

Social Media Monitoring and Surveillance for Vector Borne Diseases

This is a research project for developing a strategy and infrastructure to use social media to monitor topics for the purpose of public health surveillance in Pima County, Arizona.

These include influenza, vector borne disease, and heat related illness.

The research will in collaboration with the Pima County Health Department and it will develop and implement strategies to accomplish the following:

1. Track and monitor social media for key public health topics
2. Elicit public feedback and reporting on topics of public health relevance
3. Measure the impact and reach of social media campaigns

1. Data Collection

[Twitter](#) offers [streaming APIs](#) to give developers and researchers low latency access to its global stream of data. Public streams, which can provide access to the public data flowing through Twitter, were used in this study. Studies have estimated that using Twitter's Streaming API, researchers can expect to receive 1% of the tweets in near real-time. [Tweepy](#) -- a [Python](#) library for the Twitter API, was used to access tweet information from the Twitter Streaming API.

Keywords (vector borne diseases symptoms)

The **vector borne diseases** symptoms-related stream is to collect only tweets containing any of 64 related keywords that were suggested by clinical collaborators. Figure 1 shows the keywords used in the data collection process.

Table 1. Keywords Table

abdominal pain	activity change	acute kidney injury	altered mental status	appetite change	arthralgia (arthralgias)	back ache (back pain)	blurred vision (blurry vision)
bone pain	chills	clammy	constant throbbing headache	cough (coughing)	decreased activity	decreased appetite (lack of appetite)	decreased awareness
diarrhea	diffuse rash (rash)	dizziness	elevated lft	fatigue	fever	flaccid paralysis	headache

abdominal pain	activity change	acute kidney injury	altered mental status	appetite change	arthralgia (arthralgias)	back ache (back pain)	blurred vision (blurry vision)
hemoconcentration	high grade fever	increased agitation	increasing obtundation	joint pain	joint swelling	lethargic	leukopenia
lymphadenopathy	macular rash on upper body	menses longer than usual.	mild nausea	mild sore throat	muscle aches (muscle pain, muscle weakness)	myalgia	n/v/d
nausea diarrhea	neck pain (neck stiffness)	numbness in hands	pain in knees	pain in wrists	petechiae	photophobia	red eyes
seizure	severe headache	somnolence	sore throat	stiffness	swollen feet hands	throbbing	thrombocytopenia
transaminitis	tremors	unable to focus eyes	unable to walk	unconscious	vomiting	weakness	wide based gait

Figure 1 shows the number of tweets in our Twitter stream for all of the keywords used in data collection.

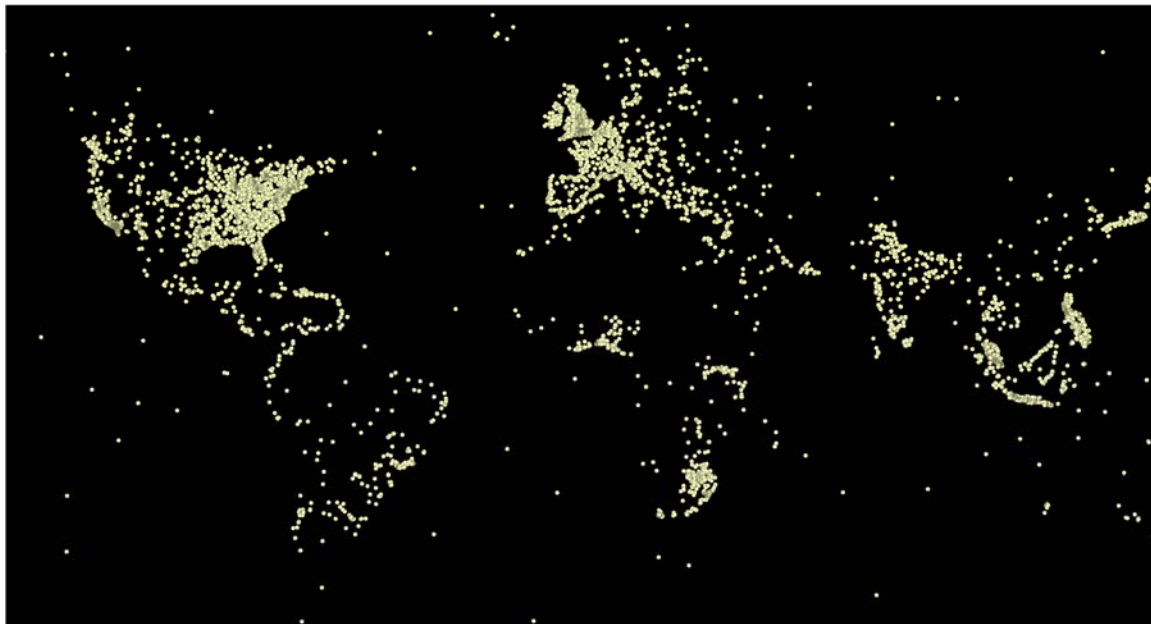
[Open Figure 1](#)

Figure 1. keywords and Frequency of Tweets

There were **61.31% (16,200,147 / 26,422,695)** tweets containing location information. However, there is considerable ambiguity in the location information, e.g., “I am somewhere on the earth” or “Moon”.

