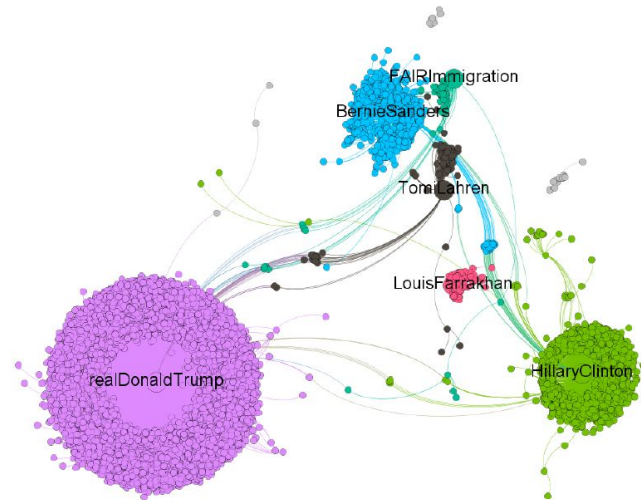
Strong Community Structure – samples of retweet networks

5% random sampled retweet networks for October 24 2016



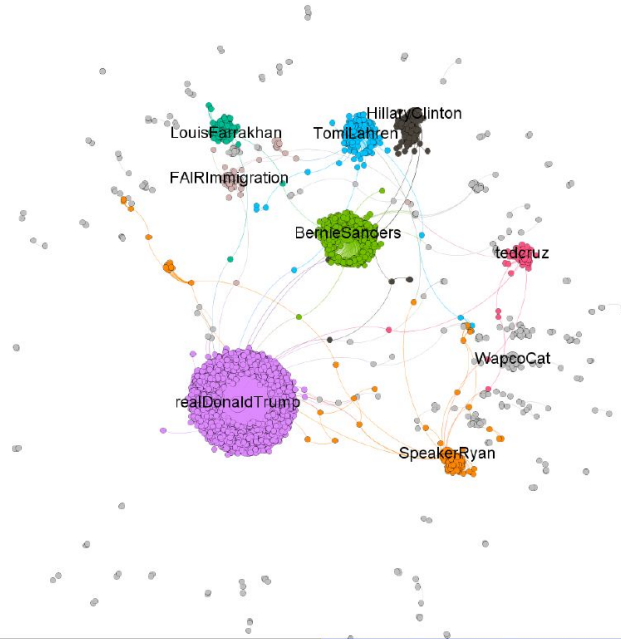
Strong Community Structure – samples of retweet networks

5% random sampled retweet networks for October 24 2016 – top 6



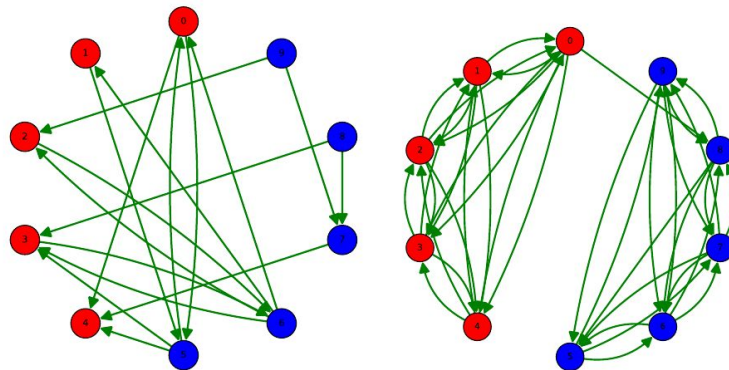
Strong Community Structure – samples of retweet networks

5% random sampled retweet networks for November 15 2016



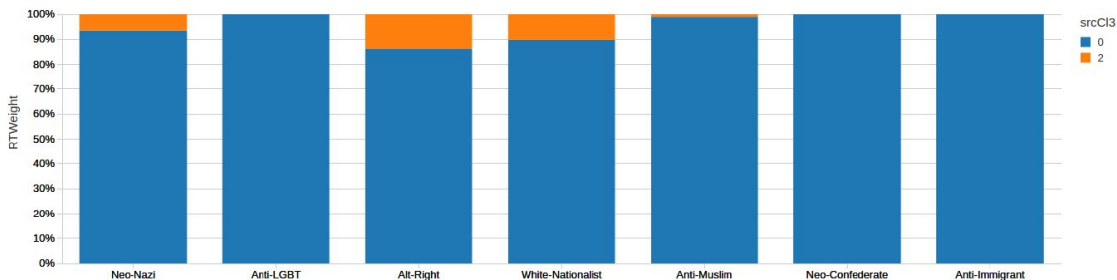
Models for Ideological Network Dynamics

- If arc $a_{i,j} = 1$ then we say i ideologically concurs with j



- Just two retweet networks out of 4, 722, 366, 482, 869, 645, 213, 696 for 9 individuals!
- We want indegree and outdegree conditioned random networks to preserve observed heterogeneity
- This is the classical *random directed configuration model* – H_0 : *apathetic retweet network*
- NEED: distributed computing using Apache Spark (fastest growing Apache project)

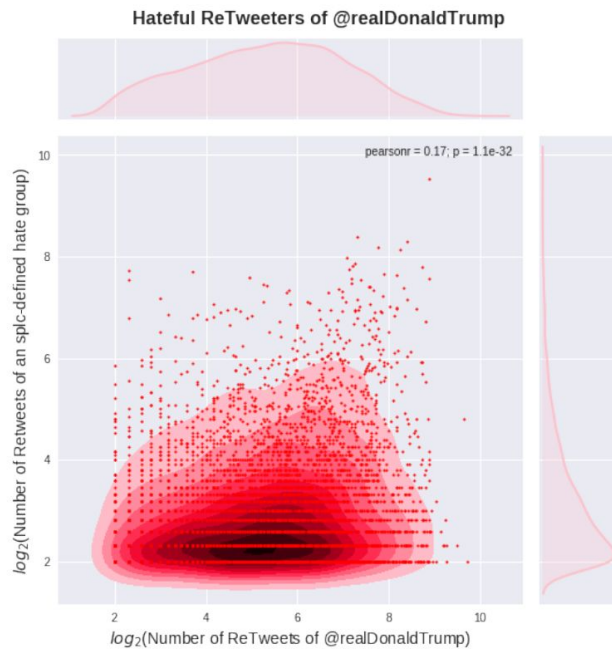
7 SPLC-defined hateful ideologies Retweet Proportions



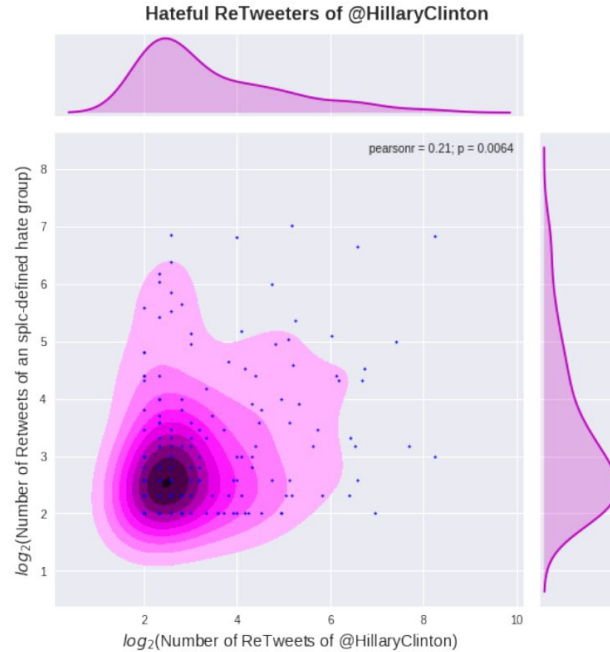
0 = Trump's cluster and 2 = Clinton's cluster

A significant proportion of retweets by leaders of seven extremist ideologies have original tweets in Trump's ideological cluster.

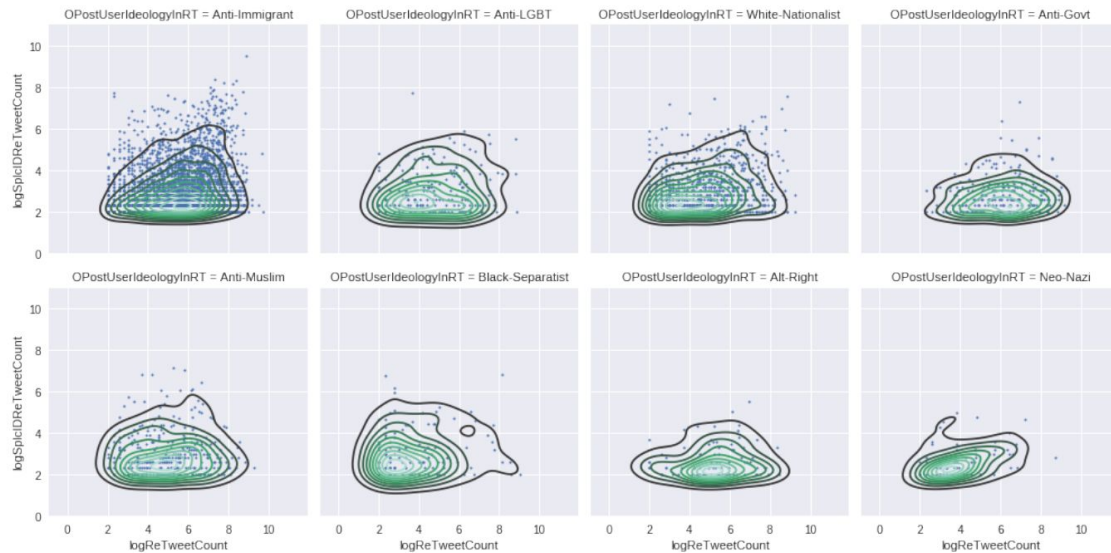
Trump's Hateful Retweeters



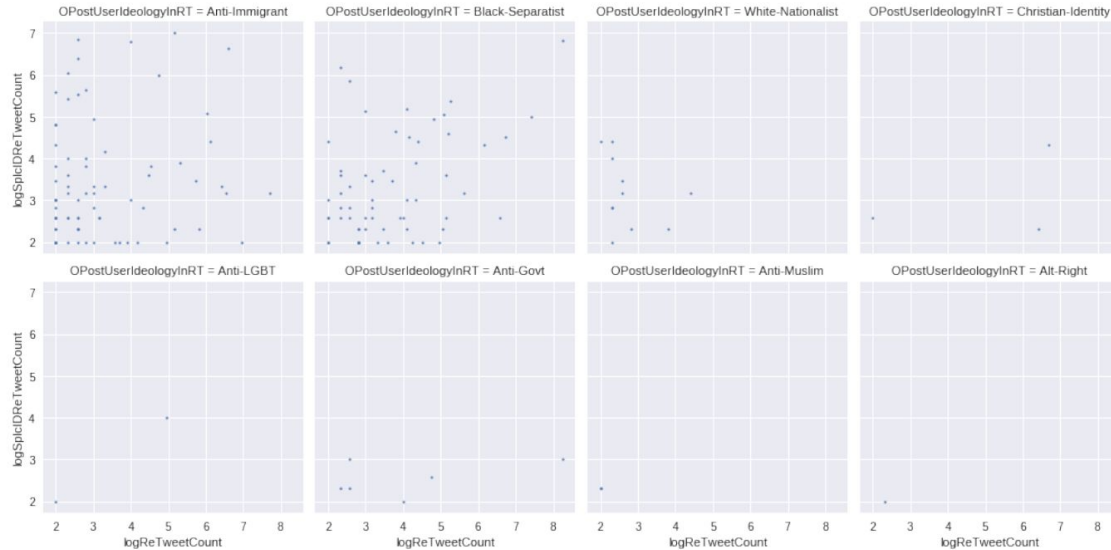
Clinton's Hateful Retweeters



Trump's Hateful Retweeters By Ideology



Clinton's Hateful Retweeters By Ideology



Chi-square tests – do NOT account for network heterogeneity

Ideology	Donald J Trump	Hillary R Clinton	Chi-Square Statistic	R ²
Alt-Right	1127 (90.7%)	116 (9.3%)	$\chi^2 = 822.30, p < .0001$	R ² =0.662
Anti-Government	1455 (89.5%)	171 (10.5%)	$\chi^2 = 1013.93, p < .0001$	R ² =0.623
Anti-Immigrant	15019 (88.6%)	1926 (11.4%)	$\chi^2 = 10116.65, p < .0001$	R ² =0.597
Anti-LGBT	1621 (88.6%)	209 (11.4%)	$\chi^2 = 1089.48, p < .0001$	R ² =0.595
Anti-Muslim	2293 (90.8%)	233 (9.2%)	$\chi^2 = 1679.97, p < .0001$	R ² =0.665
Black-Separatist	1279 (54.9%)	1049 (45.1%)	$\chi^2 = 22.72, p < .01$	R ² =0.009
Neo-Nazi	1039 (90.7%)	106 (9.3%)	$\chi^2 = 760.25, p < .0001$	R ² =0.664
White-Nationalist	5103 (89.2%)	616 (10.8%)	$\chi^2 = 3520.40, p < .0001$	R ² =0.616
Total	28992 (86.5%)	4509 (13.5%)	$\chi^2 = 18006.72, p < .0001$	R ² =0.540

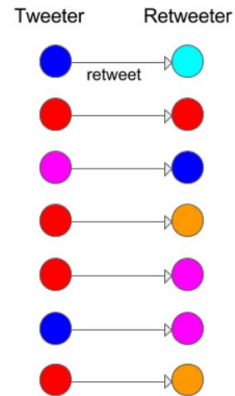
Chi-square tests – do NOT account for network heterogeneity

Restricting to retweeters who retweet at least 4 times

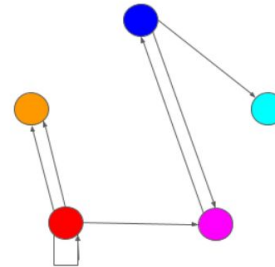
Ideology	Donald J Trump	Hillary R Clinton	Chi-Square Statistic	R ²
Alt-Right	936 (98.7%)	12 (1.3%)	$\chi^2 = 900.61, p < .0001$	R ² =0.950
Anti-Government	1388 (98.4%)	23 (1.6%)	$\chi^2 = 1320.50, p < .0001$	R ² =0.936
Anti-Immigrant	12618 (96.6%)	442 (3.4%)	$\chi^2 = 11351.84, p < .0001$	R ² =0.869
Anti-LGBT	1110 (96.0%)	46 (4.0%)	$\chi^2 = 979.32, p < .0001$	R ² =0.847
Anti-Muslim	1866 (98.8%)	22 (1.2%)	$\chi^2 = 1801.03, p < .0001$	R ² =0.954
Black-Separatist	494 (62.5%)	296 (37.5%)	$\chi^2 = 49.63, p < .001$	R ² =0.062
Neo-Nazi	692 (99.4%)	4 (0.6%)	$\chi^2 = 680.09, p < .0001$	R ² =0.977
White-Nationalist	3751 (98.0%)	76 (2.0%)	$\chi^2 = 3529.04, p < .0001$	R ² =0.922
Total	22855 (96.1%)	921 (3.9%)	$\chi^2 = 20234.71, p < .0001$	R ² =0.851

Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

Directed Retweet Edges as Two Columns

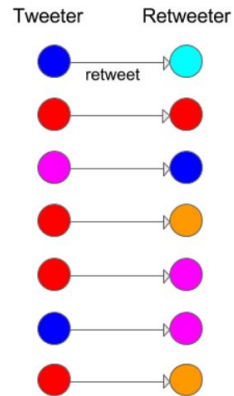


Multi-edged Self-looped Retweet Network



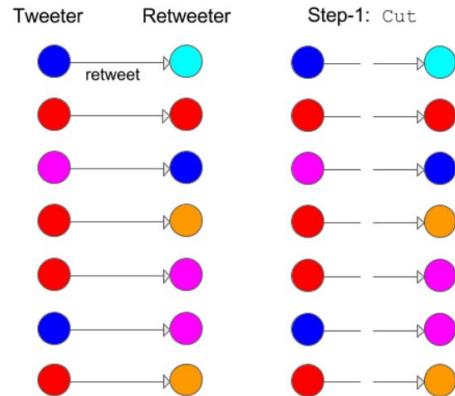
Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

Sample from Directed Multi-edged Self-looped Configuration Model



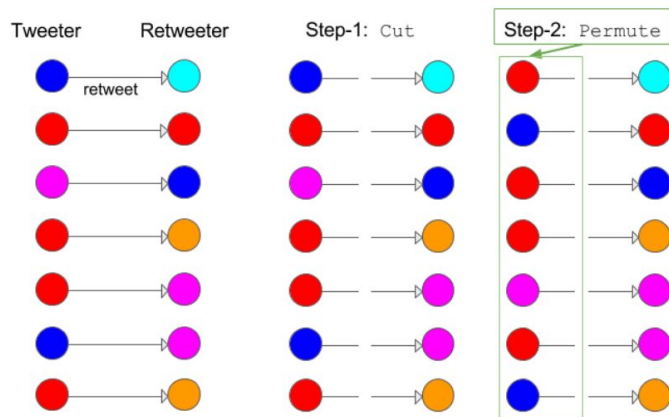
Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

