

Tweet Alert System: Development of New Real Time Event Detection and Earth Quake Reporting System

Sri Lakshmi Poornima Gorle^{#1}, P.Srinivasu Varma^{*2}

M.Tech Scholar^{#1}, Assistant Professor^{*2}

Department of Computer Science & Engineering,
Avanthi's St.Theressa Institute of Engineering & Technology,
Garividi, Vizianagaram Dist, AP, India.

Abstract

Twitter is categorized as a micro blogging service. Micro blogging is a form of blogging that enables users to send brief text message updates or micromedia data such as photographs or audio clips. In recent days twitter has received much attention from different users around the world. An important characteristic of Twitter is its real-time nature. We investigate the real-time interaction of events such as earthquakes in Twitter and propose an algorithm to monitor tweets and to detect a target event. To detect a target event, we devise a classifier of tweets based on features such as the keywords in a tweet (I.e. 'Shaking' and 'Earthquake'), the number of words, and their context. As an application, we develop an earthquake reporting system for use in Japan. Because of the numerous earthquakes and the large number of Twitter users throughout the country, we can detect an earthquake with high probability (93 percent of earthquakes of Japan Meteorological Agency (JMA) seismic intensity scale 3 or more are detected) merely by monitoring tweets. Our system detects earthquakes promptly and notification is delivered much faster than JMA broadcast announcements.

Keywords

Event Detection, Social Sensor, Location Estimation, Earthquake, Twitter

1. Introduction

Twitter, a popular micro blogging service, has received much attention recently. This online social network is used by millions of people around the world to remain socially connected to their friends, family members, and coworkers through their computers and mobile phones [1]. Twitter asks one question, "What's happening?" Answers must be fewer than 140 characters. A status update message, called a tweet, is often used as a message to friends and colleagues. A user can follow other users; that user's followers can read her tweets on a regular basis. A user who is being followed by another user need not necessarily reciprocate by following them back, which renders the links of the network as directed. Since its launch on July 2006, Twitter users have increased rapidly. The number of registered Twitter users exceeded 100 million in April 2010. The service is still adding about 300,000 users per day. Currently, 190 million users use Twitter per month, generating 65 million tweets per day.

Many researchers have published their studies of Twitter to date, especially during the past year. Most studies can be classified into one of three groups: first, some researchers have sought to analyze the network structure of Twitter [2], [3], [4]. Second, some researchers have specifically

examined characteristics of Twitter as a social medium [5], [6]. Third, some researchers and developers have tried to create new applications using Twitter [7], [8].

Twitter is categorized as a micro blogging service. Micro blogging is a form of blogging that enables users to send brief text updates or micromedia such as photographs or audio clips. Micro blogging services other than Twitter include Tumblr, Plurk, Jaiku, identi.ca, and others. Our study, which is based on the real-time nature of one social networking service, is applicable to other micro blogging services, but we specifically examine Twitter in this study because of its popularity and data volume.

This paper presents an investigation of the real-time nature of Twitter that is designed to ascertain whether we can extract valid information from it. We propose an event notification system that monitors tweets and deliver notification promptly using knowledge from the investigation. In this research, we take three steps: first, we crawl numerous tweets related to target events; second, we propose probabilistic models to extract events from those tweets and estimate locations of events; finally, we developed an earthquake reporting system that extracts earthquakes from Twitter and sends a message to registered users. Here, we explain our methods using an earthquake as a target event.

First, to obtain tweets on the target event precisely, we apply semantic analysis of a tweet. For example, users might make tweets such as “Earthquake!” or “Now it is shaking,” for which earthquake or shaking could be keywords, but users might also make tweets such as “I am attending an Earthquake Conference,” or “Someone is shaking hands with my boss.” We prepare the training data and devise a classifier using a Support Vector Machine (SVM) based on features such as keywords in a tweet, the number of words, and the context of target-event words.

After doing so, we obtain a probabilistic spatiotemporal model of an event. We then make a crucial assumption: each Twitter user is regarded as a sensor and each tweet as sensory information.

These virtual sensors, which we designate as social sensors, are of a huge variety and have various characteristics: some sensors are very active; others are not. A sensor might be inoperable or malfunctioning sometimes, as when a user is sleeping, or busy doing something else. Consequently, social sensors are very noisy compared to ordinary physical sensors. Regarding each Twitter user as a sensor, the event-detection problem can be reduced to one of object detection and location estimation in a ubiquitous/ pervasive computing environment in which we have numerous location sensors: a user has a mobile device or an active badge in an environment where sensors are placed. Through infrared communication or a Wi-Fi signal, the user location is estimated as providing location-based services such as navigation and museum guides [9], [10].

