

Wright State University

## CORE Scholar

---

Symposium of Student Research, Scholarship,  
and Creative Activities Materials

Office of the Vice Provost for Research

---

4-2020

### Anthrax Event Detection Using Twitter: Analysis of Unigram and Bigrams for Relevant vs Non-Relevant Tweets

Michele Miller

Wright State University - Main Campus, miller.1232@wright.edu

William L. Romine

Wright State University, william.romine@wright.edu

Follow this and additional works at: [https://corescholar.libraries.wright.edu/urop\\_celebration](https://corescholar.libraries.wright.edu/urop_celebration)



Part of the [Other Microbiology Commons](#), and the [Science and Technology Studies Commons](#)

---

#### Repository Citation

Miller , M., & Romine , W. L. (2020). *Anthrax Event Detection Using Twitter: Analysis of Unigram and Bigrams for Relevant vs Non-Relevant Tweets.* .

This Poster is brought to you for free and open access by the Office of the Vice Provost for Research at CORE Scholar. It has been accepted for inclusion in Symposium of Student Research, Scholarship, and Creative Activities Materials by an authorized administrator of CORE Scholar. For more information, please contact [library-corescholar@wright.edu](mailto:library-corescholar@wright.edu).

# Anthrax Event Detection Using Twitter: Analysis of Unigram and Bigrams for Relevant vs Non-Relevant Tweets

Michele E. Miller; Dr. William L. Romine; and Dr. Terry Oroszi  
Wright State University

## PROBLEM

Due to the lack of anthrax attacks in recent times, researchers have used naturally occurring events to assess their anthrax detection models, but these provide little information on how the models will perform in the context of an unannounced, intentional release of a bioterrorism agent, like anthrax (Nordin et al., 2005). Therefore it is important to develop a detection model using data surrounding real anthrax scares and events.

We develop a methodology to detect an anthrax-related event on Twitter. We describe a process to separate the tweets concerning anthrax-related events from those not related so experts can address misconceptions and fears in real-time.

## METHODS

Tweets in English containing the keyword "Anthrax" and "*Bacillus anthracis*" were collected from 9/25/2017 through 8/15/2018 using a crawling algorithm to collect tweets containing the keywords used in real-time. Data collected included the text of 204,008 tweets as well as the date and time the tweet was posted. A line graph of the total number of was created to show how tweets fluctuate daily (Signorini et al., 2011) and to quickly visualize on what days an event occurred.

A feature is a measurable characteristic or property of an observed phenomenon (Bishop, 2006). In other words, features are typically numeric phenomena used to help classify data, or in this case tweets, into the categories of interest to the researcher. Twitter specific features include hashtag (#), URL, re-tweet (RT), and at-mentions (@). These features were coded using binary coding to represent the presence or absence of the feature. Tweets were then preprocessed by removing URLs, non-ascii characters, #, @, punctuation, and retweet indicators. Capital letters would also be changed to lowercase and the words anthrax and *Bacillus anthracis* would be removed since every tweet will contain it. Tweets will then be further processed by removing single letters such as "a", "y", "s", etc, extra spaces, and stop words. Tweets were also stemmed so that only the root of the word remained. General features are features that are not unique to Twitter. These include parts of speech (POS) and n-grams. POS included adjectives, singular nouns, verbs, determiners, prepositions, personal pronouns, predeterminers, and adverbs.

To develop the gold standard, a random sample of 5000 unique tweets were annotated by three experts in terrorism as relevant versus not.

## FUTURE WORK

Logistic regression, naïve Bayes, support vector machine, and random forest will be trained using the golden standard. The remaining 195,000 tweets will then be coded as relevant versus not using the best performing algorithm.

The methods used to classify tweets as relevant versus not will then be used to classify tweets as concerning an anthrax related event or not. Topic modeling will then be used to determine topics of discussion of tweets concerning anthrax related events. The goal is to determine topics of fear in the general public so experts can address them in real time.

