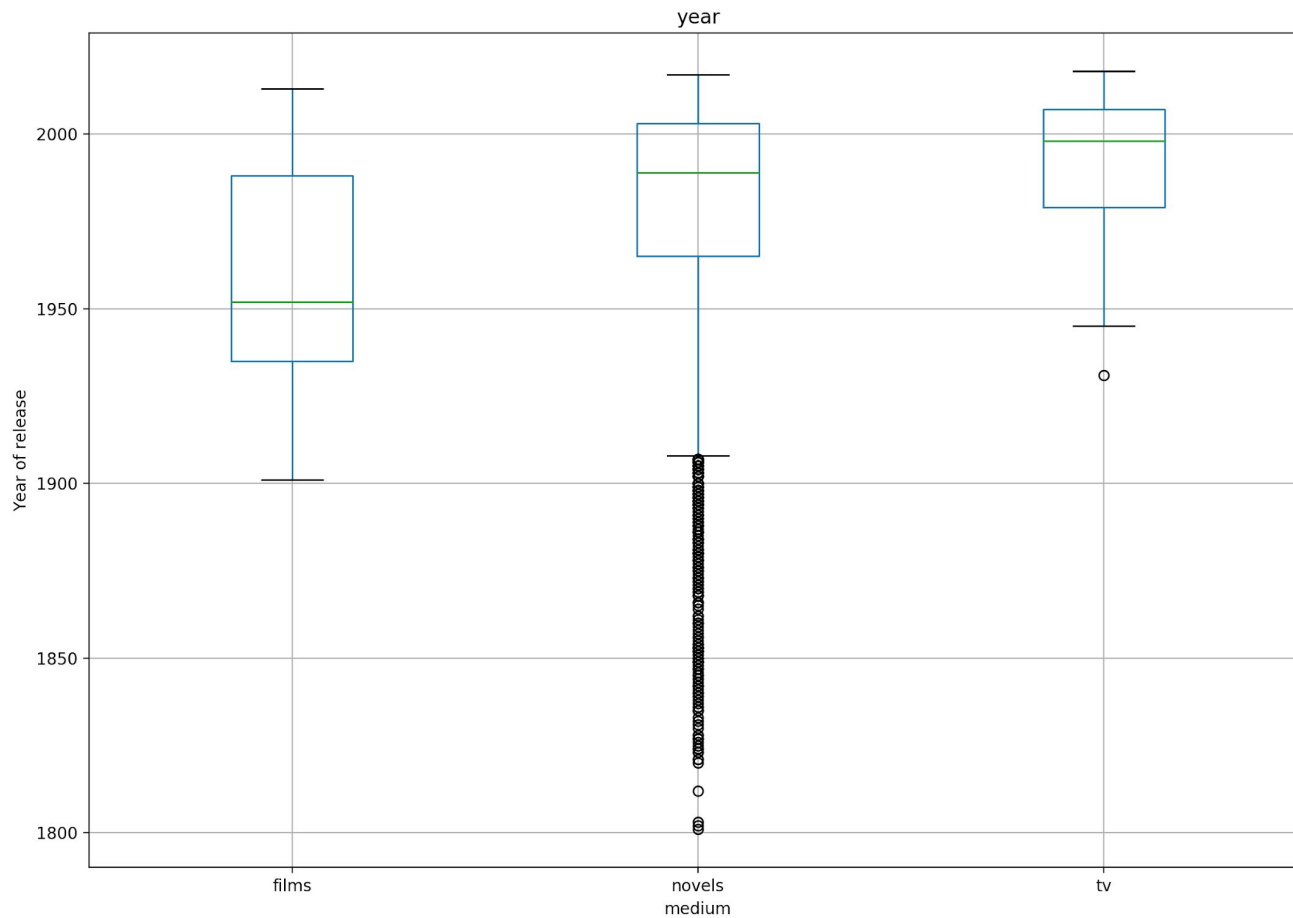


Wikipedia

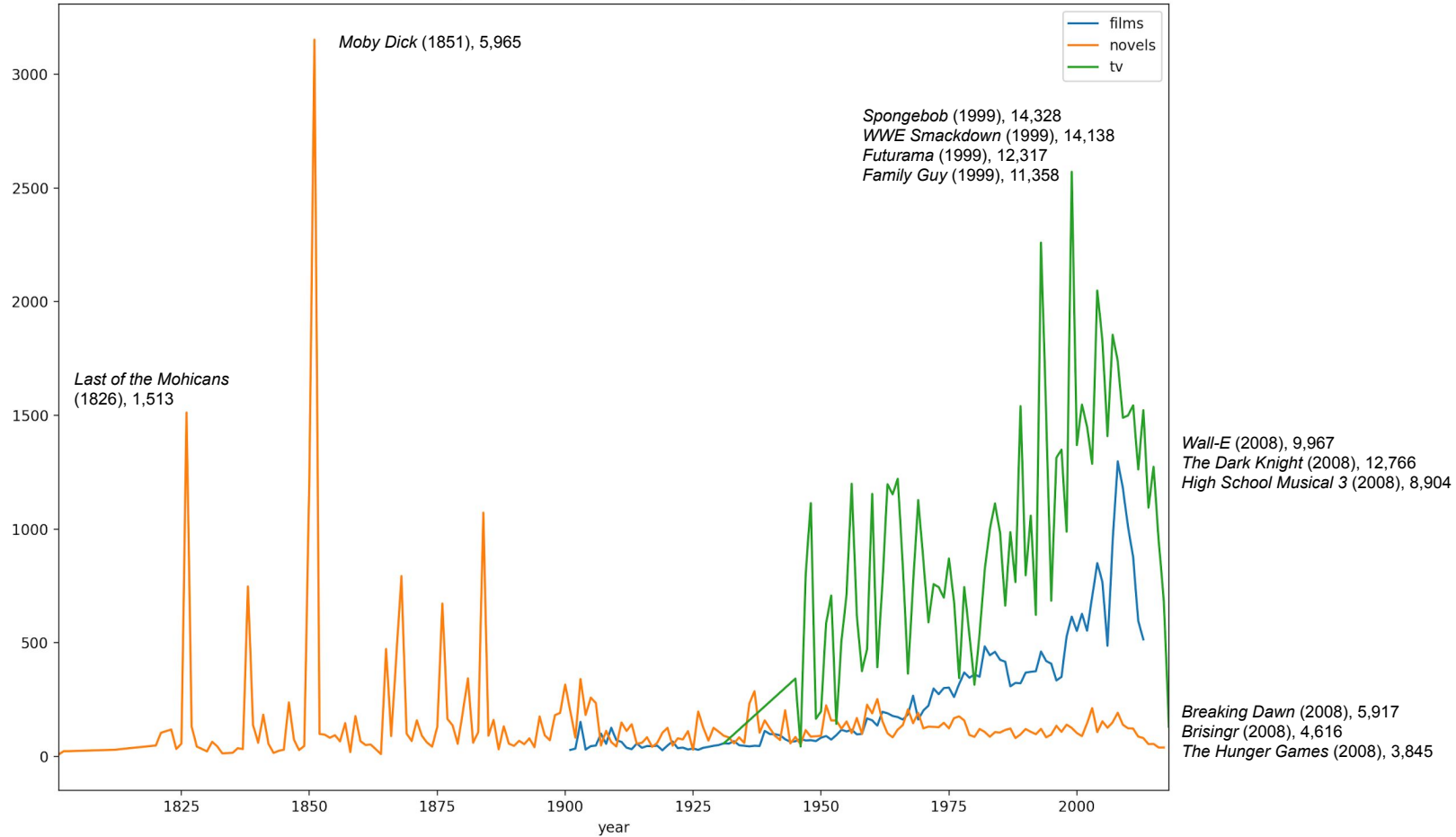
Overview of the dataset

- Complete histories of Wikipedia pages for:
 - 25,355 films (1901-2013), via en.wikipedia.org/wiki/Lists_of_American_films