2. Background Knowledge

In this section we will describe the assumptions and background knowledge that is used for developing this Tweet alert system for identifying the earthquakes.

2.1 Main Motivation

We choose earthquakes in Japan as target events, based on the preliminary investigations. We explain them in this section. First, we choose earthquakes as target events for the following reasons:

1. Seismic observations are conducted worldwide, which facilitates acquisition of earthquake information, which also makes it easy to validate the accuracy of our event detection methodology; and
2. It is quite meaningful and valuable to detect earthquakes in earthquake-prone regions.

Second, we choose Japan as the target area based on the following investigation.

Fig. 1 portrays a map of Twitter users worldwide (obtained from UMBC eBiquity Research Group); Fig. 2 depicts a map of earthquake occurrences worldwide (using data from Japan

Meteorological Agency (JMA)). It is apparent that the only intersection of the two maps, those regions with many earthquakes and large Twitter users, is Japan. Other regions such as Indonesia, Turkey, Iran, Italy, and Pacific coastal US cities such as Los Angeles and San Francisco also roughly intersect, but their respective densities are much lower than that in Japan. Many earthquake events occur in Japan and many Twitter users observe earthquakes in Japan, which means that social sensors are distributed throughout the country. We present a brief overview of Twitter in Japan: the Japanese version of Twitter was launched on April 2008. In February 2008, Japan was the No. 2 country with respect to Twitter traffic.⁵ At the time of this writing, Japan has the second largest number of tweets (18 percent of all tweets are posted from Japan) in the world. Therefore, we choose earthquakes in Japan as a target event because of the high density of Twitter users and earthquakes in Japan.

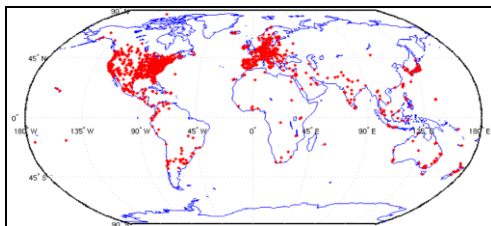


Fig. 1. Twitter User Map.

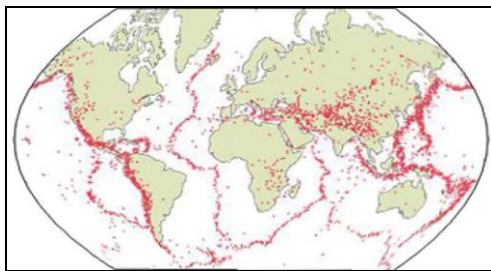


Fig. 2. Earthquake Map.

3. Real Time Event Detection Model

As described in this paper, we target event detection. An event is an arbitrary classification of a space-time region. An event might have actively

participating agents, passive factors, products, and a location in space/time [13]. We target events such as earthquakes, typhoons, and traffic jams, which are readily apparent upon examination of tweets. These events have several properties.

1. They are of large scale (many users experience the event).
2. They particularly influence the daily life of many people (for that reason, people are induced to tweet about it).
3. They have both spatial and temporal regions (so that real-time location estimation is possible).

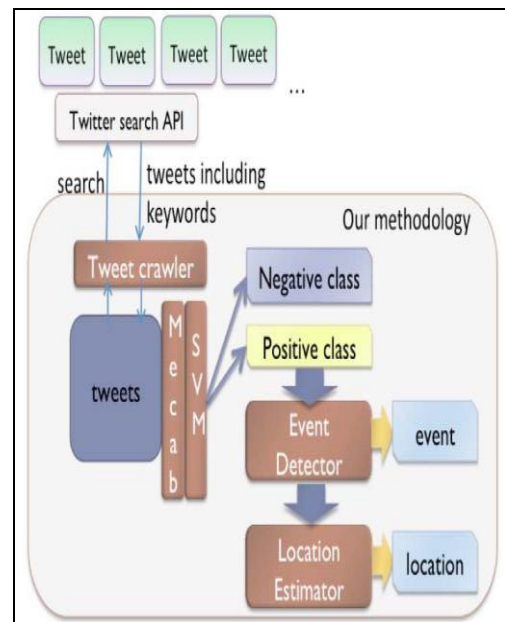


Fig. 3. Method to acquire tweets referred to a target event precisely

Such events include social events such as large parties, sports events, exhibitions, accidents, and political campaigns. They also include natural events such as storms, heavy rains, tornadoes, typhoons/hurricanes/cyclones, and earthquakes. We designate an event we would like to detect using Twitter as a target event.