## WORKS CITED

## RESULTS AND DISCUSSION

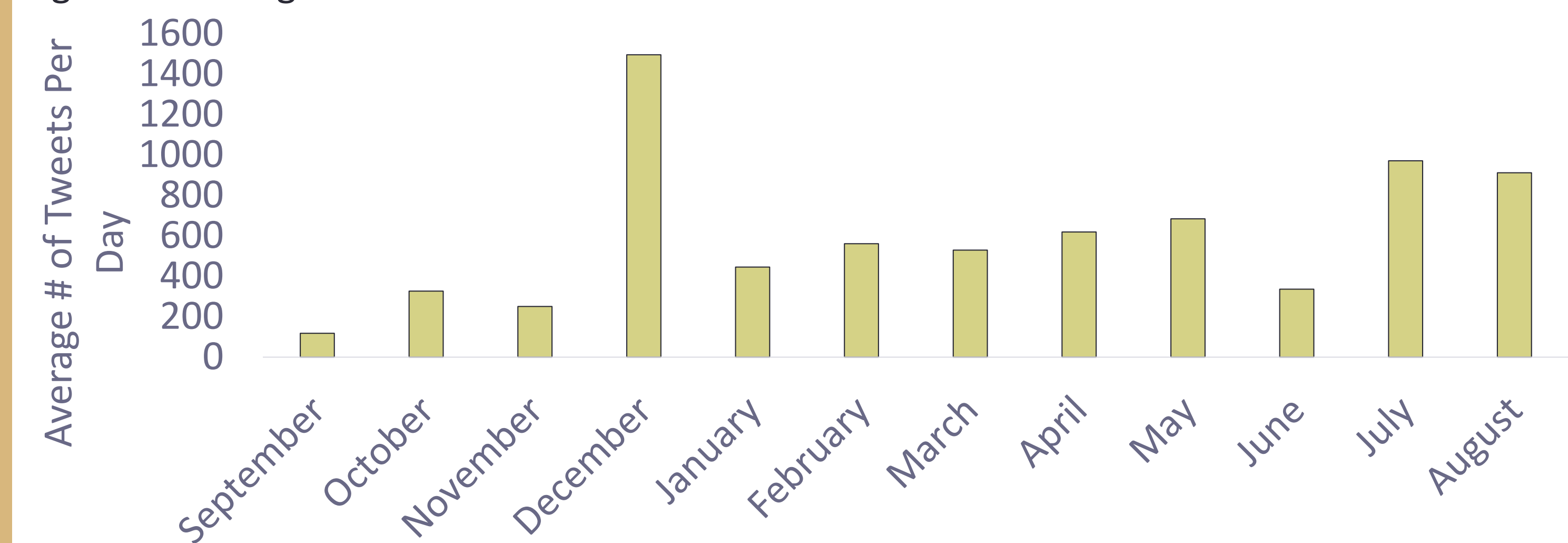
A tweet was considered relevant if the focus of the tweet was about *Bacillus anthracis*. A tweet was considered not relevant if the focus was about the band Anthrax, or something else not about the bacterium. Some examples of relevant and not relevant tweets are seen in Table 1.

Table 1. Examples of relevant versus not relevant tweets.

Relevant	Not relevant
1. RT @catoletters: Remind me again, why did DC invade Iraq? Yellow cake and Nuclear weapons? Anthrax and Bio weapons? 9/11 Saudis?	1. Put anthrax on a Tampax\nAnd slap you 'til you can't stand\nGirl, you just blew your chance\nDon't mean to ruin your plans
2. Elusys Initiates Third Clinical Safety Study of Anthim a New Anthrax Treatment; Company Successfully Completes... <a href="https://t.co/Cj5F7IN7RH">https://t.co/Cj5F7IN7RH</a>	2. Anthrax - In The End Official <a href="https://t.co/M8oSYYPC6Q">https://t.co/M8oSYYPC6Q</a>