 - 10,523 novels (1801-2017), via en.wikipedia.org/wiki/Category:American_novels
 - 2,324 tv shows (1933-2018), via en.wikipedia.org/wiki/List_of_American_television_programs

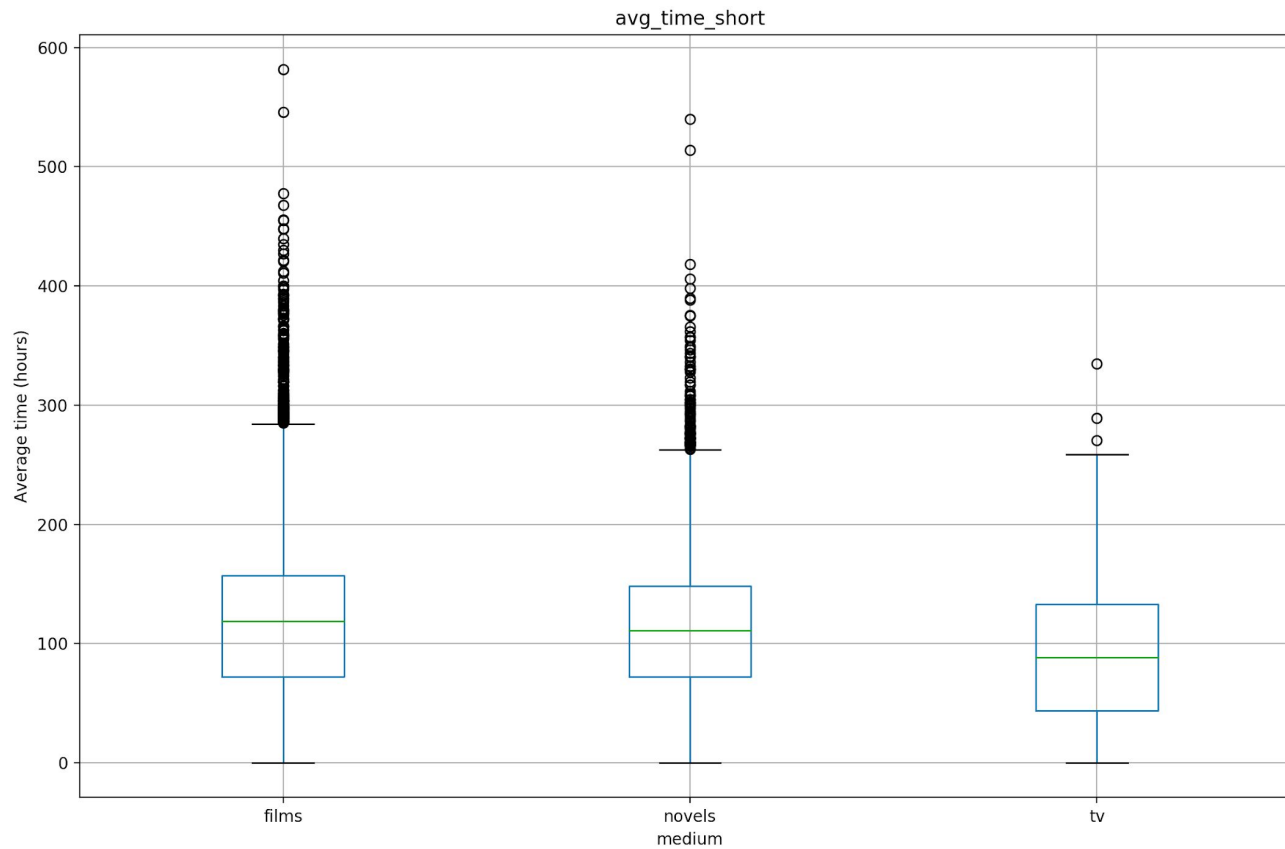
Distribution of articles by year of release



