

# Knowledge Assistant for Comment Systems

Capstone Presentation – A Case Study on HODINKEE.com

May 7, 2020

Hamza Masood

[LinkedIn](#) | Data Science – Flatiron School | [GitHub](#)



# Why is this interesting?

## To Hodinkee

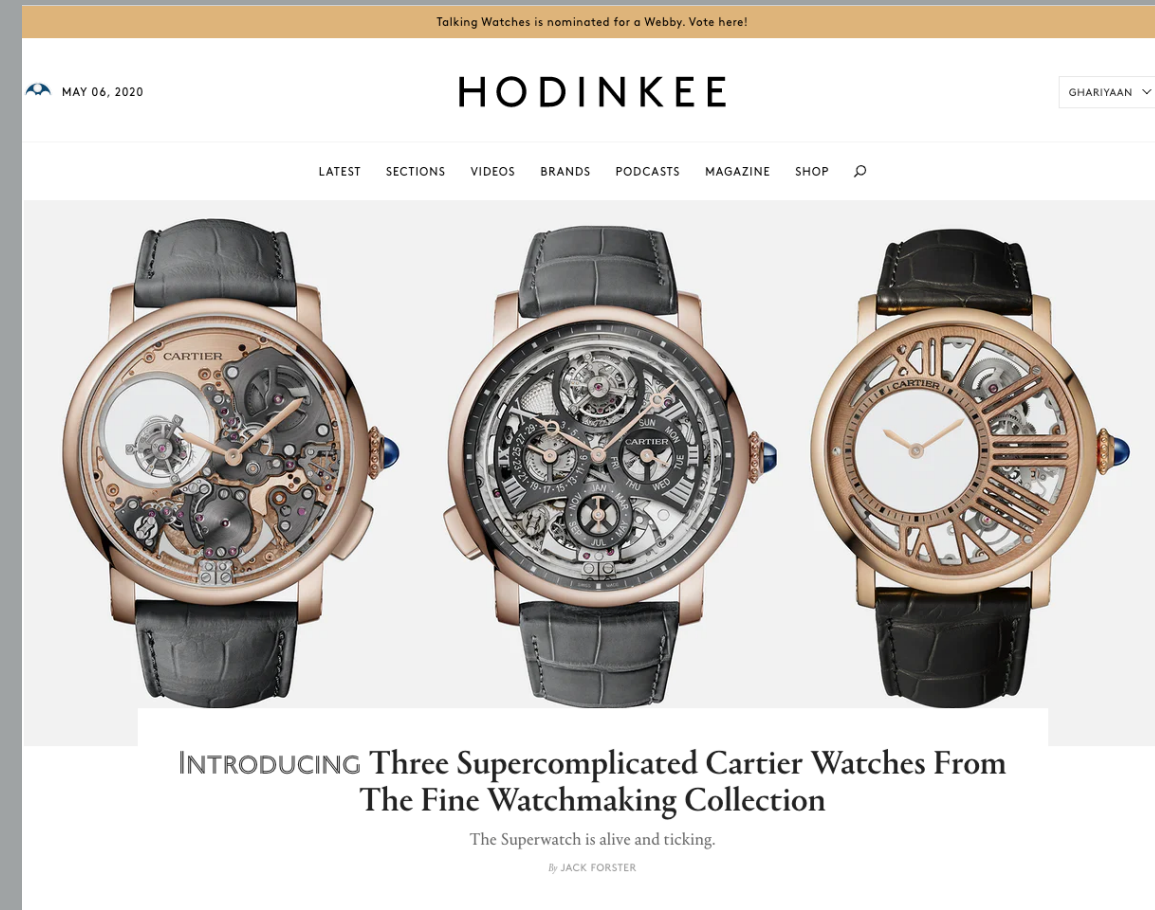
- Better user engagement → more time spent on hodinkee.com → Higher conversion rate for the Hodinkee Shop
- Build a better community than competitors

