

Network Analysis of Collaboration Structure in Wikipedia

Ulrik Brandes Patrick Kenis
Jürgen Lerner Denise van-Raaij

WWW 2009 Madrid, 20.–24. April, 2009

Long-term goal: understanding the social collaboration process in Wikipedia.

Collaboration in Wikipedia:

- ▶ group of **voluntary actors**,
- ▶ organized on the bases of **peer governance**,
- ▶ attempt to **produce** a product (articles).

Collaboration structure is emergent (not predetermined).

What distinguishes 'good' from 'bad' collaboration?

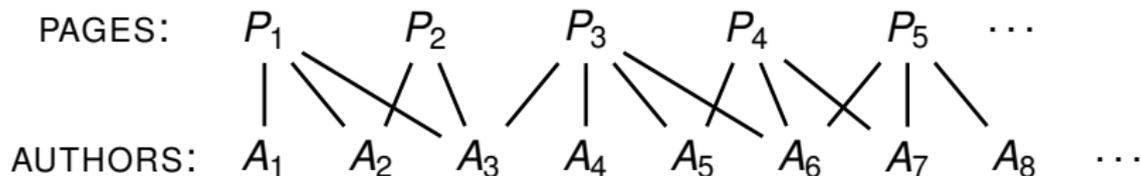
Claim: structure of collaboration network explains (to a certain extent) the quality/efficiency/effectiveness of process/product.

Collaboration networks among authors of Wikipedia.

Wikipedia's contributing users interact in several ways:

- ▶ discussion networks (from talk pages),
- ▶ leaving messages on user pages,
- ▶ blocking of users, election of administrators, ...

and **co-authoring** encyclopedic entries (content pages):



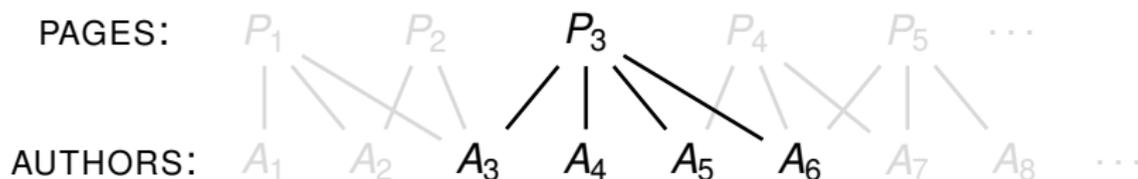
Authors of the same page are not homogeneous; their interaction reveals interesting structure \Rightarrow **edit network**.

Collaboration networks among authors of Wikipedia.

Wikipedia's contributing users interact in several ways:

- ▶ discussion networks (from talk pages),
- ▶ leaving messages on user pages,
- ▶ blocking of users, election of administrators, ...

and **co-authoring** encyclopedic entries (content pages):



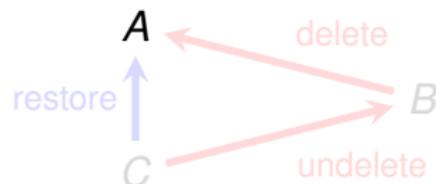
Authors of the same page are not homogeneous; their interaction reveals interesting structure \Rightarrow **edit network**.

Edit network associated with a page.

Contributing users together with time-stamped interaction:

- ▶ who **adds** which part of the text? (monadic relation)
- ▶ who **deletes** text written by whom? (dyadic relation)
- ▶ who **restores** text written/deleted by whom? (triadic rel.)

deleting and un-deleting is interpreted as **disagreement**; restoring is interpreted as **agreement**.



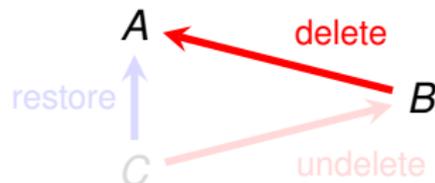
Can be computed by purely automatic means (no NLP).

Edit network associated with a page.

