

---

---

# Caring on Reddit: The Linguistic Lives of Loyalists and Vagrants

Justine Zhang  
Tianze Shi  
Skyler Seto

---

---

# Motivation

How invested are users in online communities? It depends! (Gilbert&Karahalios)

Some users are more invested than others **on reddit**.

Subreddit	% of Comments by Top 10% of users
gameofthrones	54
starcraft	66
PoliticalDiscussion	70

In a platform with multiple subcommunities, users take the opportunity to **explore**. (Tan&Lee)

# Intuitions

For a community, define two types of users...

**Loyalists**: users who spend a sizeable portion of their online life in a particular community.

**Vagrants**: users who spend most of their time in other communities, but occasionally visit.

***Why might a user belong to one particular class?***

*interest, personal investment, expertise, ...*

# Intuitions

Hence, our question:

***How do loyalists and vagrants interact within an online community?***

→ are there **linguistic differences?**

# Dataset



Reddit, duh.

Some nice things about Reddit:

- We have a lot of **data**.
- There are **multiple communities**, and users can move freely between them.
- Lots of interest-based communities which possibly **inspire different tiers of interest/investment/expertise**.

Some other things about Reddit:

- Communities are about **diverse range of things** (spoiler alert: tricky)
- We have **a lot of** data (and only 1 TB of RAM).

# Dataset (What we've looked at)

## 10 subreddits

starcraft, Android, Music, gaming, PoliticalDiscussion, gameofthrones, nfl, CFB, books, MakeupAddiction

## 5 (ish) years

2009-2013 and/or subsets thereof, **separated by months**, contingent on data size

*(... but it would be nice to look at everything)*

# Definitions

For **each month**, take all users who contribute  $> 15$  comments. (Never look at the rest again.)

For **each month and subreddit**, compute a user's activity fraction in that subreddit:

*# **unique links** commented on in subreddit / # unique links commented on in all of reddit*

**Loyalists:** Top 10th Percentile

**Vagrants:** Bottom 90th Percentile

*footnote: tweaking the definitions probably also produces a legitimate, possibly better, dataset.*

# Definitions - Depth is Confounding

Preliminary analysis:

1. Loyalists and vagrants comment at different depths.
2. Shallow comments (**level**  $\leq$  1) are linguistically different from deep ones.

Point 2 is super interesting but not our problem. Hence we split our analysis for shallow and deep comments and mostly focused on the shallow ones.

*footnote: looking at deep comments gets into the conversational behavior of users, which is interesting and hard.*



# Preliminary Statistics

Distribution of Loyalist vs Vagrant comments, numbers of Loyalist/Vagrant users

Average # of comments/unique links per loyalist, or per vagrant user.

The story: there are numerical differences, so are there content ones?

# NLP is just counting words

Hypothesized symptom of loyalty: **putting effort into posts** .

*loyalists should write longer posts than vagrants (despite writing more!)*

	Android	books	CFB	GOT	Makeup Addiction	Music	NFL	Political Discussion	Star craft
WC	<b>80, 0.635</b>	<b>89, 0.618</b>	<b>32, 0.667</b>	<b>38, 0.633</b>	<b>29, 0.604</b>	<b>59, 0.447</b>	<b>56, 0.712</b>	<b>31, 0.646</b>	<b>18, 0.75</b>
SC	<b>39, 0.619</b>	<b>42, 0.568</b>	<b>4, 0.167</b>	<b>0.0, 17, 0.567</b>	<b>2, 0.083</b>	<b>47, 0.712</b>	<b>19, 0.487</b>	<b>12, 0.5</b>	<b>4, 0.33</b>

format: uhhhhhh

# NLP is just throwing LIWC at things

Hypothesized symptom of loyalty: **personal investment** ...maybe.

Us:

*loyalists should use more personal/concrete pronouns than vagrants.*

Deja Review paper:

*vagrants are more likely to contribute personal gut reactions than loyalists ("experts")*

# NLP is just throwing LIWC at things

Hypothesized symptom of loyalty: **personal investment** ...maybe.

Results:

	Android	books	CFB	GOT	Makeup Addiction	Music	NFL	Political Discussion	Star craft
ipro	<b>10<sup>-23</sup> 95 0.754</b>	<b>10<sup>-9</sup> 92 0.639</b>	<b>0.01 42 0.875</b>	<b>10<sup>-9</sup> 51 0.85</b>	0.07 21 0.4375	<b>10<sup>-31</sup> 93 0.7045 454545 45</b>	<b>10<sup>-8</sup> 53 0.67948 717948 7</b>	<b>10<sup>-5</sup> 41 0.854166 666667</b>	<b>10<sup>-7</sup> 22 0.91666 666666 7</b>

# LIWC is pretty fun though

Other stuff:

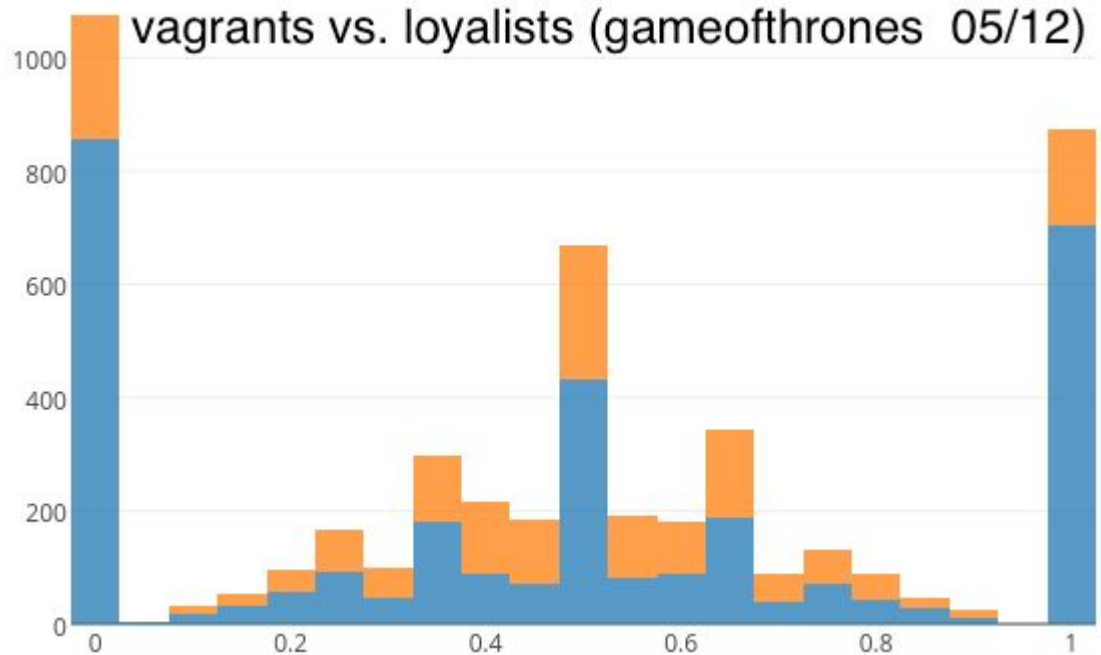
	Android	books	CFB	GOT	Makeup Addiction	Music	NFL	Political Discussion	Star craft
SW	<b>10<sup>-293</sup></b> <b>16</b> <b>0.32</b>	<b>10<sup>-26</sup></b> <b>46</b> <b>0.32</b>	<b>10<sup>-21</sup></b> <b>14</b> <b>0.29</b>	<b>10<sup>-32</sup></b> <b>13</b> <b>0.22</b>	<b>10<sup>-16</sup></b> <b>4</b> <b>0.08</b>	<b>10<sup>-18</sup></b> <b>9</b> <b>73</b> <b>0.56</b>	<b>10<sup>-8</sup></b> <b>32</b> <b>0.41</b>	<b>10<sup>-123</sup></b> <b>21</b> <b>0.44</b>	<b>10<sup>-11</sup></b> <b>0</b> <b>0</b>
Qs	<b>0.0</b> <b>40</b> <b>0.32</b>	<b>0.0</b> <b>0</b> <b>0</b>	<b>0.0</b> <b>2</b> <b>0.04</b>	<b>0.0</b> <b>2</b> <b>0.03</b>	<b>0.0</b> <b>2</b> <b>0.04</b>	<b>0.0</b> <b>27</b> <b>0.21</b>	<b>0.0</b> <b>0</b> <b>0</b>	<b>0.0</b> <b>21</b> <b>.438</b>	<b>10<sup>-238</sup></b> <b>0</b> <b>0</b>

# Is loyalty rewarded?

Hypothesis:

*communities prefer loyalists.*

Results:



# The consistency of loyalty

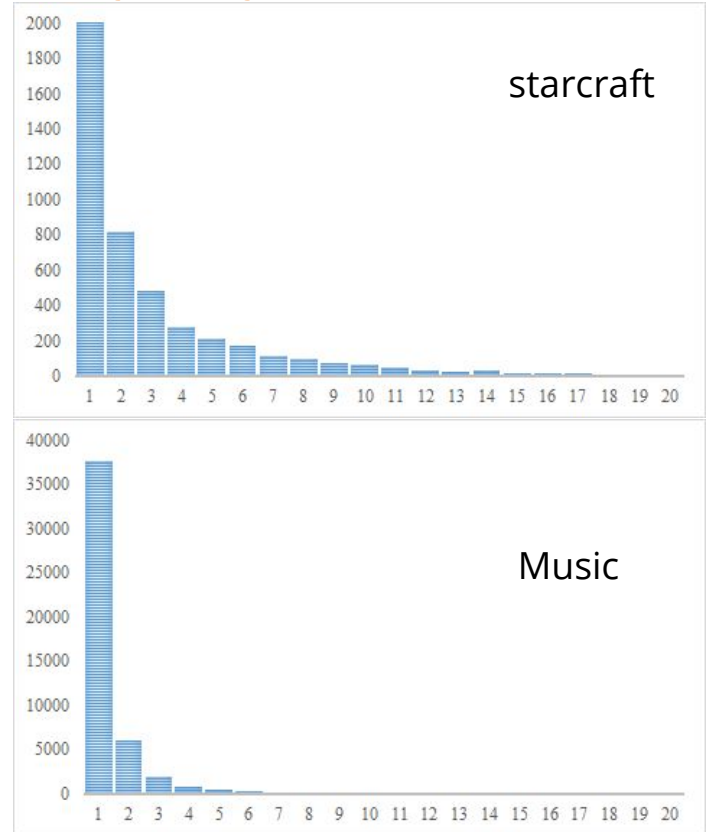
Jaccard similarity of loyalist sets between months.

