

# “I’ll be back”: Examining Restored Accounts On Twitter

Arnav Kapoor\*

Rishi Raj Jain\*

IIIT Hyderabad

IIIT Delhi

arnav.kapoor@research.iiit.ac.in

rishi18304@iiitd.ac.in

Avinash Prabhu†

Tanvi Karandikar†

IIIT Hyderabad

avinash.prabhu@students.iiit.ac.in

tanvi.karandikar@students.iiit.ac.in

Ponnurangam Kumaraguru

IIIT Hyderabad

pk.guru@iiit.ac.in

## Abstract

Online social networks like Twitter actively monitor their platform to identify accounts that go against their rules. Twitter enforces account level moderation, i.e. suspension of a Twitter account in severe cases of platform abuse. A point of note is that these suspensions are sometimes temporary and even incorrect. Twitter provides a redressal mechanism to ‘restore’ suspended accounts. We refer to all suspended accounts who later have their suspension reversed as ‘restored accounts’. In this paper, we release the first-ever dataset and methodology<sup>1</sup> to identify restored accounts. We inspect account properties and tweets of these restored accounts to get key insights into the effects of suspension. We build a prediction model to classify an account into normal, suspended or restored. We use SHAP values to interpret this model and identify important features. SHAP (SHapley Additive exPlanations) is a method to explain individual predictions. We show that profile features like date of account creation and the ratio of retweets to total tweets are more important than content-based features like sentiment scores and Ekman emotion scores when it comes to classification of an account as normal, suspended or restored. We investigate restored accounts further in the pre-suspension and post-restoration phases. We see that the number of tweets per account drop by 53.95% in the post-restoration phase, signifying less ‘spammy’ behaviour after reversal of suspension. However, there was no substantial difference in the content of the tweets posted in the pre-suspension and post-restoration phases.

## CCS Concepts

• **Human-centered computing** → **Social media**; • **Computing methodologies** → Supervised learning by classification.

## Keywords

Restored Accounts, Suspension, Twitter, User Analysis

<sup>1</sup><https://github.com/rishi-raj-jain/WIAT-Restored-Accounts-On-Twitter>

\*,† denote equal contribution.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WI-IAT ’21, December 14–17, 2021, ESSENDON, VIC, Australia

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-9115-3/21/12...\$15.00

<https://doi.org/10.1145/3486622.3493959>

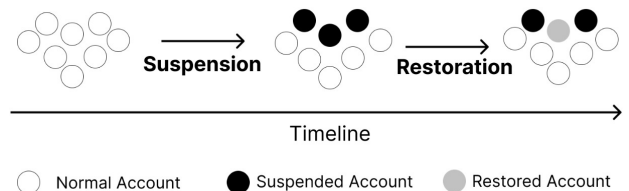
## ACM Reference Format:

Arnav Kapoor\*, Rishi Raj Jain\*, Avinash Prabhu†, Tanvi Karandikar†, and Ponnurangam Kumaraguru. 2021. “I’ll be back”: Examining Restored Accounts On Twitter. In *IEEE/WIC/ACM International Conference on Web Intelligence (WI-IAT ’21)*, December 14–17, 2021, ESSENDON, VIC, Australia. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3486622.3493959>

## 1 Introduction

Over the years, social media platforms have emerged as an effective means to voice our opinions [15]. They have become a critical part of elections around the world [27]. Twitter in particular has been used both by the common public to engage in political discourse and by politicians to reach the voters [14]. Mass misuse of these platforms during recent political events like the 2020 U.S. Presidential Election<sup>2</sup>, including voter manipulation, spam and hate [9], have led to multiple countermeasures by Twitter and other social media platforms. Twitter performs both tweet level moderation (removing a tweet for violating Twitter policy) and account level moderation (suspension of the account). The most recent high profile suspension of Donald Trump’s Twitter account - @realDonaldTrump<sup>3</sup> raised discussions around Twitter’s suspension policy.

We note that some of these suspended accounts are unsuspended (or ‘restored’) after a while and able use the platform again. Twitter policies give us an insight into the reasons for suspension as well as the means to restore the account. Common reasons for suspension by Twitter are spam, account at risk (hacked or compromised account) and abusive tweets. In some suspension cases, Twitter provides a means to restore the account by filing an appeal. In other cases, the suspension is temporary, and the account is automatically restored after a duration. The restored accounts allows us to understand the change in account activity in pre-suspension and post-restoration phases.. Henceforth, we use the term ‘restored’ to refer to any such account initially in the suspended state and later returned to the normal state. Figure 1 shows the life-cycle of a restored account.



**Figure 1: Restored accounts transition from normal to suspended state and then return to the normal state.**

<sup>2</sup>[https://blog.twitter.com/en\\_us/topics/company/2020/2020-election-update](https://blog.twitter.com/en_us/topics/company/2020/2020-election-update)

<sup>3</sup>[https://blog.twitter.com/en\\_us/topics/company/2020/suspension](https://blog.twitter.com/en_us/topics/company/2020/suspension)

Social media platforms can use this study on restored accounts to measure the change in behaviour pre-suspension and post-restoration to understand the effect of suspension. The transition of states from normal to suspended and back for a restored account raises questions about the efficacy of platform moderation. Platform moderation has to thread the fine line between infringement of freedom of speech, and enforcing platform policies and rules [11]. However, these platform moderation methods are black boxes without clear insights into what exact features or behaviour of an account can cause suspension and possible restoration. Hence, we use feature importance values on top of our classifier models to initiate discussion and form baselines to understand the factors that influence Twitter suspension. We do so by identifying the different aspects that characterise normal, suspended and restored accounts.

