

Factionalization in an online community: Combining network and linguistic information

Jean Mark Gawron
Alex Dodge

San Diego State University

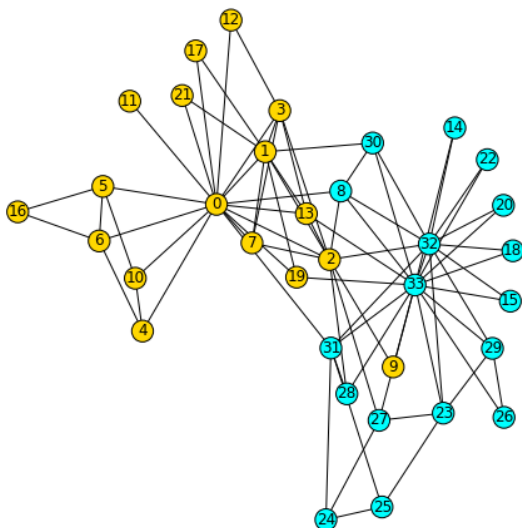
August 9, 2015

Outline

- 1 Community detection
- 2 Name dropping and information sources
- 3 Data
- 4 Conclusion
- 5 References

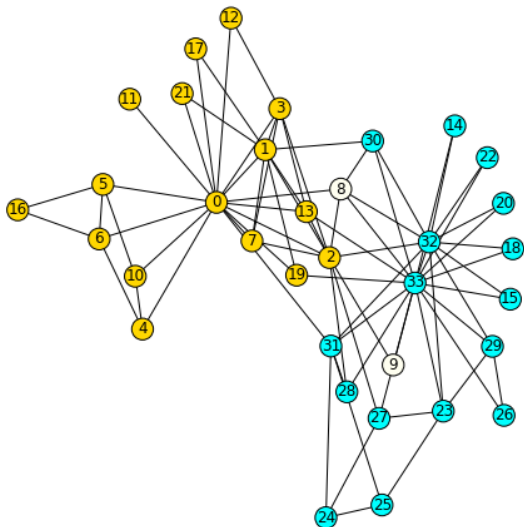
Factionalization: Karate club example (Zachary 1977)

Predict split with Newman (2006)a, Newman (2006)b



How'd we do?

Only members 9 and 8 misclassified



Why can community detection algorithms work?

- 1 Semantics of the links (Zachary: Weights reflecting number of club external activities)
- 2 Each link represents a social tie or trait or a set of social ties or a set of social traits
- 3 But community-establishing traits and ties are diverse in nature: beliefs, geographical location, family, school, dress
- 4 We proceed with a **similarity-based** approach. We investigate graphs in which links represent similarity of diverse sorts.

Expanding the sources of information

Combining ling & link info

- 1 Links between members and members, members and non members and even between non members
- 2 Language features, using a kind of graph very similar to a **shared-influence graph**
- 3 We will fold similarities in links in with language similarities, thus using a heterogeneous set of of similarity relations:
 - (a) community-internal and community-external links
 - (b) language featuresZachary used only community-internal links

Language focus

Sources

You can't measure vowel shifts with most online communities, nor are they likely to yield significant patterns with geographically dispersed groups, so we need to look for linguistic variables tied to a **discourse community** (Nystrand 1982, Swales 1990)

A key feature of many discourse communities is the existence of a shared set of **sources**, sources of beliefs and sources of practices, reflected in the names they drop and authorities they cite (both positively and negatively).

Name Dropping

Patterns of citation in Polblogs data (Adamic and Glance 2005):

Links to news media

Left

Right

Fox News

Salon

National Review

NY Times

WSJ Opinion Journal

New Republic

Washington Times

Wall Street Journal

Names of political figures

Left

Right

Donald Rumsfeld

Dan Rather

Colin Powell

Michael Moore

Zell Miller

Yassar Arafat

Tim Russert

Terry Mcauliffe

What does a link ($A \rightarrow B$) represent?

Direct Citation graph (DCG)	$A \rightarrow B$	A cites B
Co-citation	$A \leftrightarrow B$	A and B cited by C
Shared influence graph	$A \leftrightarrow B$	A and B cite C
Citer similarity graph	$A \leftrightarrow B$	A and B have similar sets of citers
Influence similarity graph	$A \leftrightarrow B$	A and B cite similar sets of influencers
Link similarity graph	$A \leftrightarrow B$	The combined sets of the citers and influencers of A and B are similar

We use the term **Source Similarity Graph** (SSG) instead of Influence Similarity Graph to accommodate both positive and negative name dropping.

