

# Detect China Online Censorship — The impact of major political events on Weibo posts

Hazel Chui  
University of Chicago  
Chicago, Illinois, USA  
hazelchui@uchicago.edu

Peihan Gao  
University of Chicago  
Chicago, Illinois, USA  
peihaan@uchicago.edu

Yulun Han  
University of Chicago  
Chicago, Illinois, USA  
yulunhan@uchicago.edu

Xin Tang  
University of Chicago  
Chicago, Illinois, USA  
tangxin@uchicago.edu

## ABSTRACT

This paper investigates the impact of censorship on online public opinion by analyzing Weibo posts made by users before and after the 20th CPC National Congress. We collected Covid-related Weibo posts and examined changes in sentiment, keywords, and topics to measure potential government censorship in China. The findings suggest that strict censorship may have been imposed prior to the Congress, while speech controls appeared to have relaxed afterwards. The study offers insights into the Chinese government’s strategic information control policies.

## KEYWORDS

Information control, Weibo, sentiment analysis, topic modeling, government censorship

## 1 INTRODUCTION

The increasing use of social media platforms in China has led to a growing need to understand the impact of censorship on online public opinion. As the Chinese government reinforces censorship measures before major political events to ensure their success, it is essential to investigate how these measures affect the public’s posts on social media platforms such as Sina Weibo (Weibo).

In this paper, we analyze the public’s posts on Weibo before and after the Twentieth National Congress of the Communist Party of China (the 20th CPC National Congress) to investigate the impact of censorship on public opinion. Our findings could provide insight into the effectiveness of censorship measures in controlling public opinion during major political events in China. By analyzing the public’s posts on Weibo before and after the 20th CPC National Congress, we can examine the atmosphere that the government desires to create before and after symbolic political events, as well as to evaluate users’ self-censorship or platforms’ filtering effects.

We use computational methods to analyze a large dataset of social media posts, and our work differs from previous studies in its focus on a specific major political event and its use of posts that have not yet been censored. This approach allows us to examine the atmosphere that the government desires to create before and after symbolic political events, as well as to evaluate users’ self-censorship or platforms’ filtering effects.

## 2 RELATED WORK

Research on online censorship in China has focused on analyzing the content of censored social media posts and identifying keywords likely to be deleted by the government. King, Pan, and Roberts (2013) found that the censorship program aims to silence collective expression, not criticism, by analyzing censored social media posts. This conclusion was supported by a separate study in which they created accounts and randomly submitted texts. Similarly, Baman, O’Conner, and Smith (2012) uncovered keywords likely to be censored by the Chinese government by analyzing deleted Weibo posts. Additionally, Tai and Fu (2020) found that even articles on seemingly non-sensitive topics and those containing pro-regime messages are also removed by censors on WeChat.

Several studies have found that the level of censorship in China may not remain constant over time. Han and Shao (2022) show that the Chinese state scales up control over citizenry complaints in response to important political events. They find that the government tightens information control before important political events. Ruan et al. (2020) have demonstrated that by studying WeChat, a popular Chinese communication app, the government alters censorship patterns before, during, and after significant events. The results suggest that censorship in China may be dynamic and responsive to political developments.

Our study aims to expand upon previous research on censorship in China in two key ways. First, instead of analyzing deleted posts, our study concentrates on Weibo posts that have yet to be censored. This approach enables us to examine the atmosphere that the government desires to create before and after symbolic political events, as well as to evaluate users’ self-censorship or platforms’ filtering effects. Second, to our knowledge, researchers have yet to study the impact of major political events on censorship in Weibo, the most widely used microblog in China. Weibo posts are an essential source for understanding the outcome of the interaction between censorship and critical political events.

## 3 METHODS AND DATA

In this project, we want to find out how the Chinese government censors social media and figure out how much it does it.

In most cases, more social stability and speech censorship might be required to keep things running smoothly when significant

political events are held in China. Weibo is China’s most often-used social media platform for people to share their thoughts and emotions towards a specific topic. Before December 2022, China put out a set of strict rules to stop pandemics during COVID-19. Some of these rules raised public outrage.

We will scrape Weibo posts mentioning keywords including “疫情” (“covid”), “隔离” (“quarantine”), and “抗疫” (“prevention and control of covid-19”) spanning one month before and after the 20th CPC National Congress (October 16th, 2022) as our main observational sample, and analyze the change of sentiments, keywords, and topics of these posts to measure the potential censorship from the Chinese government.

We have the following hypotheses: Our research will demonstrate that, during significant political events, it is much more difficult for people to disagree with policies on social media due to censorship. Our analysis will highlight the considerable decrease in posts containing politically sensitive keywords after the 20th CPC National Congress. Our findings may demonstrate how successfully censorship is used in China to control public opinion during significant political events.