This paper aims to model and understand the restored accounts and identify what differentiates them from the suspended and normal accounts. Normal accounts are the control group of active accounts which are neither suspended nor restored. We investigate different properties of all accounts - the profile properties, the content tweeted and their interaction with other accounts on the Twitter platform. We then generate features to create a classifier to predict the category of an account. We use SHAP values to get feature importance scores to interpret prediction results. The aim is to better understand what properties set restored accounts apart from normal and suspended accounts. We then focus on the entire timeline of restored accounts (all tweets posted in 2019). We look at the activities and behaviour in two time phases - pre-suspension and post-restoration. To the best of our knowledge, no previous work has focused on the category of restored accounts. We set out to answer the following Research Questions (RQs):

- RQ1 - What factors differentiate the three categories of accounts - normal, suspended and restored?
- RQ2 - What are the changes, if any, in how an account interacts with the platform post-restoration?

For RQ1, we find that specific attributes of accounts like the creation date and frequency of tweets can act as excellent indicators to distinguish between categories, while other content features are not good differentiators. This higher importance of account properties over content features suggests that most account level suspensions are not due to content tweeted but rather the associated meta-information and activity of the account. Additionally, we find that restored accounts resemble normal accounts more than suspended accounts for the majority of explored features.

For RQ2, we find that the number of tweets posted by the restored accounts in the post-restoration phase dropped by an average of 53.95% as compared to the pre-suspension phase. However, the content of the tweets by restored accounts did not change drastically.

In adherence to the FAIR Principle [25] and encouraging collaborations, we also make public the first-ever dataset on these restored accounts and provide a methodology to collect these accounts.

## 2 Background and Related Work

Suspension is one of the most common ways of platform moderation. Previous works studying suspended accounts [1, 5, 22] on Twitter have looked into multiple aspects of these accounts, primarily focusing on spam [22] and bot networks [1]. We look at

the Twitter policy for suspension to give us additional context regarding suspension on the platform. Twitter's policy states that accounts are suspended in the following cases, *'Most of the accounts we suspend are suspended because they are spammy, or just plain fake'*. The policy also says, if an account engages in abusive behaviour, like sending threats to others or impersonating other accounts, Twitter might suspend them. Twitter also gives information related to restoration of suspended accounts. In some cases, the suspension is temporary, and the account is restored to its original state after a duration. An account may also get suspended by mistake - in this case, Twitter works with the account holder to restore the account. The category of restored accounts on Twitter has not been studied in the past.

Past literature has also created models to predict suspended accounts [24] for credibility analysis. We build upon this work to focus on the interpretability of similar prediction models. The interpretable models allow us to identify why platform moderation algorithms consider certain accounts to be less credible and more likely to be suspended.

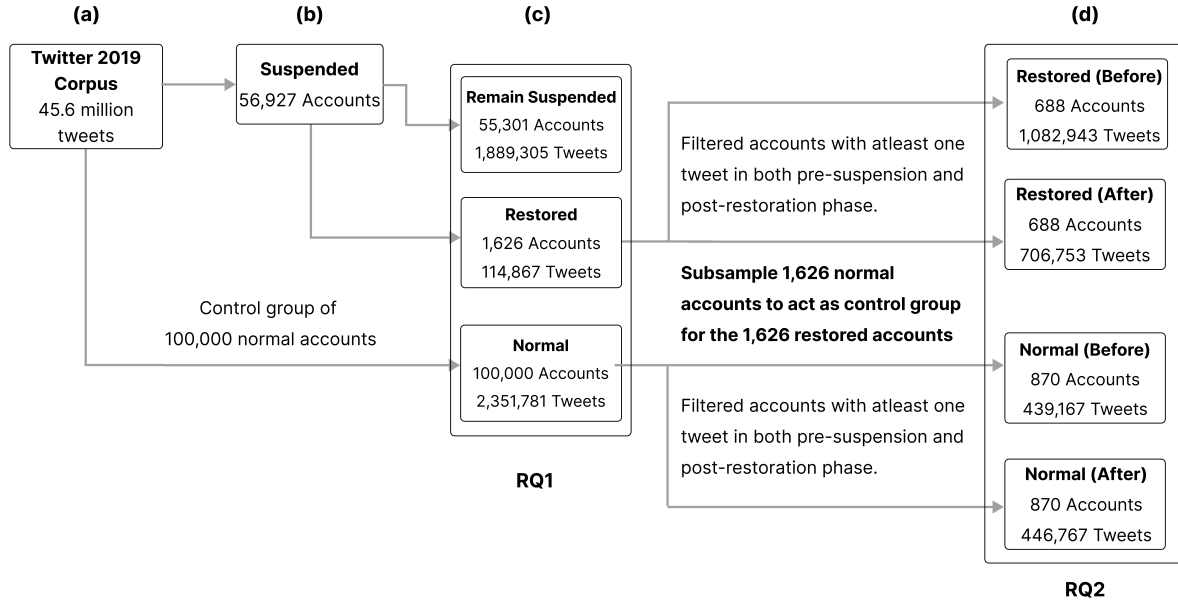
Platform moderation and suspension has been studied in the context of other social media networks; primarily Reddit [4, 20]. The focus is on community-level suspension, i.e. banning of a particular subreddit or a community of accounts, and the forced migration to either different platforms [20] or different communities (subreddits) [4]. Restored accounts on Twitter, however, do not migrate but rather return to the same state. The account returns to the same environment, which allows us to investigate the impact of suspension better as no external factors like change of community or platform is at play. These external factors in migration studies make it hard to identify if the change in account behaviour was due to the suspension event or simply a change of rules of the new community or platform. Our interpretability studies give insights into the cause of suspension and tell us about salient features of different account categories - normal, restored and suspended.

## 3 Dataset

In this section, we describe our data collection process to collect information regarding restored accounts on Twitter. The entire data collection pipeline is summarized in Figure 2.

