

Looking at Censorship on the Internet as a social process

- Pranathi Iyer

MACSS, University of Chicago

1. The social process that I look at is censorship. While censorship in itself is an extremely broad topic, I particularly wish to look at censorship of content on the internet by governmental institutions. The recent past has witnessed a visible surge in government intervention in regulation of content online. This intervention has been prevalent world over, across spheres of art, education, politics, and other forms of expression. Censorship is a huge challenge for countries, as it speaks volumes of the extent to which governments truly detach themselves from citizens' freedom of speech and expression. However, governmental interference in the form of censorship has been more pronounced in some countries than others. One avenue where this disproportionate interference across countries manifests itself, is twitter removal requests. Between 2012 and 2020, the highest percentage of removal requests came from five countries alone (Twitter, 2020). One of these five countries is India, despite being the largest democracy in the world. This solicits interest as to what could be the dynamics within the country, and hence, I further narrow down my social process to

understanding censorship of content on the internet in the Indian context, by looking at removal requests made by the government to twitter.

Why India?

India is an extremely eclectic country with 6 main religions (Pew Research, n.d.) and other religions with 21 modern languages (Government of India, n.d.), which makes its expression of speech and expression extremely diverse. This makes it extremely interesting to study, and understand how the government responds to behaviour of people of different backgrounds online, and if there are certain groups that are censored more than others.

Why twitter?

A simple google search stating “Government of India bans twitter accounts” provides a rather compelling answer to this question. While there have been other avenues such as films (BBC, 2021) , where people have faced censorship, the mass user base of twitter, along with the plethora of issues that it provides a discussion platform for, especially political discourse, makes it a suitable choice.

I choose twitter, however if feasible, it would be interesting to see how this would pan out in the space of Instagram and Facebook, where government intervention seems to be lesser, however, there is no statistical evidence supporting this.

2. What makes it hard to study censorship is the arduous task of tracking censored data, and understanding its nuances. In the context of Twitter, most government requests, if public, are murky and often difficult to trace down to the actual content of the post-- unless it has been captured online during its ephemeral time

on the web. The lack of official institutions that track censored content makes it even harder to study the process. Luckily enough, for the purpose of my research, I was able to gather and use data from the Lumen database of the Berkman Klein Center for Internet & Society at Harvard University (Lumen Database, 2017). Established in 2002, this archiving initiative attempts to capture requests of removal for content from the web. They were also kind enough to provide me with researcher's access and an authentication key, which made scraping data much easier and systematic. I use the database to acquire notices of government removal requests for twitter data, and then further access details (URLs) of accounts reported, from these notices

Why is it suitable for my social process?

As I mentioned above, the transient nature of censored content online, makes it a difficult task to analyze withheld or removed content. Consequently, there is a dearth of authentic resources--beyond newspaper articles and media--which are able to catalogue requests for removing content, along with the corresponding collection of URLs that are proposed to be infringing. The Lumen database has one such exhaustive resource of twitter removal requests made by the government of India, and I wish to use this repository to extract the specific accounts that were reported, and further my analysis from that point.

Time resolution

While the concept of censorship is a relatively old one, the internet as a medium to express ,honest and sensitive--some would say controversial-- opinion and content has erupted in the past decade or so. Platforms such as Facebook, Instagram, Reddit, and Twitter have increasingly become platforms for hosting discussions across the spectrum, right from somebody's vacation updates, to conversations around suppression of racial and religious minorities. It is this development over the past decade that makes the Lumen database apt for my research, since it captures government removal requests from 2002 to 2020, probably more accurately as time passed. Owing to several constraints on processing, storage ability, heavy throttling from the databases's API, and extremely volatile nature of twitter which has thousands of posts in a day, I stick to the timeline of August 2019 to October 2021. However, I only present a small sample of this dataset even between the specified timeline, since the number of files were too many to process. Moreover, as Brugger (Brugger, 2018) explains, the temporality of these notices is not the most well defined on the web page itself. Different elements of the website such as the website itself, the notices, and the actual request dates, are all rather intertwined, which makes using html tags to extract dates imperative . This timeline also means that for studies that wish to look at how expression, and conversations have changed since the advent of twitter, this data set might not be suitable with respect to time and spatial considerations (Brugger, 2018).

- The data is not indicative of the temporality of content for 2007-2012

- The data might not be relevant to study actual dynamics of conversation on twitter since this data only talks about accounts whose posts were withheld.
- The actual posts themselves might not even be visible, however analysis of the comment threads might provide meaningful insights.

Data transformation

I perform seven steps to transform the data which I illustrate through the table below.

