### SUMMARIZATION OF NAMED ENTITY BY TWEET SENTIMENT AND DATA ANALYTICS

#### Sharmila DeviV

Departmentof Computer Science and Engineering, College of Engineering Guindy, Anna university(India)

#### ABSTRACT

Twitter is used by millions of people to share and disseminate timely information. Each tweet is 140 characters in length. The goal of this project is to propose a novel framework for tweet segmentation in a batch mode, called HybridSeg by splitting tweets into meaningful segments. The semantic or context information is well preserved and easily extracted by the downstream applications. HybridSeg finds the optimal segmentation of a tweet by maximizing the sum of the stickiness scores of its candidate segments. The stickiness score considers the probability of a segment being a phrase in English (i.e., global context) and the probability of a segment being a phrase in English (i.e., global context) and the probability of a segment being a phrase in English (i.e., bor the latter, we propose and evaluate two models to derive local context by considering the linguistic features and term-dependency in a batch of tweets. High accuracy is achieved using named entity recognition. Sentiment analysis is used to detect named entity (players) tweets automatically. Since great attention is given to polarity of words i.e. positive or negative and calculate the quality of pitches for that players thrown in MajorLeague Baseball (MLB). Thequality of a particular pitch is evaluated as the expected number of bases conceded. Qualityis expressed as a function of various covariates including pitch velocity, pitch location, pitchtype and batted ball outcomes. The estimation of pitch quality is obtained through the use of regression model to accommodate the inherent complexity of the relationship betweenpitch quality and the associated covariates.

Keywords : Named Entity Recognition (NER), Major League Baseball (MLB).

#### **I.INTRODUCTION**

In recent years, virtual communities and networks on Internet are adopted called as social media. Social media forms such as social networking, or microblogging, allow people to create, and share any kind of information and ideas. Twitter embodies both social networking, microblogging and is a new type of social media. Twitter has become one of the most important communication channels with its ability of providing the most up-to-date and information. There are 255 million monthly active users, and 500 million tweets are sent per day. In this study, tweets are recommended according to the interests of the users. A user interest model is generated where user interests are defined by means of relationship between the user and his friends. Named entities are also extracted from tweets.

Sentiment Analysis (SA) is the most studied field now a days, is also known as Opinion Mining (OP). SA is a task of analyzing people's thoughts, point of view, attitude,towards particular product, topic, organization, etc. As electronic commerce (ecommerce) is increasing rapidly, people's review on the e-commerce website is also

251 | P a g e

# International Journal of Advance Research in Science and Engineering Volume No.07, Special Issue No.(02), March 2018 WWW.ijarse.com

increasing rapidly. Twitter is one of the most used microblogging websites these days. Detecting emotions from twitter posts automatically is a challenging task because informal nature of language is used. As great attention is given to polarity of words i.e. positive or negative, here I am using NRC word emotion association lexicon which have ten emotional categories. As NRC-10 didn't include expressions used every day in social media, here expansion of NRC word - emotion association lexicon is done for language used in social media. This expansion is done using multi-label classifiers and different word level features are extracted from the unlabeled tweets and are compared with each other.

The purpose of this investigation is the assessment of the quality of pitches thrown in MajorLeague Baseball (MLB). An initial reaction may be that this is a straightforward problem. The quality of a particular pitch is evaluated as the expected number of bases conceded. Qualityis expressed as a function of various covariates including pitch velocity, pitch location, pitchtype and batted ball outcomes. The estimation of pitch quality is obtained through the use of regression model to accommodate the inherent complexity of the relationship betweenpitch quality and the associated covariates. With the fitted model, various applications areconsidered which provide new insights on pitching and batting.

#### **II DESIGN METHODOLOGY**

Our work is also related to entity linking (EL). EL identifies the mention of a named entity and links it to an entry in a knowledge base like Wikipedia.Local linguistic features are more reliable than term dependency in guiding the segmentation process. This finding allows the use of tools developed for formal text to be applied to tweets which are noisy than formal text.Semantic meaning of tweets is preserved.The aim of this task is to look at the current moods of the players fanbases to see how they are feeling about their teams and calculate the players tweet using NRC Lexicon sentiment analysis.To calculate the probability of a strikeout for a pitcher matchup in baseball using player descriptors that can be estimated accurately from small samples.To start with the logistic regression model which has been used extensively fordescribing matchups in sports.

#### 2.1 SYSTEM ARCHITECTURE

Natural Language Processing (NLP) is a field of computer science, artificial intelligence and linguistics which helps to interact between computers and human (natural) languages. The development of natural language processing applications is challenging because computer language completely differs from human. First one is often precise, highly structured and unambiguous. From the other side human language is often not precise, ambiguous and its linguistic structure can depend on many things such as slang, social context, regional dialects, etc. Natural language processing applications or toolkits are systems which help to translate human speech into computer and vice versa. The most simple and common tasks are sentence segmentation and tokenization.