Geo-location extraction was performed by using (1) latitude and longitude mapping for the coordinates field and (2) regular expressions analysis for the location field, i.e., matching state names or postal abbreviations on the US states (e.g., Texas or TX), followed by matching city names (e.g. Dallas). We were able to identify location for approximately **9.09% (2,403,950 / 26,422,695)** tweets as one of the US states name.

Figure 2. Geographic Visualization
(2016-4-15 ~ 2016-9-20, All over the world, 0.14% 37,964 Tweets)

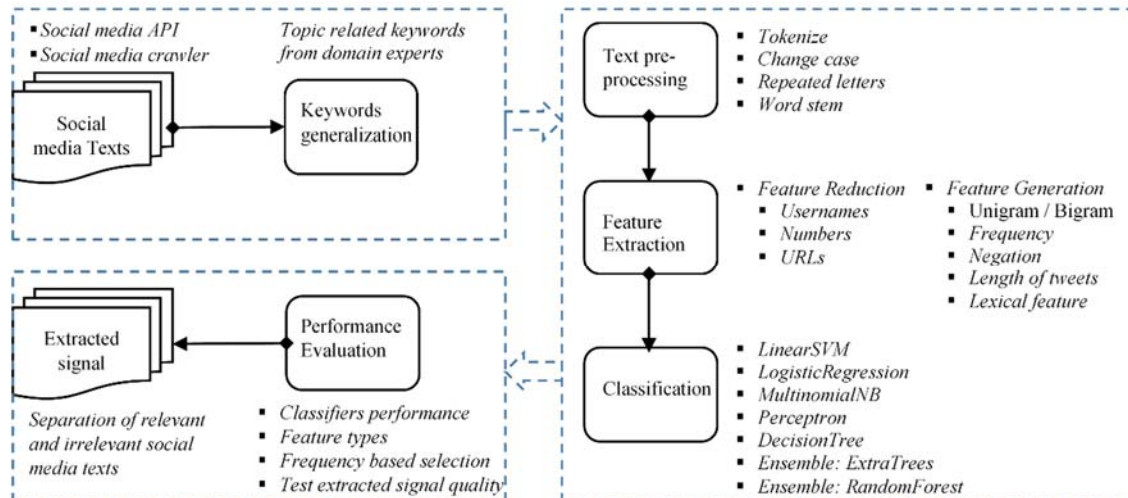


3. Natural Language Analysis

[Natural Language Analysis](#) involves natural language understanding and enabling computers to derive meaning from human or natural language input (i.e., tweets textual information).

Our approach adapts and proposes a combination of different [Natural Language Processing](#) feature generation methods along with [machine learning](#) classifiers to identify and remove various types of noise from social media text as shown in Figure 2.

Figure 3. Methodology for Signal Extraction from Social Media Text

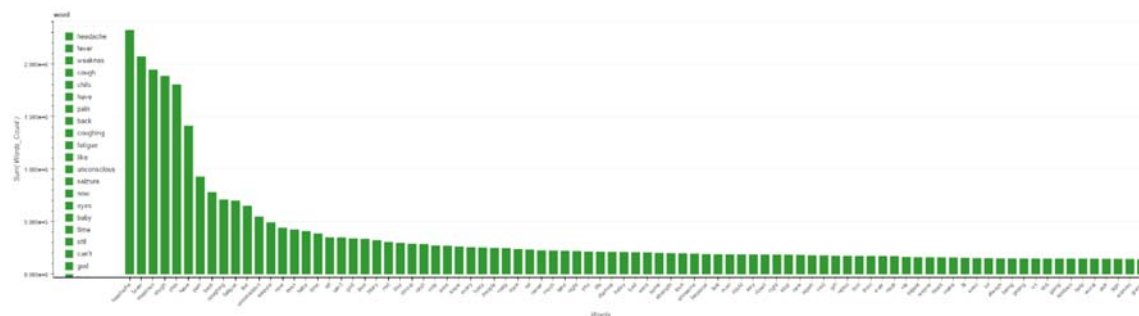


Top words used by Twitter users

We summarize the top words used by Twitter users in the dataset. "headache", "fever", "weakness", "cough (coughing)", "chills", "pain", "fatigue", "unconscious" are meaningful top words used by Twitter users.

Figure 4. Top Words Used by Twitter Users

[Open Figure 4](#)



We further analyze the lexical feature of words used by Twitter users. [Part-of-speech](#) (POS) tagging (also called grammatical tagging) was applied. POS tagging is the process of marking up a word in a text (corpus) as corresponding to a particular part of speech, based on both its **definition** and its **context**.

Based on the result of POS tagging, we summarize the top nouns and top verbs used by Twitter users in the dataset.

Figure 5. Top Nouns Used by Twitter Users

[Open Figure 5](#)