Sample from Directed Multi-edged Self-looped Configuration Model



Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

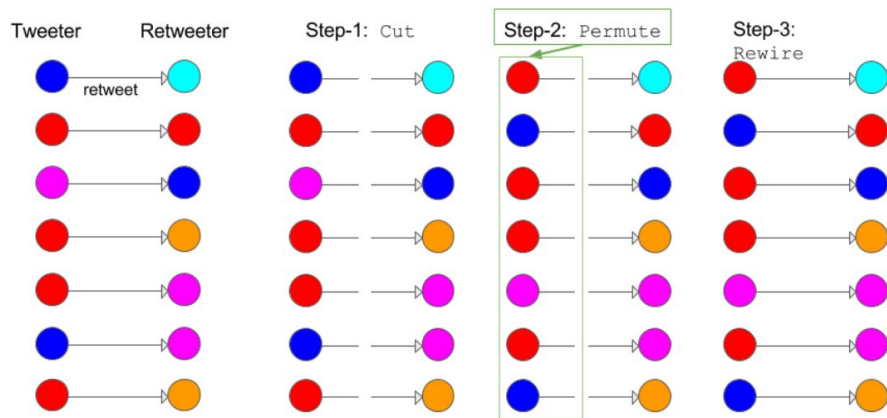
Sample from Directed Multi-edged Self-looped Configuration Model



This random permutation of row #'s of observed outbound half-edges is: $(1, 2, 3, 4, 5, 6, 7) \mapsto (7, 6, 5, 4, 3, 2, 1)$

Cut-Permute-Rewire – distributed, scalable, and fault-tolerant sampler - in pictures

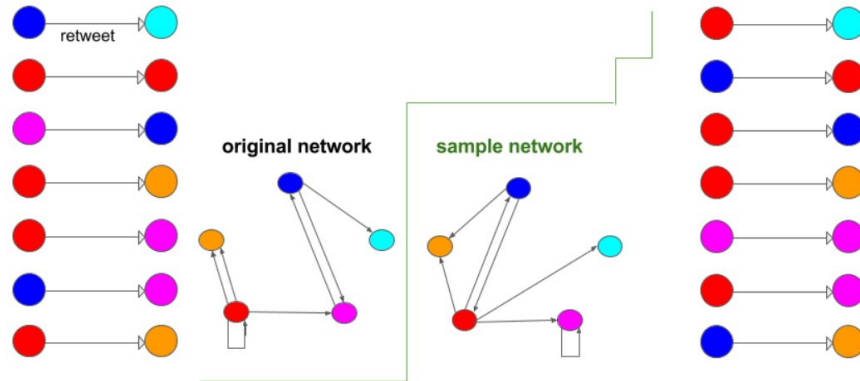
Sample from Directed Multi-edged Self-looped Configuration Model



Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

Sample from Directed Multi-edged Self-looped Configuration Model

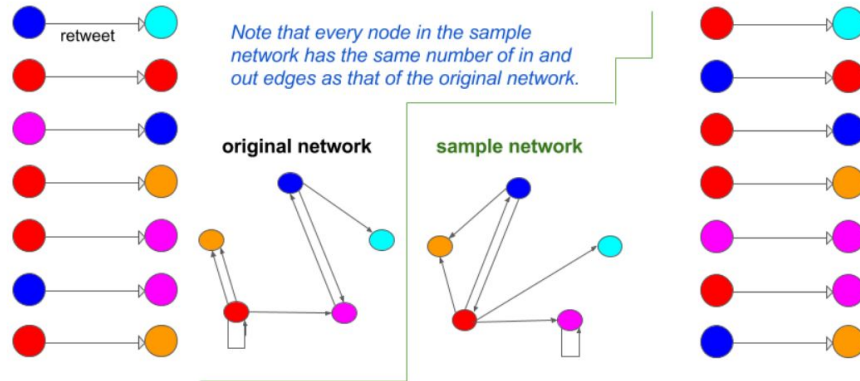
Thus, we can sample from the scalable fault-tolerant Cut-Permute-Rewire algorithm



Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

Sample from Directed Multi-edged Self-looped Configuration Model

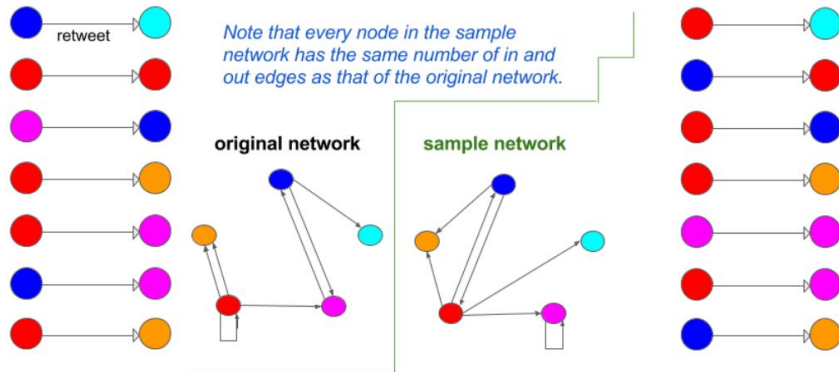
Thus, we can sample from the scalable fault-tolerant Cut-Permute-Rewire algorithm



Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler - in pictures

Sample from Directed Multi-edged Self-looped Configuration Model

Thus, we can sample from the scalable fault-tolerant Cut-Permute-Rewire algorithm

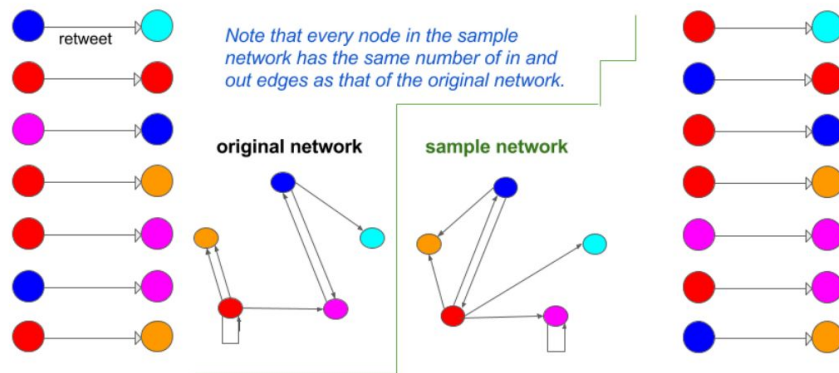


Question: What is the probability of the sample network?

Cut-Permute-Rewire – distributed, scalable, and fault-tolerant sampler - in pictures

Sample from Directed Multi-edged Self-looped Configuration Model

Thus, we can sample from the scalable fault-tolerant Cut-Permute-Rewire algorithm



Question: What is the probability of the sample network?

Answer: $1/\#\text{edges!} = 1/\#\text{retweets!} = 1/7!$

$$= 1/(7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1) = 1/(42 \times 60 \times 2) = 1/(252 \times 2) = 1/5040$$

Cut-Permute-Rewire – distributed, scalable, and fault-tolerant sampler – in English

CUTPERMUTEANDREWIRE generates sample networks from the *directed multi-edged self-looped random configuration model* (Newman, Strogatz and Watts, 2001):

- *cutting* the directed edges representing the retweets in our observed retweet network into out-bound and in-bound half edges,
- *permuting* the in-bound half edges by sorting them according to pseudo-random numbers that are generated and associated with them and
- *rewiring* the original out-bound half edges with the permuted in-bound half edges using a distributed join.

The in-degree and out-degree of each node in the observed retweet network is preserved after these three steps.

Interpret the independent and identical samples as those from the null model H_0 as the *apathetic retweet model*

Cut-Permute-Rewire — distributed, scalable, and fault-tolerant sampler – in English

`CUTPERMUTEANDREWIRE` generates sample networks from the *directed multi-edged self-looped random configuration model* (Newman, Strogatz and Watts, 2001):

- *cutting* the directed edges representing the retweets in our observed retweet network into out-bound and in-bound half edges,
- *permuting* the in-bound half edges by sorting them according to pseudo-random numbers that are generated and associated with them and
- *rewiring* the original out-bound half edges with the permuted in-bound half edges using a distributed join.

The in-degree and out-degree of each node in the observed retweet network is preserved after these three steps.

Interpret the independent and identical samples as those from the null model H_0 as the *apathetic retweet model* – “this is not reality folks!”

This is reality folks

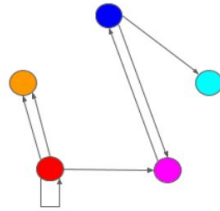


An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

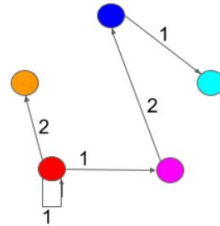
Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter

From Directed Configuration Model to Geometric Retweet Network

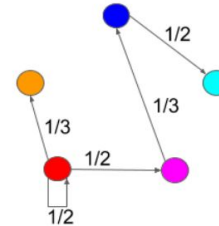
Multi-edged Self-looped
Retweet Network



Weighted Retweet Network



Geometric Retweet Network
with weights $1 / (1 + \# \text{ retweets})$

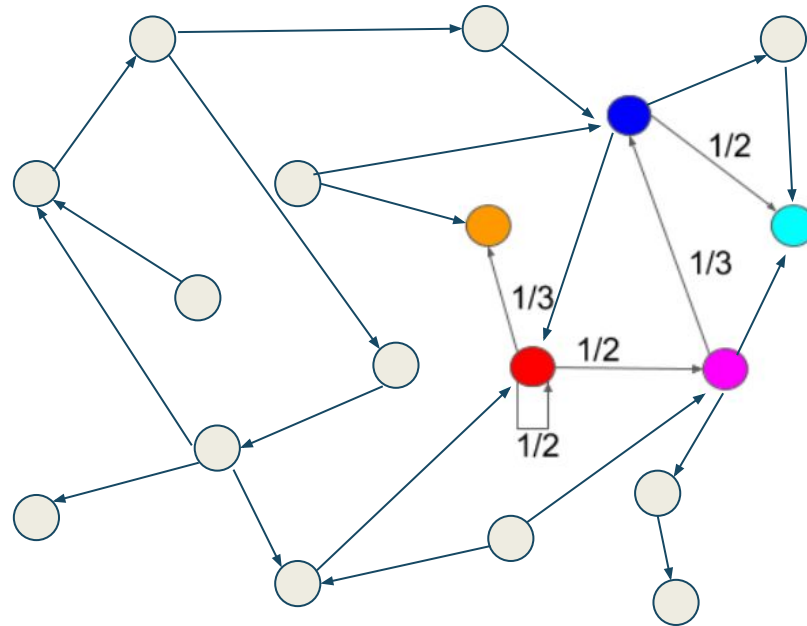


retweets \longrightarrow $1 / (1 + \# \text{ retweets})$

interpretation: In a Geometric Retweet Network, the shortest directed path from a to b is the “most retweeted path”

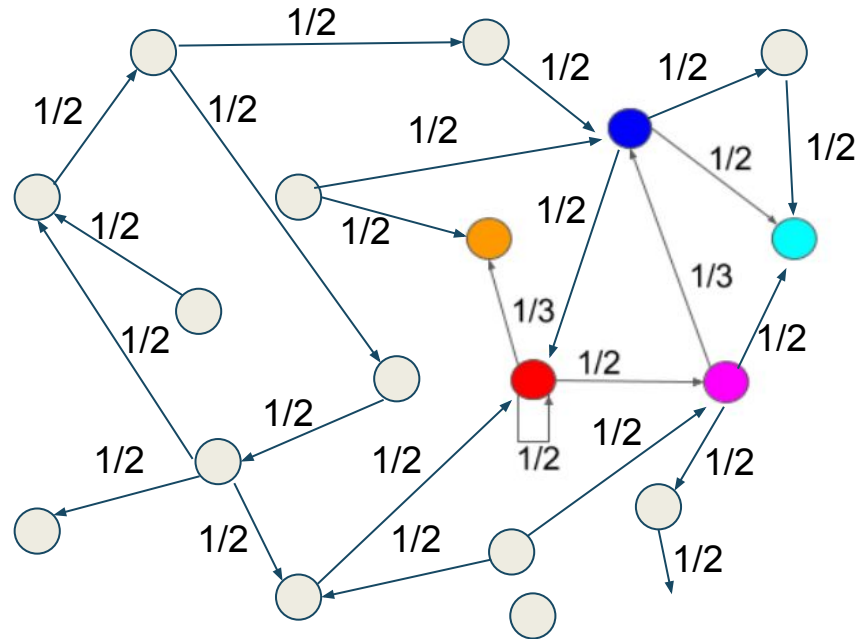
An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter



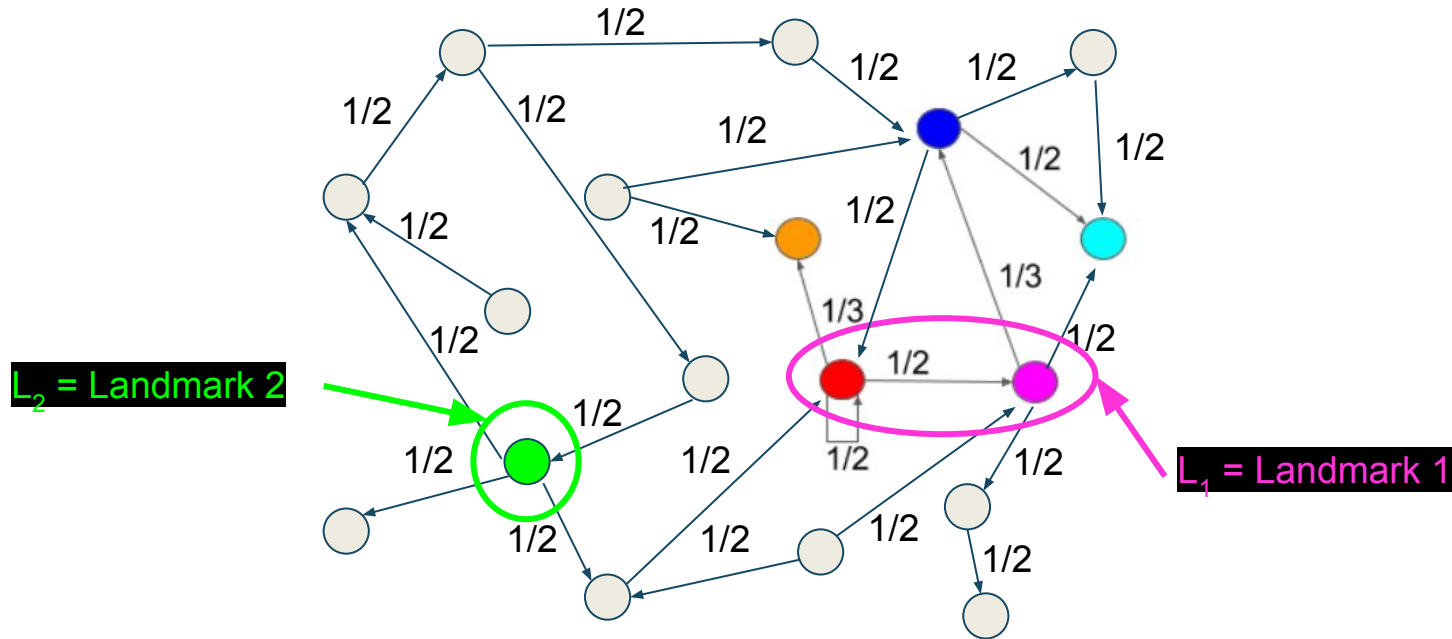
An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter



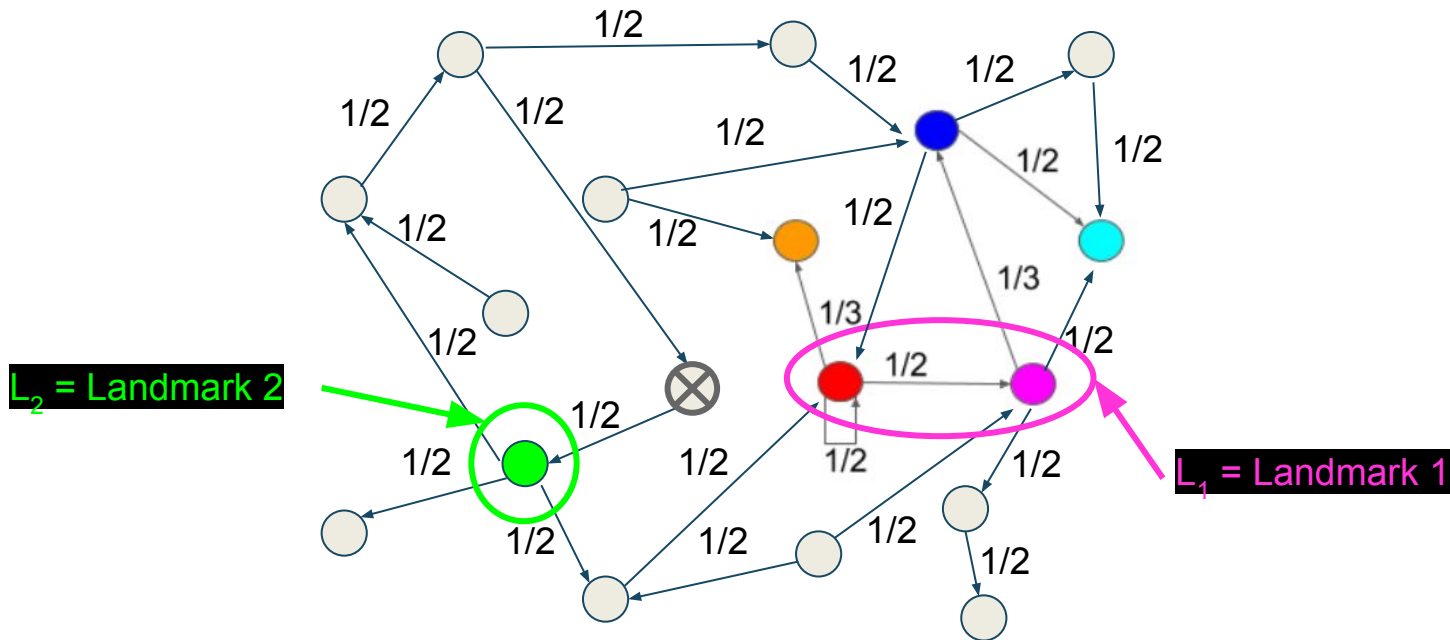
An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter



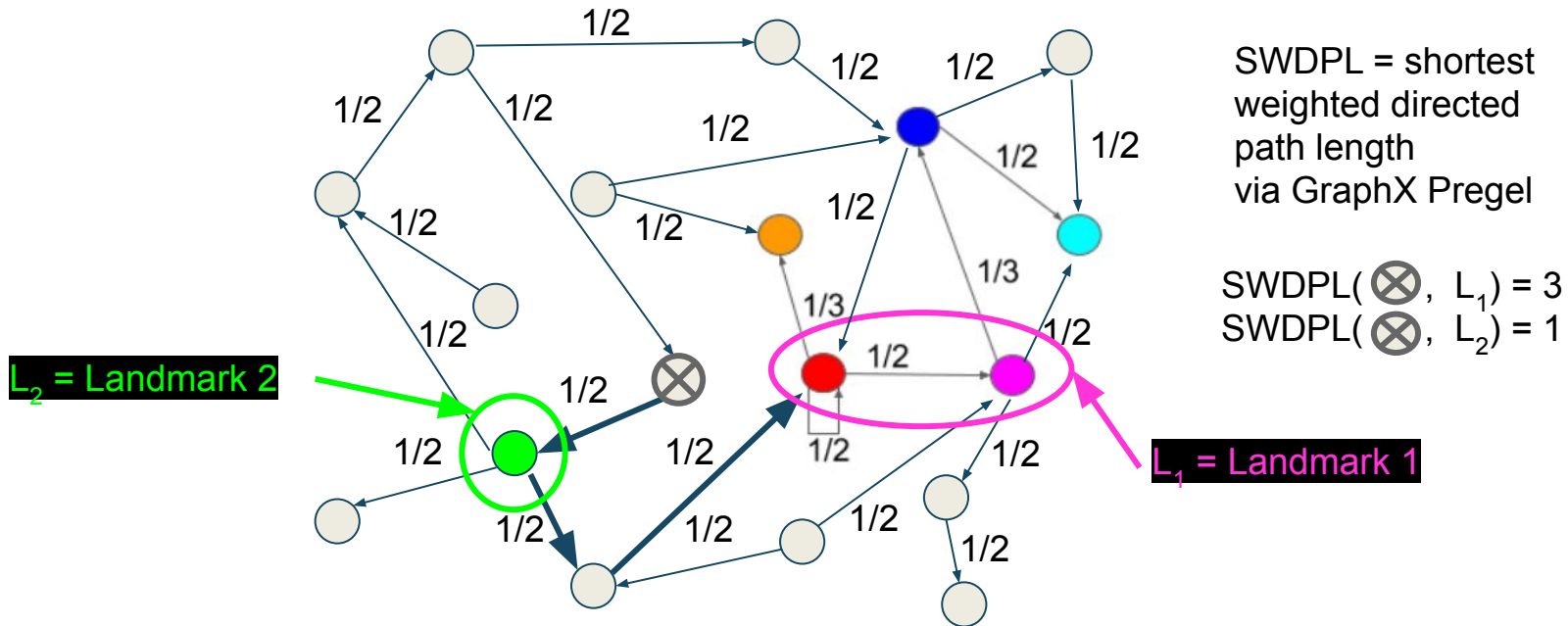
An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter



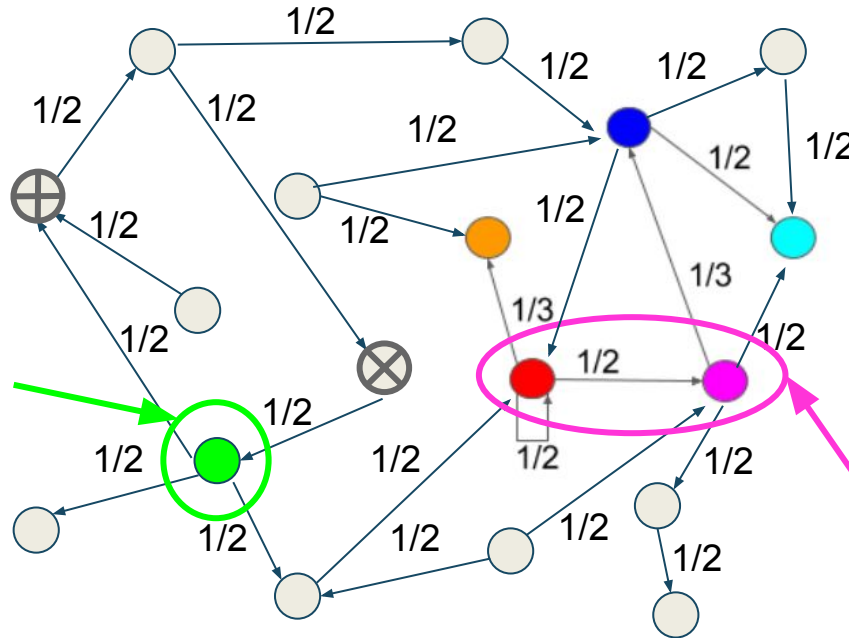
An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter



An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter



$$\text{SWDPL}(\otimes, L_1) = 3$$

$$\text{SWDPL}(\otimes, L_2) = 1$$

$$\text{SWDPL}(\oplus, L_1) = 4$$

$$\text{SWDPL}(\oplus, L_2) = 3$$

$$\text{distance}(\oplus, \otimes)$$

$$= |3-4| + |1-3|$$

$$= 3$$

$L_1 = \text{Landmark 1}$

An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter

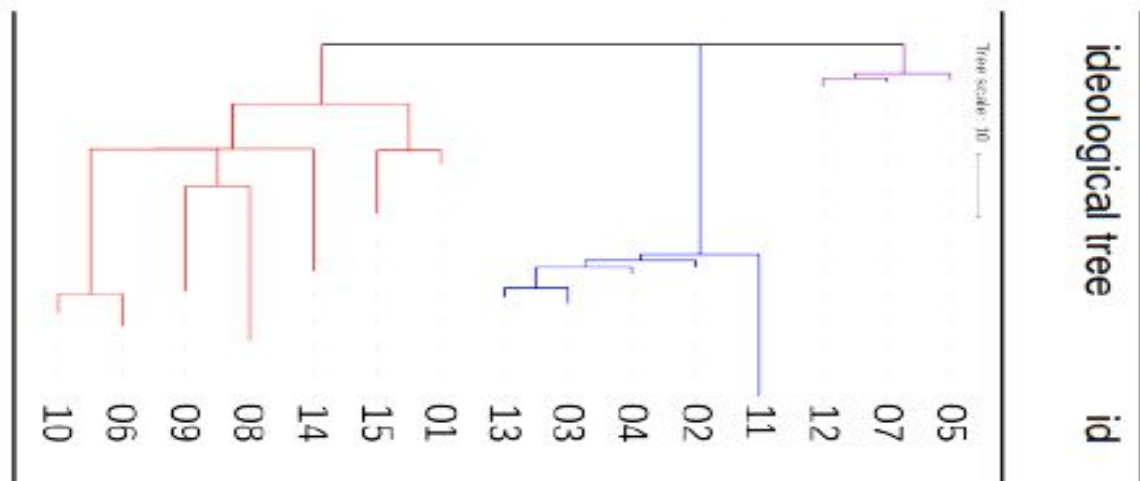
Empirical Geometric Retweet Network + distributed
multiple-sources shortest paths vertex programs
→ The “Where Am I?” Operator in Evolving *Population Ideological Trees and Forests*

- choose a set I of “influential” nodes of interest (choice is informed by the empirical out-neighborhoods and out-degrees typically)
- $I \mapsto$ most retweeted path lengths to several subsets of I
- \mapsto *Population Ideological Tree of Interest*.
- \mapsto *Population Ideological Forest of Interest* (due to multi-component retweet networks).

An Empirical Geometric Retweet Network & Most Retweeted Directed Paths — is born when distributed

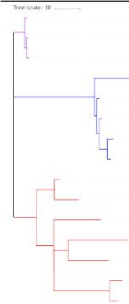
Dijkstra meets Poisson whose Expectation is Random Exponential with observed number of retweets as its mean parameter

From distance between every pair of users (based on a given set of Landmark accounts) we can obtain a retweet ideological tree of the population via Neighbor-Joining algorithm.



(Q3) Population ideological Tree & Degrees of Separation

Table 4. The top 15 groups of users according to their profiles of most retweeted path-lengths from the five politicians (DT = @realDonaldTrump, HC = @HillaryClinton, BS = @BernieSanders, PR = @SpeakerRyan, TC = @tedcruz) and eight hateful ideologies (AI = Anti-Immigrant, AM = Anti-Muslim, WN = White-Nationalist, AL = Anti-LGBT, AG = Anti-Govt, NN= Neo-Nazi, BIS=Black-Separatist, AR=Alt-Right) given by their id, frequency, percentage of population and their classification given by the ideological tree with leaf nodes as the ids.



ideological tree	id	frequency	percentage of population	Politician					Hate Group							
				DT	HC	BS	PR	TC	AI	AM	WN	AL	AG	NN	BIS	AR
	05	42853	02.005	1	1	2	4	4	5	5	7	6	4	7	7	7
	07	11481	00.537	1	2	1	4	4	5	5	7	6	4	7	7	7
	12	5868	00.274	1	1	1	4	4	5	5	7	6	4	7	7	7
	11	5972	00.279	4	2	3	5	7	8	8	9	9	7	10	10	10
	02	791286	37.016	3	1	2	4	6	7	7	8	8	6	9	9	9
	04	74126	03.468	3	1	1	4	6	7	7	8	8	6	9	9	9
	03	232093	10.857	3	2	1	6	6	7	7	9	8	6	9	9	9
	13	5173	00.242	3	1	1	6	6	7	7	8	8	6	9	9	9
	01	811586	37.965	1	4	7	4	4	5	5	7	6	4	7	7	7
	15	3892	00.182	1	4	7	1	4	5	5	7	3	4	7	7	7
	14	4011	00.188	1	4	7	4	4	1	5	3	5	4	5	7	7
	08	10460	00.489	3	5	9	1	3	3	3	5	3	6	7	9	9
	09	8069	00.377	3	3	3	3	1	4	3	3	3	6	5	6	9
	06	29997	01.403	2	3	3	3	3	5	3	3	5	5	5	3	3
	10	6257	00.293	1	3	3	4	4	5	3	3	5	4	5	3	3

(Q1) Relative frequency of retweets by any one of the hate groups or their leadership for any original tweet made by one of the politicians

Null distribution of the test statistic under the apathetic retweet network model.

Table 2. Relative frequency of retweets by any one of the hate groups or their leadership for any original tweet made by one of the politicians

Politician, observed test statistic: marginal interval for the region of acceptance at 0.001 significance level				
Donald Trump	Hillary Clinton	Bernie Sanders	Paul Ryan	Ted Cruz
0.987 : (0.6008,0.6013)	0 : (0.2708,0.2709)	0 : (0.0677,0.0682)	0 : (0.00411,0.00413)	0.0131 : (0.0024,0.0028)

(Q2) Number of unique users who retweeted a politician and a hate group at least five times each

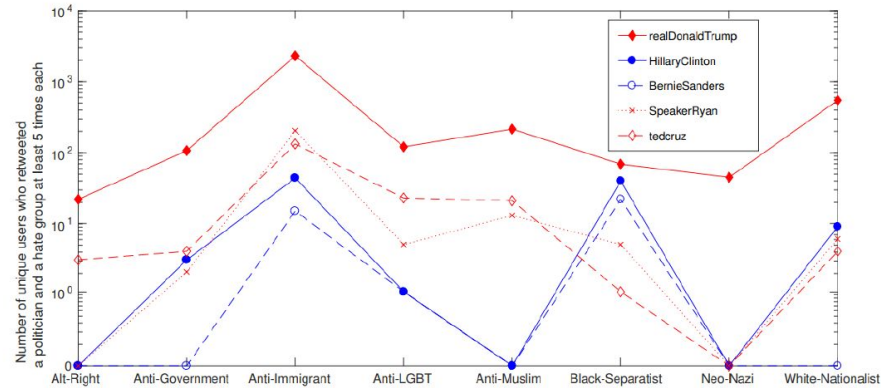


Fig. 1. Number of unique users who retweeted a politician and a hate group at least five times each (Note: The y-axis is in log-scale in powers of 10).

(Q2) Number of unique users who retweeted a politician and a hate group at least five times each

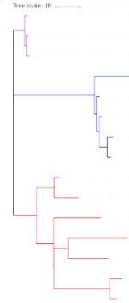
Null distribution of the test statistic under the apathetic retweet network model.

Table 3. Observed frequency of distinct users who retweeted a politician and a leader within a hate group at least 5 times each

Ideology	Politician				
	Donald Trump	Hillary Clinton	Bernie Sanders	Paul Ryan	Ted Cruz
	observed test statistic: marginal interval for the region of acceptance at 0.001 significance level				
Anti-Government	*107 : (0, 1)	3 : (0, 3)	0 : (0, 1)	*2 : (0, 1)	*4 : (0, 1)
Anti-Immigrant	*2314 : (375, 498)	°44 : (373, 492)	°15 : (369, 485)	*204 : (47, 95)	*133 : (18, 54)
Anti-LGBT	*121 : (0, 4)	1 : (0, 4)	1 : (0, 4)	*5 : (0, 3)	*23 : (0, 3)
Anti-Muslim	*215 : (0, 3)	0 : (0, 3)	0 : (0, 3)	*13 : (0, 3)	*21 : (0, 3)
Neo-Nazi	*45 : (0, 1)	0 : (0, 1)	0 : (0, 1)	0 : (0, 1)	0 : (0, 1)
White-Nationalist	*548 : (0, 12)	9 : (0, 10)	0 : (0, 10)	6 : (0, 8)	4 : (0, 7)
Black-Separatist	°69 : (653, 811)	°40 : (649, 808)	°22 : (645, 801)	°5 : (72, 128)	°1 : (28, 66)
Alternative-Right	*22 : (0, 0)	0 : (0, 0)	0 : (0, 0)	0 : (0, 0)	*3 : (0, 0)

(Q3) Population ideological Tree & Degrees of Separation

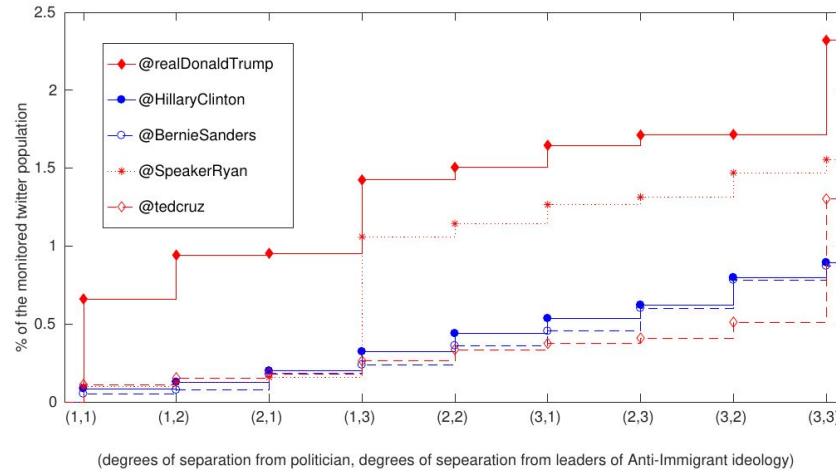
Table 4. The top 15 groups of users according to their profiles of most retweeted path-lengths from the five politicians (DT = @realDonaldTrump, HC = @HillaryClinton, BS = @BernieSanders, PR = @SpeakerRyan, TC = @tedcruz) and eight hateful ideologies (AI = Anti-Immigrant, AM = Anti-Muslim, WN = White-Nationalist, AL = Anti-LGBT, AG = Anti-Govt, NN= Neo-Nazi, BIS=Black-Separatist, AR=Alt-Right) given by their id, frequency, percentage of population and their classification given by the ideological tree with leaf nodes as the ids.



ideological tree	id	frequency	percentage of population	Politician					Hate Group							
				DT	HC	BS	PR	TC	AI	AM	WN	AL	AG	NN	BIS	AR
	05	42853	02.005	1	1	2	4	4	5	5	7	6	4	7	7	7
	07	11481	00.537	1	2	1	4	4	5	5	7	6	4	7	7	7
	12	5868	00.274	1	1	1	4	4	5	5	7	6	4	7	7	7
	11	5972	00.279	4	2	3	5	7	8	8	9	9	7	10	10	10
	02	791286	37.016	3	1	2	4	6	7	7	8	8	6	9	9	9
	04	74126	03.468	3	1	1	4	6	7	7	8	8	6	9	9	9
	03	232093	10.857	3	2	1	6	6	7	7	9	8	6	9	9	9
	13	5173	00.242	3	1	1	6	6	7	7	8	8	6	9	9	9
	01	811586	37.965	1	4	7	4	4	5	5	7	6	4	7	7	7
	15	3892	00.182	1	4	7	1	4	5	5	7	3	4	7	7	7
	14	4011	00.188	1	4	7	4	4	1	5	3	5	4	5	7	7
	08	10460	00.489	3	5	9	1	3	3	3	5	3	6	7	9	9
	09	8069	00.377	3	3	3	3	1	4	3	3	3	6	5	6	9
	06	29997	01.403	2	3	3	3	3	5	3	3	5	5	5	3	3
	10	6257	00.293	1	3	3	4	4	5	3	3	5	4	5	3	3