## To a General Audience

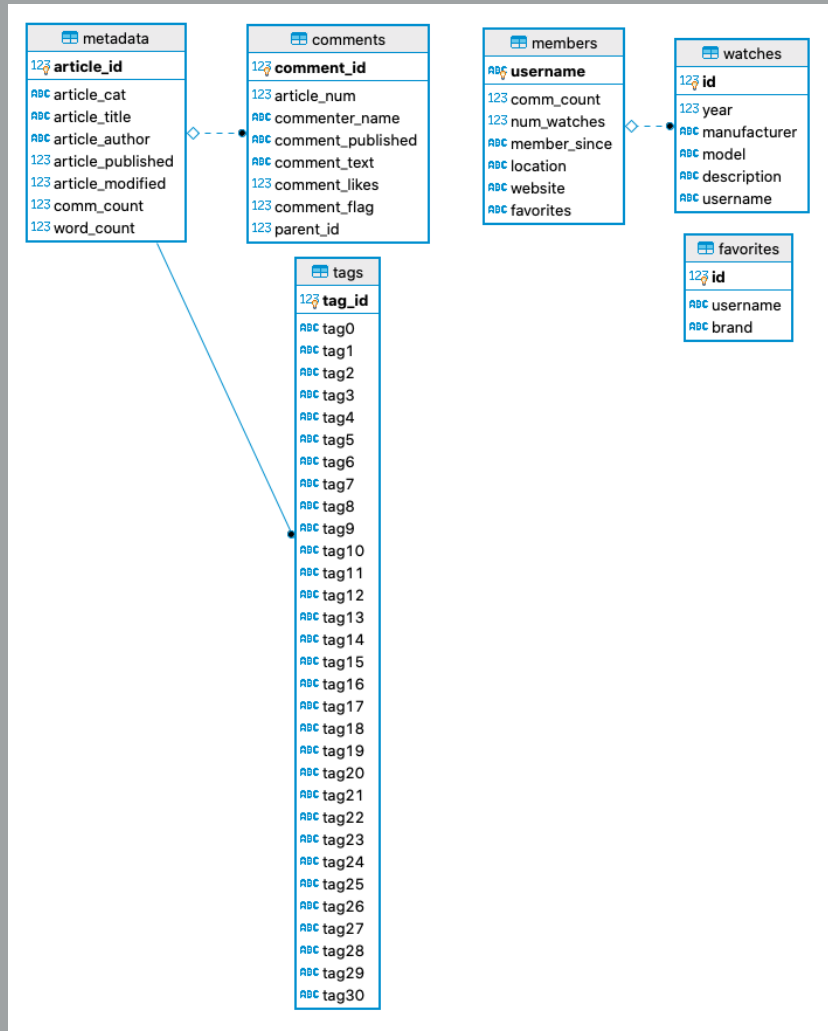
- Lots of similar online communities with highly engaged users: other online watch clubs, tech review sites, wine clubs, gaming reviews, and so on

## To Me

- I love watches
- I got to build a dataset from scratch



# Methodology Details



Website snapshot taken on Apr 5

Total scraped data >1.75GB

Python libraries used:

- pandas
- numpy
- beautifulsoup4
- sqlite
- matplotlib
- statsmodels
- datetime
- colour