	1 month	2 months	3 months		1 month	2 months	3 months
Android	0.3260	0.2575	0.2252	MakeupAddiction	0.3032	0.2026	0.1508
books	0.1546	0.1105	0.0904	Music	0.0732	0.0569	0.0480
CFB	0.3861	0.2896	0.2315	nfl	0.3764	0.2873	0.2394
gameofthrones	0.2045	0.1243	0.0883	PoliticalDiscussion	0.2979	0.2219	0.1792
gaming	0.2242	0.1684	0.1405	starcraft	0.3665	0.2775	0.2307

# The consistency of loyalty

Average number of “loyal” months

Android	3.007	MakeupAddiction	2.289
books	1.659	Music	1.415
CFB	3.341	nfl	3.326
gameofthrones	1.892	PoliticalDiscussion	2.587
gaming	2.239	starcraft	3.135





# Language Models

How vagrants and loyalists differ in language use?

Train language models and see how well they can predict the language!

# Language Models

“Training set” (set of comments used to train language models):

For each (subreddit, month), randomly choose 2,000 comments

Vocabulary: top 1,000 words for all time, others mapped to <RARE>

Predict other comments

Bunches of language models:

All comments; vagrant comments; loyalist comments;

shallow comments; deep comments

# Cross Entropy

Use the trained language models to predict other comments

Control for subreddit: LMs are subreddit-specific

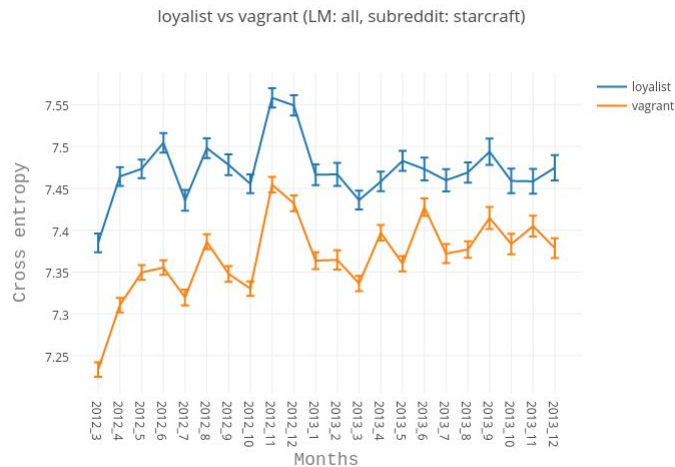
Control for time: LMs are month-specific

Control for length: take first 30 words of each comment

Control for type of users

Control for depth of comments

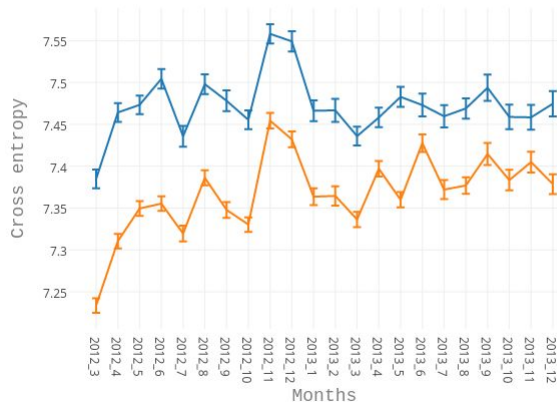
# Cross Entropy - Vagrants vs. Loyalists



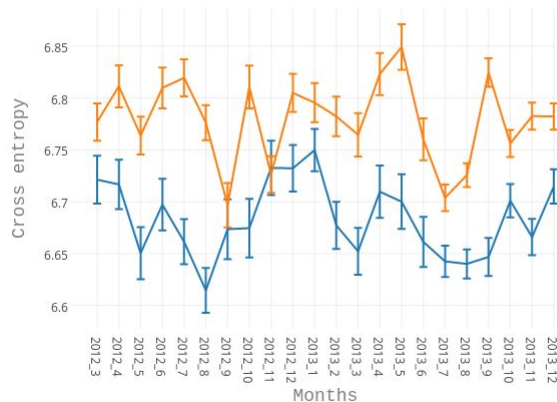
# Cross Entropy - Vagrants vs. Loyalists

But, for different subreddits?

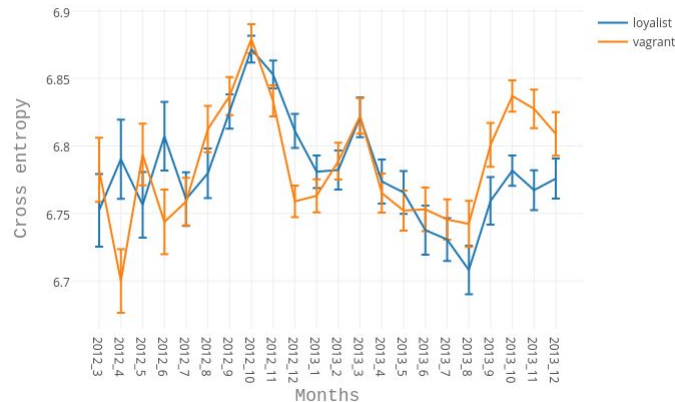
loyalist vs vagrant (LM: all, subreddit: starcraft)



loyalist vs vagrant (LM: all, subreddit: books)



loyalist vs vagrant (LM: all, subreddit: PoliticalDiscussion)