Contributing users together with time-stamped interaction:

- ▶ who **adds** which part of the text? (monadic relation)
- ▶ who **deletes** text written by whom? (dyadic relation)
- ▶ who **restores** text written/deleted by whom? (triadic rel.)

deleting and un-deleting is interpreted as **disagreement**;
restoring is interpreted as **agreement**.



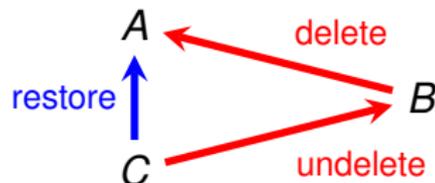
Can be computed by purely automatic means (no NLP).

Edit network associated with a page.

Contributing users together with time-stamped interaction:

- ▶ who **adds** which part of the text? (monadic relation)
- ▶ who **deletes** text written by whom? (dyadic relation)
- ▶ who **restores** text written/deleted by whom? (triadic rel.)

deleting and un-deleting is interpreted as **disagreement**;
restoring is interpreted as **agreement**.



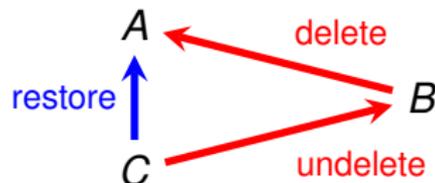
Can be computed by purely automatic means (no NLP).

Edit network associated with a page.

Contributing users together with time-stamped interaction:

- ▶ who **adds** which part of the text? (monadic relation)
- ▶ who **deletes** text written by whom? (dyadic relation)
- ▶ who **restores** text written/deleted by whom? (triadic rel.)

deleting and un-deleting is interpreted as **disagreement**; restoring is interpreted as **agreement**.



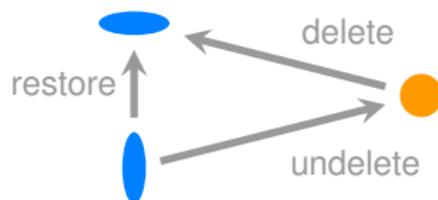
Can be computed by purely automatic means (no NLP).

Quantitative indicators of edit networks.

User **roles** and **positions**:

- ▶ providing content vs. modifying / restoring edits;
- ▶ **add/restore** vs. **delete**;
- ▶ revised  vs. revisor ;
- ▶ activity level (encoded by node size);

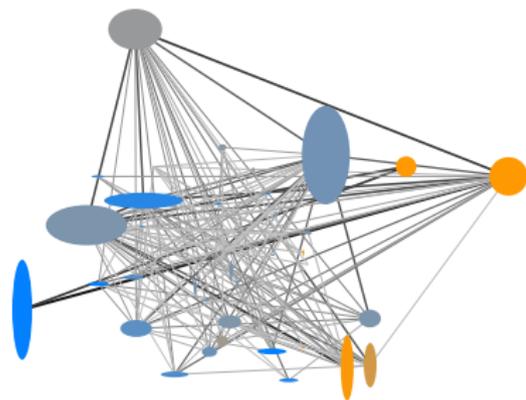
can be used to define user **reputation** etc.



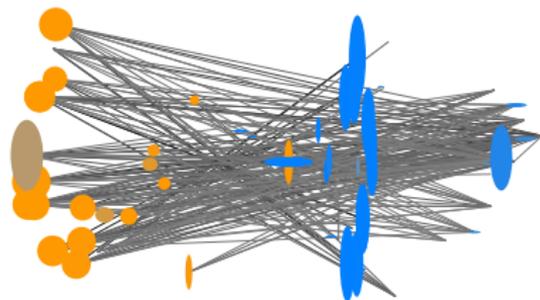
Quantitative indicators of edit networks.

Global network structure (images show only disagreement edges)

bipolarity: excess of disagreement edges between groups



no clear opposition



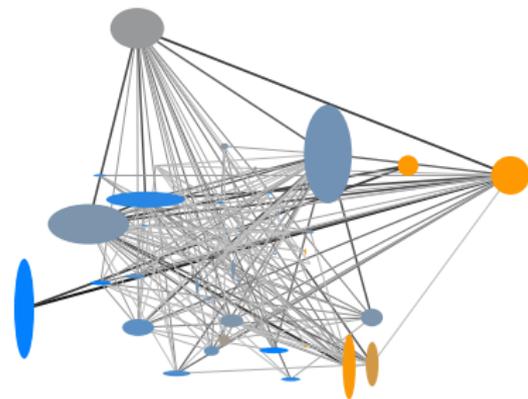
groups opposing each other

Hypothesis: bipolarity of edit networks points to controversy.