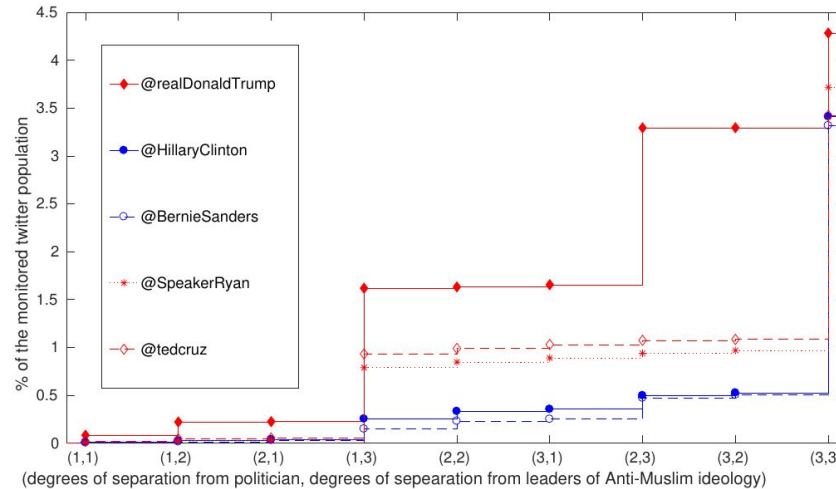
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



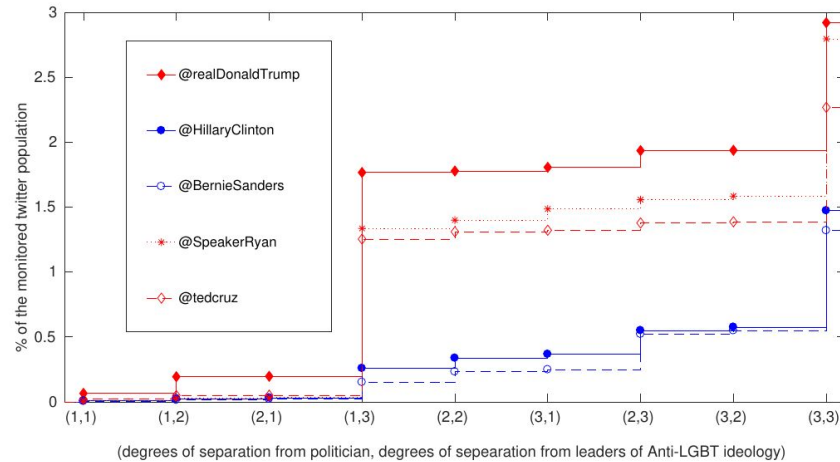
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



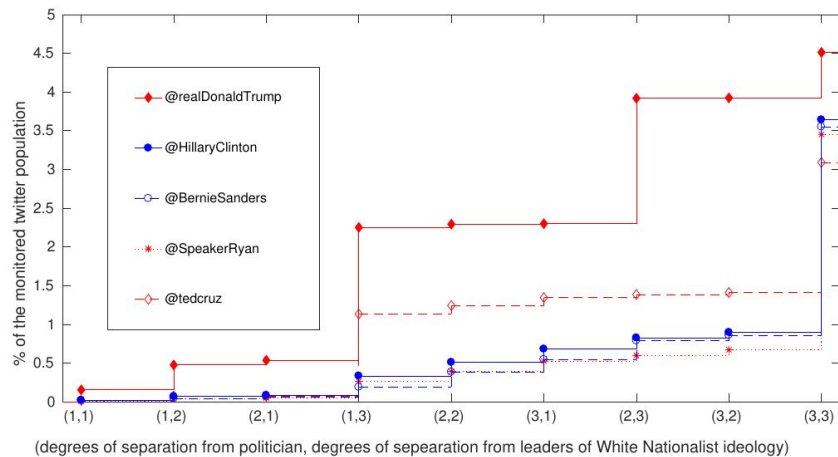
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



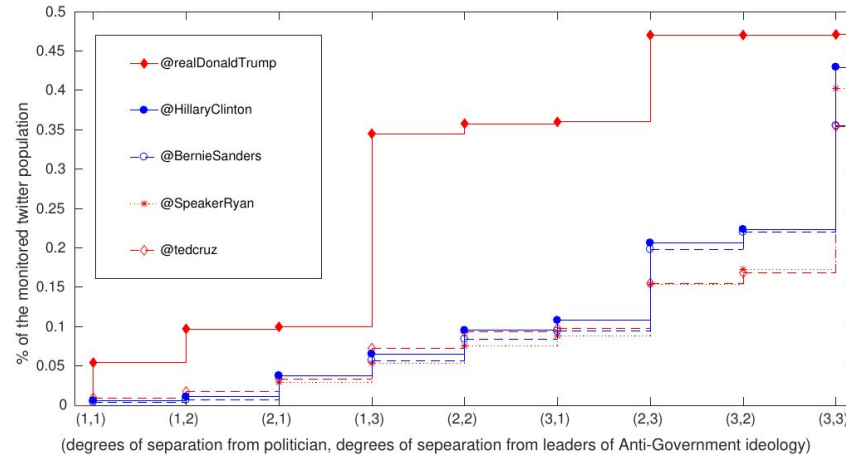
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



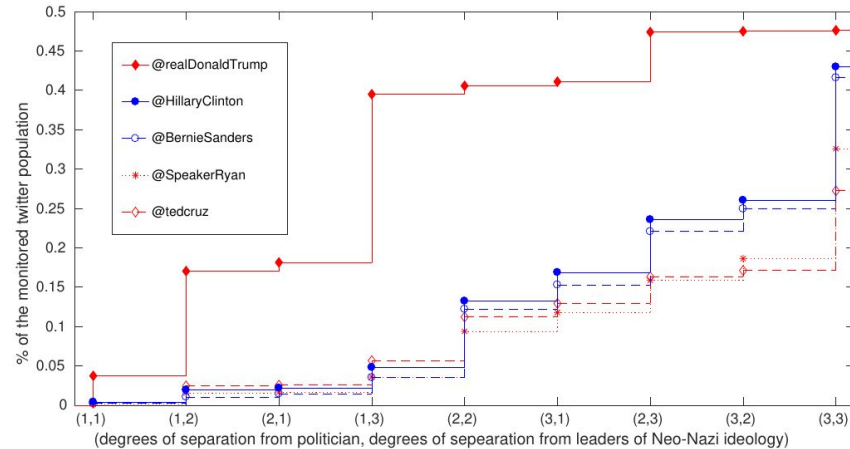
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



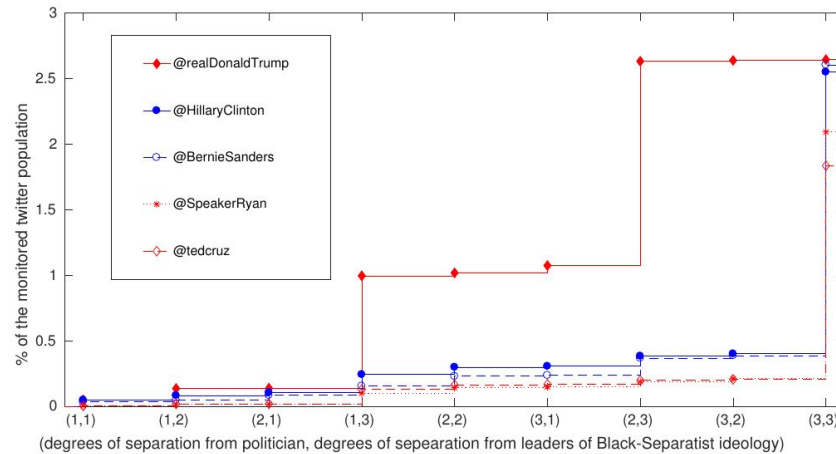
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



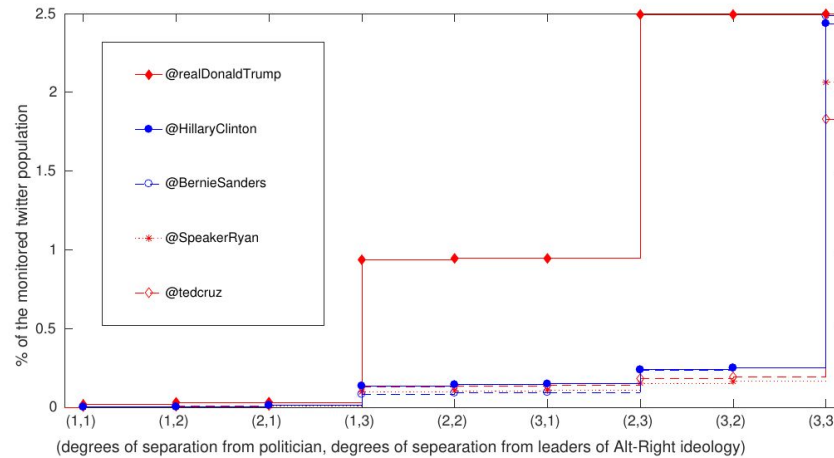
Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology —

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



Zooming into the Joint Degrees of Separation From each Politician and Hateful Ideology

Cumulative % of the monitored population who are within a given in-degree of separation from a politician and a hateful Ideology.



Significance Statement

During the 2016 US presidential election, there was significant debate on whether Donald Trump's campaign was fuelled by hate and bigotry toward minority groups. We analyzed nearly 22 million communication events on Twitter to better understand the networks of retweeters of American hate groups and five key American politicians during the late stages of the election (Donald Trump, Hillary Clinton, Bernie Sanders, Ted Cruz, and Paul Ryan). Our data reveals that Twitter users linked to various American hate groups including Anti-Government, Anti-Immigrant, Anti-LGBT, Anti-Muslim, Neo-Nazi and White-Nationalist were more strongly linked to Trump over any other politician.

Significance Statement

During the 2016 US presidential election, there was significant debate on whether Donald Trump's campaign was fuelled by hate and bigotry toward minority groups. We analyzed nearly 22 million communication events on Twitter to better understand the networks of retweeters of American hate groups and five key American politicians during the late stages of the election (Donald Trump, Hillary Clinton, Bernie Sanders, Ted Cruz, and Paul Ryan). Our data reveals that Twitter users linked to various American hate groups including Anti-Government, Anti-Immigrant, Anti-LGBT, Anti-Muslim, Neo-Nazi and White-Nationalist were more strongly linked to Trump over any other politician.

On a seemingly highly hopeful note about the "American people": Only a small fraction of those within 3 degrees of separation from @realDonaldTrump during the 9 week period are also within 3 degrees of separation from any hateful ideology!

Significance Statement

Did Trolls from Russia have an effect on our test?

Trolls := the 2,752 Twitter accounts identified by Twitter as being tied to Russia's "Internet Research Agency" troll farm

ANY GUESSES?

Significance Statement

Did Trolls from Russia have an effect on our test?

Trolls := the 2,752 Twitter accounts identified by Twitter as being tied to Russia's "Internet Research Agency" troll farm

ANSWER is NO via a non-Troll sub-graph robustness check: Out of the 12,984,331 retweets in our dataset, less than 0.1% were related to a troll account (293 were retweeted by and 12,347 were originally tweeted by a troll account) and out of 2,451,081 distinct users in our retweet network, only 172 were related to a troll account. Interestingly, *removal of these troll-related retweets from the retweet network did not alter the statistical tests.*

Generalizable Interactive Streaming-REST Design

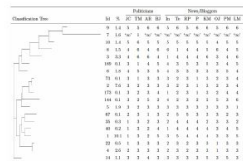
A 10 Day Design for 2017 UK Election (post-Brexit)

- **Influencers of Interest**

- **politicians:** Jeremy Corbyn (JC), Theresa May (TM), Angela Rayner (AR, Labour), Boris Johnson (BJ, Conservative)
- **journalists and bloggers:** the Independent (In), the Daily Telegraph (Te), Robert Peston (RP, journalist and author), Piers Morgan (P, journalist, tv-personality), Keven Maguire (KM, journalist), Owen Jones (OJ, left, the guardian), Paul Mason (PM, left-wing journalist (the guardian etc.)), Louise Mench (LM, previously conservative mp now blogger), and Guido Fawkes (GF, right/liberal).

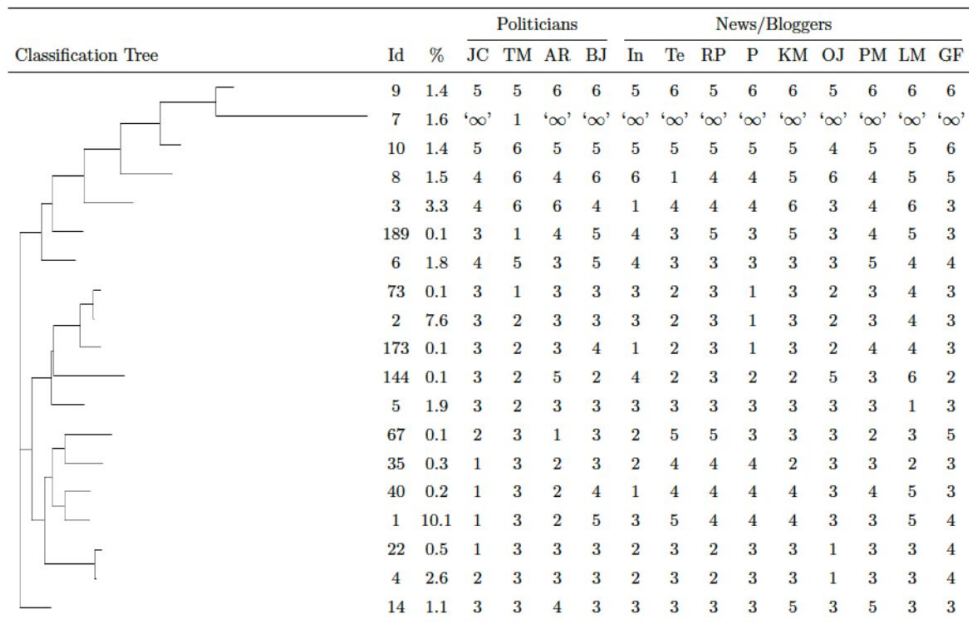


Population Ideological Tree →



Generalizable Interactive Streaming-REST Design

2017 UK Election 10 Day Design – Population Ideological Tree



Generalizable Interactive Streaming-REST Design

2017 UK Election 10 Day Design – Top 25 sorted_↓ Retweet Network Degrees

Screen Name	Out-degree	Out-nbhd	In-degree	In-nbhd
@jeremycorbyn	516833	184236	21	20
@OwenJones84	202084	78548	261	192
@Independent	195573	67341	681	22
@britainelects	130161	46921	15	14
@piersmorgan	118588	79514	157	128
@jonsnowC4	90555	53637	94	28
@paulmasonnews	74207	27358	309	222
@Telegraph	60732	29500	95	15
@LouiseMensch	53739	16916	3287	916
@Peston	48052	29552	25	8
@theresa_may	47791	31075	0	0
@faisalislam	46715	21148	101	75
@AngelaRayner	45272	15751	101	68
@DavidLammy	43043	27350	29	21
@davidallengreen	39141	15527	183	95
@bbclaurak	37683	18288	85	29
@IanDunt	36600	16069	203	157
@LordBuckethead	36436	28899	13	10
@Kevin_Maguire	36378	17015	5	1
@stephenfry	32521	26379	2	2
@Ed_Miliband	32264	23832	9	9
@MailOnline	31988	15781	594	10
@johnprescott	31906	23329	51	29
@GuidoFawkes	29033	10410	78	37
@MayorofLondon	27816	20162	44	17

All the chosen influencers, except Boris Johnson – the second most RT'd conservative MP (59) – are in top 25.

What's Happening at Project MEP Now?