**RQ1 Dataset - Identify Categories of Accounts** To identify suspended and, in turn, restored accounts, we use the 'Analysis of General Elections 2019 in India' (AGE2019) dataset [10]. This dataset has been created by collating a large corpus of Twitter accounts that tweeted during the 2019 Indian general election. We specifically chose election data for our study since manipulation and misuse of the Twitter platform is rampant during an election, leading to a more significant proportion of suspended (and consequently restored) accounts [17].

The AGE2019 dataset has about 45.6 million tweets made by 2.2 million unique accounts over six months from February 5, 2019 (two months before first polling) to June 25, 2019 (one month after election results). Some of the accounts which tweeted during this phase were suspended by Twitter. The AGE2019 dataset also released a list of 56,927 accounts (2.8% of the total 2.2 million) which were found to be suspended as of June 29, 2019. These were identified to be suspended as the Twitter API returns *code 63* for suspended accounts. The AGE2019 dataset also comes with a list of 100,000



**Figure 2:** This flowchart depicts the different components of our research approach. We used the corpus of the suspended and normal accounts from the AGE2019 dataset. We then queried all the suspended accounts, and found that 1,626 accounts were active again on the platform (referred to as restored accounts). We also fetch all the tweets made by the restored accounts in two phases: pre-suspension (1 Jan, 2019 - 29 Jun, 2019) and post-restoration (29 Jun, 2019 - 31 Dec, 2019). Further, 1,626 normal accounts were subsampled to act as the control group for the restored accounts. Tweets for this control group of normal accounts were also collected for the two phases. For restored and control group of normal accounts, we filtered the accounts with atleast one tweet in both pre-suspension and post-restoration phase.

normal accounts (*control group*) – active accounts on Twitter as of June 29, 2019.

To identify restored accounts we needed to find accounts that reverted to the normal state from the suspended state, hence we rechecked the status of the 56,927 suspended accounts. We queried the Twitter API on October 8, 2019 (around 3 months after June 29, 2019) for each of these accounts. Out of these 56,927 suspended accounts, 1,626 were active and restored to the normal state. The Twitter API no longer returned *code 63* (or any other errors) for these 1,626 accounts. Table 1 contains the summary of the different types of accounts used in our study.

Category Of Accounts	Number of Accounts	# of Tweets
Normal	100,000	2,351,781
Suspended	55,301	1,889,305
Restored	1,626	114,867

**Table 1: Statistics of normal, suspended and restored accounts identified in the context of Indian general elections 2019. (Used in RQ1)**

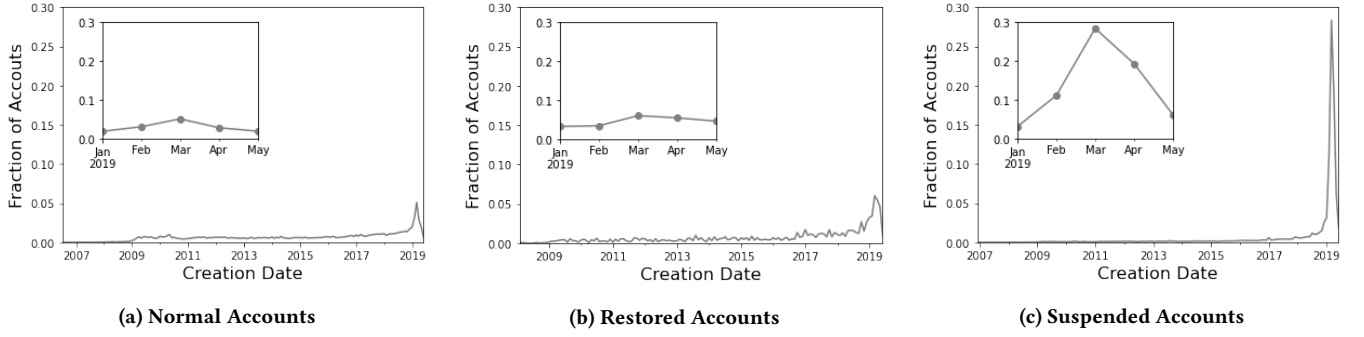
**RQ2 Dataset - Timeline of Restored Accounts** To get a more holistic view of the restored accounts, we collected tweets made by the restored accounts for the entirety of 2019 instead of just the political tweets in the election period. We crawled the entire timeline of these 1,626 restored accounts and collected a year’s

worth of data from Jan 1, 2019, to Dec 31, 2019, and got a total of about 1.78 million tweets made by these accounts. The complete timeline is a more unbiased representation of the accounts as it contains both election and non-election tweets.

We also crawled the entire timeline for 1,626 (equal to the number of restored accounts) randomly sampled normal accounts for the entirety of 2019. The normal accounts act as our control group while analysing restored accounts. We cannot collect the timelines for suspended accounts as their tweets are no longer available. We now have an extensive list of accounts and their corresponding tweets. We then split the timeline of restored accounts into two phases; pre-suspension (tweets before June 29, 2019) and post-restoration (tweets after June 29, 2019). We did this to quantify the impact of suspension and see how it affected account behaviour and its usage of the platform. We also split the control group of normal accounts along the exact dates. Additionally, we filtered out and kept only those accounts that had tweets in both the phases of pre-suspension and post-restoration. The final set of accounts and tweets are summarised in Table 2.