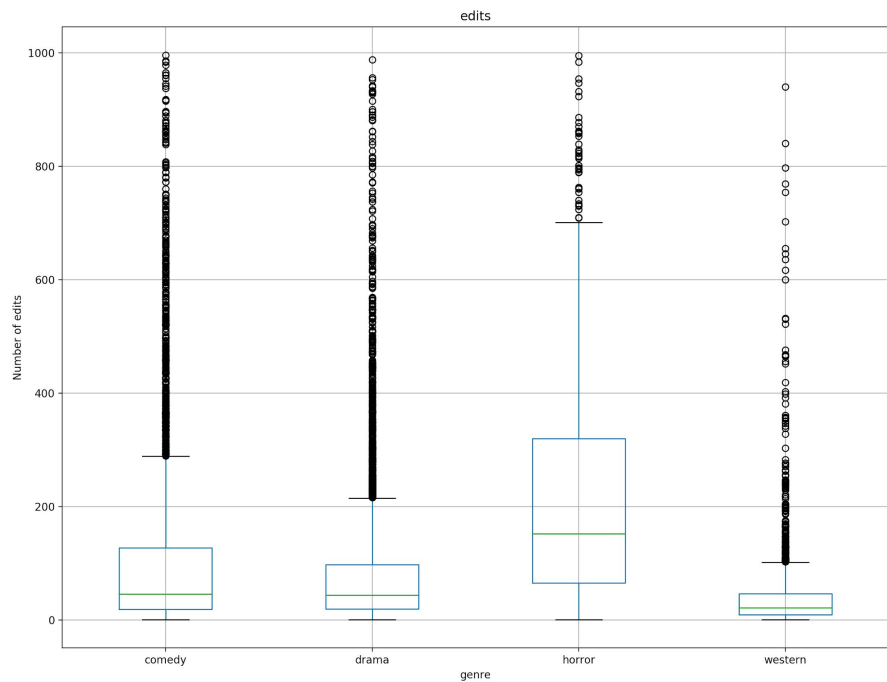
Average edits for media released in a given year



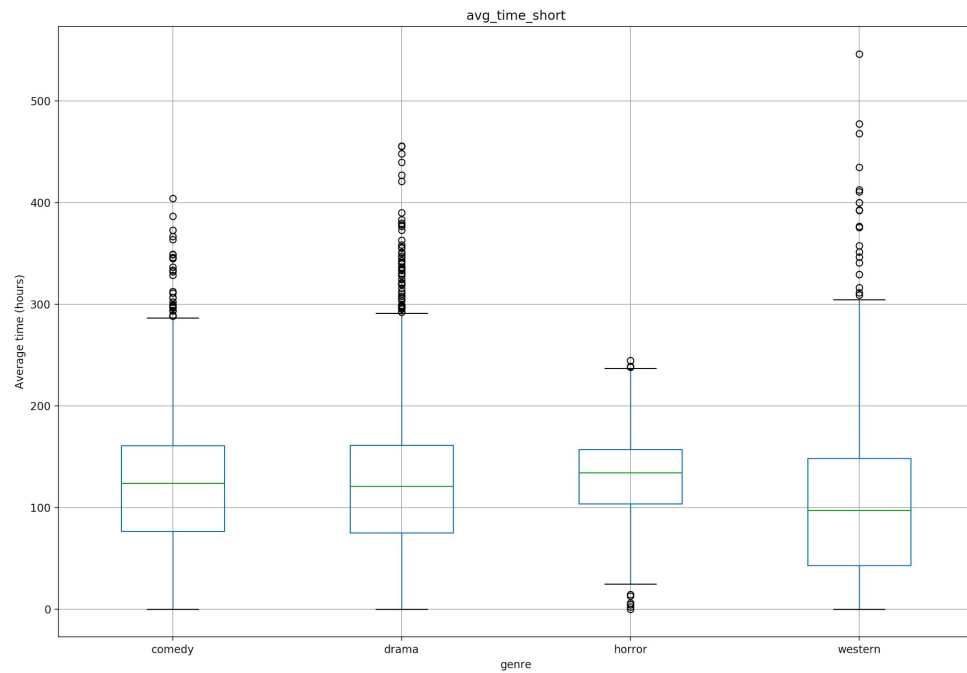
Time between edits for edits of less than one month across all mediums