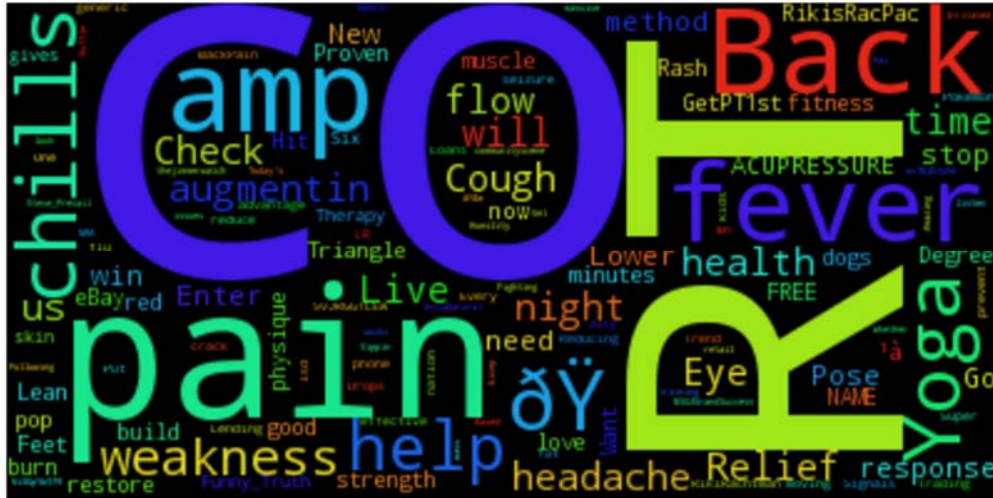
Key conversations (topics) and trends

Topic modeling is a form of text mining, a way of identifying patterns in a dataset. Latent Dirichlet Allocation (LDA) algorithms was applied in order to discover the latent topics.

Figure 8. Topic A: Symptoms

Example tweets from this topic:

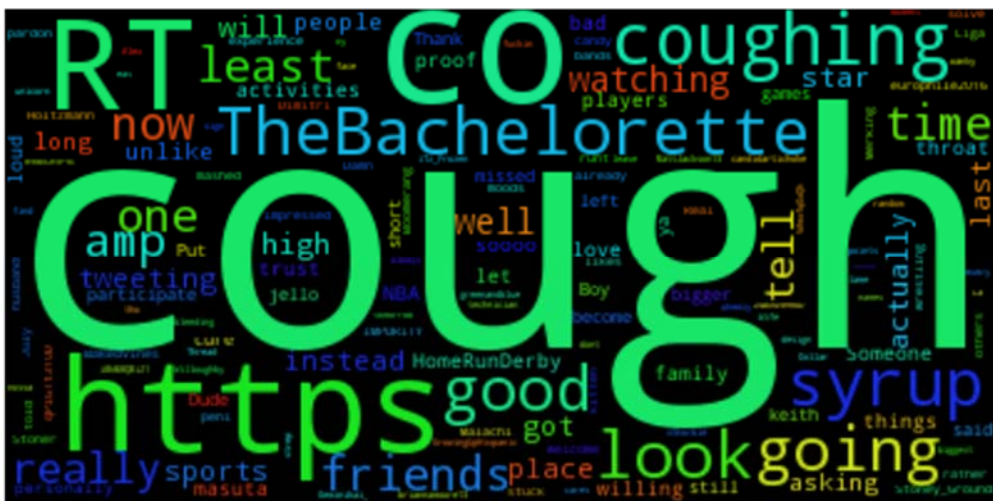
Figure 9. Topic B: News



Example tweets from this topic:

- Do you have a retail business? Have you experienced the FEVER of Pokemon-Go? Take advantage of the foot traffic... ://t.co/SE8DXpU5pG
- for more far out, headache material check out my latest film on lsd ://t.co/CAClDehimk
- Amazing Stretches To Alleviate Lower Back Pain | : ://t.co/o26X5xRYhF , , , ://t.co/R4qtN6abPK
- Cough up less for Lenovo Yoga 3 Pro - 80HE011XUS ! Bid on it! Go to eBay ://t.co/XGPjiJKWPd

Figure 11. Topic D: Irrelevant: Cough



Example tweets from this topic:



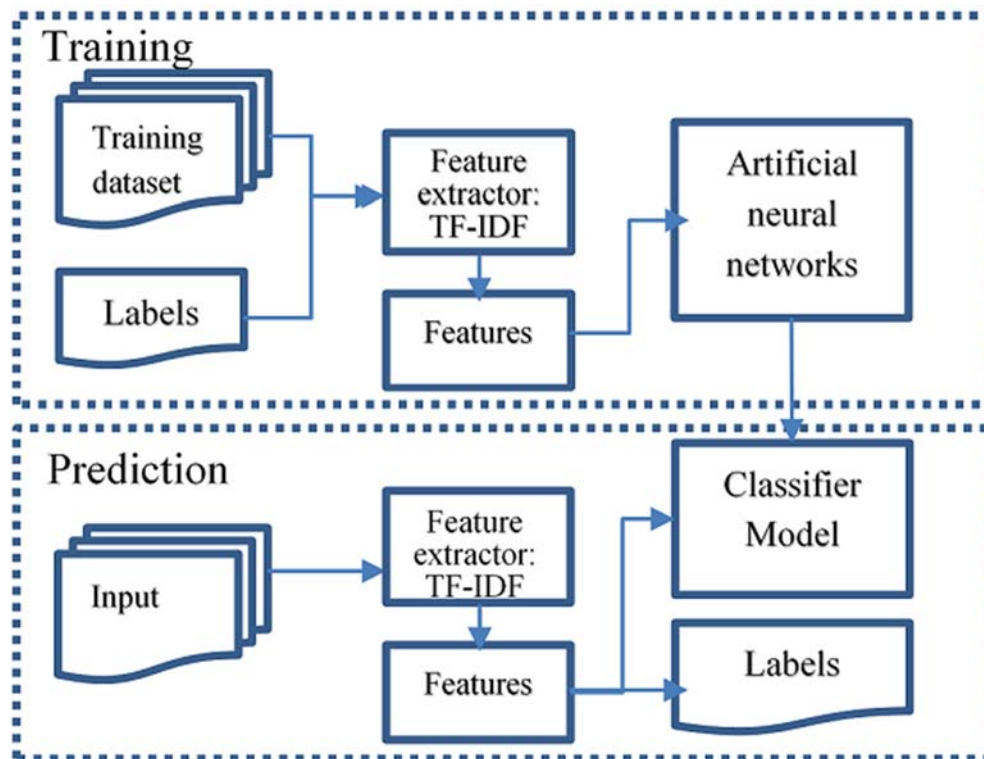
Example tweets from this topic:

- If I'm tired and have a headache she expects me to change my attitude 😊😊
- My greatest headache and favorite stress reliever is my family 💕
- Hungry and happy with a headache 😊😊
- Being in a relationship is a fucking headache...

Theme Analysis

One of the challenges we needed to address is to extract signal from the noisy Twitter dataset i.e., to distinguish tweets that are relevant to symptoms from tweets that mentioned keywords in an irrelevant context. Figure 13 shows the process used for categorizing Twitter dataset.

Figure 14. Theme Analysis



Training data

According to the result of LDA topic modeling, the **vector borne diseases** symptoms-related tweets can be classified as 4 categories:

Table 3. Tweets Categories

Label	Theme
<i>I</i>	Irrelevant: not a symptom
<i>R</i>	Relevant: disease symptoms, e.g., fever, diarrhea,
<i>N</i>	News: news, policy and government themes
<i>A</i>	Advertisement: marketing messages

Classification

[Classification](#) is applied to identify to which of a set of categories a single tweet belongs to. To this purpose, a small amount of data needs to be labeled by human operators. 15 junior and senior students from University of Arizona were invited to label **1,537 tweets** (randomly sampled from the dataset) as "*Relevant*", "*Irrelevant*", "*News*" and "*Advertisement*".

The classifier [Logistic Regression](#) was trained on this dataset. A [10 fold-cross validation](#) was executed to evaluate the [performance](#).

Table 4. Classification Performance

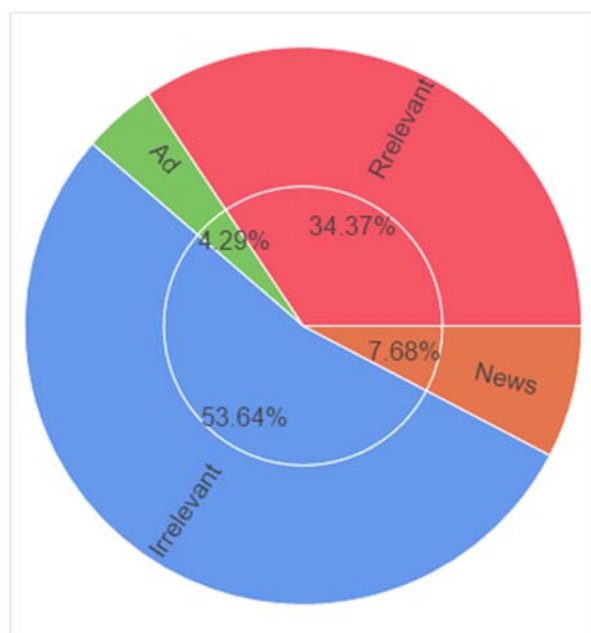
LogisticRegression	Irrelevant	Relevant	News	Advertisement
precision	0.690	0.653	0.677	0.267
recall	0.820	0.597	0.316	0.197
f1-score	0.746	0.616	0.413	0.136
accuracy	0.670			

Percentage of theme categories

The pie chart below shows the percentage of each category .

Figure 15. Theme Composition (9.09% tweets from US)

[Open Figure 15](#)



4. Time Series Analysis

[Time series](#) analysis is an illustration of the number of Tweets at successive time intervals. We could extract meaningful characteristics of the data and predict future values based on previously observed patterns.

Figure 16. Number of Tweets per Day (US)

[Open Figure 16](#)