metadata      6,501 rows

tags            6,501 rows

comments      137,495 rows

members       13,355 rows

favorites       17090 rows

watches        13490 rows

6,501 article html pages

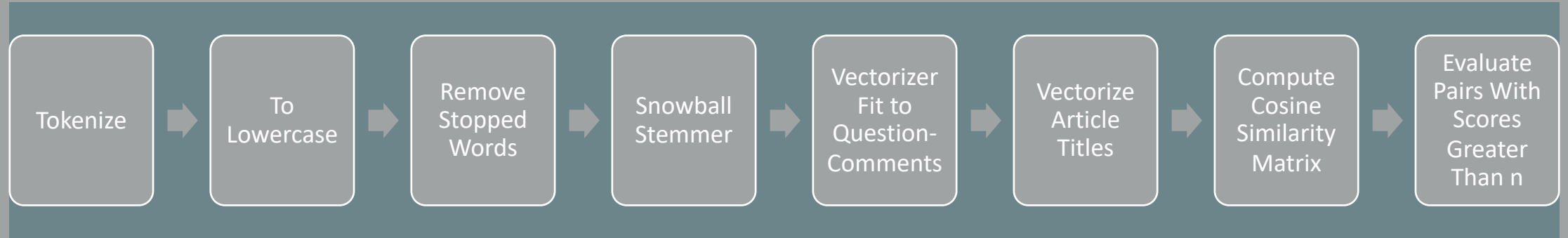
6,501 article text files

13355 user html pages

SQLite database schema diagram



# FSM



## Issues Found

- Recommended article can't be the same as the article on which the comment was written.
- Recommended article can't have a publish date greater than the comment publish date.
- Types of questions:
  - Specific (who is, what is, why is, how is)
- Not answerable by an article:
  - Replying to another comment with a question
  - contextually specific to the article itself
  - Sarcasm
  - "?!?"
  - rhetorical question answered in the same comment

Increasing Relevance



## Evaluation Scale

1. Perfect
2. Thematically Relevant
3. Related But Different
4. Irrelevant

## Results

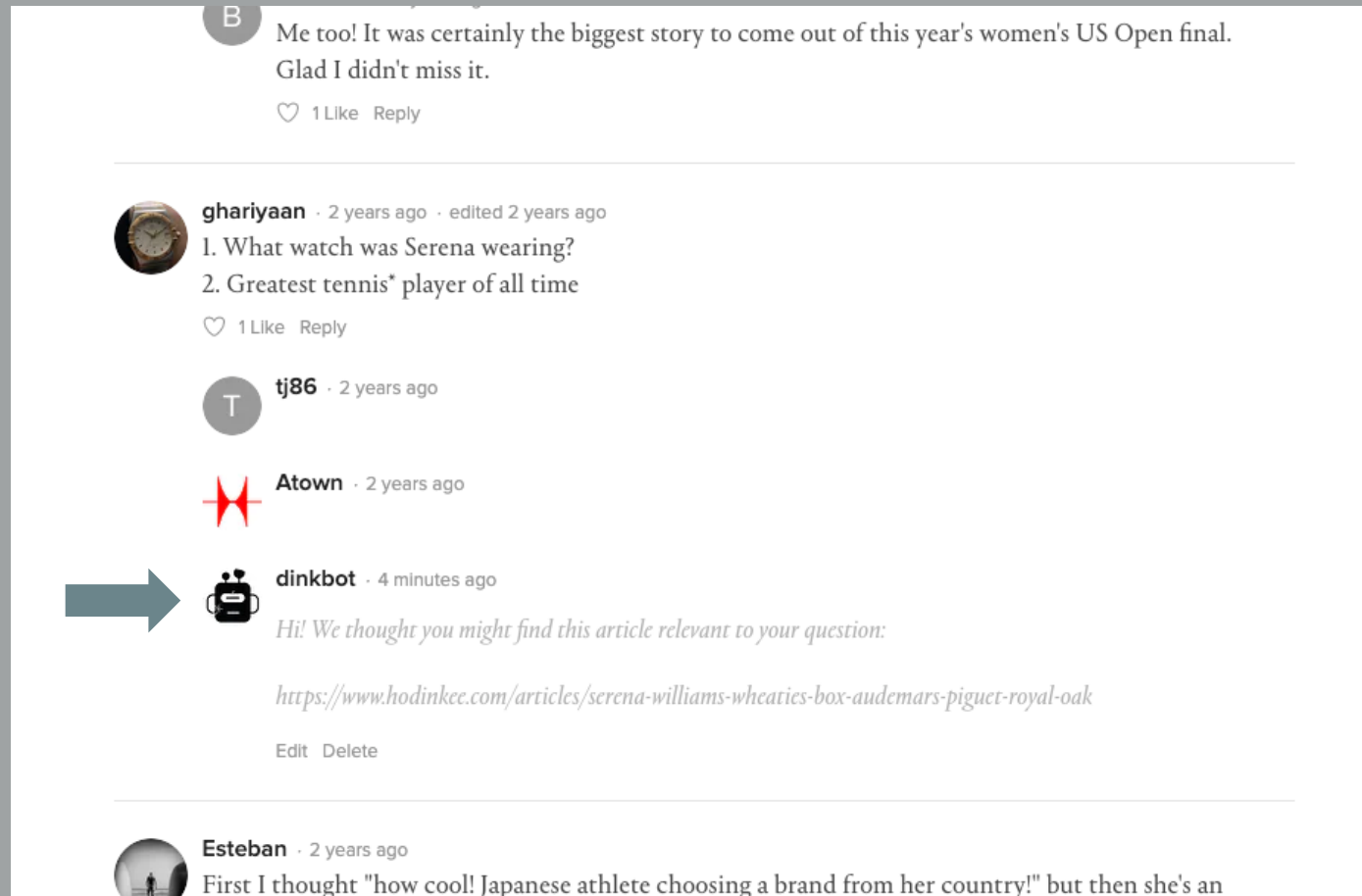
Only 49 comments with CS scores  $> 0.8$ , most were evaluated as NOT irrelevant