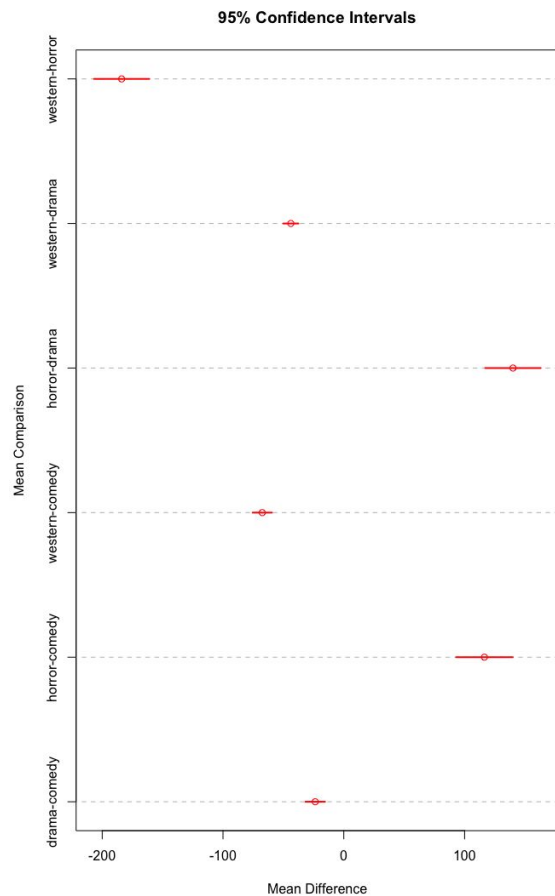
Number of edits for articles on films grouped by genre



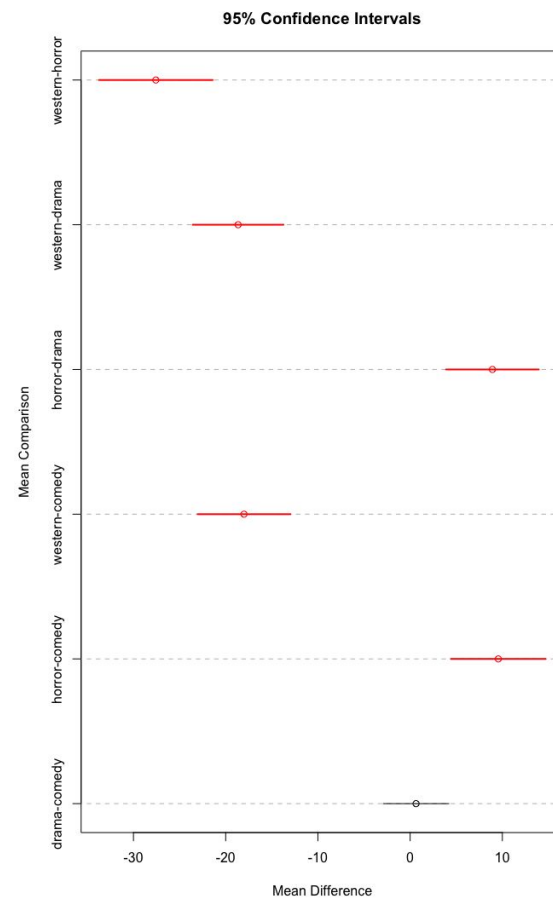
Time between edits for articles on films grouped by genre



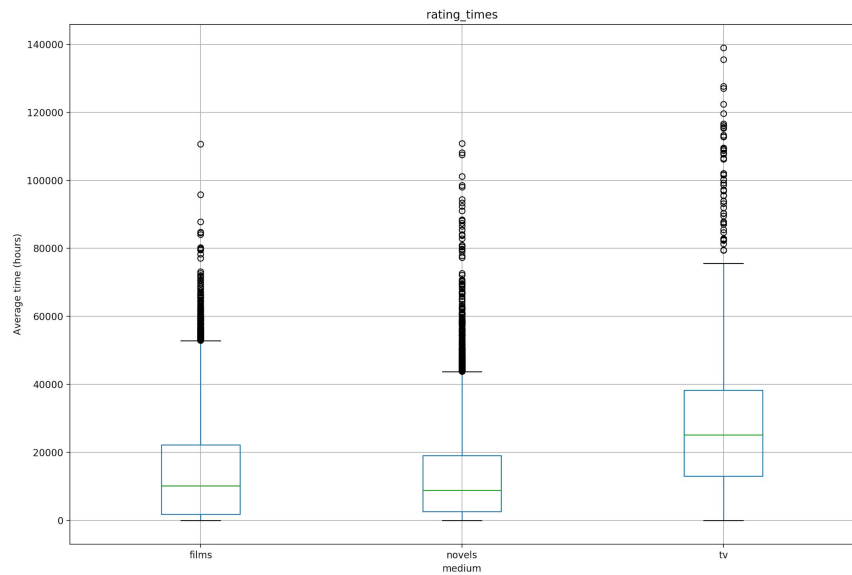
Comparison of means for number of edits by genre