Information Sources

An informal reference to an entity as a source of (mis-)information, e.g.:

An expert examination showed that the sacks contained 86.9 kilos of pure heroin, Colonel Aleksandr Kondratyev told ITAR-TASS on Saturday.

- Sources aren't necessarily other documents or people, **or even when they are, they aren't necessarily community members (out group citation)**
- Sources can be ambiguous, with references we don't know
- Source Similarity Graph allows us to draw links based on such sources (which aren't nodes in the final graph)

Assumption

Strong similarity of source citation patterns implies community affiliation.

Approximation

Extracting names

- 1 Use Stanford Named Entity Extractor (NER), see Finkel et al. (2005)
- 2 Names dropped are a feature of each user in network
- 3 Not all names are sources, so this is an approximation.
- 4 But similarity in name-usage-patterns may still be a community indicator (locations for religious or geographically based communities).

Summary

Find communities based on **feature similarity**

- 1 Name-use similarity
- 2 Community-External links
- 3 Community-Internal links

Because we use a similarity-based approach, we don't really care whether a feature refers to a community member or not (names may or may not be other community members; hyperlinks may or may not be community internal)

Puppygate

How Sci-Fis Hugo Awards Got Their Own Full-Blown Gamergate
By Katy Waldman

What on Earthsea is happening with the 2015 Hugo Awards? On Saturday, nominations for the prestigious science fiction and fantasy prizes were announced. As usual, the finalists were determined by ballot; any member of the 2014, 2015, or 2016 WorldCons . . . could vote. And yet the names and works that rose to the top provoked a tsunami of controversy. That's because a group of rightwing activists managed to game the selection process, proposing a fixed slate of nominees and feverishly promoting it. Since small margins are sufficient to secure Hugo nods, what emerged was what many are calling a strange, ideologically driven, and unrepresentative sample of fiction.

<http://www.slate.com/>

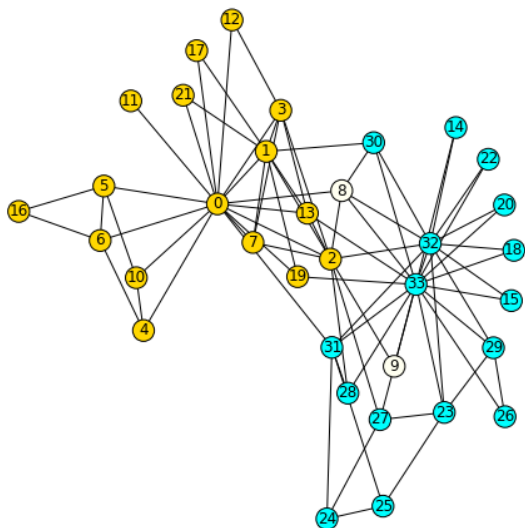
A larger conflict being played out in a smaller community

The Guardian

If all this still seems like a storm in a dragon-shaped teacup, then look at the bigger picture. This is just the latest skirmish in a culture war that has been raging, in one form or another, for at least 30 years. It was called identity politics in 1980s, political correctness gone mad in the 1990s, and Gamergate last summer. It . . . is what emboldened Nigel Farage to claim that the handpicked, balanced audience at the challengers debate on Thursday was riddled with leftwing bias.

▶ The guardian

Goal: Predict the split this community is undergoing



Observation

Faction within a small community are often reenactments of larger social conflicts

Sad puppies

*Conservative: anti-political correctness
Hurray for good old-fashioned xenophobic
bug-killing science fiction. Identify with
gamers in Gamergate*

Social justice warriors

*Liberal and politically correct
Hurray for diversity in characters, writers,
and social themes. Identify with game
critics in Gamergate*

Consequence

Because external forces play a role in the split, this is exactly the sort of situation where traits other than community internal links will come into play in predicting the split.