# Digging into the questions





# Ideas for the future



- Live site deployment
- Article corpus expanded from Hodinkee to the Watchville publisher family and beyond

Questions?

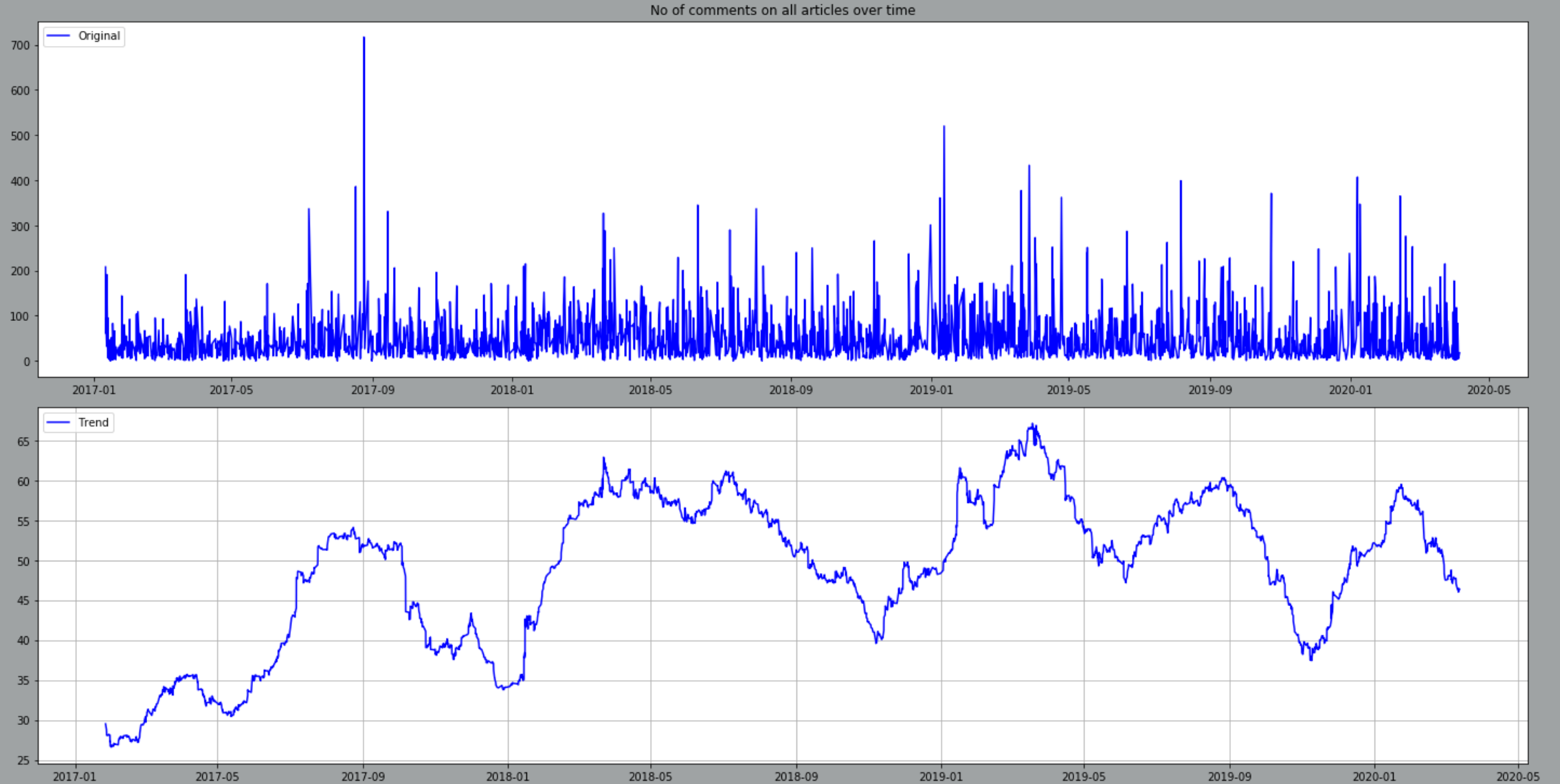
# Sources & Kudos

- [Hodinkee.com](http://Hodinkee.com)
- Erin Hoffman
- Greg Damico
- Brian McGarry