- Working with Theologists at UU to bring field ethnographic domain expertise into monitoring and analysis systems around SE 2018 Election – Towards Twitter Societal Conversational Health Metrics
- To build dynamic “Where Am I?” Operators over Dynamic Population Ideological Trees and Forests (akin to proximity networks:
<https://piratepeel.github.io/proximitynetwork.html>)
- Data Science boot-camps for researchers:
<https://lamastex.github.io/360-in-525/>

- Customizable dynamically adaptable set of set of “landmark” accounts to define the desired notion of diversity in the population ideological forests
- “Where Am I?” Operator for a kind of “ideological weather report” that can be done by any Citizen “towards participatory democracy in the big data age!”
- Live Research on:
Meme Evolution Programme
 - <https://bit.ly/2OTiUH9>

The End

Many thanks to:

- Databricks Academic Partners Programme and AWS Educate & Cloud Computing Credits for Research
- Research Chair in Mathematical Models of Biodiversity (for mathematical theorizing) held jointly by:
 - 1 Veolia Environnement
 - 2 French National Museum of Natural History, Paris, France and
 - 3 Centre for Mathematics and its Applications, Ecole Polytechnique, Palaiseau, France.
- Code Contributors: Ivan Sadikov, Akinwande Atanda and Joakim Johansson
- The Transmission Process: A Combinatorial Stochastic Process for the Evolution of Transmission Trees over Networks, Raazesh Sainudiin and David Welch, *Journal of Theoretical Biology*, Volume 410, Pages 137–170, 2016 10.1016/j.jtbi.2016.07.038
- Seeded by Hate? Characterizing the Twitter Networks of Prominent Politicians and Hate Groups in the 2016 US Election, Kumar Yogeeswaran, Kyle Nash, Rania Sahioun and Raazesh Sainudiin, 2017 <http://lamastex.org/preprints/2017HateIn2016USAElection.pdf>
- See: Project MEP for more information: <http://lamastex.org/lmse/mep>

PART - II

“Polarised State of the Swedish Political Twitterverse”

Uppsala Summer Math Camp Report 2019



UPPSALA
UNIVERSITET

Summer Math Camp 2019

Polarised State of the Swedish
Political Twitterverse Around the 2018
General Election

Agnes Davíðsdóttir

Albert Nilsson

Amela Mehic

Andreas Lindgren

Claes Fälth

Johannes Graner

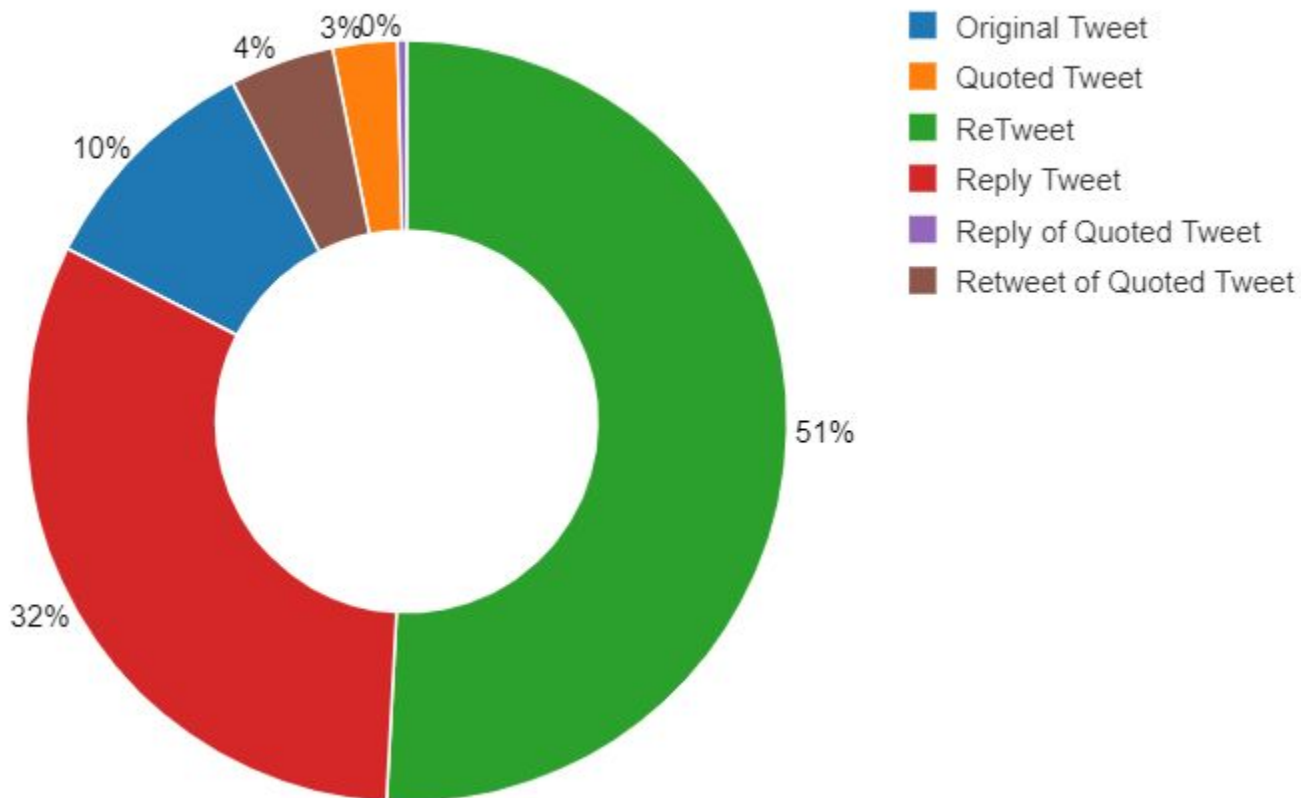
Magdalena Fischerström

Supervised by: Raazesh Sainudiin & Tilo Wiklund



UPPSALA
UNIVERSITET

Twitter interactions



↻ Aron Flam Retweetade

 **Peter Sellei**
@PeterSellei

Antar att ni då fryser bistånd till Palestina tills öppna och demokratiska val sker samt att mänskliga rättigheter respekteras?
Ordföranden i ert systerparti är inne på sitt 15:e år som president efter att ha blivit vald för en 4-årig mandatperiod.

 **Margot Wallström** ✓ @margotwallstrom · 4h
Vår demokratioffensiv innebär också mer av vår feministiska utrikespolitik. I den rådande globala polariseringen är det viktigare än någonsin att Sverige fortsätter vara en stark röst för jämställdhet och åtnjutande av mänskliga rättigheter för alla. #Priosutrikespol

10:22 fm · 28 aug. 2019 · [Twitter for iPhone](#)

45 Retweets **197** gilla-markeringar

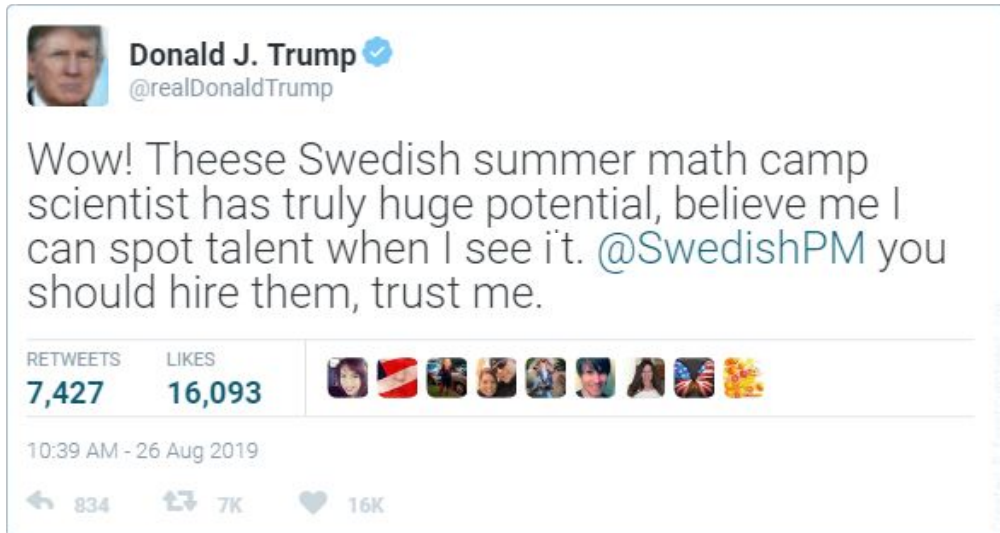
 **Eric Danell** @eric_danell · 3h
Svarar @PeterSellei och @konsensuseliten
Finns det mer hyckleri än hos regeringspartierna med kringfjäskande C och L samt V. Varför ställer inte C och L krav på regeringen bl.a. om biståndet till Palestina. Låt inte regeringen härja fritt inom utrikespolitiken! C och L, sluta fjäska för Löfven o co! Visa lite kurage!

   1 



UPPSALA
UNIVERSITET

Data



- Meme evolution programme Sverige, Raazesh Sainudiin
- Mattias Gardell, field ethnography, department of theology Uppsala
- Simon Lindgren, digital sociology, department of sociology Umeå

- 91 million tweets



UPPSALA
UNIVERSITET



VAL
2018

What data do we want?

- Structure of the Swedish Twitter network around the 2018 general election
- Swedish Twitter users
- Not bots or other spam



UPPSALA
UNIVERSITET

Cleaning Data

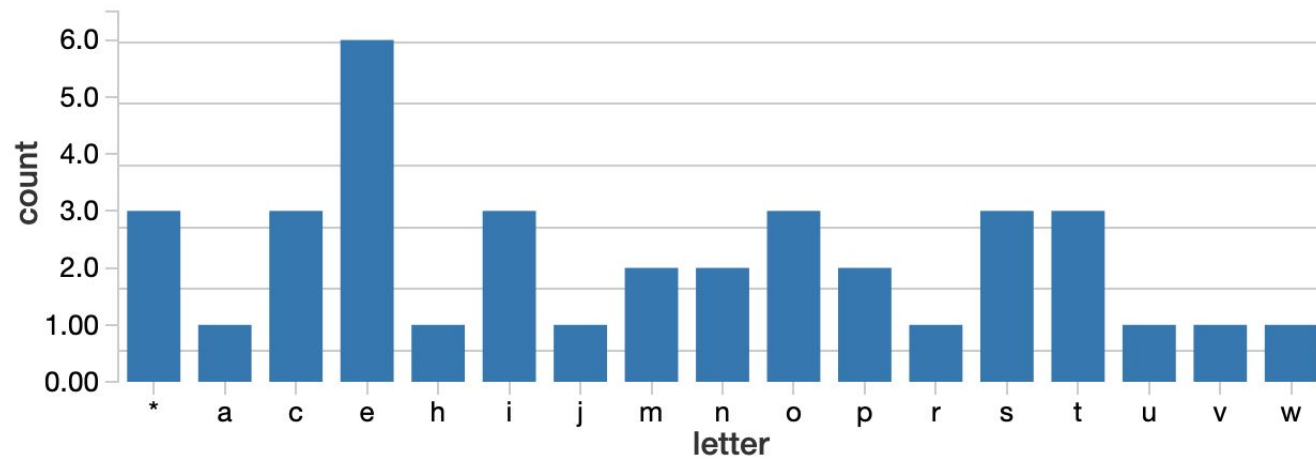
- Language filter
 - exclude non-Swedish Twitter users
- High frequency filter on retweets





- Every language has a unique letter frequency
- Swedish letter frequency compared to user letter frequency
- Remove punctuations, emojis, white space, etc ✨🎉🧶 조선글
- Remove punctuationsemojiswhite spaceetc조선글
- non-Swedish letter is mapped to *

Language Filter





Tolerance of Distance to Swedish

- A “Swedish word” in our sense:
ajhbksdnmäaösfm sakfpåslfafebgh
mtk
- Result: Tolerance 5% gave max
distance 0.485

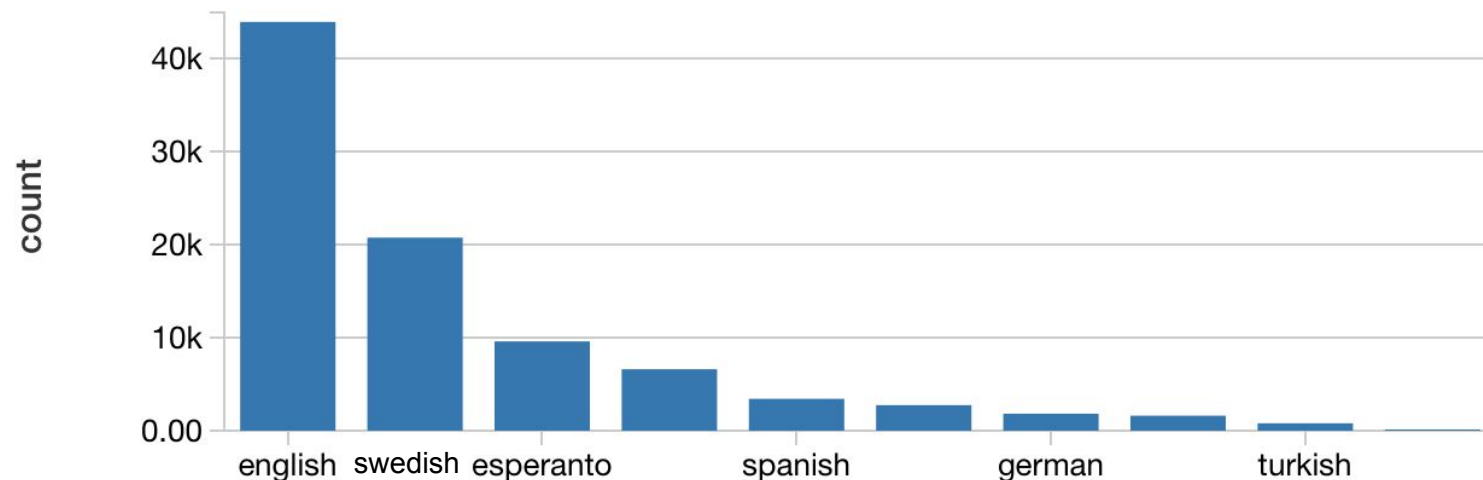
Distributions

- Swedish letter frequency
- Users tweet length



- We made sure that language was closer to Swedish than English
- Result: Reduced data from 91 to 29 million data points

Language filter



Sample of the language distribution before cleaning



High Frequency Filter

- Frequent retweets
 - Anyone retweeting the same person more often than 1000 times over the 8 months was taken out of the data (more than 4 times a day on average)



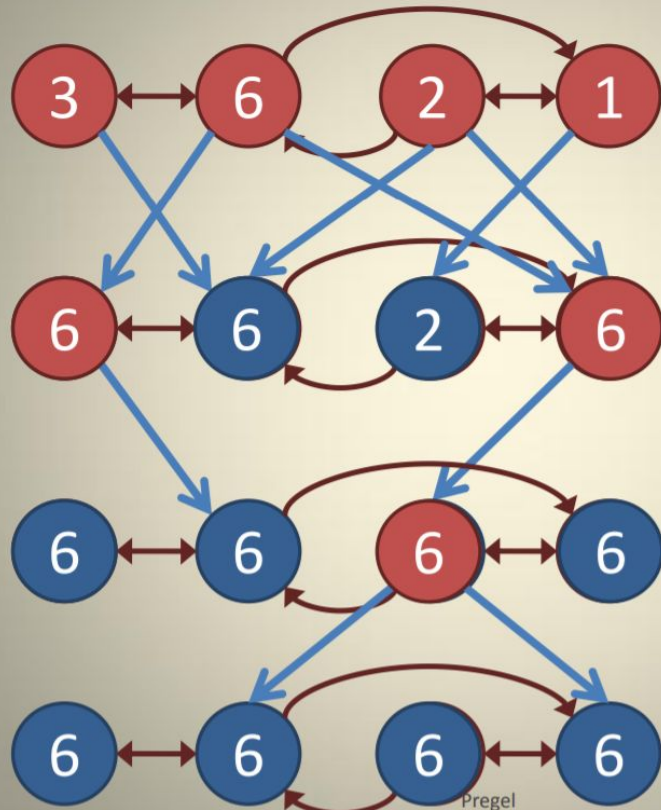


Distributed Vertex Programming

- With the clean data, we want to group similar users for analysis
 - This is easily done if we have a graph
- Problem: Computing things on large graph
 - Solution: Only pass messages between neighbours!