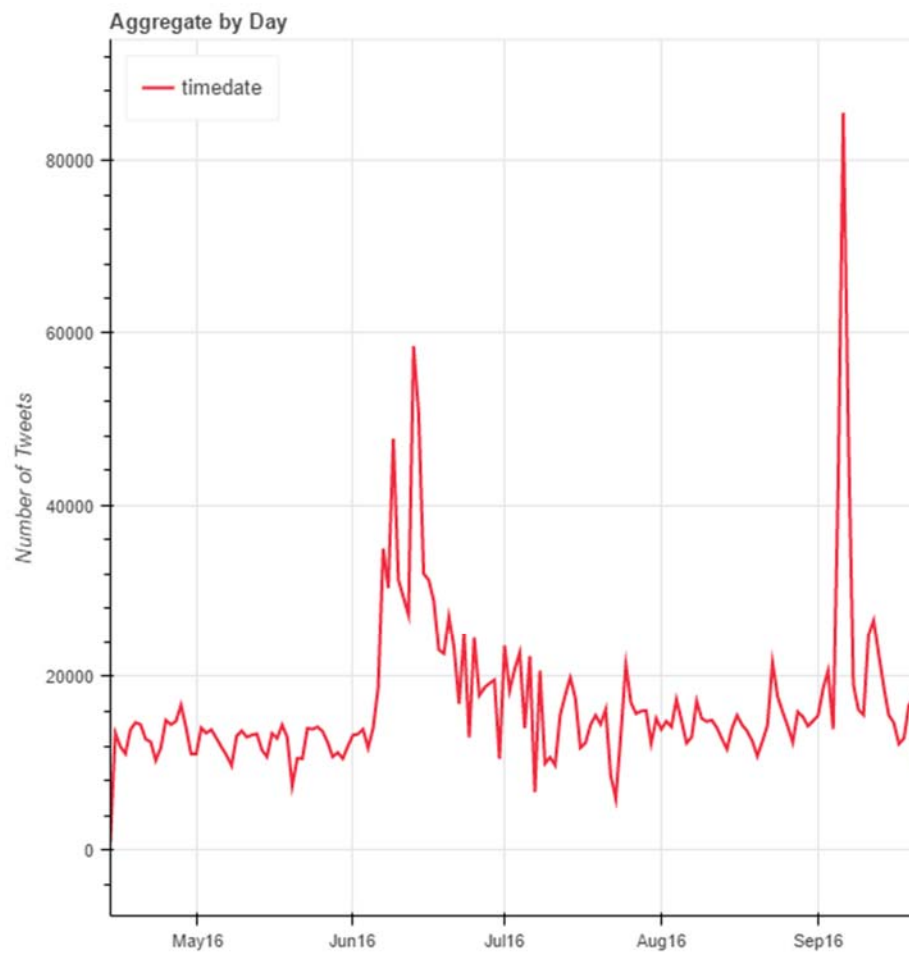


Figure 17. Number of Tweets per Hour (US)

[Open Figure 17](#)

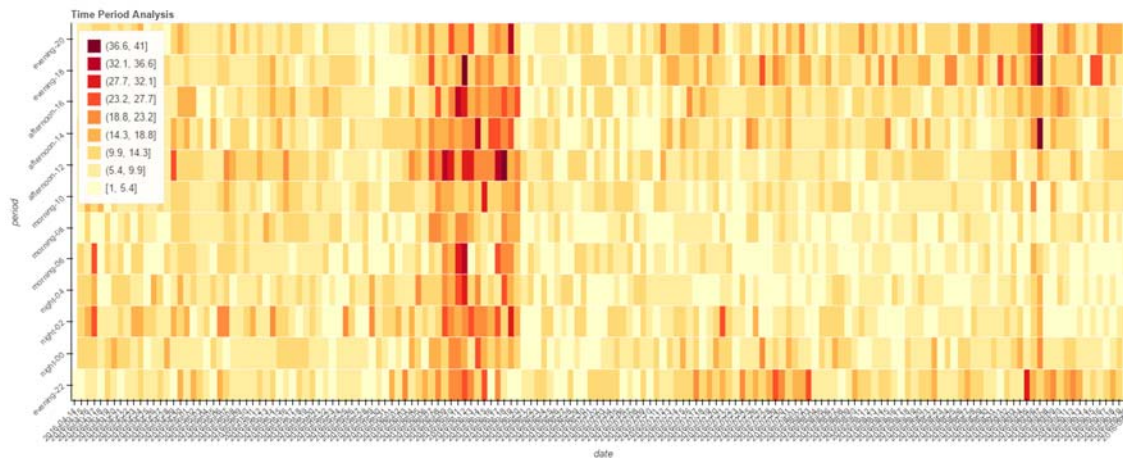
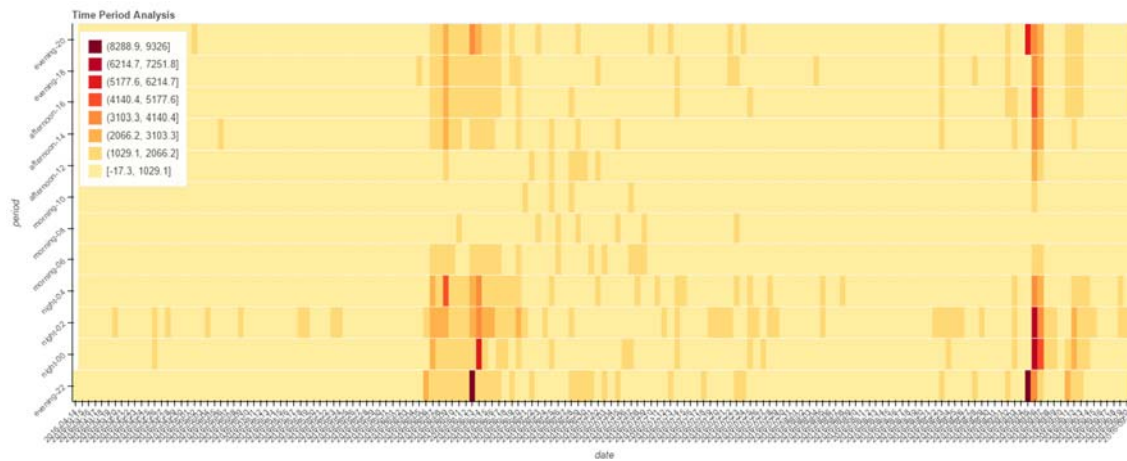


Figure 18. Number of Tweets per Hour (World)

[Open Figure 18](#)



5. Mapping and Visualization

Figure 19. Mapping Timeline

[Open Figure 19](#)