3.1 Semantic Analysis of Tweets

To detect a target event from Twitter, we search from Twitter and find useful tweets. Our method of acquiring useful tweets for target event detection is portrayed in Fig. 3.

Tweets might include mention of the target event. For example, users might make tweets such as “Earthquake!” or “Now it is shaking.” Consequently, earthquake or shaking might be keywords (which we call query words). However, users might also make tweets such as “I am attending an Earthquake Conference.” or “Someone is shaking hands with my boss.” Moreover, even if a tweet is referring to the target event, it might not be appropriate as an event report.

For instance, a user makes tweets such as “The earthquake yesterday was scary.” or “Three earthquakes in four days. Japan scares me.” These tweets are truly descriptions of the target event, but they are not real-time reports of the events. Therefore, it is necessary to clarify that a tweet is truly referring to an actual contemporaneous earthquake occurrence, which is denoted as a positive class. To classify a tweet as a positive class or a negative class, we use a support vector machine [14], which is a widely used machine-learning algorithm. By preparing positive and negative examples as a training set, we can produce a model to classify tweets automatically into positive and negative categories.

TABLE 1
SVM Features of an Example Sentence

Feature Name	Features
Features A	7 words, the fifth word
Features B	I, am, in, Japan, earthquake, right, now
Features C	Japan, right

We prepare three groups of features for each tweet as described below.

Features A (statistical features): the number of words in a tweet message, and the position of the query word within a tweet.

Features B (keyword features): the words in a tweet.

Features C (word context features): the words before and after the query word.

We can give an illustrative example of these features using the following sentence.

“I am in Visakhapatnam, earthquake right now!”

(Keyword: earthquake)

For this example, Features A, B, C are presented in Table 1. To process Japanese texts, morphological analysis is conducted using Mecab, which separates sentences into a set of words. For English, we apply standard stop-word elimination and stemming. We compare the usefulness of the features in the discussion in later sections. Using the obtained model, we can classify whether a new tweet corresponds to a positive class or a negative class.

3.2 Tweet as a Sensory Value

We can search the tweet and classify it into a positive class if a user makes a tweet about a target event. In other words, the user functions as a sensor of the event. If she makes a tweet about an earthquake occurrence, then it can be considered that she, as an “earthquake sensor,” returns a positive value. A tweet can therefore be regarded as a sensor reading. This crucial assumption enables application of various methods related to sensory information.

Fig. 4 presents an illustration of the correspondence between sensory data detection and tweet processing. The motivations are the same for both cases: to detect a target event. Observation by sensors corresponds to an observation by Twitter users. They are converted into values using a classifier.

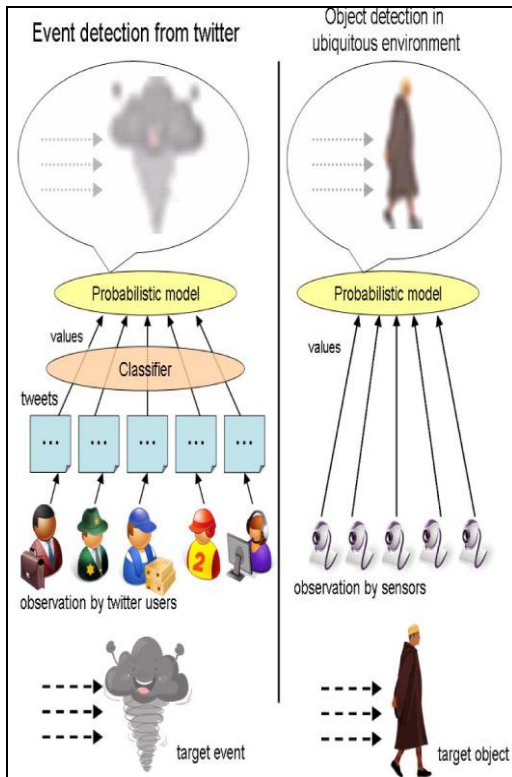


Fig. 4. Correspondence between event detection from Twitter and object detection in a ubiquitous environment.