Distributed Vertex Programming

Example. Finding largest state in strongly connected graph:



Blue Arrows
are messages

Blue vertices
have voted
to halt

- Every vertex begins with an *initial state*
- Vertices *send a message* to their neighbours
- Each vertex *updates* its state based on incoming messages
- A vertex can choose to halt, not participating in the next iteration
- This is called a Pregel program



Clustering

- We want to create clusters of retweeting users
- First we need a directed graph of retweet network
 - Vertices: Users and their unique ID
 - Edges: If A retweets B, include edge from B to A
 - We get a *directed, multi-edged, looped graph*



Clustering

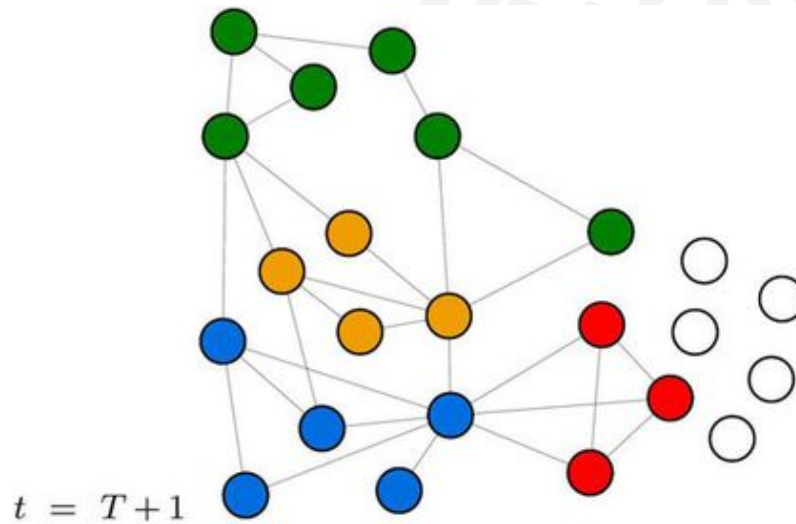
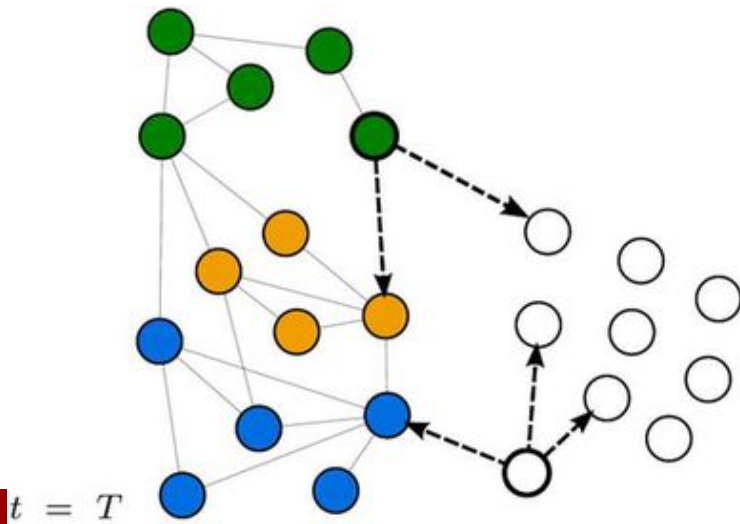
component	size
528	1157438
159429982	1873
15261401	493
15005510	170
913991205159096320	157
53614956	130
73332929	123
69549657	86
319579913	85
40227371	84
25949039	82
148152741	77
94012801	77
460893708	71
20087934	71
46806220	70
15021968	62
161418081	59
35569691	57
954735119805292544	57

- We only consider users who appear in the largest connected component of the graph
 - Computing the connected components is done by a Pregel program
- 89.9% of users are in the largest component



Clustering

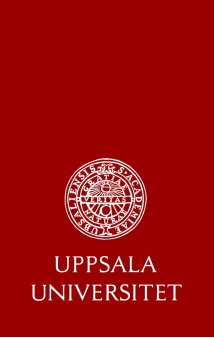
- We cluster users based on who they retweet
 - Pregel program implementing *label propagation*
 - Initial state: ID of the vertex
 - Sent Message: Current state
 - Update: Take mode of incoming messages





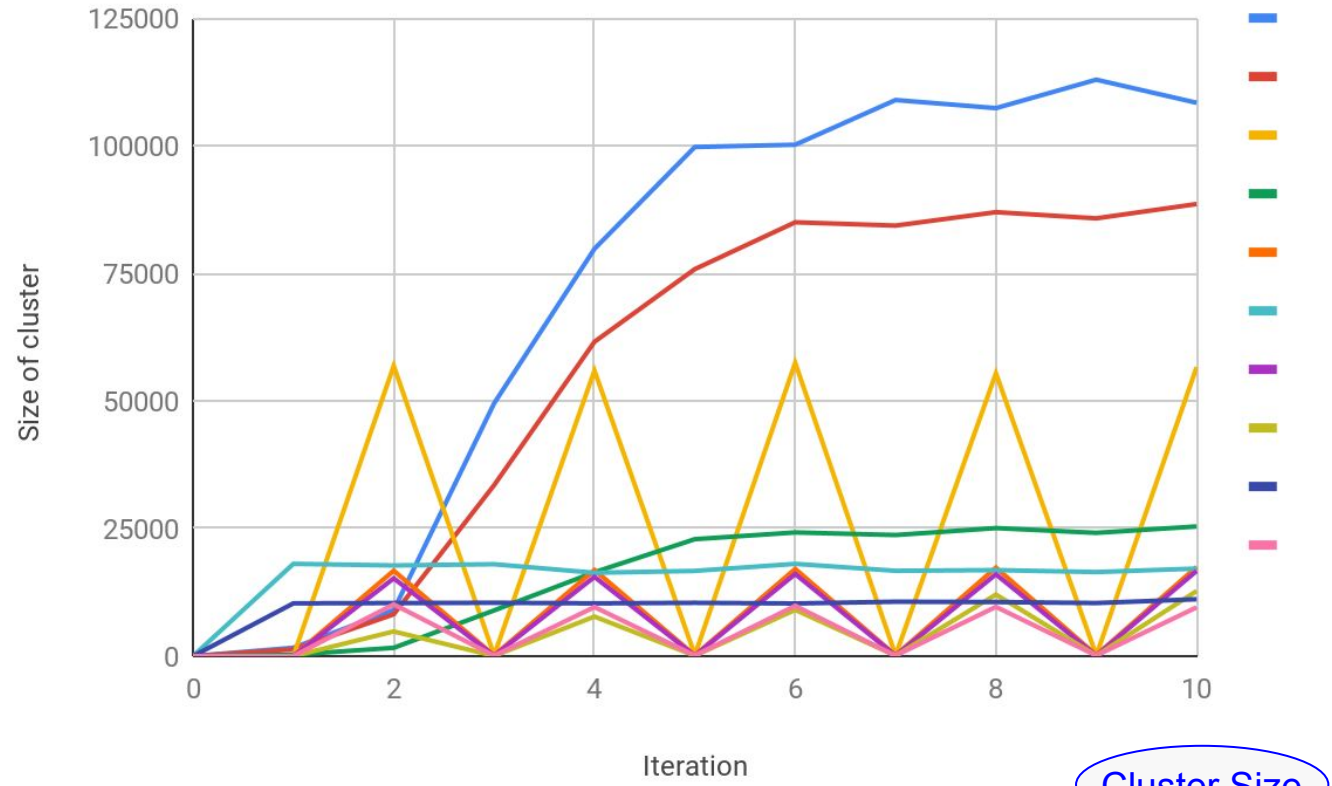
Clustering

- Pros
 - Relatively cheap
 - No information about graph necessary
- Cons
 - Convergence not guaranteed
 - Can put all vertices in the same cluster



Clustering

- After 10 iterations, the three Swedish clusters have settled



Swedish!

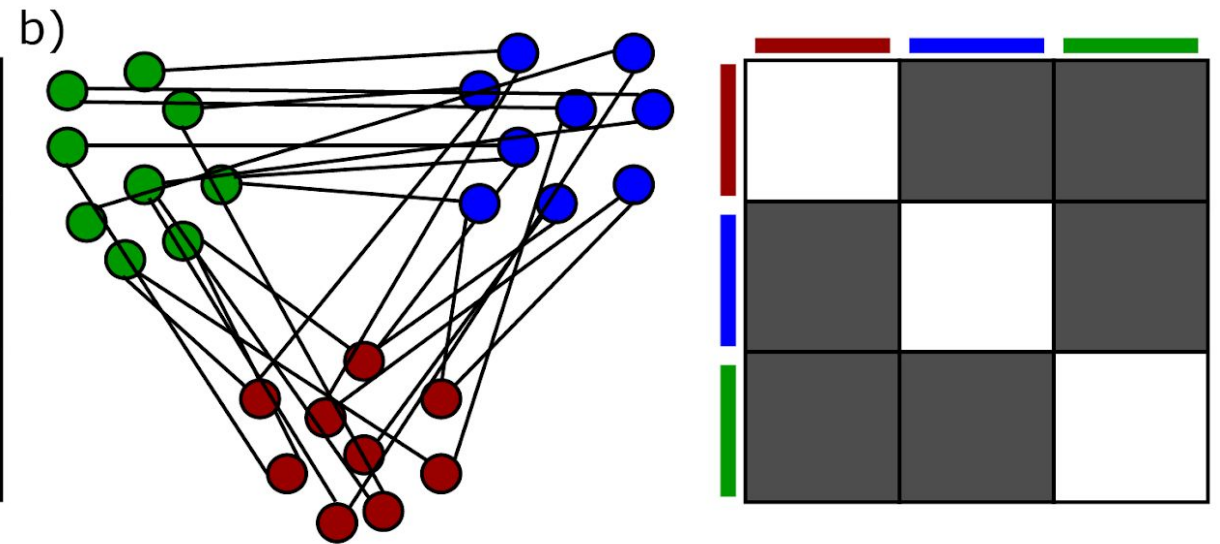
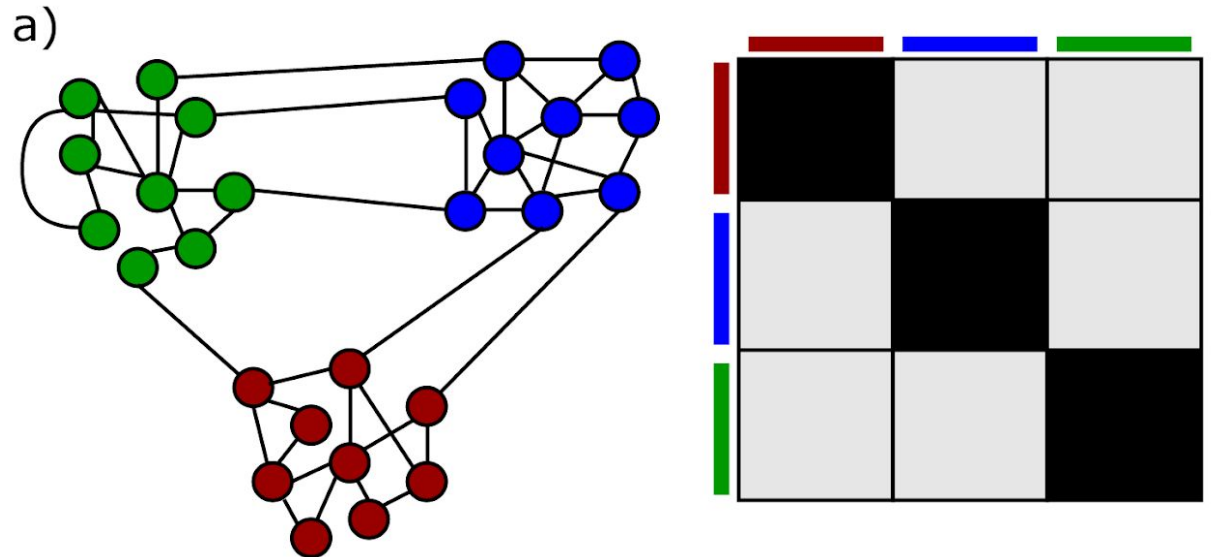
Cluster Size

label	size1	size2	size3	size4	size5	size6	size7	size8	size9	size10
343663197	1641	9033	49519	79794	99747	100188	108949	107359	112934	108408
3048723709	1330	8197	33493	61579	75819	84996	84346	86978	85763	88576
3285105132	1	56869	167	55939	202	57465	144	55412	177	56632
434315852	281	1630	8957	16436	22913	24232	23726	25084	24152	25405
700434386160852992	1	16737	138	16874	149	17082	157	17339	194	17475
1283934055	18081	17769	18005	16330	16700	18076	16718	16855	16474	17152
832135844706148353	1	15242	55	15530	74	16092	83	16037	86	16708
22558580	1	4837	43	7756	78	9068	79	12067	95	12764
2982376457	10345	10398	10460	10340	10448	10319	10661	10597	10402	11149
2297996020	2	10143	75	9618	126	9857	100	9654	118	9562



Twitter Interactions Between Clusters

- Build network based on retweets or reply tweets
- Form clusters on the network with label propagation
- Test the clustering with unseen retweets and replies
- Stochastic block model
- Probability for edge, p within clusters, q between clusters



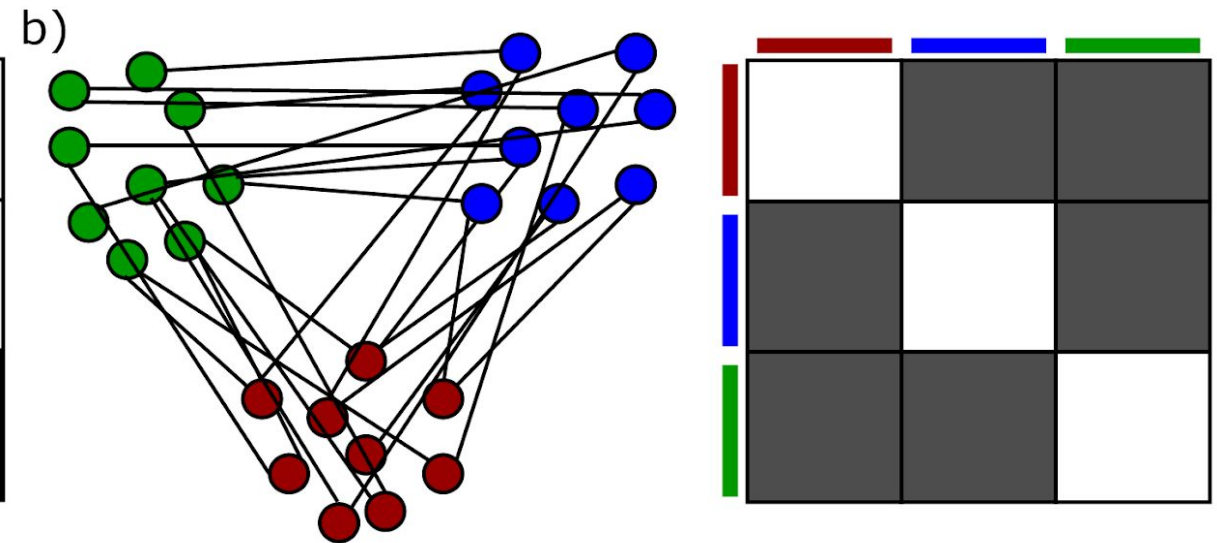
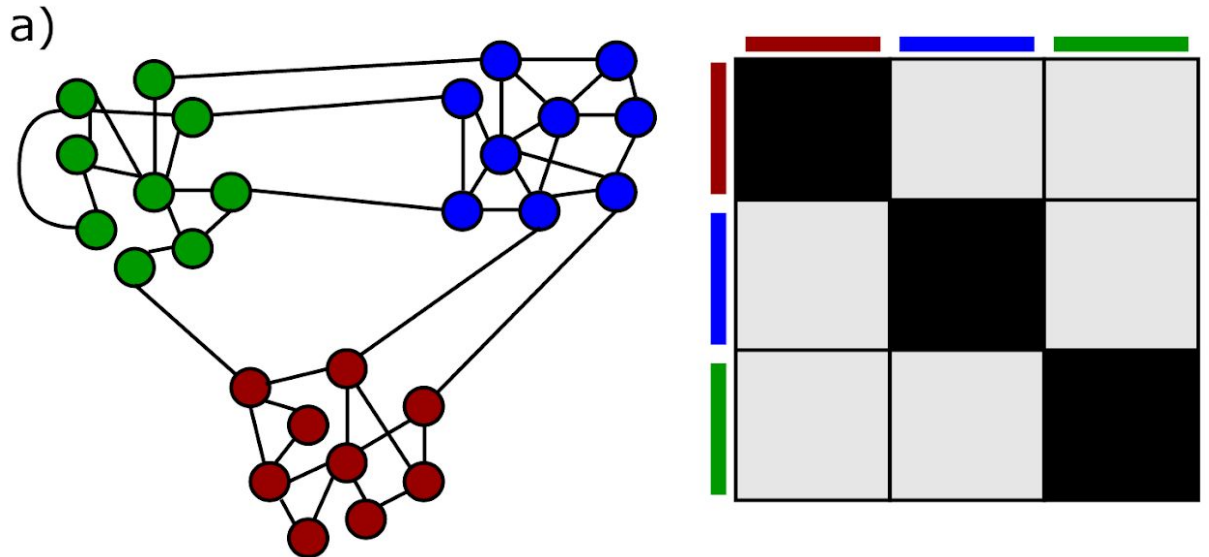
Twitter interactions between clusters

Retweet Network

Tweet Type	Within	Between
Retweet	67.6%	32.4%
Reply Tweet	45.4%	54.6%
Random connection	0.38%	99.62%

Reply Network

Tweet Type	Within	Between
Retweet	1.49%	98.51%
Reply Tweet	9.79%	90.21%
Random connection	0.00057%	99.99943%