	Normal	Restored
Number of Accounts	870	688
# of Tweets Pre-Suspension	439,167	1,082,943
# of Tweets Post-Restoration	446,767	706,753

**Table 2: Statistics of normal and restored accounts for all tweets made in 2019. (Used in RQ2)**



**Figure 3: Creation timeline of (a) normal accounts, (b) restored accounts, and (c) suspended accounts. Suspended accounts had a far larger proportion of accounts created in March-April 2019 as compared to the normal and restored accounts.**

## 4 RQ1 - Comparing Account Categories

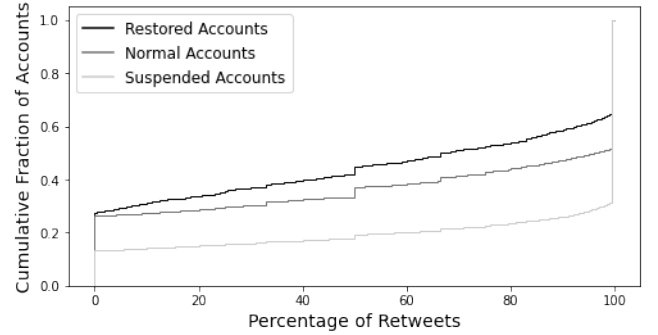
To address RQ1, we investigate the properties of restored accounts and identify features that distinguish between the three categories of accounts (normal, suspended and restored). We also conduct comparisons between these categories. To understand these categories of accounts comprehensively, we use three dimensions of analysis in our study. These are the account properties [4.1], the content tweeted [4.2] and interaction with other accounts on Twitter [4.3]. We finally model features using the above properties and combine them to create a classifier to predict the account type. The focus while creating the classifier is on creating an interpretable model to shed some light on the relatively opaque moderation methodologies and identify features that distinguish normal, suspended and restored accounts.

### 4.1 Account Properties

One of the established aspects of suspended accounts is that they are very short-lived compared to normal accounts [22], which led us to look at account creation dates. We found a stark difference between the creation dates of suspended and normal accounts. We found that more than 47.60% of the suspended accounts were created in the two months preceding the election (March and April 2019). In comparison, only 7.90% of the normal accounts were created in this period (observe Figures 3a and 3c).

Similar to the trend observed in normal accounts, only 11.50% of restored accounts were created in March and April 2019. This shows that restored accounts resemble normal accounts when it comes to account creation patterns (observe Figures 3a and 3b). The sharp peak in March 2019 for suspended accounts suggests that there was a bulk creation of accounts before the election that ended up being suspended (observe Figure 3c).

We now look at each account's proportion of retweets to tweets to see what an account does more – retweeting or tweeting original content. We plot cumulative distributions (Figure 4) of the percentage of retweets. Suspended accounts use retweets more than normal and restored accounts. About 70% of all suspended accounts only retweet and have zero unique tweets. On the other hand, the percentage of accounts that only retweet are much lower for normal (48%) and restored (35%). Overall, we see a sharp distinction in the proportion of retweets, with suspended accounts using retweets far more frequently as opposed to normal and restored accounts.

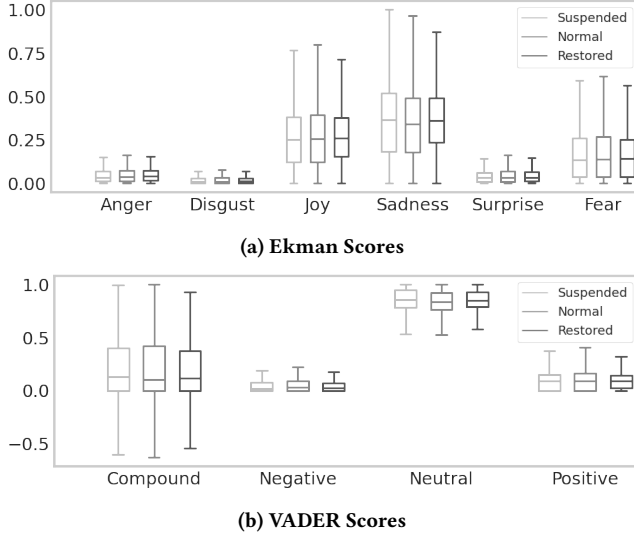


**Figure 4: CDF plots of percentage of retweets. Suspended account retweet more than normal and restored accounts.**

We further analyse other account properties such as no. of followers, no. of following, favorite count, etc. as part of our interpretable classifier in Section 4.4. The exhaustive list of account properties analysed is listed in Table 5. Most of these features had low predictive power as evident in 6.

### 4.2 Content Characteristics

In this section, we examine the content of the tweets made by the accounts. We apply language detection to filter and keep only English tweets for content analysis. We looked at two properties – Ekman emotion scores and VADER sentiment scores. We use these properties to further analyse the differences between the classes of accounts. Ekman emotions provide scores for six basic emotions [7]. We used emotion recognition models [6] to compute the Ekman scores depicted in Figure 5a. We notice that across the emotions, there was no significant difference between the three groups – normal, suspended and restored. To verify this lack of significant difference, we also looked at the VADER sentiment scores for the three groups. VADER (Valence Aware Dictionary and sEntiment Reasoner) [12] is a rule-based sentiment analyser. The VADER Scores shown in Figure 5b also show no significant difference across the groups, indicating that the content characteristics were not a factor in distinguishing the category of an account.



**Figure 5: Box plots for Ekman and VADER scores for normal, restored and suspended accounts. Observe how all values almost remain the same across classes.**

### 4.3 Network Effects

In this section, we look at how suspended and restored accounts interact on the platform, and detect communities to identify group level characteristics. We use the Leiden algorithm to form the retweet and user mentions communities, and further observe how the interaction between accounts varies intra-community vs inter-community.