We’ll use sentiment analysis to compare the evolution of posts’ sentiments and display the trend of sentiment scores in a line graph. The keywords for posts will be displayed as word clouds, with the size of the words indicating how frequently they occur. To analyze content changes, we will then use topic modeling. Our expected results will show that censorship measures significantly impact the public’s ability to express dissenting opinions on social media during major political events. Specifically, our analysis will show a significant decrease in posts containing keywords related to sensitive political topics after the National Congress.

## 4 EVALUATION AND ETHICS

We have decided to evaluate this project by 1. comparing sentiment scores before and after (the 20th CPC National Congress) to see if the difference is significant or if the scores were positive before the 20th CPC National Congress and negative after that; 2. comparing the word clouds for the two time periods and analyzing the relatively bigger words in each of the graphs; 3. analyzing the topic modeling results for the two time periods and determining if the topics the models identified make sense and if the topics can be informative. 4. deciding if the methods are feasible. For example, if we can find a sentiment measuring function to analyze the sentiment for this specific dataset reasonably.

There would be biases in our project. For one thing, our result might be influenced by “解封” (the lift of lockdown), which happened approximately two months after the 20th CPC National Congress. There is a chance that during the month after the 20th CPC National Congress, there have already been some changes in the policies China chose to implement. For another, we might not be able to detect significant results if the sample size that we get is not big enough due to different reasons e.g., unable to scrape a sufficient number of Weibo posts because of rate limits.

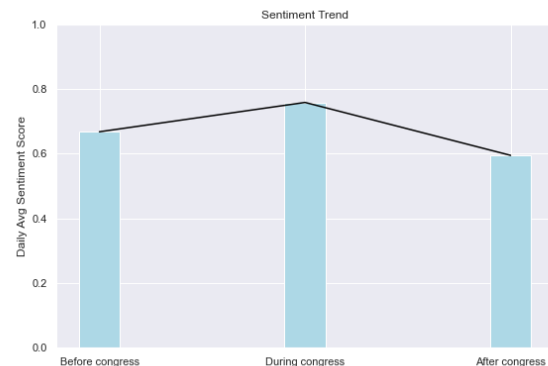
## 5 RESULTS

### 5.1 Sentiment Analysis

We compute the daily average sentiment scores, which are adjusted for post popularity, and fit three lines based on the trend of the sentiment scores using linear regression, separately for the three periods - the month before the 20th CPC National Congress started, the week of the 20th CPC National Congress, and the month after the Congress ended - to determine how the posts’ sentiments change as the days get closer to the 20th CPC National Congress and further away from the event.

More specifically, we carry out the sentiment analysis illustrated below:

Based on the default customer comments training set of the SnowNLP Python library, we compute a sentiment score for each post. This library pre-processes simplified Chinese texts, including tokenization, part-of-speech tagging, sentence splitting, and keyword extraction. When the API is called, SnowNLP performs text pre-processing on the posts before performing sentiment analysis. Each word in the message is assigned a value between 0 and 1, with 0 designating a negative sentiment and 1 a positive one. So, it generates a sentiment score for the tested post after averaging the dummy values of all the terms in the post. The daily average sentiment score is derived using the weights of the amount of popularity on each day since it is likely that posts on Weibo with varied levels of popularity will have a greatly varying influence on social media. The total of retweets, comments, and likes for each post serves as an indicator of its level of popularity.



**Figure 1: Weibo Posts’ Sentiments across the Congress**

The daily average sentiment scores for the three time periods—before, during, and after the Congress—are 0.67, 0.76, and 0.6, respectively, as shown in the graph above. It is consistent with the hypothesis that posts’ sentiment would increase in the days leading up to the 20th CPC National Congress and peak during that time. When the Congress is over, the attitudes resurface and are considerably more unfavorable than in the month leading up to the event. The average emotion trends point to the likely censoring of unfavorable posts before and during the significant event. Instead, when the emphasis is not on maintaining the atmosphere of stable public opinion, censorship vanishes after the event.

The daily time-series shift of sentiment is clearly visible in the graph below. We use linear regression to fit lines for the three eras by plotting the daily sentiment average by day. The level of the line after the 20th CPC National Congress is lower than the fitted line of the sentiment score prior to the Congress. Also, when the Congress's opening day draws near, the mood improves and peaks on that day before suffering a sharp decrease until the event's conclusion.