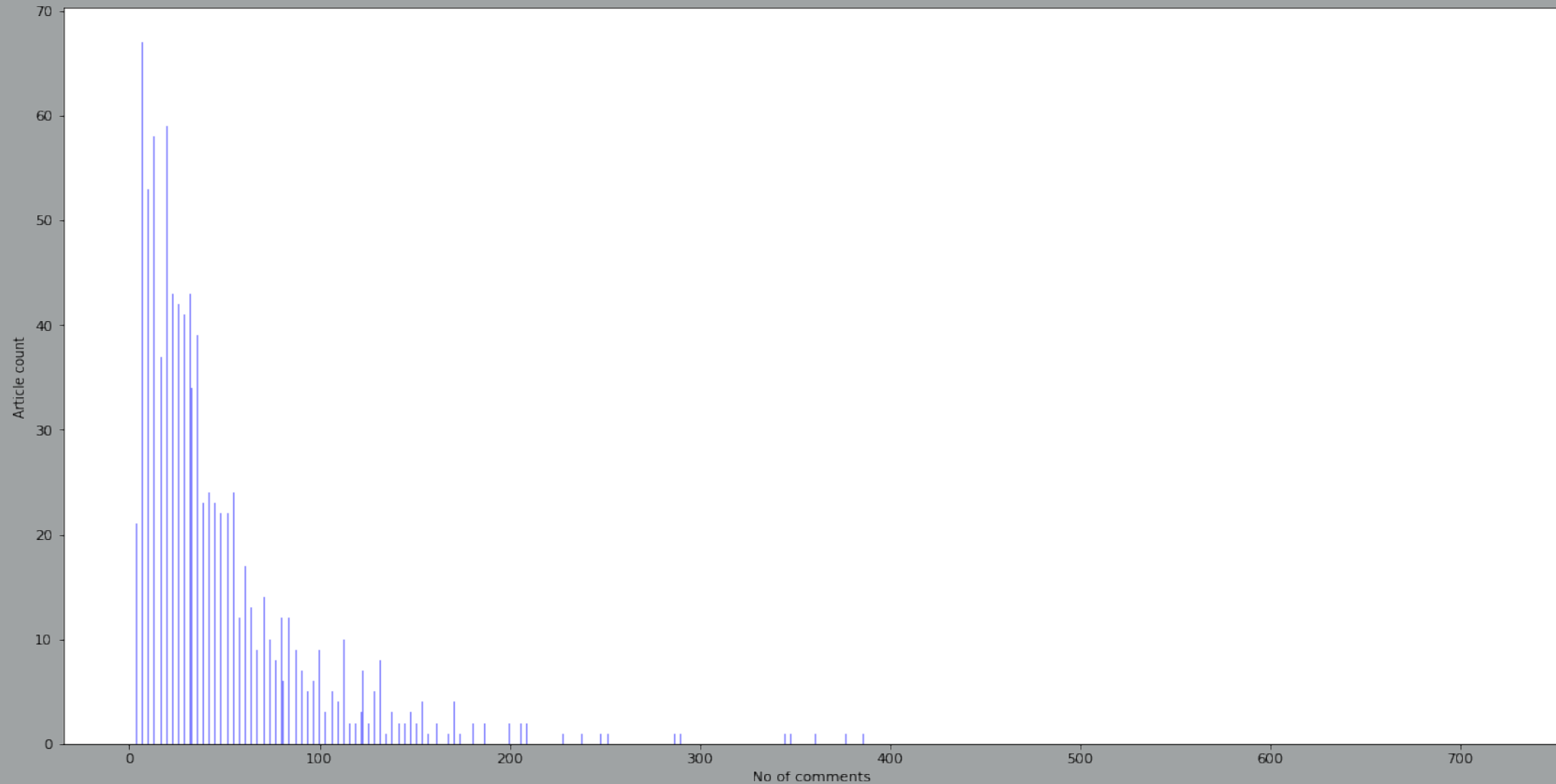




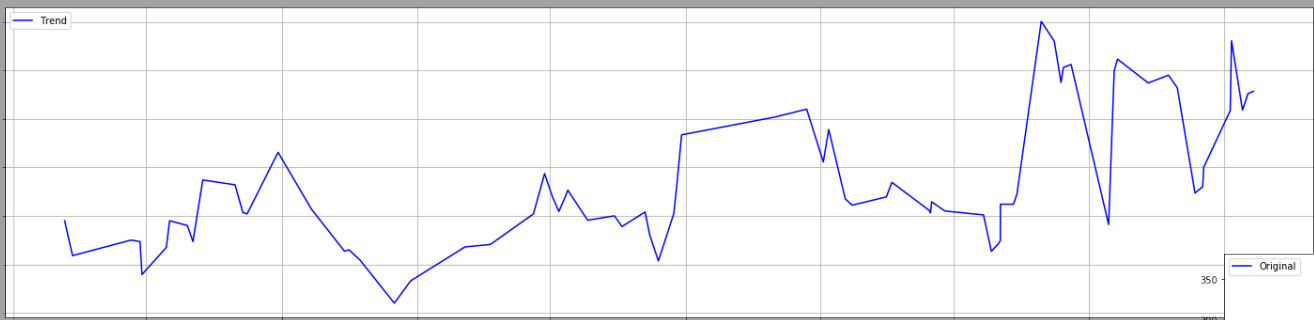