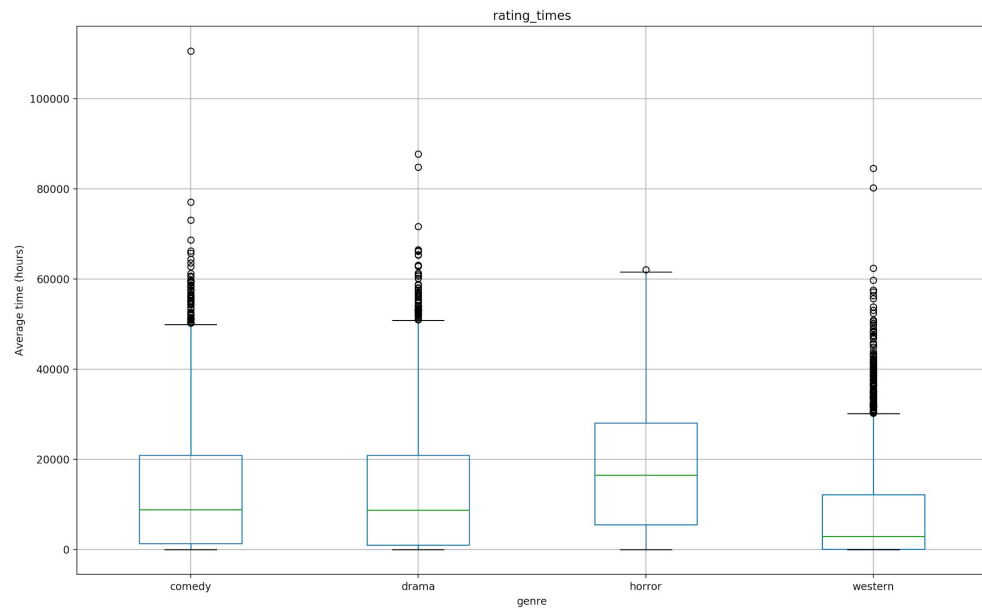
Comparison of means for time between edits by genre



Time for articles to receive a new quality rating



Time for articles to receive a new quality rating



Measuring community coherence

Louvain community detection

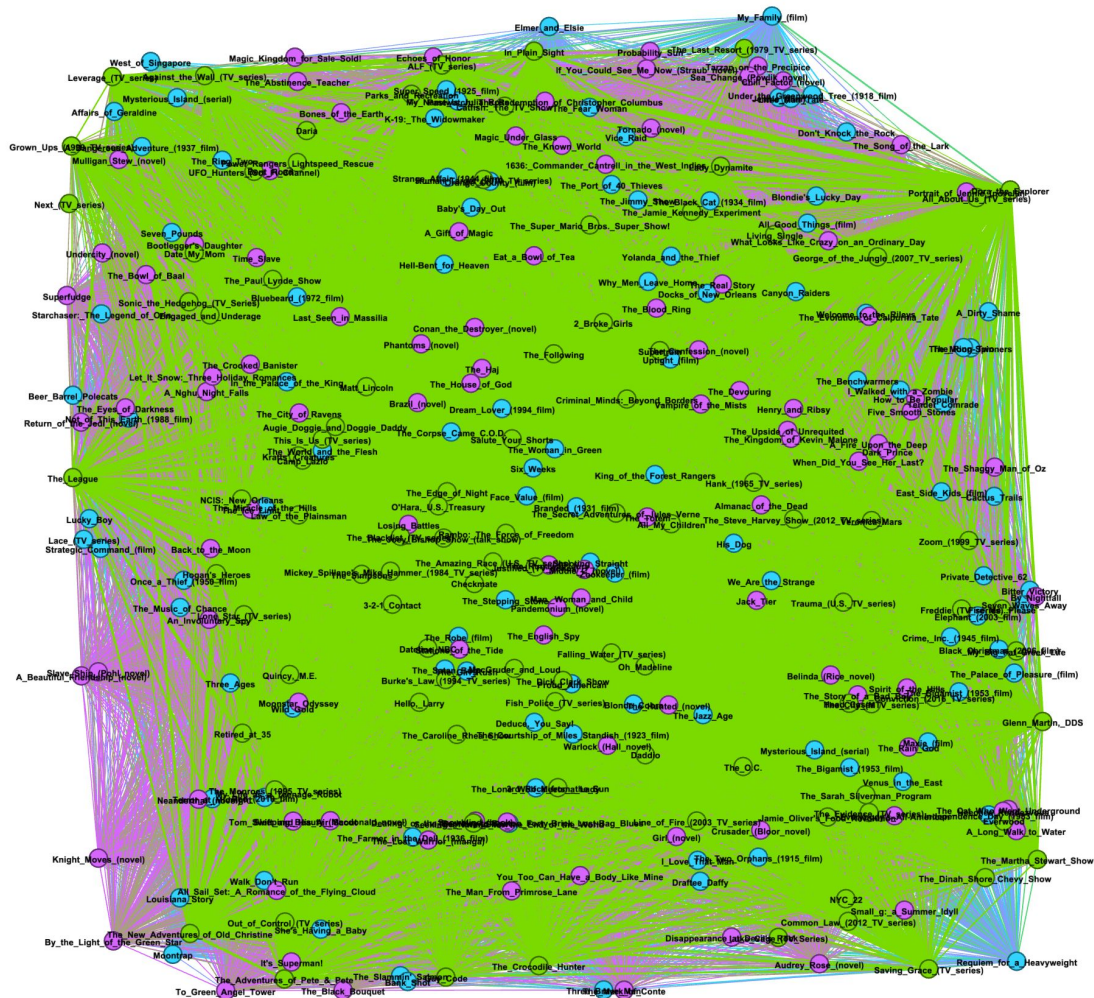
- Unsupervised
- Modeled on user behavior