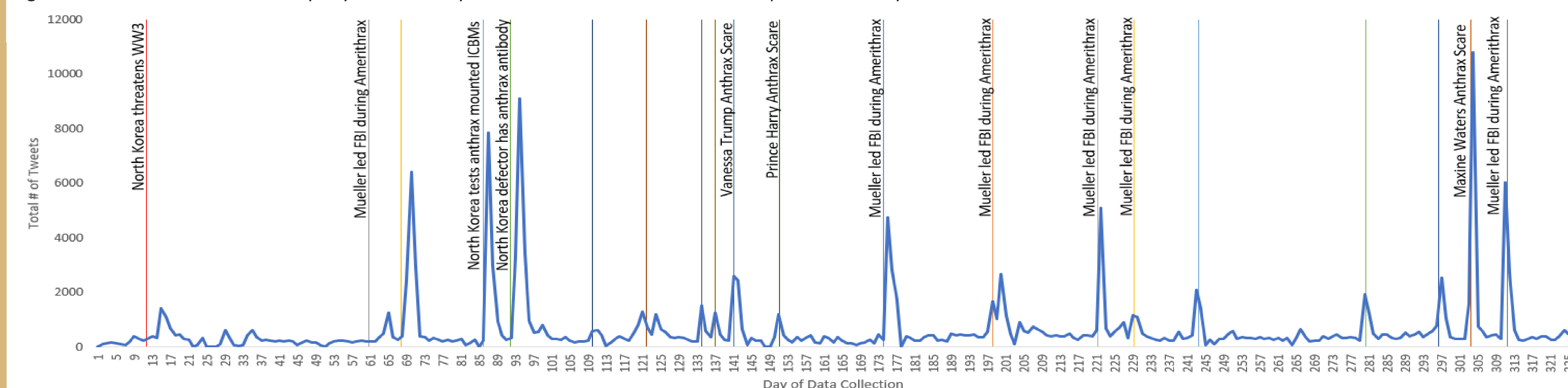
The average number of tweets for each day was found for all tweets. December had the most and September had the least. September most likely had the fewest tweets since no events occurred that month. December had the most because there were concerns over North Korea having anthrax and people discussed errant reports by Brian Ross concerning Saddam Hussein and anthrax.

Figure 1. Average number of tweets for each month.



Using a crawling algorithm, spikes in tweets due to anthrax related events were detected in real-time. Events included discussions of how Mueller handled Amerithrax, concerns over North Korea having anthrax, hippos having anthrax, people receiving packages containing powder, and Anthrax the band announcing tour dates.

Figure 3. Total number of tweets by day over 323 days of data collection. Vertical lines represent the day an anthrax related event occurred



In the gold standard, there were twice as many relevant tweets than non-relevant. Using the gold standard, the logistic regression algorithm performed best at classifying the tweets as relevant versus not (Table 2).

Figure 2. Average number of tweets for each month.