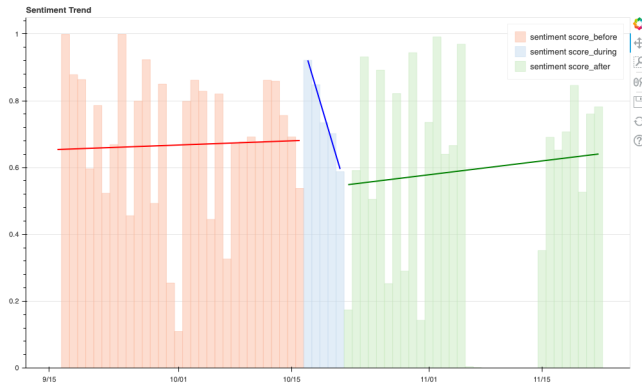


Figure 2: Weibo Posts' Average Daily Sentiment Trend

## 5.2 Topic Modeling

We also tried topic modeling to get a sense of the topics for the time period before, during, and after the 20th CPC National Congress. Topic modeling is a method that requires careful consideration of the goals and context of the analysis. To be specific, we need to set the relevance metric of topic modeling as well as summarize each topic based on the words related to them.

After removing stop words including numbers and unicode, we tokenized the text and got 5 topics for each of the three time periods (before, during, after the 20th CPC National Congress). It is important to note that a good measure of the distance between topics is one that reflects these qualities: If the topics are coherent and distinct, then they should be far apart from each other in the topic space. On the other hand, if the topics are not coherent or overlap with each other, then they may be close together in the topic space. Therefore, a good result of topic modeling in terms of the distance of the topics is one that produces coherent and distinct topics that are far apart from each other in the topic space. As a result, we only chose to show meaningful topics that do not overlap with others. Figure 3, Figure 4 and Figure 5 are some of the results we got.

We set 0.2 as the relevance metric for the three time periods. We found that, for the time before the event, apart from words related to covid e.g. “感染” (infect), “健康” (health), “肺炎” (pneumonia), other words are related to things that are almost hard to summarize. It might indicate topic modeling is not very successful in our case, but we can still use it for reference. For example, some are about celebrities (“鹿晗”) and others are about being calm (“情绪稳定”) and laughing hard (“爆笑”). For the time after the event however, negative words begin to appear, e.g. “杀” (kill), “侵犯” (violate), “毁” (destroy).

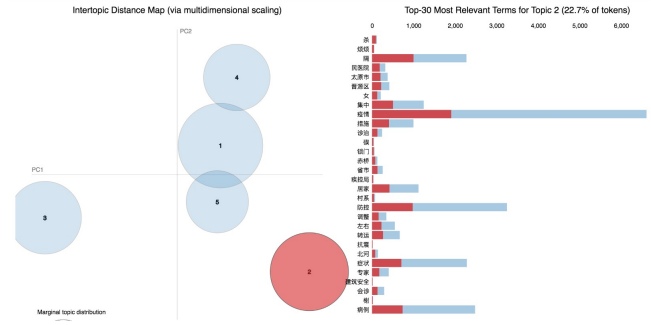


Figure 3: Topic Modeling Result Before the Congress

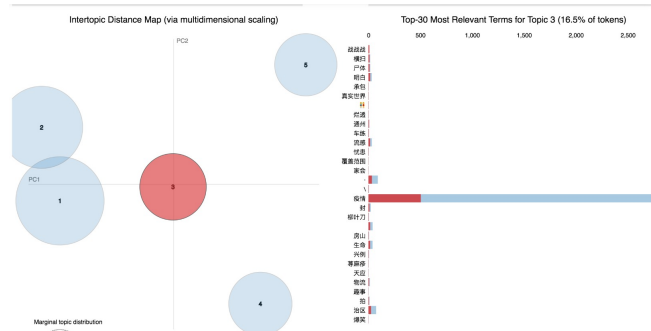


Figure 4: Topic Modeling Result During the Congress

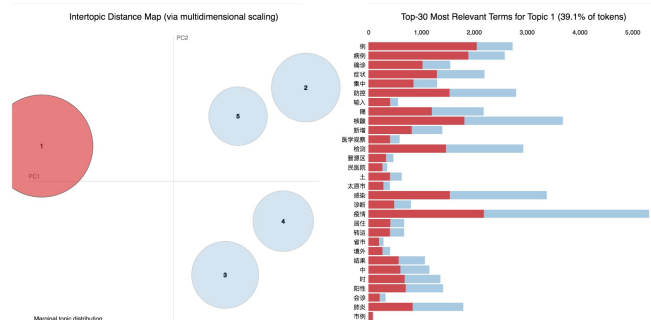


Figure 5: Topic Modeling Result After the Congress

In all, our results could lead to the conclusion that before and during the political event, positive words are frequently used in the content of Weibo. In contrast, after the event, negative words begin to appear.