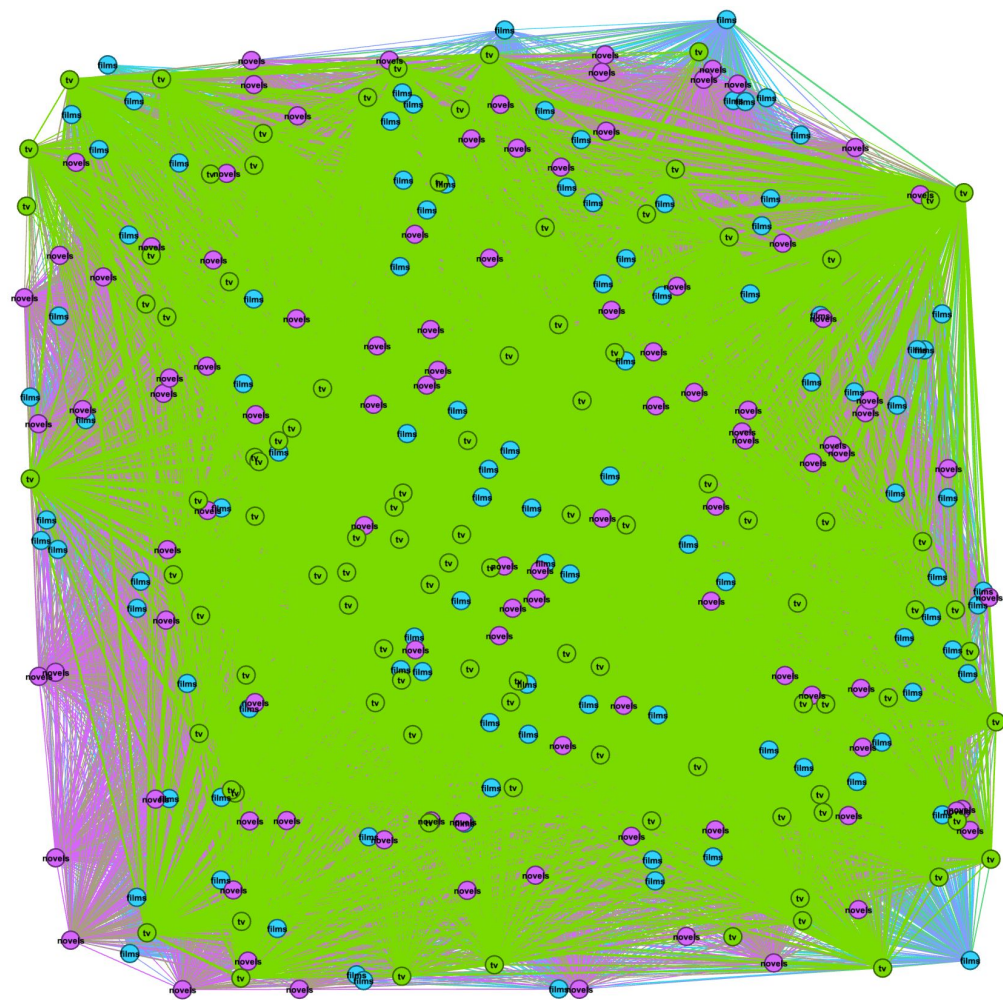
Logistic regression

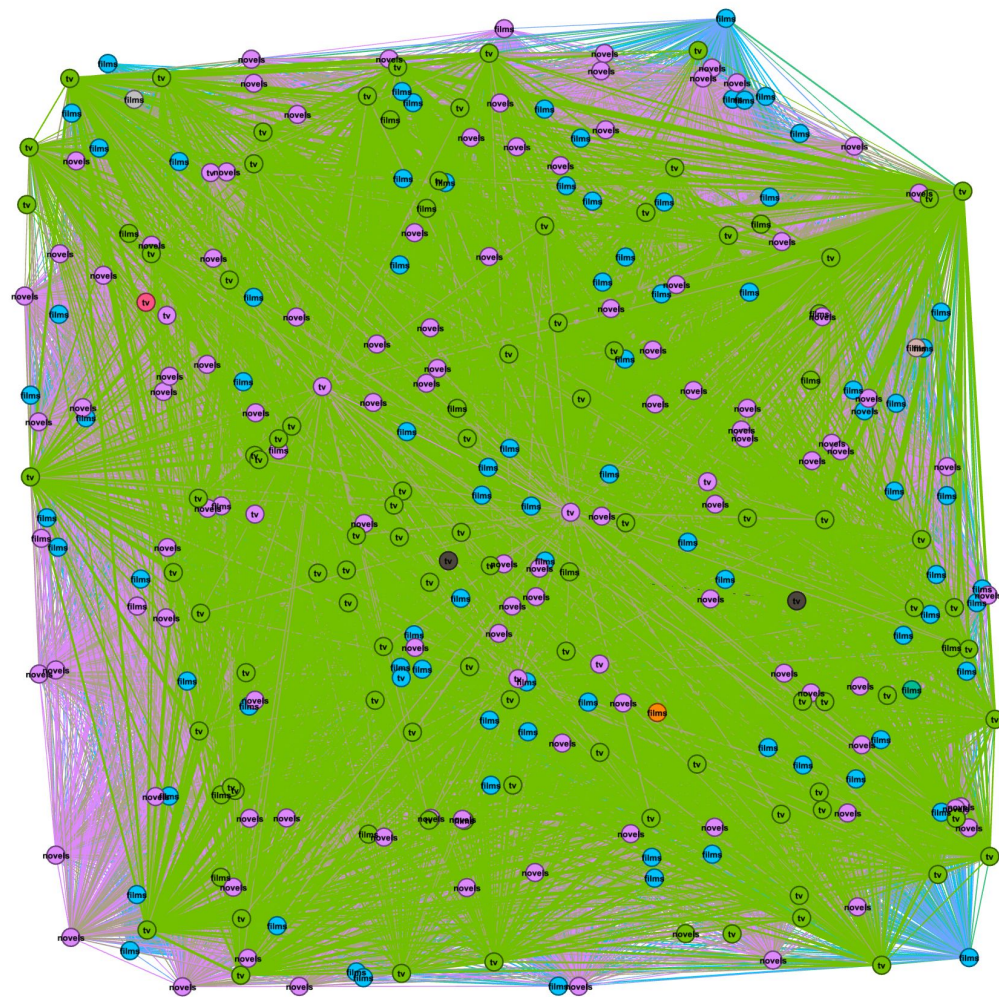
- Supervised
- Modeled on revision semantics