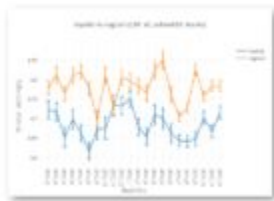
# Cross Entropy - Vagrants vs. Loyalists

Relation to consistency of royalty (avg. number of loyal months)



Android

3.007



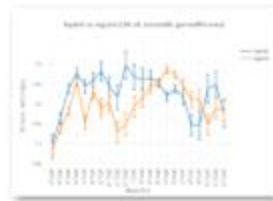
books

1.659



CFB

3.341



gameofthrones

1.892



gaming

2.239



MakeupAddiction

2.289



Music

1.415



nfl

3.326



PoliticalDiscussion

2.587



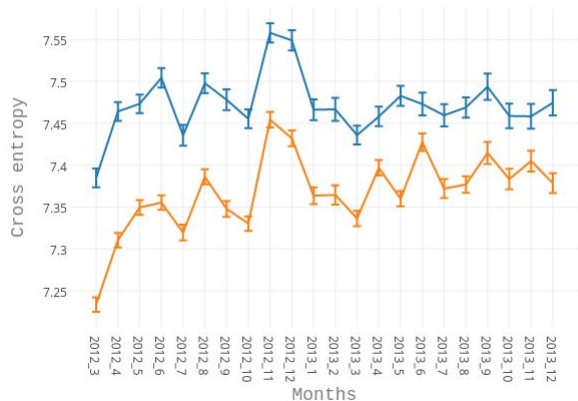
starcraft

3.135

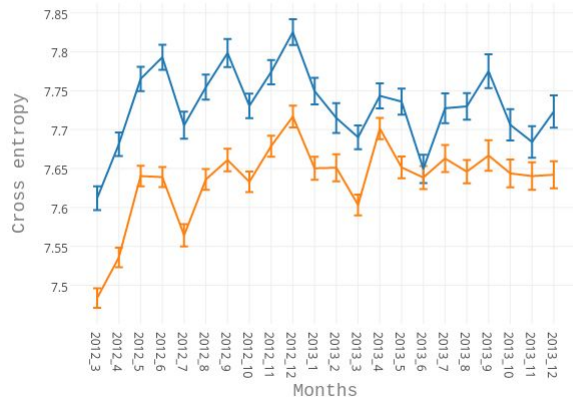
# Cross Entropy - Vagrants vs. Loyalists

At different depth

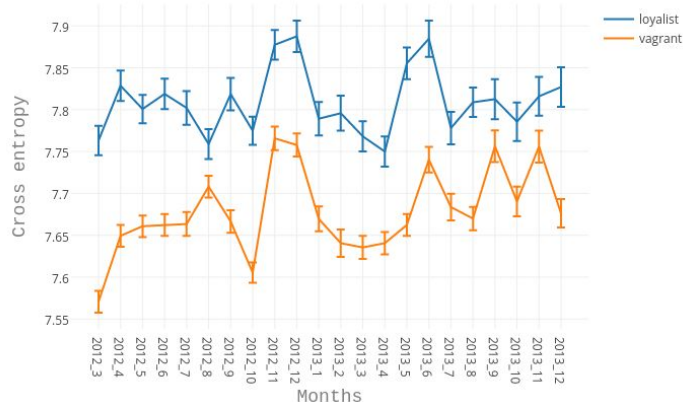
loyalist vs vagrant (LM: all, subreddit: starcraft)



loyalist vs vagrant (LM: shallow, subreddit: starcraft)

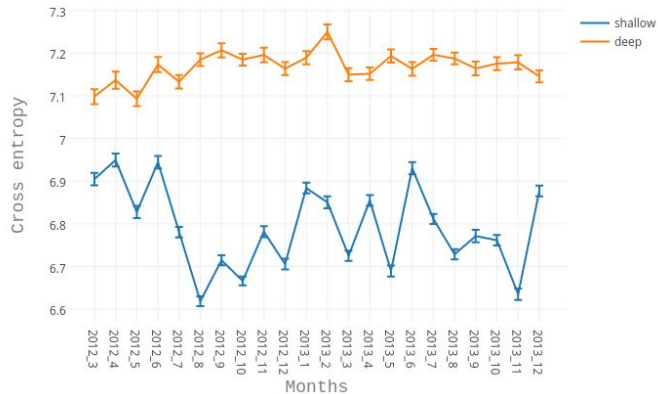


loyalist vs vagrant (LM: deep, subreddit: starcraft)



# Cross Entropy - Shallow vs. Deep

shallow vs deep (LM: all, subreddit: Music)