A community is a subset of nodes within the network such that connections between the nodes are denser than connections with the rest of the network [19]. While several factors may show an account’s involvement in a certain set of groups/activities, explicit user interaction is visible in retweets and user mentions. An interaction is defined either as an account X mentioning account Y (mentions network), or an account X retweeting a tweet by account Y (retweet network). The nodes of such a network are accounts, and the directed edges represent retweets or user mentions. We used the Leiden algorithm [23] to identify a set of well-connected communities that are guaranteed to be locally optimally assigned. In general, the Leiden algorithm when applied, iteratively creates communities such that the intra-community connections are guaranteed to be more concentrated when compared to inter-community. To measure the strength of division of a network into communities, we use the modularity score. Networks with high modularity have dense connections between the nodes within communities but sparse connections between nodes in different communities. We construct two networks – the user mention network and retweet network. These networks are constructed for all three categories of accounts.

In the retweet network, we observe a greater modularity for suspended accounts (0.746) as compared to normal (0.681) and restored accounts (0.638) (observe Table 3). This indicates a greater separation between the communities of suspended accounts. We also see a large number of suspended retweet communities being formed, hence denoting a larger spread of the information and

Retweets Network	Normal	Suspended	Restored
Modularity	0.681	0.746	0.638
Number of Communities	7,190	11,748	299
Largest Community Size	15.88 %	11.55 %	15.96 %

**Table 3: Retweets network statistics of normal, suspended and restored accounts. Largest community size for the suspended class is smaller than the largest community size for restored or normal class.**

User Mentions Network	Normal	Suspended	Restored
Modularity	0.550	0.575	0.546
Number of Communities	2,488	590	117
Largest Community Size	14.82 %	24.96 %	17.38 %

**Table 4: User mentions network statistics of normal, suspended and restored accounts. The modularity is comparable across categories of accounts. In contrast to retweet network (Table 3), suspended accounts have the largest community size in the user mention network.**

possibly a means to engage more accounts into the agenda via retweets.

For the user mention network, we see that all categories of accounts have similar modularity (observe Table 4). Similar modularity indicates that the network structure is similar with respect to separation across communities and interaction within communities.

### 4.4 Classifier

We used an interpretable classifier to further understand the role of the three aspects – account properties, content tweeted, and interactions within the network. We combine features from these aspects for training the classifier to see which features play a crucial role in differentiating between the three classes. Features used from each aspect are shown in Table 5. The dataset used to train the classifier contains 55,301 suspended accounts, 100,000 normal accounts and 1,626 restored accounts. There is a severe class imbalance, with the minority class (restored accounts) being 1.626% of the majority class (normal accounts).

Machine learning models are negatively affected when trained on datasets with class imbalance [3, 26]. Class imbalance refers to the issue where the number of samples from a particular class is much lower than samples from other classes. Many methods deal with class imbalance using data sampling methods. Two such methods are Random Over Sampling (ROS), which duplicates samples from the minority class, and Random Under Sampling (RUS), which eliminates samples from the majority class. Both these methods can cause issues like bias the model [8]. ROS, by oversampling may cause the model to overfit on the training data and RUS, by undersampling, may remove critical information from the dataset which is not ideal for our use case [16]. A heuristic over-sampling algorithm, Synthetic Minority Over-Sampling Technique (SMOTE), produces artificial samples from the minority class by interpolating existing instances which are very close together.  $K$  intraclass nearest neighbours are found for each sample in the minority class

**Account Properties**

Time since account creation, No. of friends, No. of followers,  
Total no. of likes, Ratio of friends to followers, Name len in chars,  
Bio len in chars, Screen name len in chars,  
Screen name len in words, Bio len in words, No. of tweets,  
Avg. no. of tweets per hour, Avg. tweet gap, Median tweet gap,  
Avg. tweet len in words, Avg. tweets len in chars,  
Proportion of retweets to total tweets (RT\_rate\_proportion),  
Rate of capital chars

**Tweet Content**

Median Anger, Mean Anger, Median Disgust, Mean Disgust,  
Median Fear, Mean Fear, Median Joy, Mean Joy,  
Median Sadness, Mean Sadness, Median Surprise, Mean Surprise,  
Median Compound, Mean Compound, Median Negative,  
Mean Negative, Median Neutral, Mean Neutral,  
Median Positive, Mean Positive

**Network**

No. of hashtags, No. of unique hashtags, No. of mentions,  
No. of unique mentions, No. of retweets, Avg. Hashtag count,  
Avg. unique hashtags count, Avg. mentions, Avg. unique mentions

**Table 5: Three categories of features (Account Properties, Tweet Content and Network) used to train the classifier.**

and synthetic samples are found in the direction of some or all of the nearest neighbours [2].

We used several classifier models such as Random Forest Classifier, Gradient Boosting Classifier, XGB Classifier and LGBM Classifier in combination with the above-mentioned sampling techniques. We use the best performing combination (XGB Classifier with SMOTE) for further analysis. The F1 scores for pairwise separation of the three classes are shown in Table 6.

#### 4.5 Explainability and Feature Importance

The Shapley value is a solution concept of fairly distributing gains and costs to several actors working in a coalition [21]. In the context of the explainability of machine learning models, the Shapley value allows us to identify which features contribute more to the final prediction of the model. However, the complexity constraints of using Shapley values makes it not feasible on larger datasets. To solve this problem, we used Shapley Additive Explanations, or SHAP values [18], which capture the contribution of each feature based on local explanations and principles of game theory. It uses approximate algorithms to predict the Shapley values. The SHAP plots in Figure 6 represent the 20 most important features of our model sorted by the mean absolute value of the SHAP values of each feature. Here, SHAP values are used to interpret and understand which features play an important role in determining the output of each classifier.

*Restored vs Normal Classifier:* When distinguishing between restored and normal accounts, the SHAP analysis from Figure 6a shows us that higher average tweet length in characters and higher number of friends are signs of restored accounts. On the other hand, normal accounts are characterised by higher time since account creation, higher retweet rate proportion and higher average tweet

length in words. Out of the top 20 features, 60% are account property based, 30% are tweet content based, and 10% are network based.