Louvain

- Network representation of editor behavior
 - Nodes represent articles, weighted edges represent the number of common editors between any two articles
- Structure
 - Graphs generated through random selection of articles
 - To test coherence by medium: 50 graphs, each with 100 nodes from each medium
 - To test coherence by period and genre: 150 graphs (50 representing each medium), each with 300 nodes from the same medium

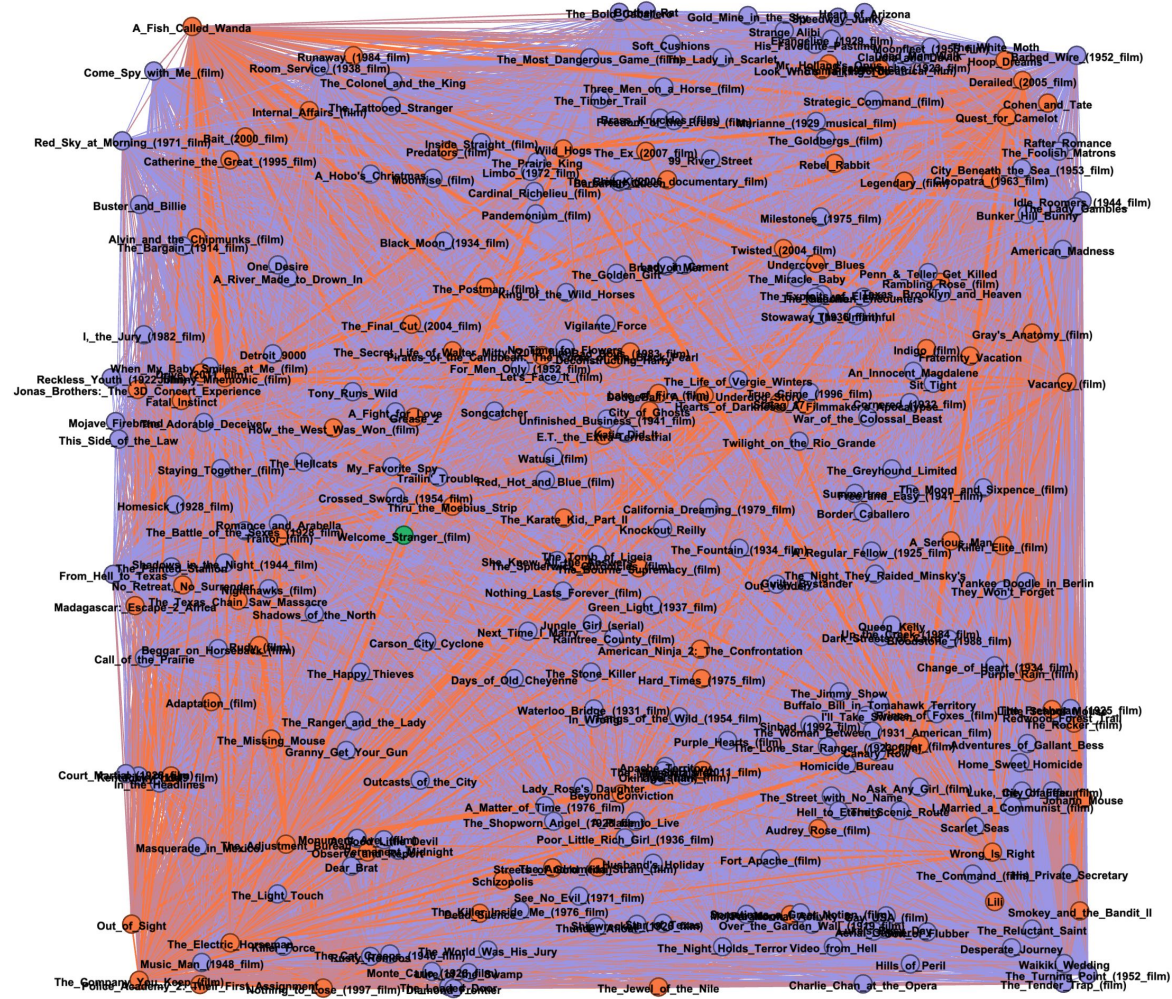






Louvain by medium

- Only clusters containing >50 nodes were considered for analysis
 - 47.8% of the total clusters across all graphs.
 - 3.02 large clusters per graph
- Average purity of 92.5% with 5.3% standard deviation
 - 33.1% of all clusters were majority novels
 - 33.1% of all clusters were majority films
 - 33.8% of all clusters were majority television shows
- User behavior is extremely confined by medium



Louvain by period

- Films

- Average purity of 78.6% (9.2% stdev)
- 50% were majority early period (71% average purity), 50% were majority late period (85% average purity)
- Strong division based on period: Editors divide equally as either early- or late-period and stay more strictly confined to editing media from their period of choice

- TV

- Average purity of 70.8% (11.2% stdev)
- 60.7% were majority early period (70.3% average purity), 39.3% were majority late period (71.6% average purity)
- Division based on period: Editors cluster more strongly around early period and tend to edit only either early- or late-period TV shows

Louvain by period

- Novels
 - Average purity of 57.9% (6.4% stdev)
 - 72.9% were majority early period (59.2% average purity), 27.1% were majority late period (54.4% average purity)
 - Little to no division based on period: Editors contribute equally to early- and late-period articles, with slightly stronger clustering around early period novels

Louvain by genre

- Clustering only using articles tagged as either dramas or comedies
- Films
 - 57.5% purity
 - 90.0% majority drama
 - This split matches the division between comedies and dramas in the dataset: 55.8% of films in the set are tagged by genre. This suggests little to no division of editor behavior by genre for films.
- Television
 - 76.4% purity
 - 100% majority comedy
 - This split also matches the division between comedies and dramas; 77.3% of TV articles in this set are comedies. Again, little to no division of editor behavior by genre.