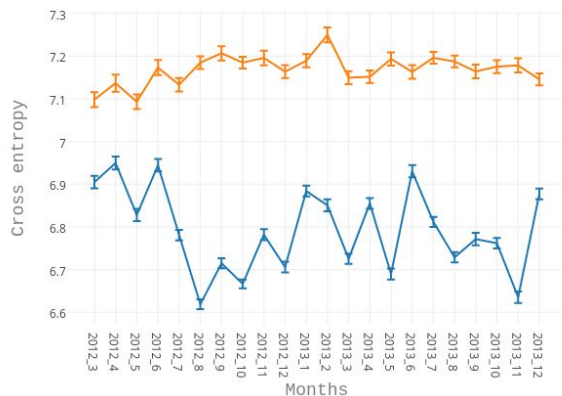




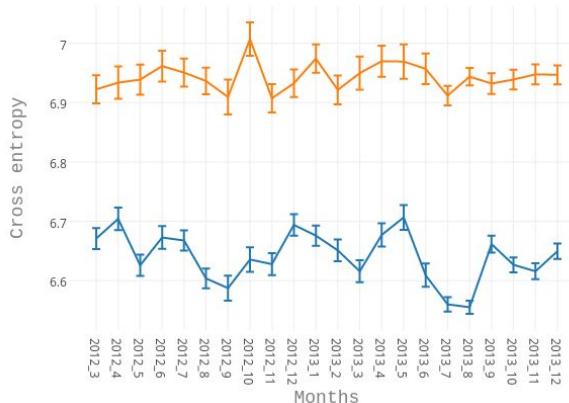
# Cross Entropy - Shallow vs. Deep

Different subreddits?

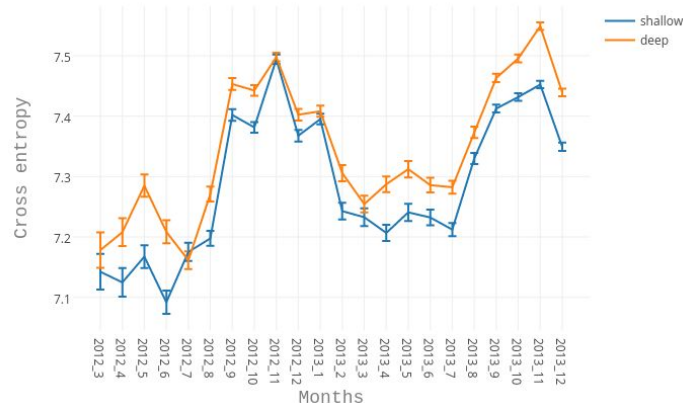
shallow vs deep (LM: all, subreddit: Music)



shallow vs deep (LM: all, subreddit: books)



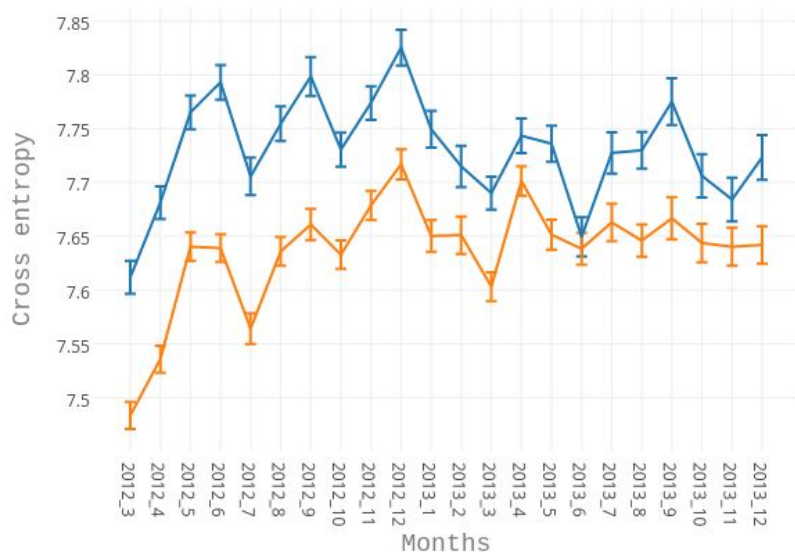
shallow vs deep (LM: all, subreddit: CFB)



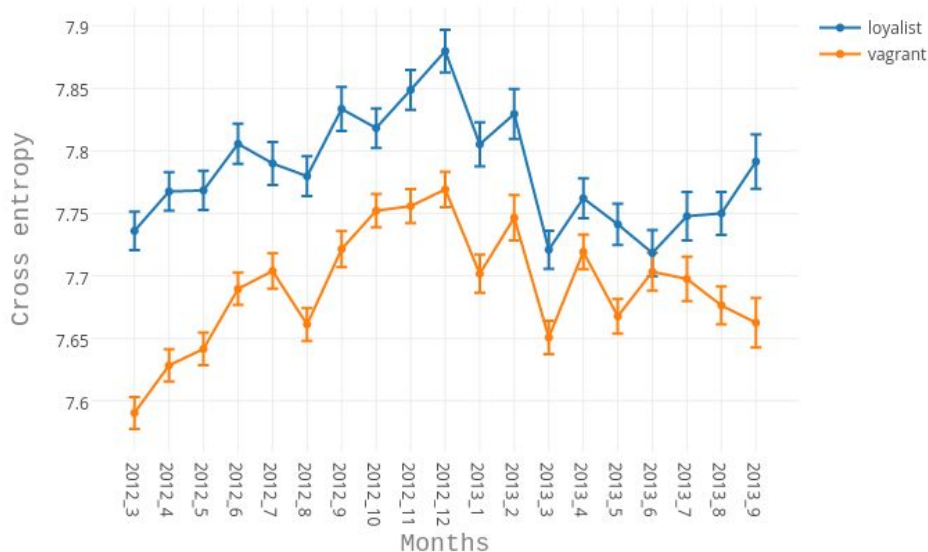
# Whose language defines the future?