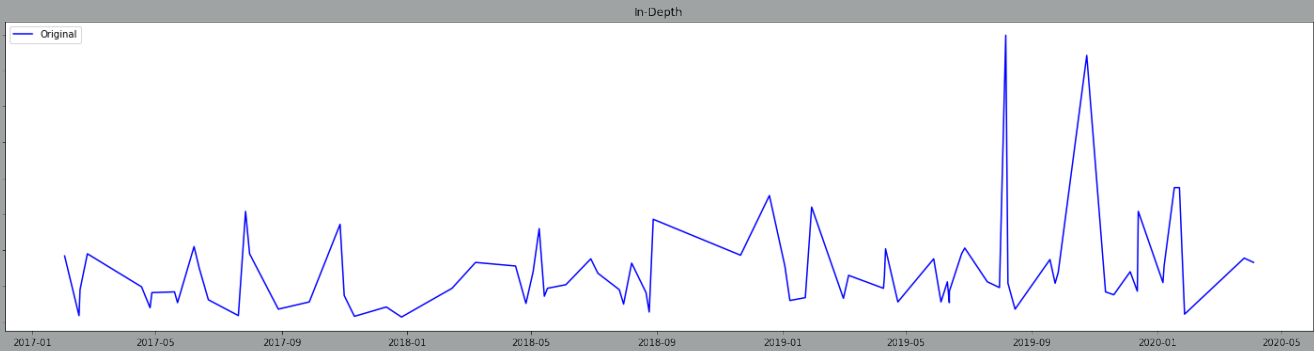
# All comments over time: seasonal, cyclical trend



