

# Social Media Content Analyst Tool: An NLP based system to identify informative tweets in Emergency Events

**Hardik Sonetta**

sonetta@uwindsor.ca  
School of Computer Science  
University of Windsor

**Jalpa Deepak Patel**

patel2fu@uwindsor.ca  
School of Computer Science  
University of Windsor

## 1. Overview

- **Motivation and Description:**

Information is cascaded on social media platforms and it has been proven to be the best medium for raising situational awareness for various causes and events. It helps us to disseminate information regarding current incidents including natural disasters. In this study, proposed work will help us to understand the authenticity of information diffused on the occurrence of disasters through social media networks and portray the involvement of humanitarian organizations and news agencies in this social network.

The information diffused on social networks during the time of disasters can greatly help in locating the exact source of the disaster and aid in efficiently allocating resources to bring relief to affected individuals. But, at the same time, substantial amount of irrelevant or less informative data is generated which may hinder in analysing the disastrous situation. So, this motivated us in identifying and eliminating the false information which may be of no use to the social network. Our work would aid in identifying such information and will lead to cascading of the data that has some social value and awareness about the disaster in the social network.

- **Project Outcome and Potential Users:**

While a huge amount of information is being generated during the time of disasters, there are a lot of users whose post is misunderstood for it to be related to the disaster situations. For this project, we will utilize data from widely known social media platform known as Twitter, which is used by masses to express thoughts and to cascade information on real life events and situations. With the help of this data, it will enable us to analyze and identify which tweets hold valuable information describing the event and eliminate the tweets that hold no information. This project can be used by disaster relief organizations for efficient allocation of resources and by end-users of the platform to identify and cascade useful information.

- **Methodology:**

As Twitter generates overwhelming number of tweets during massive disasters affecting humankind and infrastructure on a very large scale, there are substantial number of tweets

that does not hold much valuable information. Such tweets are classified as “false” tweets, which does not give any information on locating the source of the disaster as when and where the disaster occurred and how many people got affected of it. Thus, it hinders the efficient allocation of resources to relieve the affected people.

This section describes the methods used in detection of “False tweets” in a social network using word embedding and text classification using Recurrent Neural Network (RNN).

#### **i. Data Gathering:**

For this study, there are two methods that can be used to extract and collect the relevant data. First method is to make use of the open source Twitter Application Programming Interface (API) called as Tweepy. This approach, despite being free, limits the number of tweets that can be extracted for our project. Hence, along with using Tweepy, we are utilizing open source data available on Kaggle which has over 7000 tweets containing real and false tweets addressing the disaster situation.

#### **ii. Data Preprocessing:**

As twitter is widely known and accepted by masses as a platform for expressing thoughts and ideas. The platform is used by diverse audience ranging from different locations, demographics and culture. Thus, freedom of expressing thoughts can take many shapes and forms leading to the problem of ambiguity and disorientation of language used. This contributes as a major challenge, working with the real-world social media content. To solve this technique, we are going to use various data preprocessing techniques to remove slangs, stop words, punctuations and spelling corrections which will lead to accurate analysis of social network data and will contribute in taking informed decisions.

#### **iii. Data Analysis and System Modelling:**

The last component of methodology is analysing the tweets and modelling a Recurrent neural network to classify the tweets based on the authenticity of information it holds about the disaster. A recurrent neural network (RNN) is a class of artificial neural networks where connections between nodes form a directed graph along a temporal sequence<sup>[1]</sup>. Before passing the data under study into our model, we will be converting each tweet into a word embedding. Each tweet will be converted into a 100-Dimensional word vector. As there are many pre-trained word embeddings models, we will be use Glove model to represent our tweet in a multi-dimensional vector space.

- **Literature Review**

The study conducted by Shu et al. says that there are a lot of fake news detecting systems but most of them contain either linguistic or social context features. They propose a way to have linguistic and social features with dynamic features for detecting fake news<sup>[2]</sup>. There is

abundance of information from tweets during a disaster, but not all tweets originally posted or re-tweeted are informative. The study of the tweets being informative or not during the time of crisis is conducted by Madichetty et al<sup>[3]</sup>. The research done by Kim et al. shows the graphical representation of the tweets and retweets during a storm with each node representing a user and the edges represent the relationship between the Twitter Users<sup>[4]</sup>. A recent study done by Mohammadrezaei et al addresses the security issues in social networks by classifying if an account is Fake or not by implementing graph analysis on social Networks<sup>[5]</sup>.