Procedure: Blog data

Data	Number of sites
Social Justice	102
Puppies	73
External Right	27
External Left	24
	227

Features Proper names and hyperlinks

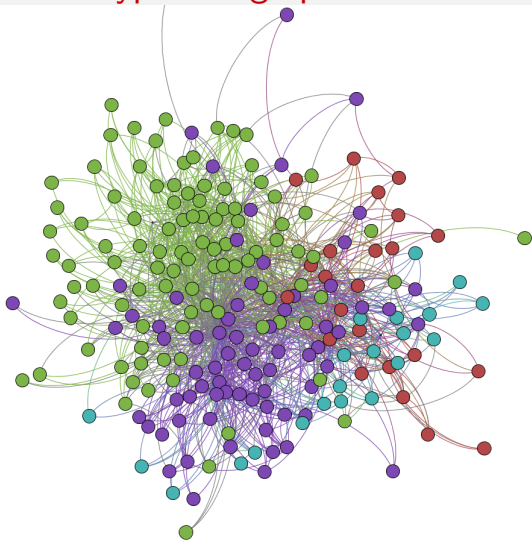
Method Unsupervised clustering

Evaluation Recovery of hand annotated puppy/SJ communities

Hypothesis A Better to use link information and linguistic features than only one of the two

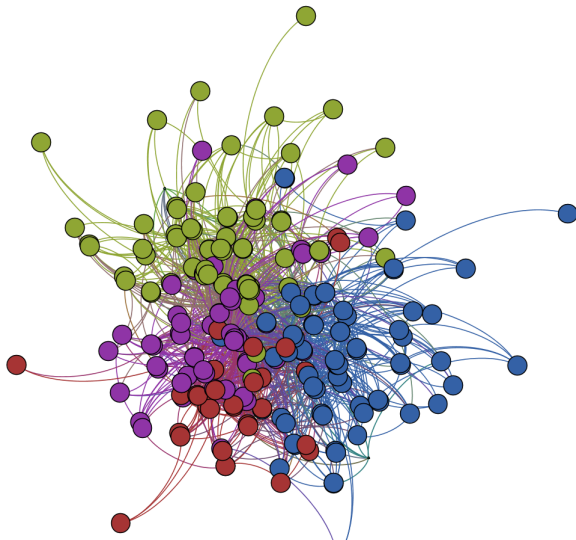
Hypothesis B Better to use external links than internal links only

The hyperlink graph: 4 web communities



- 1 Social Justice
- 2 Puppies
- 3 External right
- 4 External left

Communities a la Newman: 4 different communities



Systems

Dimensionality reduction	Run Newman's community discovery algorithm on the hyperlink graph (4 communities used as features) Run SVD on term/doc matrix
Graph-based clustering	Kmeans
Evaluation	Adjusted Mutual Information (Vinh et al. 2010)

Results

Hyps		AMI
	All sites	.128
	SF community only	.176
Names		
	All sites	.154/.060
	SF community only	.131/.013
Names + Hyps		
	All sites	.213/.011
	SF community only	.212/.046

Conclusion

- 1 Hypothesis A validated: links + ling info gave the best prediction of the community split.
- 2 Hypothesis B not validated. External links did not help, possibly because there wasn't enough structure among them, possibly because they **were** ambiguous [both sides cited slate.com article]
- 3 We have shown that using name features combined with graph-based community features helps unsupervised clustering of the Puppygate sites, supporting the intuition that names are useful community-discriminating features.
- 4 In this case, names alone actually did worse than link info alone. Speculation: This is due to data sparsity.

References I

Adamic, Lada A., and Natalie Glance. 2005.

The political blogosphere and the 2004 us election: divided they blog.
In *Proceedings of the third international workshop on Link discovery*,
36–43. ACM.

Finkel, Jenny Rose, Trond Grenager, and Christopher Manning. 2005.

Incorporating non-local information into information extraction
systems by gibbs sampling.

In *Proceedings of the 43rd Annual Meeting of the Association for
Computational Linguistics (ACL 2005)*, 363–370. Association for
Computational Linguistics.

References II

Newman, Mark EJ. 2006a.

Finding community structure in networks using the eigenvectors of matrices.

Physical review E 74(3):036104.

Newman, M.E.J. 2006b.

Modularity and community structure in networks.

Proceedings of the National Academy of Sciences 103(23):8577–8582.

Nystrand, Martin. 1982.

What writers know: The language, process, and structure of written discourse.

Academic Press New York.

References III

Swales, John. 1990.

Genre analysis: English in academic and research settings.
Cambridge University Press.

Vinh, Nguyen Xuan, Julien Epps, and James Bailey. 2010.

Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance.
The Journal of Machine Learning Research 11:2837–2854.

Zachary, W. W. 1977.

An information flow model for conflict and fission in small groups.
Journal of Anthropological Research 33:452–473.