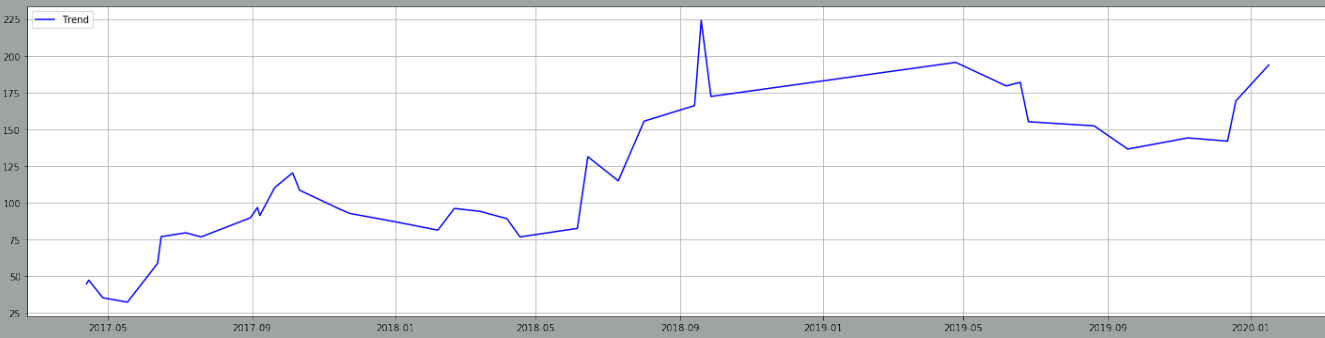
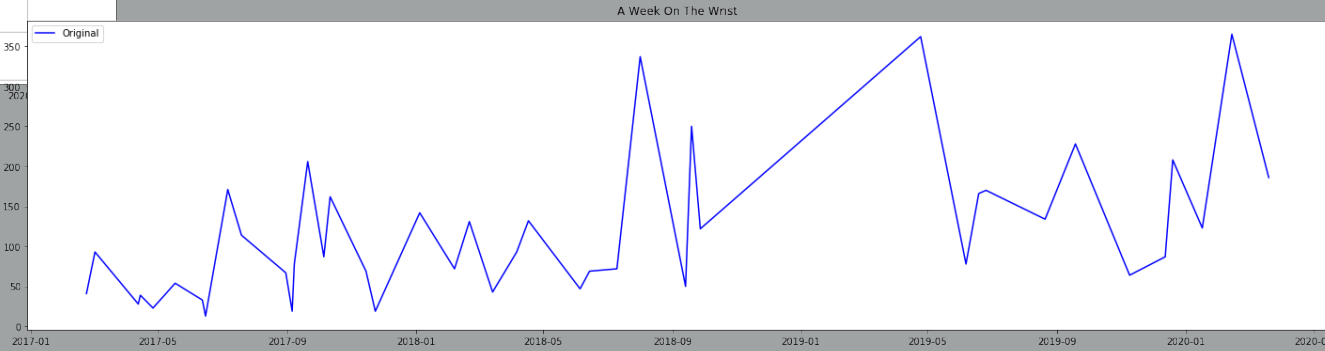
# Distrubution of comment count over articles