Patterns are similar

loyalist vs vagrant (LM: shallow, subreddit: starcraft)



loyalist vs vagrant (LM: shallow in future 3 month, subreddit: starcraft)



# Where does the difference come from?

Content words? (word frequency ranked 1,001-5,000)

**vagrants vs. loyalists**



# Where does the difference come from?

Function words / Stop words? (word frequency ranked within 100)

**vagrants vs. loyalists**



# Behind Users' Activity Change

Users change their activity level:

within a consecutive 6-month window

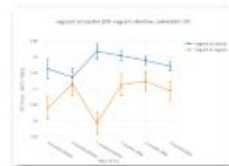
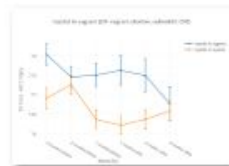
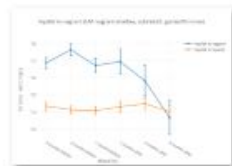
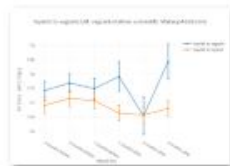
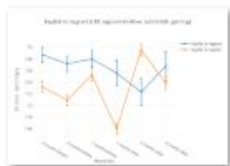
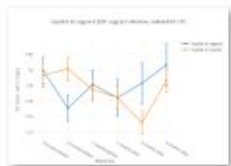
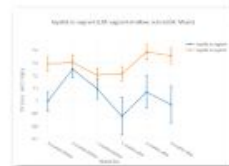
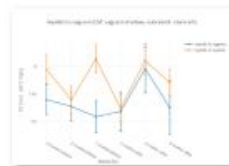
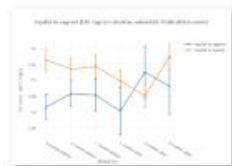
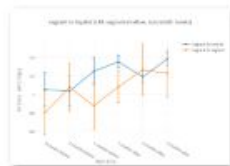
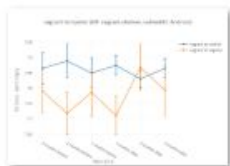
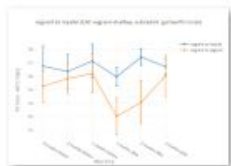
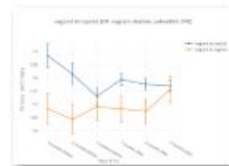
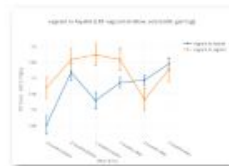
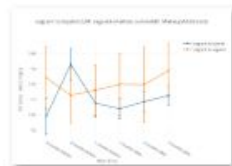
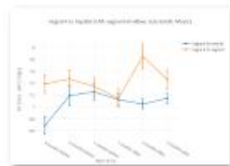
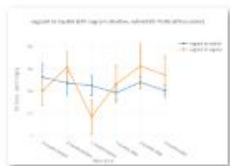
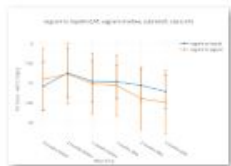
patterns like: *v v v l l l* or *l l l v v v*

pair them with: *v v v v v v* *l l l l l l*

Control for: subreddit, time, depth (all shallow), length (first 30 words)

# Behind Users' Activity Change

Results are noisy, hard to interpret (blue: changing users; orange: non-changing)



# Fightiest Words

What words distinguish loyalists from vagrants?

*investment, expertise*

For each **month and subreddit**, compute fightin' word z-scores for vocabulary split between loyalists and vagrants (we focused on **shallow** comments here).

**Fightiest Words:**  $\{w: \text{z-score} > k\}; \{w: \text{z-score} < -k\}$

We take  $k=3$ .

# Fightiest Words

**Fightiest Words:** {w: z-score > 3}; {w: z-score < -3}

## Examples:

*gameofthrones*, March 2013:

loyalist: *asos, chapters, narrative, barristan, stoneheart*

vagrant: *emilia, interview, watching, disney, doctor*

*Music*, October 2012

loyalist: *arcade, collective, quadrophenia, pitchfork, underrated*

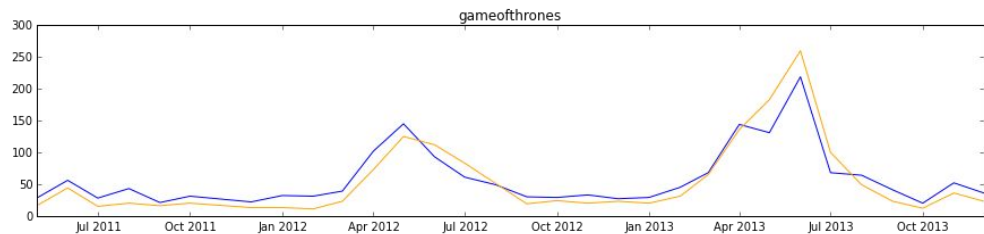
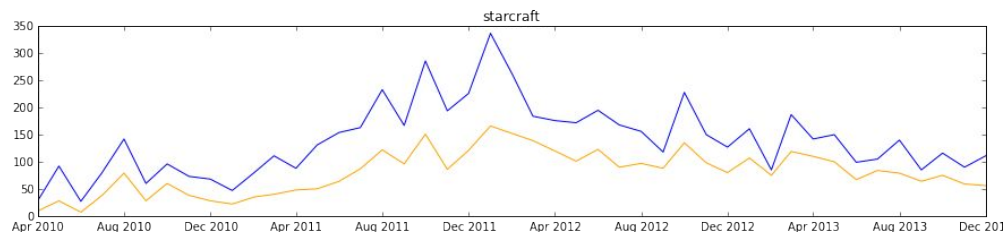
vagrant: *nicki, groovespark, remix, stupid, videos*



# Is loyalty fightier?

Mean number of fighting words per subreddit (over months):

subreddit	loyalists	vagrants
<b>starcraft</b>	<b>139</b>	<b>80</b>
gameofthrones	58	53
<b>books</b>	<b>42</b>	<b>32</b>
<b>Android</b>	<b>95</b>	<b>68</b>
<u>MakeupAddiction</u>	<u>105</u>	<u>87</u>
<u>PoliticalDiscussion</u>	<u>44</u>	<u>34</u>
gaming	316	326



# Fightiest Words over time

**How consistent are fightin' words over time?**

Hypothesis 1 (mine):

*loyalist words are more consistent, because loyalists have clearly-defined takes on the subject matter that are stable.*

Hypothesis 2:

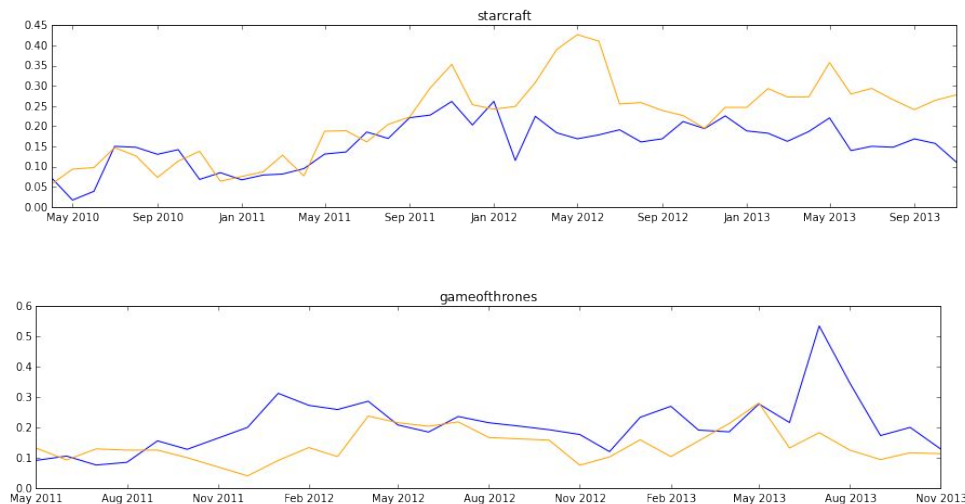
*loyalists words are less consistent, because they diligently follow changes in the subject matter.*

# Fightiest Words over time

For each month  $m$ , compute the **jaccard similarity** between fightiest words in  $m$  and  $m+dt$ .

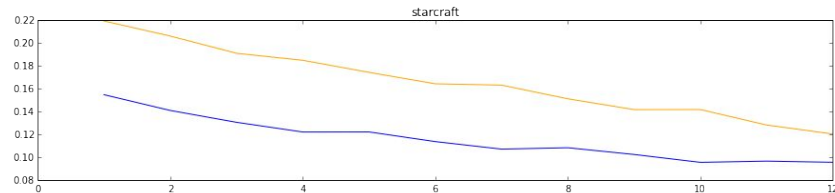
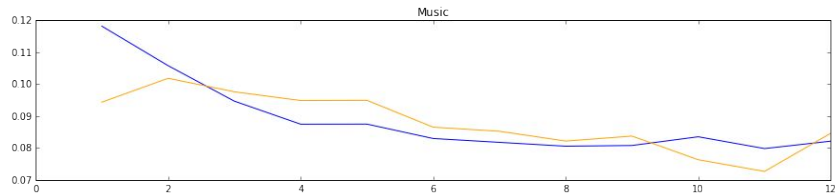
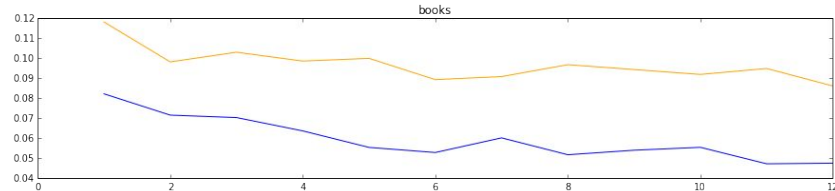
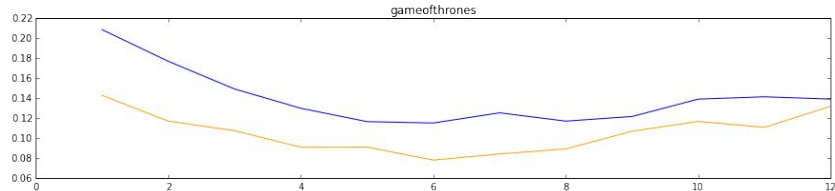
**Mean Jaccards,  $dt=3$ :**

subreddit	loyalists	vagrants
<b>starcraft</b>	<b>0.15475</b>	<b>0.21943</b>
<b>gameofthrones</b>	<b>0.20875</b>	<b>0.14281</b>
<u>books</u>	<u>0.08215</u>	<u>0.11820</u>
<b>Android</b>	<b>0.20075</b>	<b>0.16077</b>
<u>gaming</u>	<u>0.22155</u>	<u>0.19515</u>