*Suspended vs Normal Classifier:* When distinguishing between suspended and normal accounts, the SHAP analysis from Figure 6b shows us that higher time since account creation, higher name length (in characters) and higher screen name length (in characters) are signs of normal accounts. When it comes to suspended accounts, a higher number of average tweets per hour and higher average tweet length (in characters) are signs of suspended accounts. The SHAP analysis also shows us that out of the top 20 features, 70% are account property based, 20% are network based, and 10% are tweet content based.

*Suspended vs Restored Classifier:* When distinguishing between suspended and restored accounts, the SHAP analysis from Figure 6c shows us that higher time since creation, higher number of mentions and higher number of tweets are signs of restored accounts. Suspended accounts have higher retweet rate proportion, higher average unique hashtags count, higher average number of tweets per hour and higher compound sentiment are signs of suspended accounts. The SHAP analysis also shows us that out of the top 20 features, 45% are account property based, 30% are tweet content based, and 25% are network based.

Overall, we find that account properties often have the highest SHAP value and thus affect the model’s output the most. Most notably, time since account creation and retweet rate proportion (i.e. ratio of number of retweets to tweets made by the account) have the highest impact on the outputs of the three classifiers. We find that content and network features are low in number or have very low SHAP values across all classifiers. This suggests that the reason behind suspending or restoring an account is not based on the tweet content or network properties, but rather the account properties. The lack of importance of the content features is in line with our findings in Section 4.2.

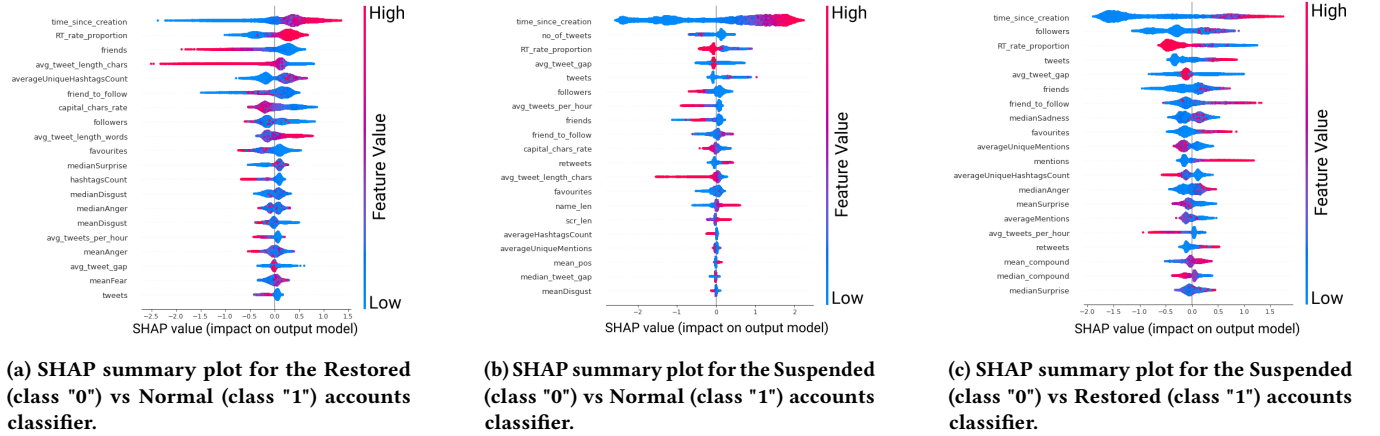
Comparison	F1-Score
Restored vs Normal	0.71242
Suspended vs Normal	0.84008
Suspended vs Restored	0.72282

**Table 6: F1 scores for pairwise prediction of account categories.**

#### 5 RQ2 - Timeline of Restored Accounts

Aim of RQ2 was to understand what the impact of suspension was on these restored accounts. Did a ban (albeit wrong or temporary) affect the way they interacted with the platform when they came back? To critically understand the impact of the suspension event, we looked at the entire timeline of the restored accounts and split it into two phases – pre-suspension and post-restoration. We also created a control group of accounts from the set of normal accounts to compare the change in behaviour. To get a better understanding, we looked at two aspects – the amount of activity on the platform and the content produced. We study the difference in these aspects across two phases (pre-suspension and post-restoration).

**Changes in activity levels** - In order to attribute the change in a property of an account to suspension, there must be a measurable



**Figure 6:** Here, we show the SHAP summary plots generated for the three classifier models. The features in each subfigure are ranked by mean absolute value of the SHAP values for each feature. The feature value is color coded; red means that the feature value is high and blue means that the feature value is low. The X-axis gives the SHAP value for each feature value. A positive SHAP value means that the corresponding feature value is pushing the output towards class "1" and a negative SHAP value means that the corresponding feature value is pushing the output towards class "0".

difference in that property pre-suspension and post-restoration. The properties we look at are the number of tweets posted per day per account (activity on the platform) and the different content characteristics of the tweets (Ekman emotions and VADER scores) created by the accounts. Now, to model a property, we are concerned with two aspects – absolute value (i.e. what the current state is) and trend (i.e. moving upwards or downwards). Hence, we fit a simple straight line to any property. The y-intercept of the line gives us the absolute value while the slope of the line gives us the trend. We perform a regression discontinuity analysis [13] to measure the difference in properties between the pre-suspension and post-restoration phases. We use two linear models:

$$y_t = \alpha_0 + \beta_0 t \quad (t < 0) \quad (1)$$

$$y_t = \alpha_1 + \beta_1 t \quad (t > 0) \quad (2)$$

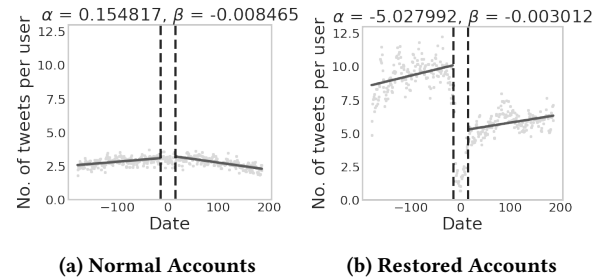
where  $t$  is the date ranging from -179 to 186 ( $t = 0$  represents June 29, 2019, i.e. the date as of which the accounts were suspended).  $y_t$  represents the statistic we are modelling (tweets posted per account, Ekman and VADER scores). Equation 1 gives us the equation of the line pre-suspension and equation 2 gives us the equation of the line post-restoration. These models assume that we can approximate the various data points as a straight line defined by  $\alpha_0$  and  $\beta_0$  pre-suspension and  $\alpha_1$  and  $\beta_1$  post-restoration.

Since, we are interested in studying the change in behaviour before and after suspension, we define two additional terms:  $\alpha$  and  $\beta$  to represent this change.  $\alpha$  is the change in the y-intercept post-restoration (ie.  $\alpha = \alpha_1 - \alpha_0$ ) and  $\beta$  is the change in the slope post-restoration (ie.  $\beta = \beta_1 - \beta_0$ ).  $\alpha$  represents the dip or rise in the absolute value caused by suspension, while  $\beta$  represents the change in long-term trend. We further exclude data from a grace period before and after suspension to account for the bursty behaviour occurring in the days around suspension.

**Amount of activity:** Figures 7a and 7b show the total number of tweets per day normalized by the number of accounts in that class before and after suspension. The normal accounts saw

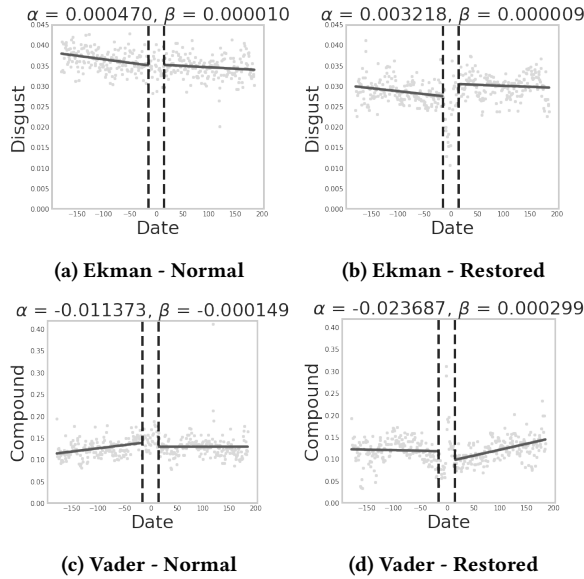
an average increase of around 0.15 tweets per day ( $\alpha = 0.154817$ ) after the suspension date which represents a percentage change of around +5.5% of the Mean Value Before Suspension (henceforth referred to as MVBS). On the other hand, the restored accounts saw an average decrease of about 5 tweets per day ( $\alpha = -5.027992$ ). This represents a percentage decrease of 53.95% of the MVBS. This significant decrease shows that suspension does indeed have an effect on the number of tweets posted by restored accounts. We do not observe any significant change in the long-term trend ( $\beta$ ).

**Content pushed:** Figures 8a, 8b, 8c and 8d show the average VADER and Ekman score respectively for each day before and after suspension, for normal and restored accounts. Across the six basic Ekman emotions and the four VADER sentiment classes we observe little change in behaviour. For brevity, the figure contains only Disgust (Ekman basic emotion) and Compound sentiment (VADER sentiment class), but the same was observed across all Ekman emotions and VADER sentiment scores.



**Figure 7: Activity Levels:** The daily number of tweets tweeted by (a) normal accounts and (b) restored accounts. Date "0" (29th June 2019). A grace period of 15 days before and after suspension are represented by the dotted vertical lines. On top of each subplot, we report the coefficients associated with the suspension policy ( $\alpha$  and  $\beta$ ).

On observing the figures, we can see that there is a negligible difference in  $\alpha$  and  $\beta$  values before and after the suspension for both normal and restored accounts. This is in line with our observations from section 4.5 and section 4.2 which show that content features do not differ much between the classes.



**Figure 8: Content Levels: The daily Ekman (Disgust) and VADER (Compound) scores of tweets made by normal and restored accounts. The grace period and coefficients have the same meaning as Figure 7. The pre-suspension and post-restoration properties for restored accounts are similar for both Ekman (b) and Vader (d).**

## 6 Discussion

Restored accounts provide a new dimension to study the effect of suspension on Twitter. The most recent (9th August, 2021) case of Rep. Marjorie Taylor Greene’s one-week suspension<sup>4</sup>, highlights the importance of studying restored accounts. We can gauge the impact of suspension by looking at changes in tweet activity and behaviour pre-suspension and post-restoration. Twitter can harness such analysis to understand the effectiveness of suspension on their platform and come up with customised suspension policies and duration of suspension to maximise effectiveness. Our work on the interpretable prediction models provides a baseline to understand the features that can impact suspension and restoration on Twitter. Restored accounts are a relatively small proportion of total accounts. To overcome this limitation we plan to scale our data collection and examine a larger more varied corpus of accounts. We also plan to expand this work to other platforms as the concept of suspension is ubiquitous across social networks.

## References

[1] Noor Abu-El-Rub et al. 2019. BotCamp: Bot-driven Interactions in Social Campaigns. In *The World Wide Web Conference on - WWW '19*. ACM Press, San Francisco, CA, USA, 2529–2535.