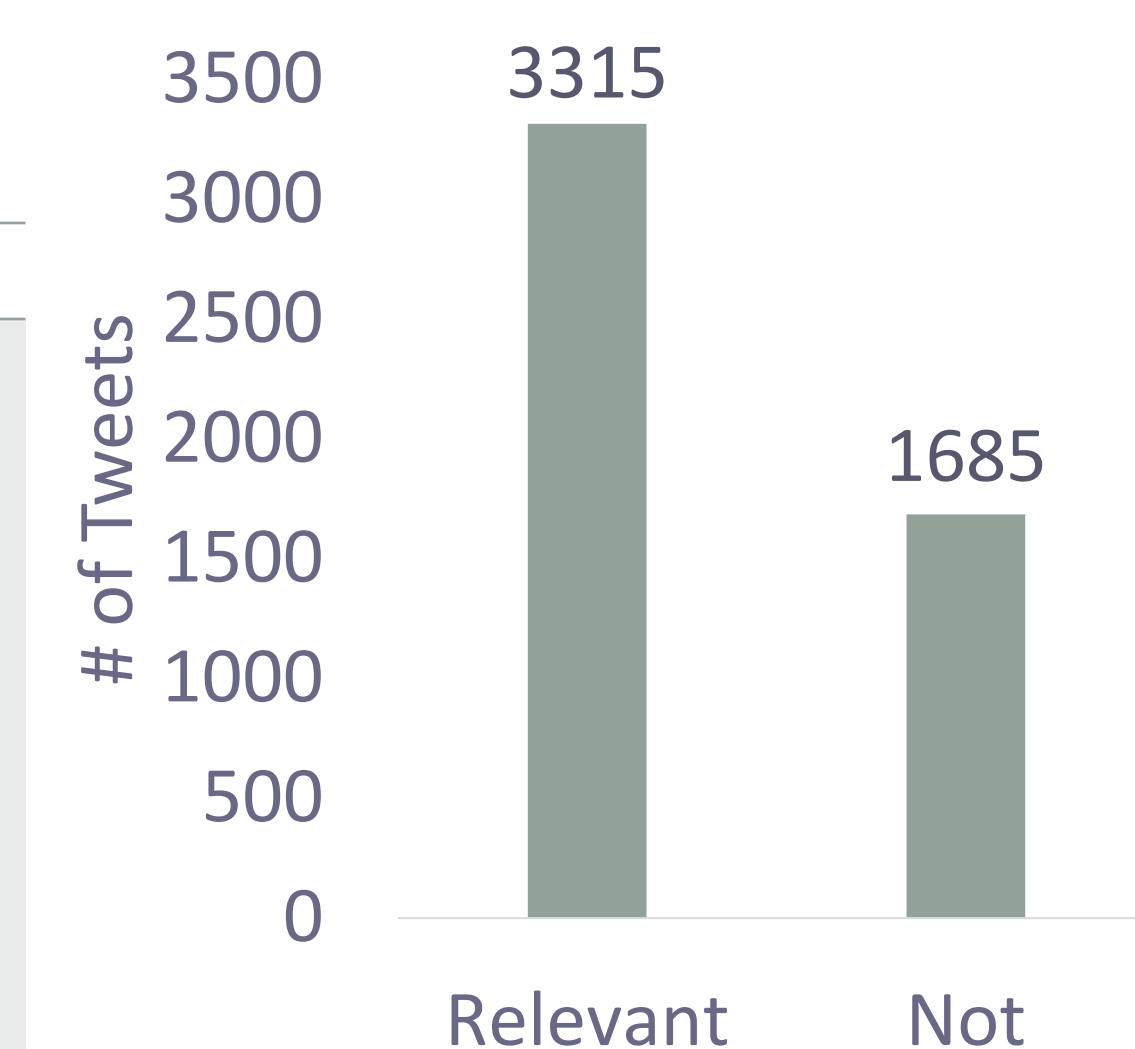


Table 2. Performance of four machine learning algorithms using the gold standard and ten-fold cross validation

	Precision	Recall	F1
SVM	0.71	0.72	0.70
Random Forest	0.79	0.80	0.79
Naïve Bayes	0.80	0.72	0.72
Logistic Regression	0.81	0.81	0.81

Logistic regression was used on the gold standard to classify tweets as relevant or not. The top ten unigrams and bigrams were found for the relevant tweets and not-relevant tweets. All the relevant n-grams concerned anthrax related events while the top n-grams for the not-relevant tweets concerned the metal band anthrax (Table 3).

Table 3. top ten n-grams for relevant and not-relevant tweets in the gold standard.

Relevant				Not			
Count	Unigram	Count	Bigram	Count	Unigram	Count	Bigram
468	mueller	261	north korea	241	slayer	53	slayer lamb
309	korea	48	Korea defector	156	god	52	public enemy
274	north	46	maxine waters	152	testament	42	scott ian
255	attack	40	meghan markle	147	lamb	35	god behemoth
240	case	38	whitey bulgar	117	metal	27	charlie benante
216	investigation	37	white powder	102	video	24	john bush
126	sent	36	prince harry	98	tour	23	god testament
116	fbi	33	biological weapon	86	live	18	testament
114	mail	32	innocent man	82	behemoth	17	napalm death
113	trump	29	korea begin	75	time	17	iron maiden