## 2. Project Goal

The goal of the project is to develop a system that identifies the tweets generated at the time of disaster holds any valuable information about the disaster or not. This is achieved by classifying and visualizing the tweets of various users that are relevant to a catastrophe and to differentiate the tweets that are misinterpreted to have information about the disastrous situation.

## 3. Project Team

Name	Responsibility	Comment
Hardik Sonetta	Data Extraction and Consolidation	<ul style="list-style-type: none"> <li>- Working with twitter API and writing a script to extract relevant tweets</li> <li>- Collecting data from open source</li> <li>- Consolidating data collected across various origins</li> <li>- Modelling RNN for identifying authentic tweets</li> </ul>
Jalpa Patel	Data Preprocessing and Cleaning	<ul style="list-style-type: none"> <li>- Performing Exploratory data analysis</li> <li>- Handling missing values and preprocessing data</li> <li>- Using Glove embedding model to convert tweets in multi dimensional vectors</li> </ul>

## 4. Schedule and Milestones

Milestones	Description	Milestone Criteria	Planned Date
M0	Problem identification	Identify task related to social network data analysis	2020-01-22
M1	Targeted Social media platform	Identify a social media platform for analysis	2020-01-29
M2	Data source	Identify various data sources for data extraction	2020-02-5

M3	Data gathering	Design and script a data gathering module	2020-02-12
M4	Data consolidation	Consolidate data in a single file and validate with input standards	2020-02-19
M5	Exploratory Data Analysis	Visualize raw data and note down remarks	2020-02-26
M6	Data Cleaning and Pre-Processing	Clean data and validate with input standards	2020-03-4
M7	Generate word embeddings	Generate and store word embeddings	2020-03-11
M8	Modelling System	Create baseline model	2020-03-18
M9	Documentation and Presentation	Create presentation for project delivery	2020-03-25

## 5. Communication and Reporting

The internal communication between team members regarding work distribution and reporting of work progress will be done through blackboard group tools. Informal communication method will be Instant Messaging application (WhatsApp). The work progress will be updated internally by written and verbal communication methods among team members. Weekly scrum meetings will be scheduled to review the achieved and to be achieved milestones.

## 6. Delivery Plan

Sr. No	Deliverable name	Description	Deliverable date	Owner
1.	Project proposal	Document describing project goals and timeline	2020-01-31	Hardik Sonetta Jalpa Patel
2.	Data gathering script	Script file to scrape data	2020-02-12	Hardik Sonetta
3.	Consolidated data file	File consisting of data from various sources	2020-02-19	Hardik Sonetta
4.	Preprocessed data file	Cleaned and pre-processed data file	2020-03-4	Jalpa Patel
5.	Project document	Final project document	2020-03-25	Hardik Sonetta Jalpa Patel

## 7. References

- [1] [https://en.wikipedia.org/wiki/Recurrent\\_neural\\_network](https://en.wikipedia.org/wiki/Recurrent_neural_network)
- [2] Shu K, Mahudeswaran D, Wang S, Lee D, Liu H. FakeNewsNet: A Data Repository with News Content, Social Context and Spatialtemporal Information for Studying Fake News on Social Media. 2018;(September).
- [3] Madichetty S, Sridevi M. Detecting Informative Tweets during Disaster using Deep Neural Networks. *2019 11th Int Conf Commun Syst Networks, COMSNETS 2019*. 2019;2061:709-713.
- [4] Kim J, Bae J, Hastak M. Emergency information diffusion on online social media during storm Cindy in U.S. *Int J Inf Manage*. 2018;40(February):153-165.
- [5] Mohammadrezaei M, Shiri ME, Rahmani AM. Identifying Fake Accounts on Social Networks Based on Graph Analysis and Classification Algorithms. *Secur Commun Networks*. 2018;2018.