## 5.3 Word Cloud

To gain insight into the most commonly used words in Weibo posts during different periods, we generated three-word cloud visualizations. The first word cloud depicts the most frequently used words in Weibo posts before the 20th CPC National Congress. It reveals that the predominant tone of these posts was positive. For example, the phrase “中国人受新冠影响全球最小”(Chinese were least affected by Covid in the world) appeared most frequently in Weibo



Figure 6: Word cloud (before the Congress)

posts. This phrase suggests that the Chinese government had effectively controlled the pandemic, resulting in the minimal negative impact on its citizens. Another positive message found in the word cloud is “辽宁好网民抗疫在行动” (Good Liaoning netizens are fighting against Covid), which indicates that people were actively participating in the fight against the pandemic, thereby fostering a positive and harmonious atmosphere on Weibo before the 20th CPC National Congress.

However, this appraising attitude towards the government and people’s proactive reactions to the pandemic in Weibo posts may have been the result of information censorship by the Chinese government, which aimed to maintain a stable political environment, at least online, before the 20th CPC National Congress.



Figure 7: Word cloud (during the Congress)

The second word cloud illustrates the most commonly used words in Weibo posts during the 20th CPC National Congress. Interestingly, the words that appeared most frequently were largely uninformative. For example, “郑州疫情” (epidemic in Zhengzhou, a city in China), “大同疫情” (epidemic in Datong, another city in China), and “居家隔离” (home quarantine) were the most frequently used terms during the 20th CPC National Congress. These words did not provide much substantive information about the pandemic.

Following the 20th CPC National Congress, the third word cloud reveals a shift towards more negative news and language. For example, the phrase “新冠口服药被要求即刻下架” (Covid oral medication was immediately taken down) suggests that the medication may have been problematic. Additionally, words associated



Figure 8: Word cloud (after the Congress)

with criminal activity, such as “扫黑除恶” (crackdown on crimes) and “剪破与黑恶势力交织的关系网” (cutting off the network of mafia-like relationships), were frequently used in Weibo posts after the 20th CPC National Congress. This shift towards more negative language may reflect a loosening of censorship by the Chinese government, allowing for the expression of negative opinions and concerns on Weibo.

## 6 CONCLUSION

To determine how censorship affects public opinion, this paper examines Weibo posts made by users before and after the 20th CPC National Congress (the 20th CPC National Congress). In this study, we examine public posts on Weibo before and after the 20th CPC National Congress using computational approaches.

As our main observational sample, we scrape Weibo posts mentioning the keywords “疫情” (“covid”), “隔离” (“quarantine”), and “抗疫” (“prevention and control of covid-19”) for one month before and after the 20th CPC National Congress (October 16, 2022). We then analyze the change in sentiments, keywords, and topics of these posts to gauge the possibility of Chinese government censorship.

To summarize, our findings indicate that the sentiments expressed in posts rose and peaked in the days preceding the 20th CPC National Congress. After the Congress, the attitudes return and are far more negative than in the month before the event. Our findings from topic modeling suggest that positive phrases are often employed in Weibo posts before and during political event. In contrast, negative words start to surface after the occurrence. From the perspective of word clouds, we find that more encouraging news has emerged before the National Congress (e.g., China is the nation least impacted by Covid-19 globally), which may be an indication of rigorous government censorship. Following the National Congress, there were more unfavorable reports and harsh language (such as anti-crime and evil), indicating that speech controls had started to loosen.

The study shed light on how effectively Chinese censorship methods manage public opinion during important political events, and we discovered that there may be censorship when a significant political event is just around the corner.

## 7 APPENDICES

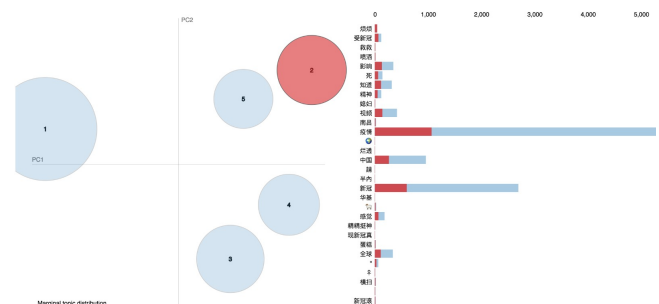




**Figure 12: Topic Modeling Result After the Congress 1**



**Figure 13: Topic Modeling Result After the Congress 2**



**Figure 14: Topic Modeling Result After the Congress 3**

## REFERENCES

- [7] github. (2020). GitHub. Retrieved from [https://github.com/eruong/weibo\\_corpus](https://github.com/eruong/weibo_corpus)