# Fightin' Words Have Different Lifespans

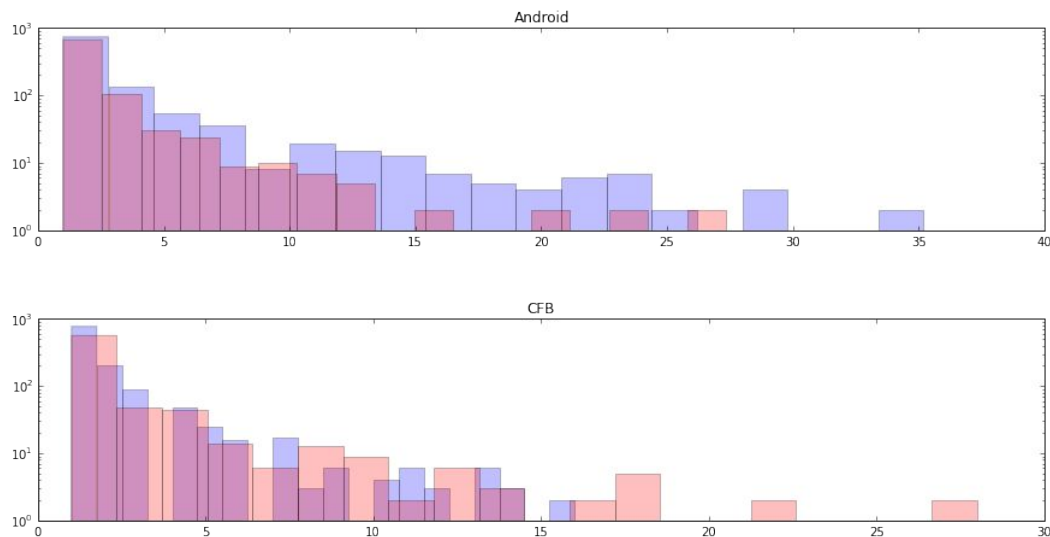
Jaccards vs  $dt$ :



# Loyalists are more loyal to their fightin' words

Average number of months a word is fightin':

subreddit	loyalists	vagrants
<b>starcraft</b>	<b>3.08</b>	<b>3.05</b>
gameofthrones	2.86	2.47
<b>books</b>	<b>2.33</b>	<b>2.19</b>
<b>Android</b>	<b>3.44</b>	<b>2.83</b>
nfl	2.64	2.74
CFB	2.01	2.55



# Longest-living fightin' words (nsfw: language)

## **starcraft:**

loyalist: teamliquid (33), gsl (26), protoss (22)

vagrant: twitter (28), dota (20), bot (20)

## **gameofthrones:**

loyalist: asos (26), affs (24), speculation (18)

vagrant: fuck (17), hodor (15)

## **Music:**

loyalist: album (44), radiohead (26), beatles (21)

vagrant: fuck (28), shit (14), dubstep (9)

## **MakeupAddiction:**

loyalist: brush (19), palette (18), blush (15)

vagrant: beautiful (22), eyebrows (20), gorgeous (19)

# Fine print

Fightiness and usage are slightly different; a word can become a fighting word if the other side uses it *less*. On average they're the same so I thought of them as such, and this is certainly true of the longest-living fightin' words. Letting the two cases be "straightforward" and "gah", vagrant fightin' words were more likely to be "straightforward" than loyalist fightin' words.

In light of this, what the fighting words analysis really tries to say is that one side's language use is more/less distinctive than the others, and how it is distinctive is preserved over time in different ways.

It's probably worth coupling FW with actual word frequencies. I tried and didn't find anything interesting, probably cause I wasn't asking the right questions. Here's a query that yielded small p values: is the lifespan of a FW related to the magnitude of the change in frequency between the month it became a FW and the month previous? Yes for loyalists, for some subreddits; not really, for almost all vagrants.

Also note that while relatively infrequent, a word can be fighty for loyalists and vagrants at different points in time. On average it wasn't clear to me who was introducing such words (loyalists or vagrants) so I didn't pursue this analysis further.

Who introduces such words, anyways? I have no idea about this and plenty of ideas about confounds.

# Conclusions and future work

Loyalists and vagrants behave differently, but precisely how depends on the subreddit.

We can make up nice anecdotes to explain the effect direction in individual subreddits. **Can we characterize subreddits by the behaviour of their loyalists and vagrants?**

Loyalists and vagrants both seem to define some community roles in their language use. **Can we couple an analysis of community network structure with loyalist and vagrant roles?**