# Categories with growing engagement



In-Depth



A Week On The Wrist



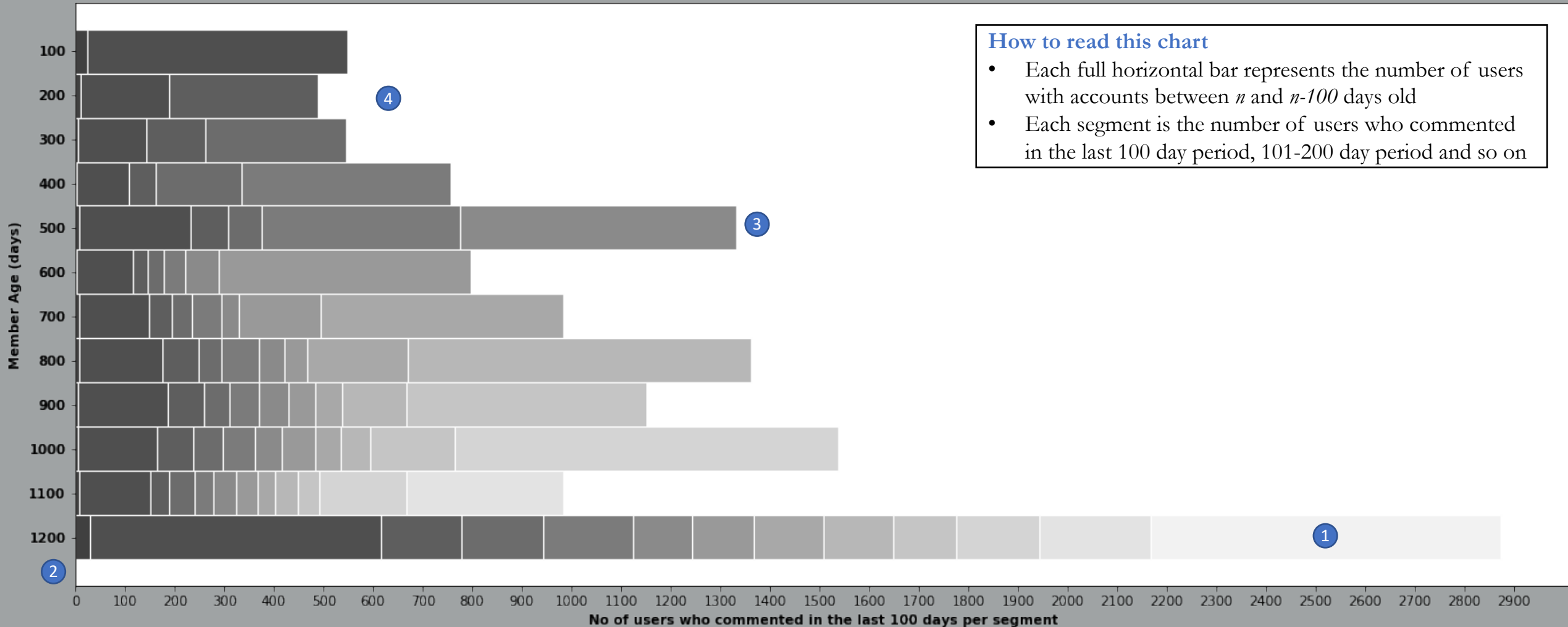
# Most prolific commenters (normalized by account age)

Rank	Username	Prolificness Score	Comment Count	Account Age (days)	No. of Watches in Profile
1	Shatners	1.38	721	521	0
2	ICH	1.38	1392	1009	1
3	Gav	1.35	499	370	5
4	GreatScot	1.34	127	95	0
5	JackForster	1.11	1324	1190	0
6	Bside	0.98	1171	1190	0
7	Boman	0.97	625	644	5
8	wkf	0.90	1072	1190	12
9	PaulMiller	0.85	1008	1190	0
10	ripwatch	0.81	720	886	0
11	Oliver_H	0.77	27	35	5
12	Yev	0.77	703	917	5
13	CynicalBastard	0.76	626	825	0
14	TheOmegaMan	0.75	853	1131	10
15	Putito	0.73	136	187	1
16	AJ117	0.72	244	340	0
17	ThicknessMatters	0.71	762	1070	0
	...	...	...	...	...
20	Cole	0.63	271	429	9
	...	...	...	...	...
82	BenClymer	0.27	318	1190	3
	...	...	...	...	...
124	CaraBarrett	0.22	264	1190	2
	...	...	...	...	...
133	ghariyaan	0.21	131	613	7

Prolificness is calculated as  $\frac{\text{Comment Count}}{\text{Account Age in days}}$

Included some Hodinkee Staff and myself for reference  
(I'm surprised and disappointed my score isn't higher 😊)

# User age and engagement pyramid



- 700 of the oldest users on the site haven't commented in 1,200 days
- There is a continuously engaged cohort of users across all account ages

- Users picked up around December 2018 continue to be more active
- Fewer new users have signed up in the last year than in the previous two years

**Key Takeaways:** Most users who signed up didn't remain engaged and fewer users have been signing up over time