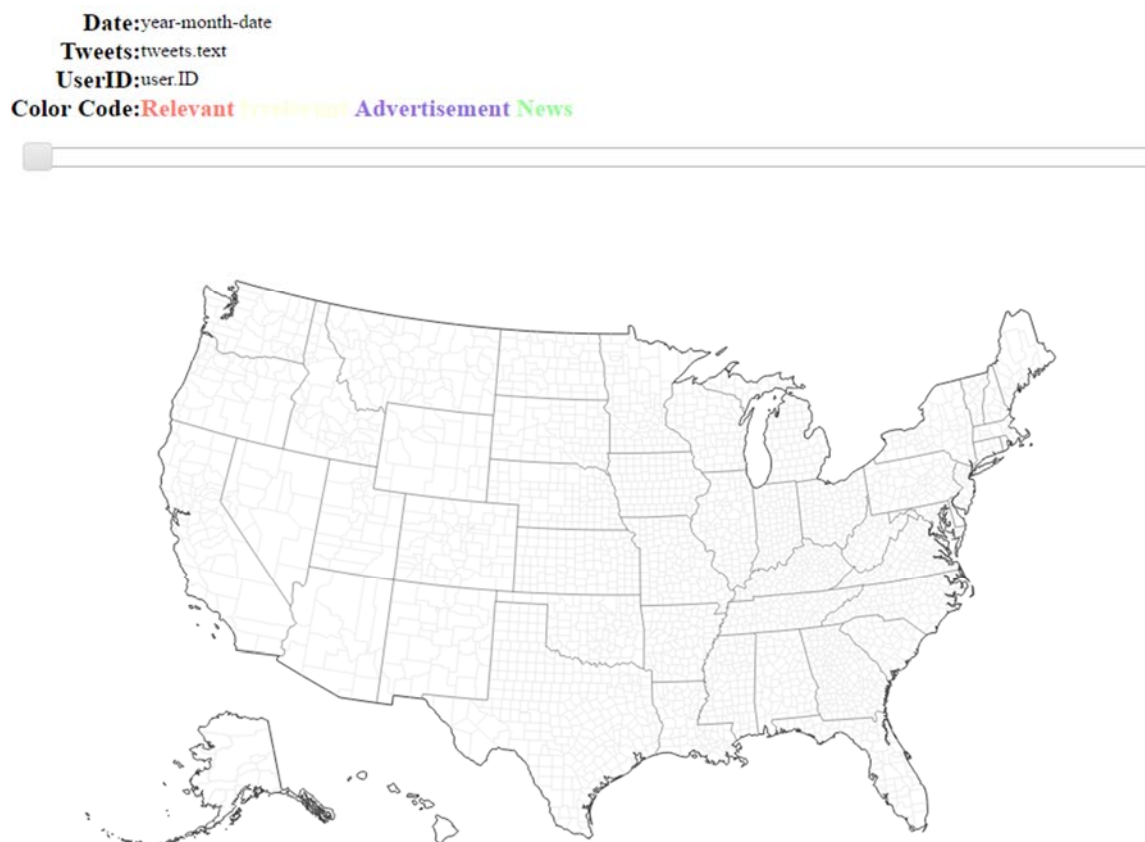
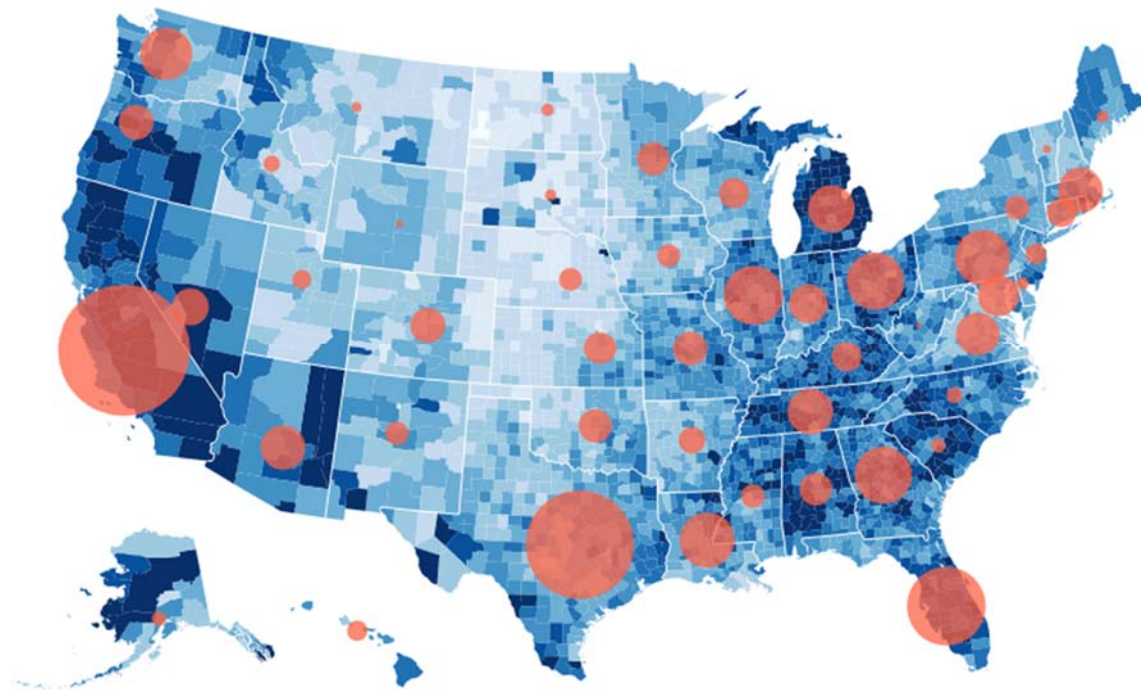