A tweet can be associated with a time and location: each tweet has its post time, which is obtainable using a search API. In fact, GPS data are attached to a tweet sometimes, such as when a user is using an iPhone. Alternatively, each Twitter user makes a registration on their location in the user profile. The registered location might not be the current location of a tweet. However, we infer it that a person is probably near the registered location. Some tweets include place names in those bodies. Some researchers describe their efforts to extract place names from tweets as a part of Named Entity Recognition [15], [16]. However, the performance derived from those efforts remains insufficient for practical use (precision ranges from 0.6 to 0.8). For the present study, we use GPS data and the registered location of a user. We do not use tweets

for spatial analysis if a location is not available; however, we use the tweet information for temporal analyses.

4. Implementation Modules

Implementation is the stage where the theoretical design is automatically converted into practically by dividing this into various modules. We have implemented the current application in Java Programming language with Front End as JSP, HTML and Back End as MYSQL data base. Our proposed application is divided into following 6 modules. They are as follows:

- 1) Tweet Collection Module
- 2) Crawling Tweets From Twitter Module
- 3) Twitter Search API Module
- 4) Filtering Tweets Using Machine Learning Module
- 5) Semantic Analysis On Tweets Module
- 6) Earthquake Reporting System Module

1. Tweet Collection Module

In this module, we develop our system by posting tweets by the users. It is necessary to collect tweets referring to an earthquake from Twitter. This process includes two steps: crawling tweets from Twitter and filtering out tweets that do not refer to the earthquake. For crawling and filtering tweets, we recommend using script programming languages.

2. Crawling Tweets from Twitter Module

To collect tweets or some user information from Twitter, one must use the Twitter Application Programmers Interface (API). Twitter API is a group of commands that are necessary to extract data from Twitter. Twitter has APIs of three kinds: Search API, REST API, and Streaming API. In this section, we introduce Search API and Streaming API, which

are necessary to crawl tweets from Twitter. We explain REST API later because REST API is necessary to extract location information from Twitter information. Additionally, it is known that Twitter API specifications are subject to change. When using Twitter API, it is necessary to know the latest details and requirements. They are obtainable from Twitter API documentation.

3. Twitter Search API Module

The Twitter Search API extracts tweets from Twitter, including search keywords or those fitting other retrieval conditions, in chronological order. It is possible to use language, date, location and other conditions as retrieval conditions.

Some points must be considered when using Twitter Search API:

- It is possible to collect tweets posted only during the prior five days. It is not possible to search tweets posted six days ago.
- It is only possible to collect the latest 1500 tweets at one time. (Technically speaking, it is possible to access one page with a request and track pages back to the 15th page.
- One page includes 100 tweets at most. Therefore it is possible to acquire the latest 1500 tweets at one time.)
- One is limited to API requests.

4. Filtering Tweets Using Machine Learning Module

We collected data from tweets including keywords related to earthquakes, such as earthquake, shake. Those tweets include not only tweets that users posted immediately after they felt earthquakes, but also tweets that users posted shortly after they heard earthquake news, or perhaps they misinterpreted some sense of shaking from a large truck passing nearby. When the seismic activity reached its peak, the graph of tweets invariably

showed a peak. However, when the graph of tweet counts showed a peak, the seismic activity did not necessarily show a peak. Some "false-positive" peaks of the graph of tweet counts arise from mistakes by users or some news related to earthquakes. Therefore, we must filter tweets to extract those posted immediately after the earthquake. We designate tweets posted by users who felt earthquakes as positive tweets, and other tweets as negative tweets. Here, we describe the creation of a classifier to categorize crawled tweets into positive tweets and negative tweets, using Support Vector Machine: a supervised learning method.