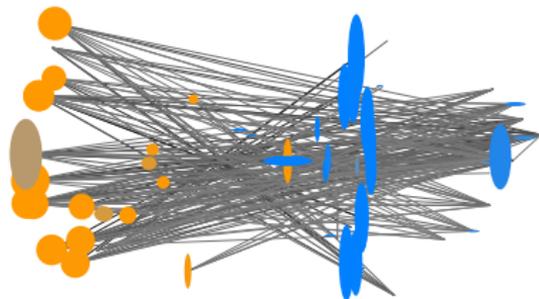
Quantitative indicators of edit networks.

Global network structure (images show only disagreement edges)

bipolarity: excess of disagreement edges between groups



no clear opposition



groups opposing each other

Hypothesis: bipolarity of edit networks points to controversy.

Statistical validation.

Hypothesis: bipolarity of edit networks points to controversy.

Select pages listed as `controversial issues`.

Use `featured` articles as control group (rather than arbitrary articles which are too small).

Randomly select 60 articles from each category.

Results: avg. bipolarity of controversial articles (0.72 ± 0.022) is significantly higher than avg. bipolarity of featured articles (0.60 ± 0.022).

Drawbacks so far.

Static analysis (ignoring time) of a dynamic network.

- ▶ which edit precedes which?
- ▶ is network structure only the *consequence* of past edits or also the *cause* of future edits?

Article labels (featured/controversial) might not be the best indicators for collaboration quality.

- ▶ not the opposite of each other; not disjoint;
- ▶ controversy is often caused by topic rather than by authors;
- ▶ controversy might be either vandalism-induced or opinion-induced.

Both issues are addressed in ongoing and future work.

Dynamic analysis of edit networks. (ongoing and future work)

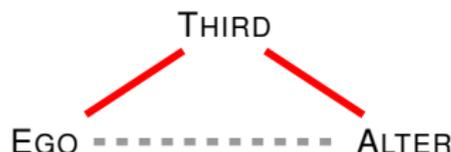
Edit networks are **event networks** (time-stamped interaction).

Model the probability and type of future events dependent on (specific aspects of) the network of past interaction.

- ▶ Does network position/role of contributors influence longevity of edits? (might point to implicit **reputation**)
- ▶ Are ties biased dependent on relations to third actors?

structural balance

(“the enemy of my enemy is my friend”)



See: Brandes, Lerner, and Snijders. “Networks Evolving Step by Step: Statistical Analysis of Dyadic Event Data,” *Proc. ASONAM 2009*, to appear.

Dropout hazard as a different aspect of collaboration quality. (ongoing and future work)

Fact: some active Wikipedians stop contributing at some time.

Question: is the collaboration network (or certain aspects of it) responsible for these dropouts?

Characteristics that increase dropout hazard can be interpreted as sources of frustration, pointing to bad collaboration.

Loss of contributors is also likely to cause stagnation in the further development of Wikipedia.

See: Brandes, Kenis, Lerner, and van Raaij. "Is Editing More Rewarding Than Discussion? A Statistical Framework to Estimate Causes of Dropout from Wikipedia," *Proc. Webcentives'09*.

Conclusion.

Edit network: interaction from page editing.

- ▶ Reveals user roles and positions
 - ▶ content providers, modifiers, watch dogs, . . .
 - ▶ author reputation (persistent edits vs. deleted edits)
- ▶ Global structure is correlated with quality labels.

Ongoing and future work

- ▶ Longitudinal analysis of edit networks (taking time-ordering of edit events into account).
- ▶ Analyze alternative outcome var's (e. g., dropout hazard).
- ▶ Model the co-evolution of different types of interaction (e. g., edit and discussion networks).