Figure 20. Mapping Number of Tweets aggregated by US States

[Open Figure 20](#)



6. Flights and Travel Analysis

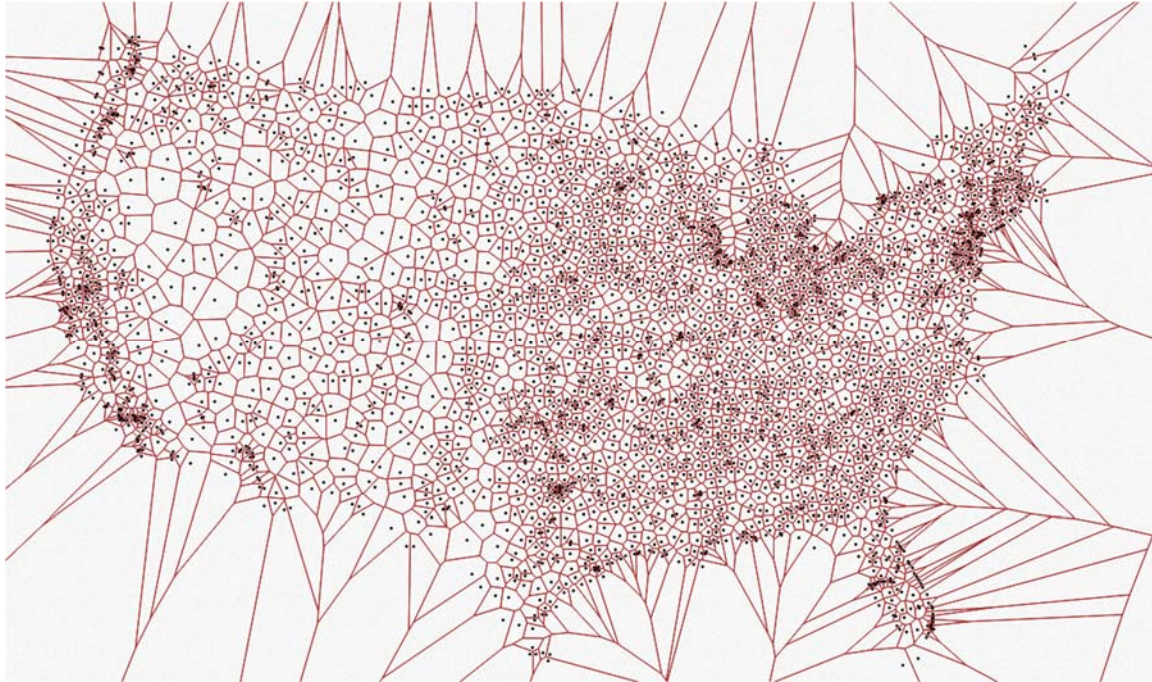
Airports and Flights

According to the literature review, with the outbreaks of vector borne diseases, the number of cases among travelers visiting or returning to the US from affected areas will likely increase. These imported cases could result in local spread of the virus in other parts of the US.

We visualized all US airposts and Voronoi algorithm was applied. [Voronoi diagram](#) is a partitioning of a plane into regions based on distance to points in a specific subset of the plane. For each seed there is a corresponding region consisting of all points closer to that seed than to any other. These regions are called Voronoi cells.

Figure 21. All Airports in US

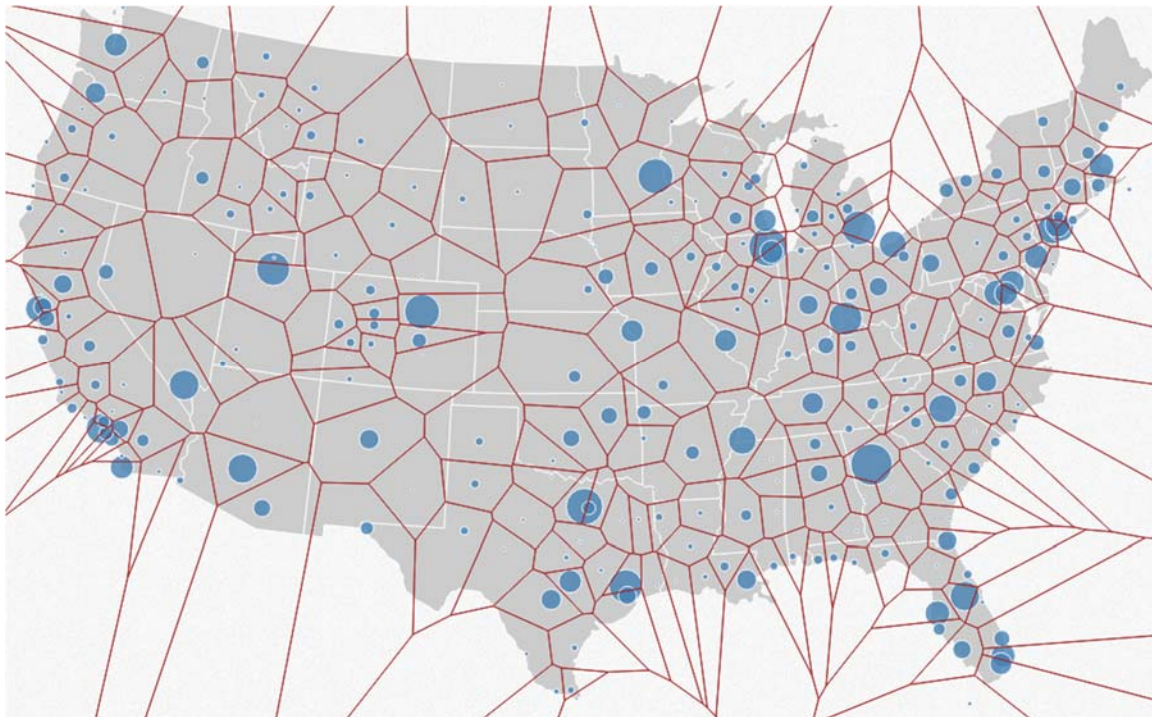
[Open Figure 21](#)