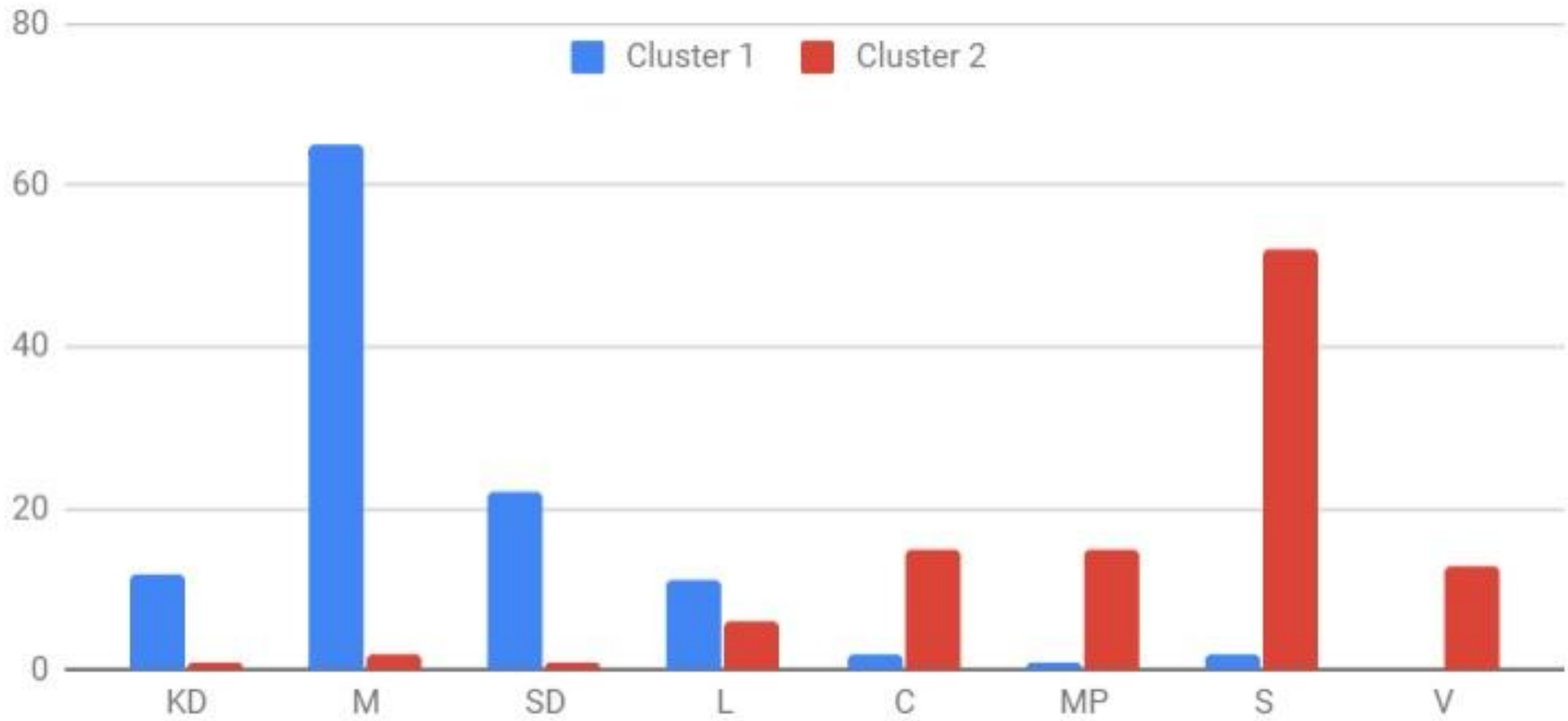




UPPSALA
UNIVERSITET

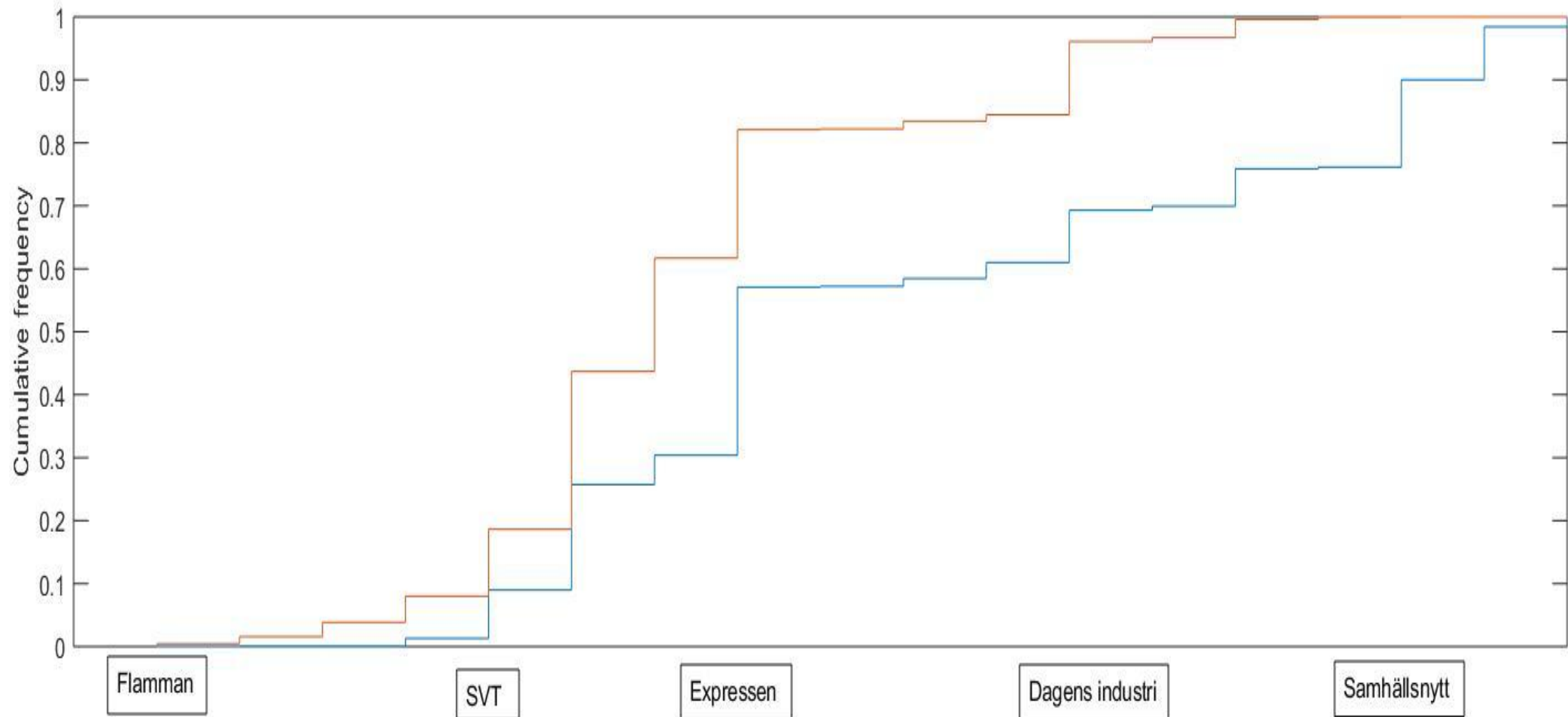
Exploring the Clusters





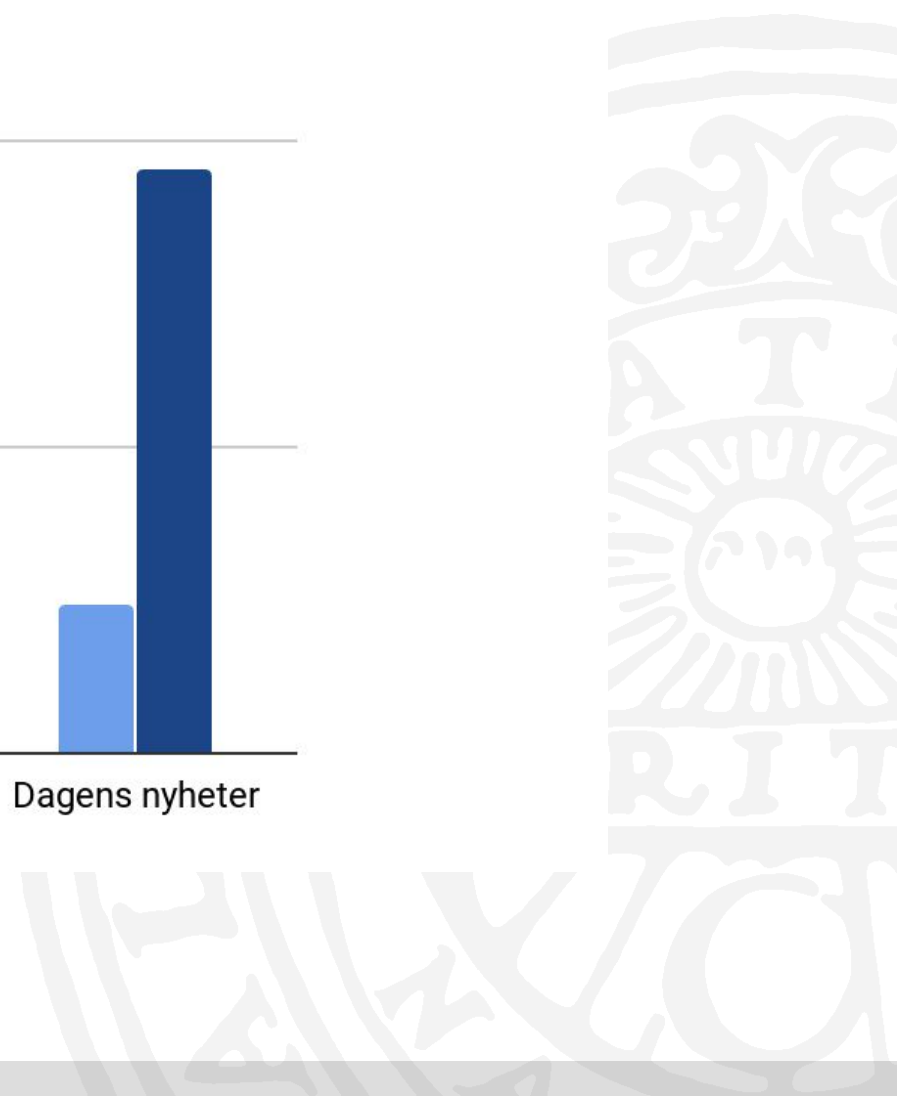
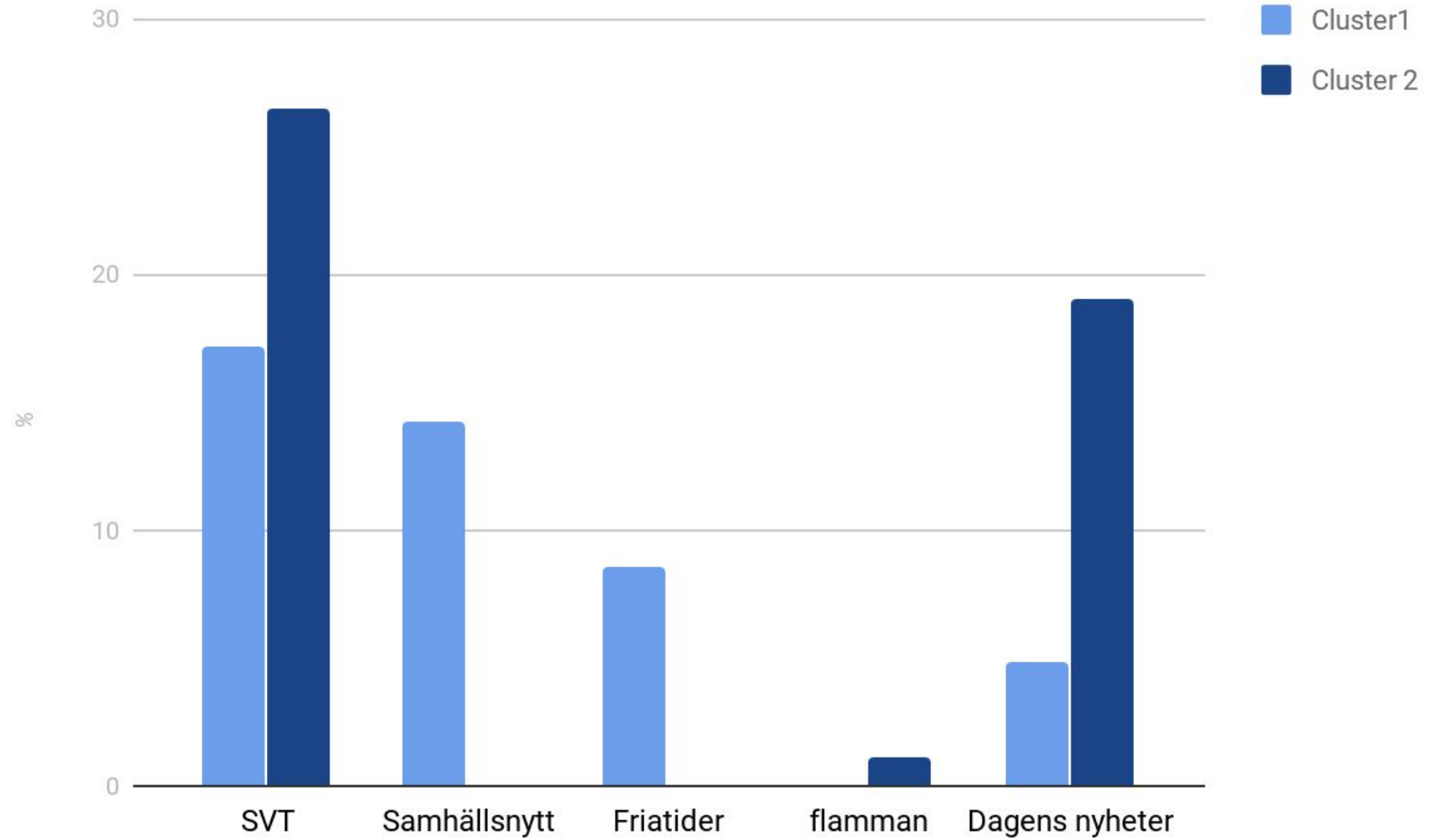


Kolmogorov–Smirnov Test



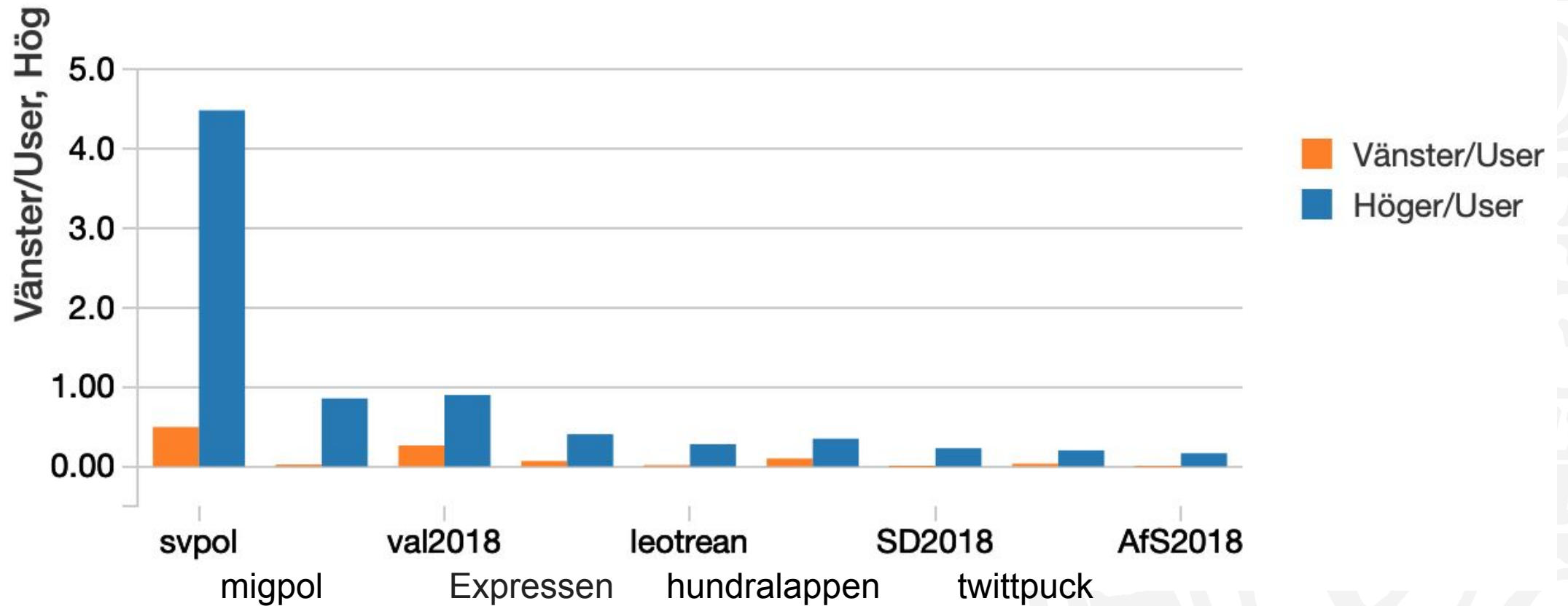


UPPSALA
UNIVERSITET





Hashtag-Distribution





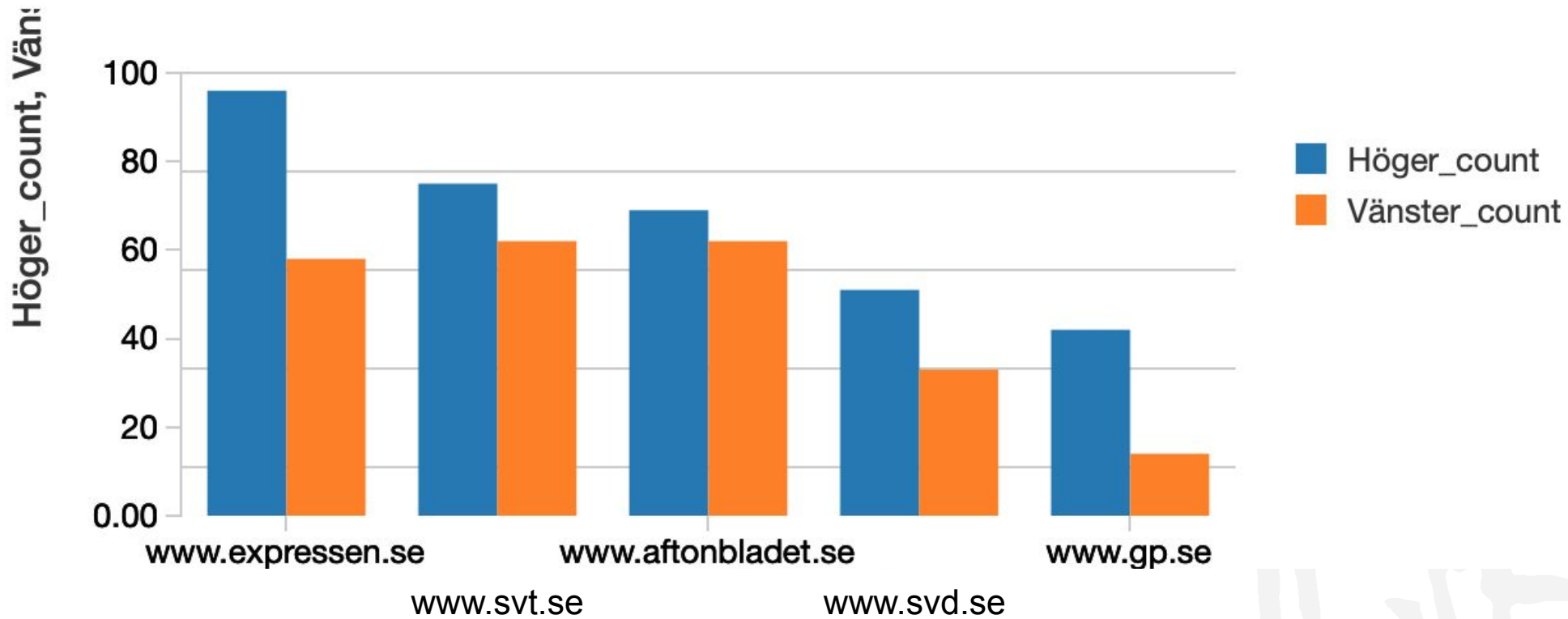
Hashtag Distribution in Retweets

- H_0 : Is the hashtag distribution of the clusters different from the global distribution?
- Sampled random subgraphs 1000 times with same size as the cluster for each cluster
- Total variation distance from sampled graphs to global hashtag distribution
- Null hypothesis was rejected for all 3 clusters with 0.1 % significance