Logistic regression: Accuracy

- Pairwise medium comparisons
 - 93.3% accurate between novels and films
 - 95.2% accurate between novels and TV
 - 96.6% accurate between TV and film
- Pairwise period comparisons (between sets split by the median release year)
 - 86.0% accurate for films
 - 82.0% accurate for television
 - 76.3% accurate for novels
- Pairwise genre comparisons (between drama and comedy)
 - 88.9% accurate for television
 - 88.4% accurate for films

Distinctive features by medium

Medium specific
Genre specific
Of note

- As the intersection of each set of pairwise comparisons
 - Television: lasting, **cancellation**, rancher, **cancels**, **premiering**, decadence, and resale
 - Novels: **publishers**, **publishing**, paperboy, and novelty
 - Films: **silents**, **stars**, **directing**, **mayhem**, and **westerners**

Between each medium

- Novels compared to television
 - TV: **stars**, liveliest, direly, completely, productively, hostess, clearance, variously, scheduling, **starry**, and **producers**
 - Novels: **authoress**, **spoiling**, **novel-ettes**, trimester, **adapting**, thrilling, **booksellers**, deli, harebrained, **suicides**, **presser**, **bookcase**, considerably, and message
- Television compared to films
 - TV: **reruns**, site, revolvers, dirge, thirties, lift, **stubble**, amazons, boldest, focusing, and **syndicating**
 - Films: bangs, festive, **spoiling**, rotter, tabor, **republicanism**, weaned, schemer, **romanticism**, permeate, **sales**, **picturing**, and **screen**
- Novels compared to films
 - Novels: **finality**, repudiate, harden, **serials**, **reprints**, beta, **bestselling**, locusts, **illustrating**, torched, prosecute, reciprocal, **novelization**, omnipotence, nebulous, and **sordidness**
 - Films: forcing, **archives**, temples, billet, monograph, **melodramatic**, **dramatic**, beeswax, distributes, digresses, comer, uncritically, **documented**, and tomboy

Distinctive features by period

- Novels (1989)

- Early period: pseudonyms, **serials**, **illustrating**, **adapting**, originate, **reprints**, spying, **editor**, swiftness, rulers, **westerners**, drugstores, repudiate, chickens, **publishers**, coping, scientific, stover, femme, and **adventures**
- Late period: ween, scheduling, announcements, **fantasies**, positives, torched turvy, harden, ailed, officials, entertains, graphically, **wed**, **writhe**, virtue, cabaret, synod, measly, **adulterer**, and homeland

- Films (1950)

- Early period: **silents**, preserves, **profitable**, domesticated, **melodramatic**, beeswax, censured, twerp, survivor, codebooks, remainder, overheated, reelection, fecal, dome, smilingly, **corporatism**, crisply, congressman, and printed
- Late period: **engulfing**, **devouring**, externally, theatrics, festive, bulletproof, spoiling, **erotically**, globetrotters, thrilling, furbelows, **prostitutes**, rocked, **tomboy**, cultists, heisting, rotter

Distinctive features by period

- Television
 - Early period: runaways, **starry**, **releasing**, regularly, including, syndrome, **programing**, appears, chill, **releases**, yeast, hag, **broadcasters**, dusk, answer, beck, **policemen**, **fatherhood**, seeped
 - Late period: **premieres**, **realizations**, millions, fond, channeling, beep, follows, relative, comas, seconded, whim, **defends**, **killed**, anyhow, **premiering**, willful, islanders, lookalikes, and beached

Distinctive features by genre

- Films

- Dramas: **dramatic**, **melodramatic**, **dramatics**, reigning, **thrilling**, **feelings**, informative, novelist, prisoners, activate, formidable, **lovesick**, gibbon, holes, possum, pennies, wardens, understandably, someplace, and wellbeing
- Comedies: comer, **tickle**, gawker, **sequence**, functional, nip, canes, superb, originally, harebrained, opposed, reinforced, watery, **heroic**, **manslaughter**, **serious**, doves, thirds, turning, and convincingly

- TV

- Dramas: **dramatic**, **produces**, **seriously**, italics, historians, offbeat, newscaster, outraged, inaccuracy, part, laurel, teaming
- Comedies: site, comers, **upbringing**, **filming**, **youngest**, players, allegations, hammering, west, showcase, songwriter, **episodic**, brace, one, frontier, hind, washed, sheen, headache, and them

Further work

- Tracing distinctive features
 - How does usage of distinctive features for genres and mediums change over time?
- Topic modeling
 - Topic model the history of each document, compare topics for high-frequency and low-frequency pages to understand differences between high activity and low activity pages
- Intra- and extra-Wikipedia comparisons
 - How do editing habits on Wikipedia compare to academic writing in journals?
 - How do other social domains on Wikipedia compare to cultural domains?