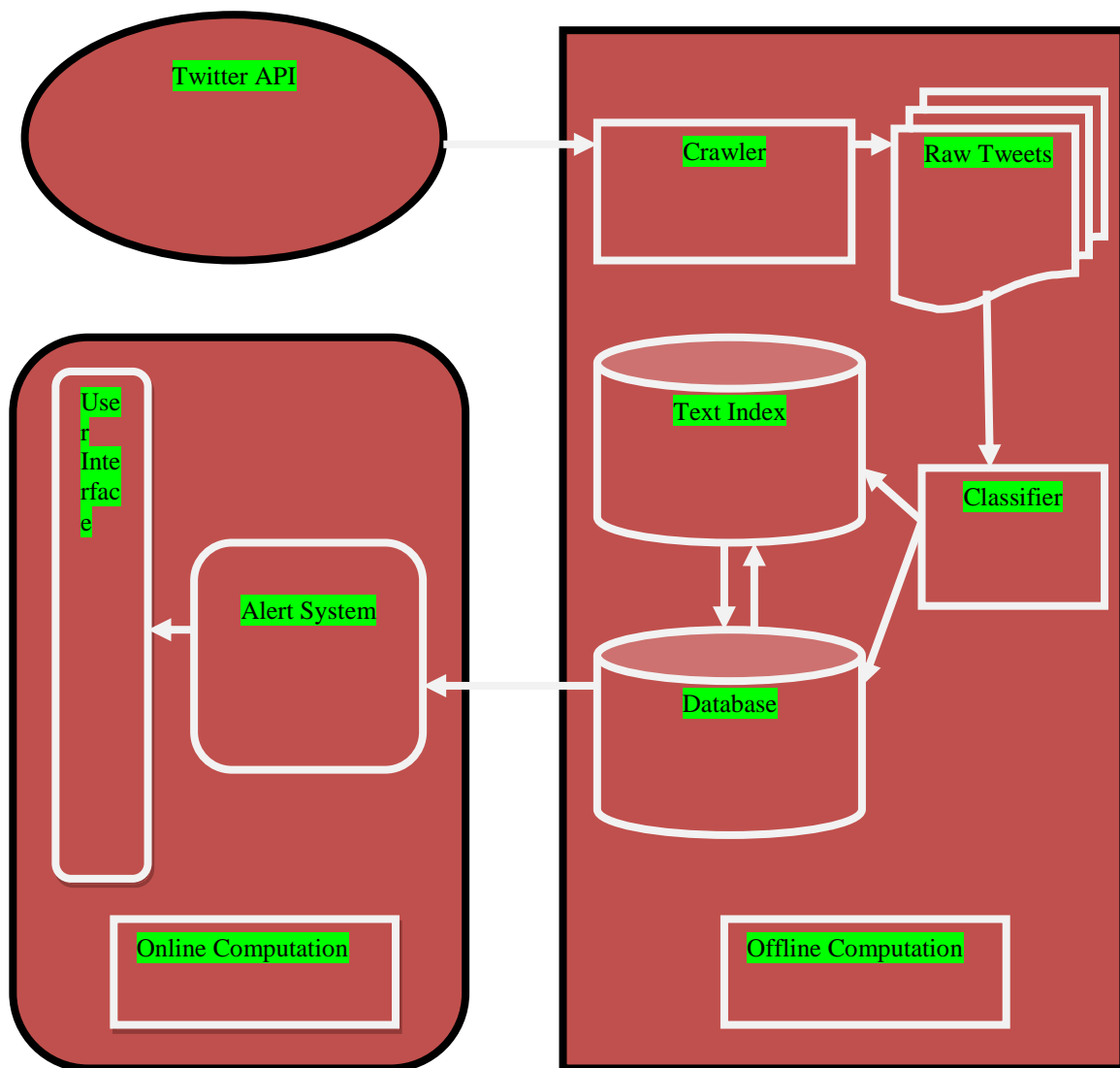
5. Semantic Analysis on Tweets Module

Semantic Analysis on Tweet Search tweets including keywords related to a target event Example: In the case of earthquakes "shaking", "earthquake" Classify tweets into a positive class or a negative class Example: "Earthquake right now!!" ---positive "Someone is shaking hands with my boss" --- negative Create a classifier Semantic Analysis on Tweet Create classifier for tweets use Support Vector Machine (SVM) Features (Example: I am in Japan, earthquake right now!) Statistical features (7 words, the 5th word) the number of words in a tweet message and the position of the query within a tweet Keyword features (I, am, in, Japan, earthquake, right, now) the words in a tweet Word context features (Japan, right) the words before and after the query word.

6. Earthquake Reporting System Module

In this module, the users will be alerted if the earthquake occurs based on their location and the tweets. Effectiveness of alerts of this system Alert E-mails urges users to prepare for the earthquake if they are received by a user shortly before the earthquake actually arrives.