<sup>4</sup><https://edition.cnn.com/2021/08/10/tech/twitter-marjorie-taylor-greene/index.html>

- [2] Kevin W. Bowyer et al. 2011. SMOTE: Synthetic Minority Over-sampling Technique. *CoRR* abs/1106.1813 (2011). arXiv:1106.1813 <http://arxiv.org/abs/1106.1813>
- [3] Mateusz Buda et al. 2017. A systematic study of the class imbalance problem in convolutional neural networks. *CoRR* abs/1710.05381 (2017). arXiv:1710.05381 <http://arxiv.org/abs/1710.05381>
- [4] Eshwar Chandrasekharan et al. 2017. You Can’t Stay Here: The Efficacy of Reddit’s 2015 Ban Examined Through Hate Speech. *Proceedings of the ACM on Human-Computer Interaction* 1, CSCW (Dec. 2017), 1–22.
- [5] Farhan Asif Chowdhury et al. 2020. On Twitter Purge: A Retrospective Analysis of Suspended Users. In *Companion Proceedings of the Web Conference 2020*. ACM, Taipei Taiwan, 371–378.
- [6] Niko Colneric and Janez Demsar. 2020. Emotion Recognition on Twitter: Comparative Study and Training a Unison Model. *IEEE Transactions on Affective Computing* 11, 3 (July 2020), 433–446.
- [7] Paul Ekman. 1992. An argument for basic emotions. *Cognition and Emotion* 6, 3–4 (May 1992), 169–200.
- [8] Victor González-Bárceñas, Erendira Rendón, Roberto Alejo, Everardo Granda-Gutiérrez, and Rosa Valdovinos. 2019. Addressing the Big Data Multi-class Imbalance Problem with Oversampling and Deep Learning Neural Networks. 216–224.
- [9] Shashi Gupta. 2019. What makes Twitter so vulnerable to controversies in India - ET BrandEquity.
- [10] Saurabh Gupta et al. 2020. #IVoted to #IGotPwned: Studying Voter Privacy Leaks in Indian Lok Sabha Elections on Twitter. (2020).
- [11] Jeffrey W. Howard. 2019. Free Speech and Hate Speech. *Annual Review of Political Science* 22, 1 (2019), 93–109. arXiv:https://doi.org/10.1146/annurev-polisci-051517-012343 <https://doi.org/10.1146/annurev-polisci-051517-012343>
- [12] C. Hutto and Eric Gilbert. 2014. VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. *Proceedings of the International AAAI Conference on Web and Social Media* 8, 1 (May 2014), 216–225.
- [13] Guido W. Imbens and Thomas Lemieux. 2008. Regression discontinuity designs: A guide to practice. *Journal of Econometrics* 142, 2 (2008), 615–635. The regression discontinuity design: Theory and applications.
- [14] Andreas Jungherr. 2016. Twitter use in election campaigns: A systematic literature review. *Journal of Information Technology & Politics* 13, 1 (2016), 72–91. <https://doi.org/10.1080/19331681.2015.1132401>
- [15] Kawaljeet Kaur Kapoor et al. 2018. Advances in Social Media Research: Past, Present and Future. *Information Systems Frontiers* 20, 3 (June 2018), 531–558.
- [16] M. Kubat. 2000. Addressing the Curse of Imbalanced Training Sets: One-Sided Selection. *Fourteenth International Conference on Machine Learning* (06 2000).
- [17] Huyen Le, G. R. Boynton, Zubair Shafiq, and Padmini Srinivasan. 2019. A Post-mortem of Suspended Twitter Accounts in the 2016 U.S. Presidential Election. In *Proceedings of the 2019 IEEE/ACM International Conference on ASONAM* (Vancouver, British Columbia, Canada) (ASONAM ’19). Association for Computing Machinery, New York, NY, USA, 258–265. <https://doi.org/10.1145/3341161.3342878>
- [18] Scott M Lundberg et al. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems* 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 4765–4774.
- [19] Filippo Radicchi et al. 2004. Defining and identifying communities in networks. *Proceedings of the National Academy of Sciences* 101, 9 (March 2004), 2658–2663.
- [20] Manoel Horta Ribeiro et al. 2020. Does Platform Migration Compromise Content Moderation? Evidence from r/The\_Donald and r/Incels. *arXiv:2010.10397 [cs]* (Oct. 2020). arXiv: 2010.10397.
- [21] L. S. Shapley. 1953. 17. A Value for n-Person Games. In *Contributions to the Theory of Games (AM-28), Volume II*, Harold William Kuhn and Albert William Tucker (Eds.). Princeton University Press, 307–318.
- [22] Kurt Thomas et al. 2011. Suspended accounts in retrospect: an analysis of twitter spam. In *Proceedings of the 2011 ACM SIGCOMM (IMC '11)*. Berlin, Germany, 243–258.
- [23] V. A. Traag et al. 2019. From Louvain to Leiden: guaranteeing well-connected communities. *Scientific Reports* 9, 1 (2019), 5233.
- [24] Svitlana Volkova and Eric Bell. 2017. Identifying Effective Signals to Predict Deleted and Suspended Accounts on Twitter Across Languages. *Proceedings of the International AAAI Conference on Web and Social Media* 11, 1 (May 2017), 290–298.
- [25] Mark D. Wilkinson et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data* 3, 1 (Dec. 2016), 160018.
- [26] Kaixiang Yang et al. 2020. Hybrid Classifier Ensemble for Imbalanced Data. *IEEE Transactions on Neural Networks and Learning Systems* 31, 4 (2020), 1387–1400. <https://doi.org/10.1109/TNNLS.2019.2920246>
- [27] Ekaterina Zhuravskaya et al. 2020. Political Effects of the Internet and Social Media. *Annual Review of Economics* 12, 1 (Aug. 2020), 415–438.