# Using Linking Language to Discover Text Communities and Predict Group Fracture

Jean Mark Gawron
Alex Dodge
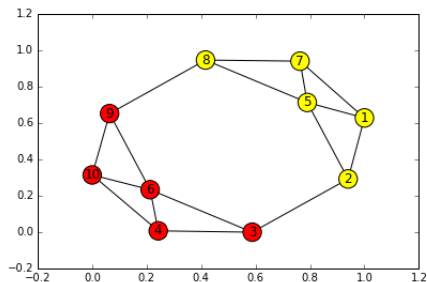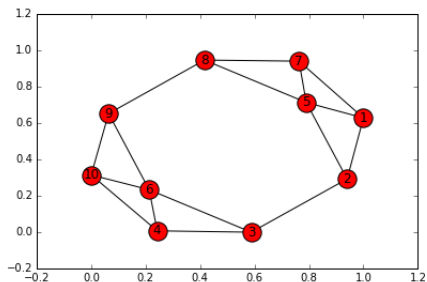
San Diego State University

August 2, 2016
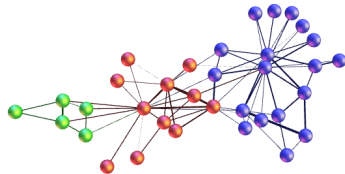
# Outline

# What is community detection?
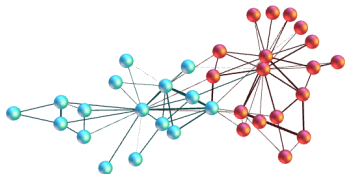
# Factionalization

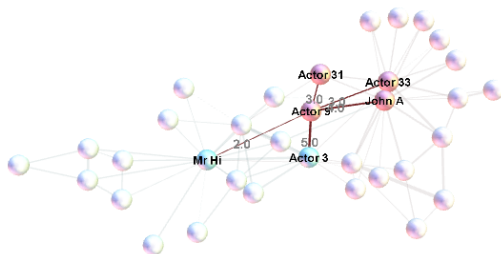Karate club example (Zachary 1977)



| Purity | 100% |
|--------|------|
| AMI    | .693 |

# The awful truth

Actor 9 misbehaved.



1. Actor 9 has strong connections to both communities (weights shown)
2. Decisive information may not be contained in the graph.

# Data

### Gamergate refresher

In August 2014, a number of posts attacking prominent women in the video gaming industry appeared in social media devoted to gaming. The women targeted included clutural critic Anita Sarkeesian (creator of the website "Feminist Frequency"), Zoe Quinn (co-creator of Depression Quest), and Brianna Wu (game developer and journalist). The attacks were inspired in part by the release of video games such as Depression Quest exploring darker real-life themes and challenging the traditional role of video-gaming as pure escapism; participants often saw themselves as reacting against a strain of "political correctness" that had "infected" the gaming world. Personal attacks followed, including "doxing" (publication of personal contact details of Quinn and Wu), accusations of trading sex for good news coverage, and death threats.

# Sad Puppies

A virtual copy of Gamergate within the **much** smaller media market of science fiction books.

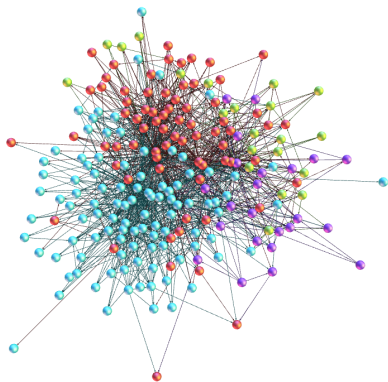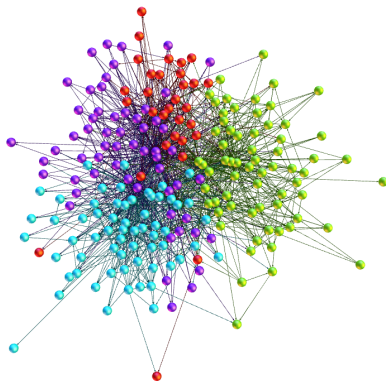| Sad puppies | Conservative: anti-political correctness Hurray for good old-fashioned xenophobic bug-killing science fiction. Identify with gamers in Gamergate |
|---|---|
| Social justice warriors | Liberal and politically correct Hurray for diversity in characters, writers, and social themes. Identify with game critics in Gamergate |

NB: $\frac{1}{3}$ of all link edges are between communities.

# A much harder problem



True                           Proposed

Purity   75%
AMI      .177

# Summarizing

1. Community detection can be genuinely revealing in social networks when the links contain the relevant social information.