Clusters	Interval	Obs
Right-wing	[0.042,0.044]	0.149
Left-wing	[0.078,0.082]	0.541
Sport	[0.261,0.284]	0.744



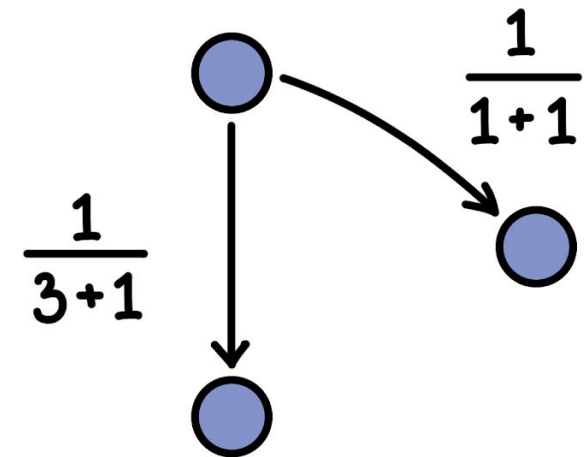
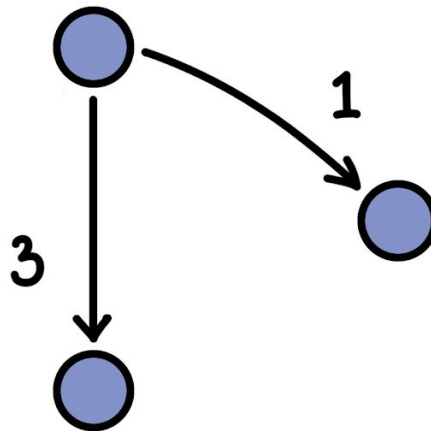
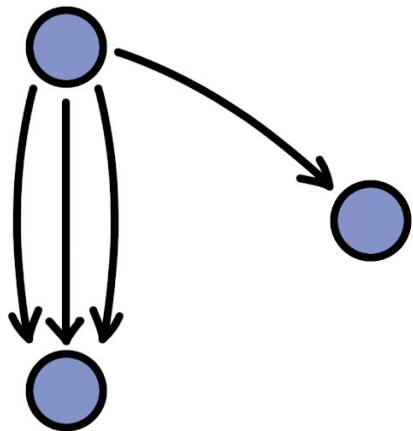
Retweets - URLs





Shortest Weighted Path

- Distributed Dijkstra's algorithm
- Weights - Number of retweets
- Landmarks - Key users
- Degrees of separation





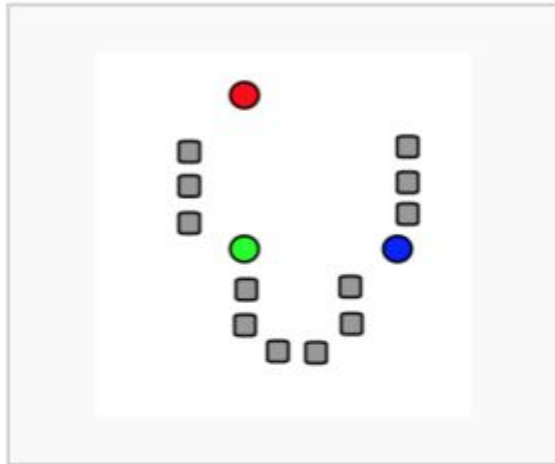
Shortest Weighted Path

freq	percentage	cs	hanifbali	sdriks	katjanouch	socialdemokrat	strandhall	jsjostedt
84371	37.9	37.9	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4
32732	14.7	52.6	∞	∞	∞	∞	∞	∞
8995	4.04	56.7	3	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4
8312	3.73	60.4	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4	3
7914	3.56	64.0	≥ 4	≥ 4	2	≥ 4	≥ 4	≥ 4
5988	2.69	66.7	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4	2
5742	2.58	69.3	2	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4
4192	1.88	71.2	≥ 4	≥ 4	3	≥ 4	≥ 4	≥ 4
3598	1.62	72.8	3	≥ 4	≥ 4	≥ 4	≥ 4	3
3005	1.35	74.1	≥ 4	≥ 4	≥ 4	3	≥ 4	≥ 4
2892	1.30	75.4	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4	1
2885	1.30	76.7	≥ 4	≥ 4	≥ 4	≥ 4	3	3
2866	1.29	78.0	≥ 4	≥ 4	≥ 4	≥ 4	3	≥ 4
2515	1.13	79.1	1	≥ 4	≥ 4	≥ 4	≥ 4	≥ 4
2429	1.09	80.2	≥ 4	≥ 4	≥ 4	2	3	≥ 4
2280	1.03	81.3	≥ 4	≥ 4	1	≥ 4	≥ 4	≥ 4
2240	1.01	82.3	≥ 4	3	2	≥ 4	≥ 4	≥ 4
2107	0.947	83.2	≥ 4	2	2	≥ 4	≥ 4	≥ 4
1854	0.834	84.0	2	≥ 4	2	≥ 4	≥ 4	≥ 4
1720	0.773	84.8	3	3	3	≥ 4	≥ 4	≥ 4

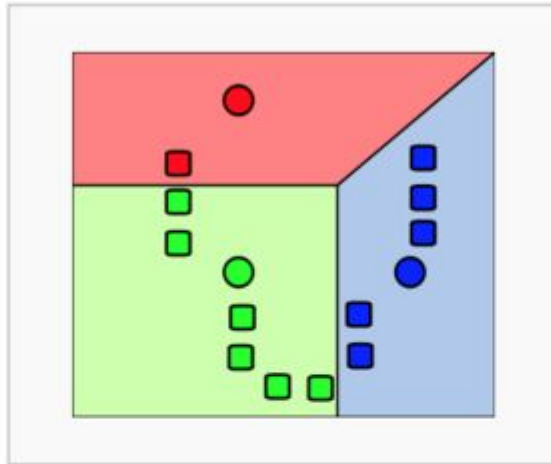


K-Means

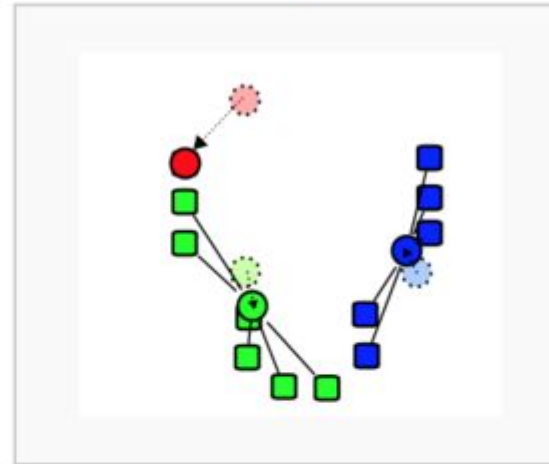
Demonstration of the standard algorithm



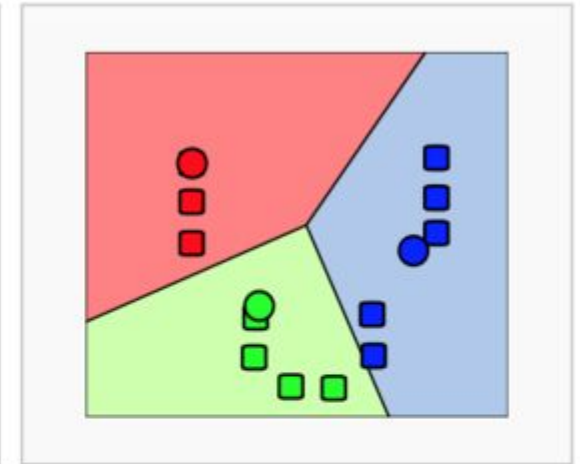
1. k initial "means" (in this case $k=3$) are randomly generated within the data domain (shown in color).



2. k clusters are created by associating every observation with the nearest mean. The partitions here represent the [Voronoi diagram](#) generated by the means.



3. The **centroid** of each of the k clusters becomes the new mean.



4. Steps 2 and 3 are repeated until convergence has been reached.



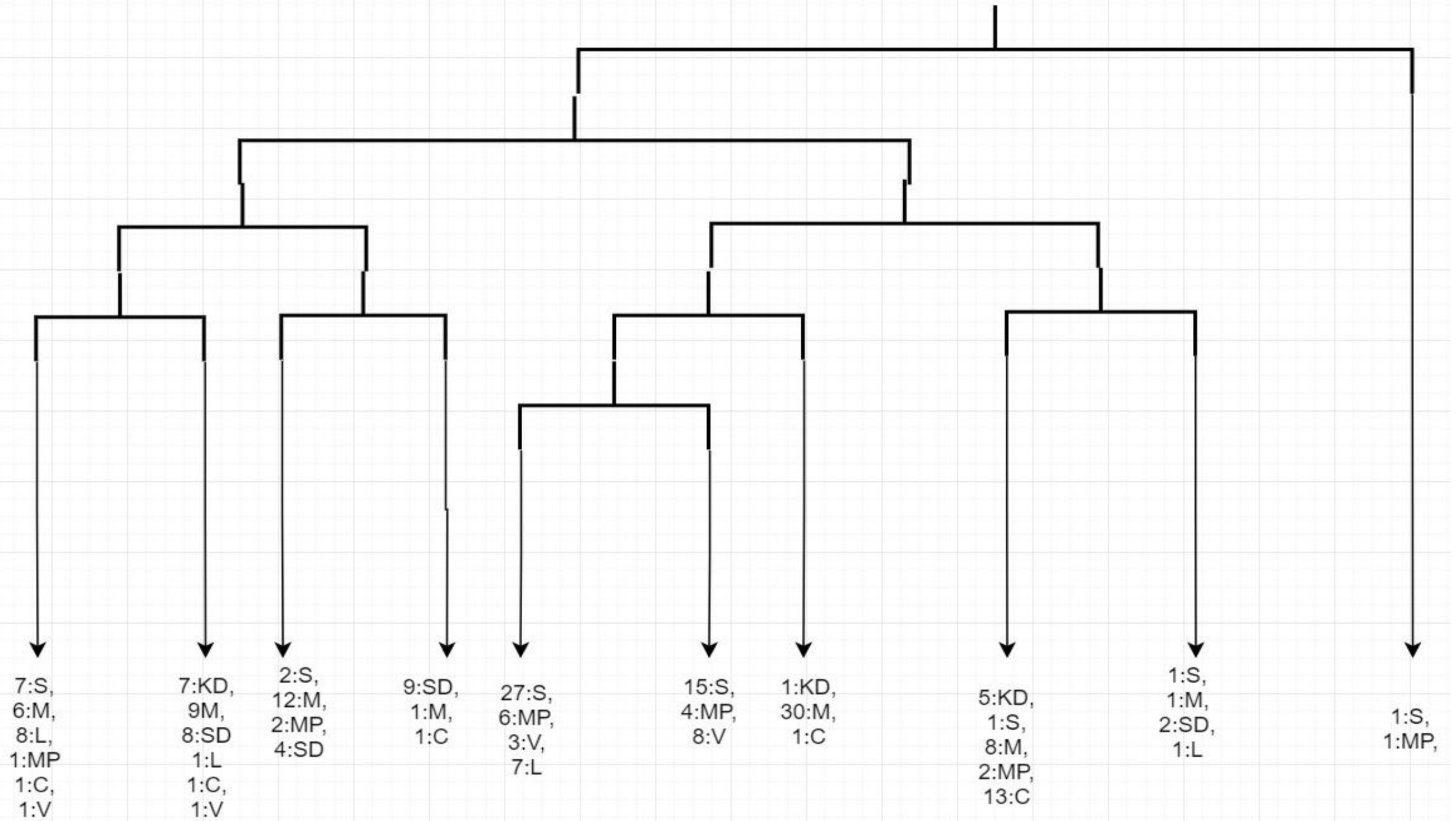
UPPSALA
UNIVERSITET

Bisecting K-Means

Hanif Bali(M) Isabella Lövin(MP) Jonas Sjöstedt(V)
Paula Bieler(SD) Jeff Ahl(AFS)

Nazispotting katjanouch







Swedish Twitter

Overall Conclusions:

- Clustering captures political affiliation
- The clusters use different news sources
- Different distributions of hashtags and URLs
- We intended to also look closer at the content of the tweet-text itself but the method we had available didn't give meaningful results



UPPSALA
UNIVERSITET

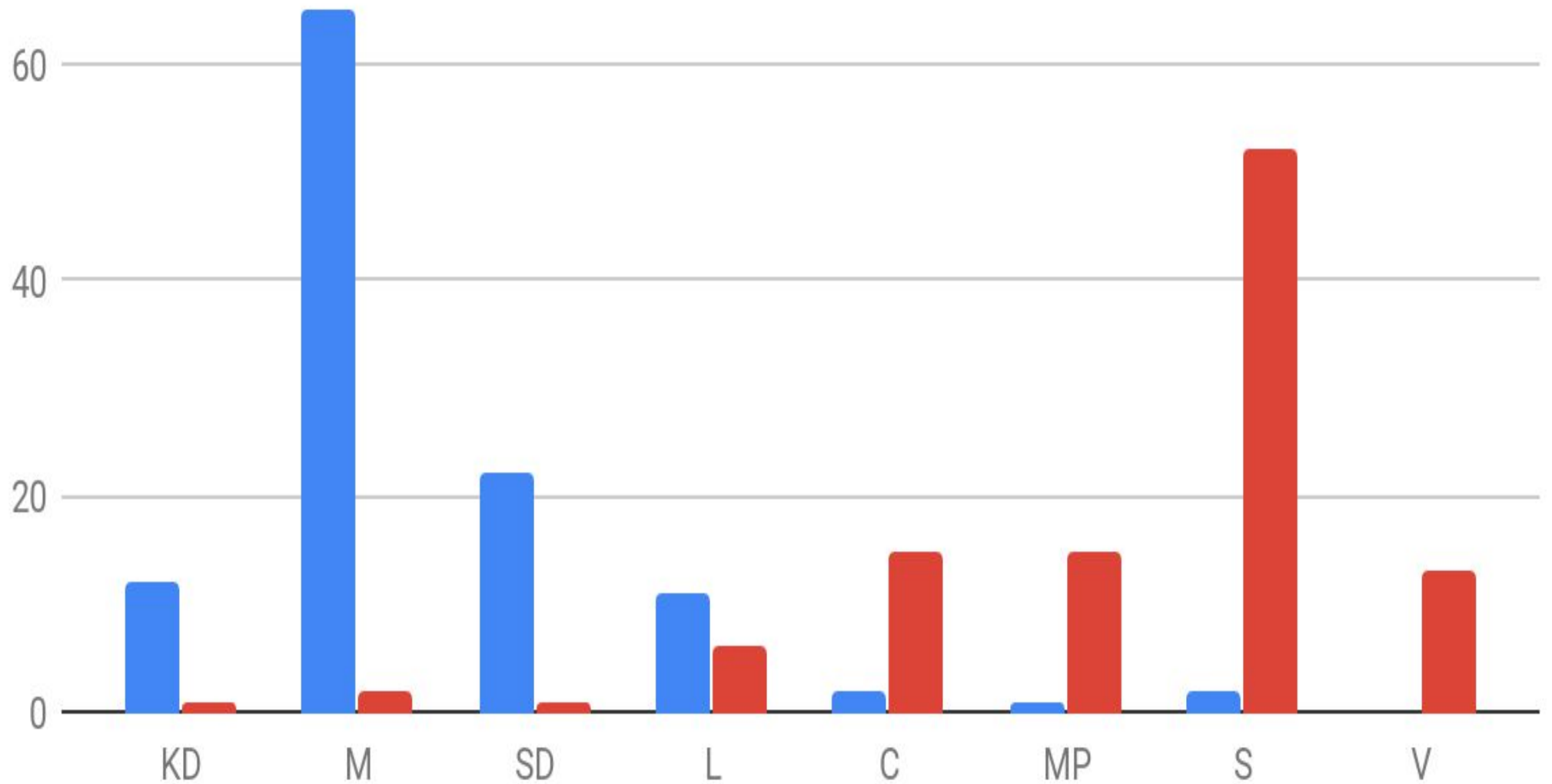
Acknowledgements

- Databricks in USA provided free cloud computing
- Combient AB in Sweden donated 4 NUCs
- Tilo has worked hard on on-premise computing support





Members of Parliament





Kolmogorov–Smirnov Test

$$D_{n,m} = \sup_x |F_{1,n}(x) - F_{2,m}(x)|$$

$$D_{n,m} > c(\alpha) \sqrt{\frac{n+m}{nm}}$$