Then airports with only one flight were filtered out. We can hover to see flight routes associated with each airport.

Figure 22. Mapping Airports and flights

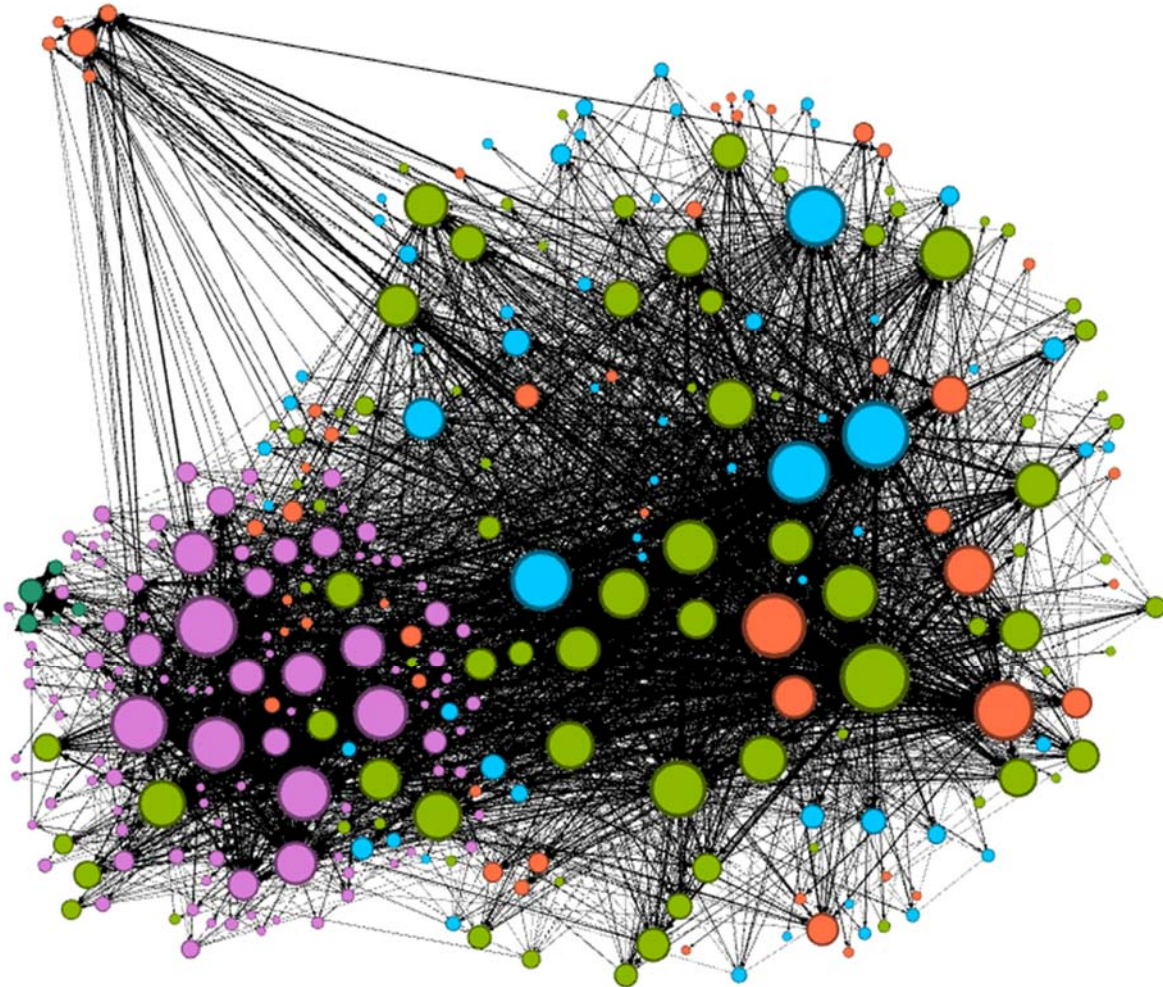
[Open Figure 22](#)



Flights Network Analysis

Figure 23. Flights Network Analysis

[Open Figure 23](#)



7. Twitter Dataset Visualization for Specific Vector (mosquito)

We represented "mosquito" related tweets as a full representation with nodes for Users, Tweets, Hashtags, Urls and Media.

Twitter Type		
User		(43.3%)
Tweet		(40.5%)
Link		(8.38%)
Hashtag		(4.61%)
Media		(3.21%)