2. That information may be both complex and heterogeneous (Jackson 2008).

3. Useful clustering will not be possible unless the right mix of information is represented in the links.

# A useful goal: Ways of adding new information

1. We adopt a **similarity-based** approach.

   *A general framework providing ways of combining information of diverse kinds into one graph, suitable for community discovery.*

2. Appropriate for this data: combining linguistic informatiom with link information. Information about faction membership in the language:

   *This is just one little battle in an ongoing culture war between artistic free expression and puritanical bullies who think they represent* real *fandom. (Correia 2015)*

# Combining types of similarity



Similarity graphs of (b)

(a)

(b)

co-cit

co-ref

sum

# Zachary's karate club as a similarity graph



Purity    100%
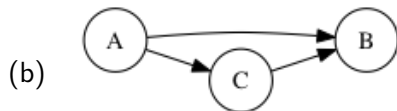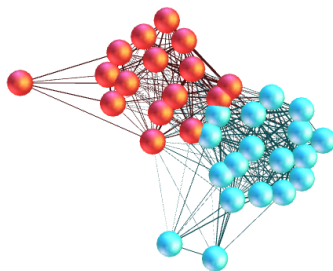AMI    1.00

# Polblogs data of Adamic and Glance



Links                                    Sim

# Polblogs numbers

|         |          | Purity | AMI    | Communities |
|---------|----------|--------|--------|-------------|
| Louvain | Link     | 94.76  | **.559** | 9           |
|         | CoCit    | 87.70  | .456   | 4           |
|         | CoRef    | 82.3   | .254   | 5           |
|         | CoCitRef | 94.76  | **.562** | 9           |
| Newman  | Link     | 95.25  | **.727** | 2           |
|         | CoCit    | 87.32  | .511   | 2           |
|         | CoRef    | 82.4   | .267   | 3           |
|         | CoCitRef | 95.17  | **.724** | 2           |

# Language focus

1. Unsupervised clustering problem: What kind of language is most likely to signal our connectedness to others?

2. Social function of language in the foreground, Ideational function takes a backseat.

# Linking Language: Bibliometric Inspiration

## Small (1973)

1. Co-citation: Two authors are co-cited when they both cited by a third. Define the co-citation **strength** of two authors/papers.
2. Co-reference: Two authors cite a third. Define co-reference strength.
3. Small builds a similarity graph, observes that these are different (only weakly correlated kinds of similarity)

# Linking language

It is not hard to imagine linguistic generalizations of both co-citation and co-reference. What might be called generalized co-citation holds between two individuals when some third individual refers to both of them in speech or in a text. What might be called generalized co-reference holds between two speakers or writers when they both refer to some third individual. We will refer to the Noun Phrases (NPs) used to make such linking references as **source NPs**.

## 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)

# Approximation: Proper names

### 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.

# Name Dropping

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

| News organizations | Right | Left |
|---|---|---|
| | Salon | Fox News |
| | NY Times | National Review |
| | New Republic | WSJ Opinion Journal |
| | Wall Street Journal | Washington Times |

| People | Right | Left |
|---|---|---|
| | Dan Rather | Donald Rumsfeld |
| | Michael Moore | Colin Powell |
| | Yassar Arafat | Zell Miller |
| | Terry Mcauliffe | Tim Russert |

# Similarity components

$$Sim(i, j) = CoCit(i,j) + CoRef(i,j) +$$
$$LingCoRef(i,j) + LingCoCit(i,j)$$
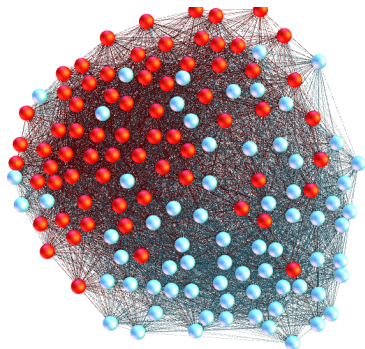
# Community Discovery System



**Sites** = input text data. **Ling**: named-entity recognition (*ner*), feature
selection (*feats*), SVD-reduction (*svd*). Community Discovery: similarity
graph construction (*sm*), **Community Discovery (cd)**.
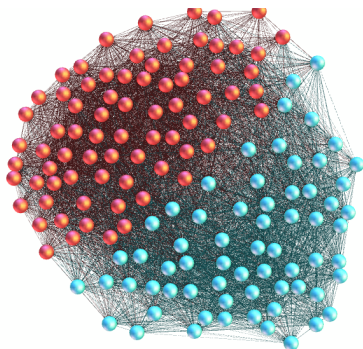
# Results

|           | AMI            |
|-----------|----------------|
| Links     | $.177 \pm 0$   |
| Ling      | $.127 \pm .023$|
| LingLinks | $.309 \pm .039$|

# Graph info still insufficient



Actual classes                    Newman communities

# Conclusion

1. Links + Ling info gave the best prediction of the community fracture.
2. Unsupervised clustering: Adding more features does **not** always improve your. The features need to be relevant to the task, and these features are.
3. The similarity graph we used did not contain the information needed to separate the factions:
   1. Co-citation (reference resolution)
   2. eXtracting all source NPs, icnluding group and anti-group references (including racial slurs)
   3. Meme recognition (*Je suis Charlie*, Je ne suis pas Charlie, je suis Ahmed).

# 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.

Correia, Larry. 2015.
> Sad puppies update: The nominees announced and why i refused my nomination.
> *monsterhunternation.com* April 4.

# References II

Finkel, Jenny Rose, Trond Grenager, and Christopher Manning. 2005.
Incorporating non-local information into information extraction
systems by gibbs sampling.
In *Proceedings of the 43nd Annual Meeting of the Association for
Computational Linguistics (ACL 2005)*, 363–370. Association for
Computational Linguistics.

Jackson, Matthew O. 2008.
*Social and economic networks*.
Princeton, NJ: Princeton University Press.

Small, Henry. 1973.
Co-citation in the scientific literature: A new measure of the
relationship between two documents.
*Journal of the American Society for information Science*
24(4):265–269.

# References III

Zachary, W. W. 1977.
An information flow model for conflict and fission in small groups.
*Journal of Anthropological Research* 33:452–473.