7. System Architecture Diagram



8. Conclusion

In this paper, we investigated the real-time nature of Twitter, devoting particular attention to event detection. Semantic analyses were applied to tweets to classify them into a positive and a negative class. We regard each Twitter user as a sensor, and set the problem as detection of an event based on sensory observations. Location estimation methods such as particle filtering are used to estimate the locations of events. As an application, we developed an earthquake reporting system in Java technology using search API for identifying the keywords based on location and the event occurred time, which is a novel approach to notify people promptly of an earthquake event.

9. References

- [1] M. Sarah, C. Abdur, H. Gregor, L. Ben, and M. Roger, "Twitter and the Micro-Messaging Revolution," technical report, O'Reilly Radar, 2008.
- [2] A. Java, X. Song, T. Finin, and B. Tseng, "Why We Twitter: Understanding Microblogging Usage and Communities," Proc. Ninth WebKDD and First SNA-KDD Workshop Web Mining and Social Network Analysis (WebKDD/SNA-KDD '07), pp. 56-65, 2007.
- [3] B. Huberman, D. Romero, and F. Wu, "Social Networks that Matter: Twitter Under the Microscope," ArXiv E-Prints, <http://arxiv.org/abs/0812.1045>, Dec. 2008.
- [4] H. Kwak, C. Lee, H. Park, and S. Moon, "What is Twitter, A Social Network or A News Media?" Proc. 19th Int'l Conf. World Wide Web (WWW '10), pp. 591-600, 2010.
- [5] G.L. Danah Boyd and S. Golder, "Tweet, Tweet, Retweet: Conversational Aspects of Retweeting on Twitter," Proc. 43rd Hawaii Int'l Conf. System Sciences (HICSS-43), 2010.
- [6] A. Tumasjan, T.O. Sprenger, P.G. Sandner, and I.M. Weppe, "Predicting Elections with Twitter: What 140 Characters Reveal About Political Sentiment," Proc. Fourth Int'l AAAI Conf. Weblogs and Social Media (ICWSM), 2010.
- [7] P. Galagan, "Twitter as a Learning Tool. Really," ASTD Learning Circuits, p. 13, 2009.
- [8] K. Borau, C. Ullrich, J. Feng, and R. Shen, "Microblogging for Language Learning: Using Twitter to Train Communicative and Cultural Competence," Proc. Eighth Int'l Conf. Advances in Web Based Learning (ICWL '09), pp. 78-87, 2009.
- [9] J. Hightower and G. Borriello, "Location Systems for Ubiquitous Computing," Computer, vol. 34, no. 8, pp. 57-66, 2001.
- [10] M. Weiser, "The Computer for the Twenty-First Century," Scientific Am., vol. 265, no. 3, pp. 94-104, 1991.
- [11] V. Fox, J. Hightower, L. Liao, D. Schulz, and G. Borriello, "Bayesian Filtering for Location Estimation," IEEE Pervasive Computing, vol. 2, no. 3, pp. 24-33, July-Sept. 2003.
- [12] T. Sakaki, M. Okazaki, and Y. Matsuo, "Earthquake Shakes Twitter Users: Real-Time Event Detection by Social Sensors," Proc. 19th Int'l Conf. World Wide Web (WWW '10), pp. 851-860, 2010.
- [13] Y. Raimond and S. Abdallah, "The Event Ontology," <http://motools.sf.net/event/event.html>, 2007.
- [14] T. Joachims, "Text Categorization with Support Vector Machines: Learning with Many Relevant Features," Proc. 10th European Conf. Machine Learning (ECML '98), pp. 137-142, 1998.
- [15] X. Liu, S. Zhang, F. Wei, and M. Zhou, "Recognizing Named Entities in Tweets," Proc. 49th Ann. Meeting of the Assoc. for Computational Linguistics: Human Language Technologies (HLT '11), pp. 359-367, June 2011.
- [16] A. Ritter, S. Clark Mausam, and O. Etzioni, "Named Entity Recognition in Tweets: An Experimental Study," Proc. Conf. Empirical Methods in Natural Language Processing, 2011.

10. About the Authors



Sri Lakshmi Poornima Gorle is currently pursuing her 2 Years M.Tech (CSE) in Department of Computer Science and Engineering at Avanthi's St.Theresa Institute of Engineering & Technology, Garividi, Vizianagaram District. Her area of interests includes Networks, Information Security.



P.Srinivasu Varma is currently working as Assistant Professor, in Department of Computer Science and Engineering at Avanthi's St.Theresa Institute of Engineering & Technology, Garividi, Vizianagaram District. His research interests include Networks Security & Information Security, Data Mining.