Final Objective: to get access to details of accounts (URL) against which the government of India took action

Step 1	I scrape the first 5 search pages of removal requests by the Indian government on Lumen's database to get access to URLs of notices from August 2019 to October 2021.
Step 2	From the URLs scraped above, I use html tags to access specific URLs of pdfs associated with these notices which have URLs of accounts that were reported.
Step 3	I access dates of each of the pdfs by accessing Epoch dates from their urls and convert them to YYYY-MM-DD format which I'll use as filenames for the pdfs.
Step 4	I clean these dates by stripping off the whitespaces..
Step 5	I use a curl command with necessary arguments using the subprocess library and bash in python, to download these pdfs onto my system and save them with their dates as filenames
Step 6	Once I have all pdfs on my system, I write a python script to extract URLs of the required twitter accounts from these documents..
Step 7	I finally save the date and the account URLs of the accounts whose content is withheld, in a dataframe

The change in digitality is two fold, the first is that the notices themselves are being accessed using the metadata of the webpage. Secondly, the documents with account details are no longer just pdfs online, they are also pdfs accessible on my local system. In some sense, the born-digital document has been reborn. Furthermore, getting access to the links to the accounts and the posts which are withheld--not visible--gives access to the comments on the post. This opens up a whole new realm of trying to understand what the post could be about from the comments. In some sense this is similar to a true palimpsest (Bailey, 2007) where the original data or meaning is removed, and we can make sense of it from what remains on the web.

Gaps

- One of the biggest gaps in this dataset could be its spatial consideration. Posts are redacted on the internet almost everyday, and even with a track of the removal requests, the best we can get is the comments made on the post, and not the post itself.
- The fact that account details are being retrieved from downloaded documents means that this data can be inaccessible if documents are damaged or truncated themselves.
- The size of this specific archive itself is extremely small owing to constraints mentioned above. More data would have to be scraped off using the same code, that can then be used for more meaningful research.

- Owing to several constraints, the current archive also does not allow to have the macroscopic picture that big data and archives of this kind can provide (Graham et al, 2016).

Ethical Implications

- Data of this kind can be extremely sensitive since on some level, it captures which content the government of a country wanted removed from the online platform itself.
- Account details of people also have personal identities of individuals whose accounts were reported by the government and must be handled with utmost care. This relates to the topic of understanding who is being researched, and looking at digital data beyond the text that appears on screen (Lamborg, 2019)
- To some extent this is resolved since several individuals oppose being reported by the government, and have been public about their identities online. However, this does not change the fact that the data must be handled carefully.
- Lastly, the database itself is not easily accessible to the public without an API key or researcher credentials. This should make researchers extremely cautious about where they use this data and what form they present this in. However, non-researchers can access this data, and it is only the speed and scale at which they might be able to access it will differ.
- I choose to not disclose any of the personal details of these accounts publicly, thereby trying to secure the identities of people.

Future scope of research:

The URLs from the archive can be used to access censored posts, and if the post is withheld, analysis of the comment section can provide meaningful insights into the kind of content that the government tends to censor.

I believe that studying censorship can provide significant insights--perhaps not surprising-- into how governments think, and how accommodative countries are of dissent and creation without constraints. Akin to several other researchers in today's times, I wish to understand online censorship and its repercussions, specifically in the Indian context. This archive is only a stepping stone towards my larger goal of looking at censorship as a social process over time

References

Bailey, G. (2007). Time perspectives, palimpsests and the archaeology of time. *Journal of Anthropological Archeology*.

BBC. (2021). Bollywood: Filmmakers cry foul over censorship proposals. *BBC*.

<https://www.bbc.com/news/world-asia-india-57676214>

Brugger, n. (2018). *The archived web*. MIT Press.

el al, G. (2016). *The Joys of Big Data for Historians*.

Government of India. (n.d.). *The Indian Linguistic Space*. Education Ministry.

https://www.education.gov.in/en/sites/upload_files/mhrd/files/upload_document/languagebr.pdf

Lomborg, S. (2019). Ethical considerations for web archives and web history research. In N. Brugger, & I. Milligan (Eds.), *The SAGE handbook of web history*. Sage UK. Credo
http://proxy.uchicago.edu/login?url=https://search.credoreference.com/content/entry/sageukweb/ethical_considerations_for_web_archives_and_web_history_research/0?institutionId=170

Lumen Database. (2017). *Lumen Database*. Lumen Database.

<https://www.lumendatabase.org/>

Pew Research. (n.d.). *Key findings about religious composition of India*. PewResearch.

<https://www.pewresearch.org/fact-tank/2021/09/21/key-findings-about-the-religious-composition-of-india/>

Twitter. (2020). *Removal requests*. Twitter. Retrieved October 18, 2021, from

<https://transparency.twitter.com/en/reports/removal-requests.html#2020-jul-dec>