# International Journal of Advance Research in Science and Engineering Volume No.07, Special Issue No.(02), March 2018 WWW.ijarse.com



Fig 1. Architecture diagram

The system architecture in Figure1 describes the overall system architecture of the proposed HybridSeg framework segment tweets in batch mode. The entire architecture is divided into different major parts. The first portion is the retrieval module: Tweets are collected from twitter API. Filter the hashtag, URL, Punctuation and Number. Remove the stop words and perform stemming process. Scrape atbat and pitch data is available on the Gameday website (http://www.mlb.com/mlb/gameday/) through XML files. The second portion is the sentiment and emotion identification the aim is to look at the current moods of the players fanbases to see how they are feeling about their teams and calculate the tweets using NRC Lexicon sentiment analysis here lexicon is a set of english word and associates with eight basic emotion (anger, fear, anticipation, trust, surprise, sadness, joy and disgust) and two sentiment (positive and negative). The final step is to take the combined players dataframe and plot the sentiment analysis as percents of total tweets. The third portion is sub event detection here find the quality of pitches in baseball. Quality is expressed as a function of various covariants including pitch velocity, pitch type, pitch location and batted ball outcomes. The fourth portion is summarization.

#### **III IMPLEMENTATION AND RESULTS**

This section discusses the experimentation to efficiently evaluate the accuracy and performance of the proposed system against the existing approach.

#### 3.1 DATA SET

The tweets are collected fromtwitter API. The players data is available on the gamedaywebsite (http://www.mlb.com/mlb/gameday/) through XML files. At the Gameday website,Informationis organized into five tables: pitch, atbat,runner, action and po. Within each of these tables are variables that

### International Journal of Advance Research in Science and Engineering Volume No.07, Special Issue No.(02), March 2018 www.ijarse.com

provide information on the pitch, batsman, runners, and more. Some variables overlap between tables, where the pitch and at-bat tables are by far the largest. The at-bat and pitch tables containall of the variables relevant to our investigation.



#### **3.2 EXPERIMENTAL RESULTS**



Figure 2 shows the snapshot of the combined dataframe and plot the sentiment analysis as percents of total tweets.





#### **IV.PERFORMANCE ANALYSIS**

The accuracy of NER is evaluated by Precision (P), Recall (R) and F measure.

#### 4.1 Precision

Precision is the proportion of classified named entities that belong to the target type. It is defined by equation 4.1.

$$P = \frac{TP}{TP + FP}$$

# International Journal of Advance Research in Science and Engineering Volume No.07, Special Issue No.(02), March 2018 WWW.ijarse.com

#### 4.2 Recall

Recall measures the proportion of named entities of a given type which has been correctly classified is given by equation 4.2.

$$R = \frac{TP}{TP + FN}$$

#### 4.3 F - measure

F-measure combines the effect of the metrics as formulated in equation 4.3.

$$F_1 = 2 \cdot P \cdot R/(P+R)$$

The F-measure scores of locations and organizations indicate that the selected core attributes represent good classification criteria for identifying name entities.

#### 4.4 ROUGE

ROUGE is suitable for measuring the quality of data. To evaluate the metric on a continuous time period [T0,T]. Calculate the integral of the metric over the period, which is given by

 $\int_{T}^{t} metric(t)dt.$ 

#### V. CONCLUSION

Players are modeled using strikeoutand ground ball rate descriptors that can be estimatedreliably using data for a single platoon configuration for asingle season. We used a constrained three-term logit modelto show that the log5 formula provides an accurate model forstrikeout probability and that small changes to the log5 coefficientsmight be used to improve the accuracy of the modelfor the LHP versus RHB platoon configuration. We alsoshowed that a batter/pitcher ground ball rate interactionvariable is highly significant when added to the three-termlogit model. This interaction variable has a strong physicaljustification and adjusts the predicted strikeout probabilitybased on the relative ground ball versus fly ball tendencies of the batter and pitcher. The models were used to show thatbatters are responsible for most of the variance in the predictedstrikeout probability. The method employed to extend the log5 model to include additional variables can easily beadapted for other application areas. This paper has focused on the development and evaluation of low-dimensionalmodels for strikeout probability and anatural next step is to assess the utility of these models for prediction.

#### REFERENCES

 Dehong Gao, Wenjie Li, Xiaoyan Cai, Renxian Zhang, and You Ouyang, "Sequential Summarization: A Full View of Twitter Trending Topics," IEEE Transactions on Audio, Speech, and Language Processing, Vol. 22, No. 2, pp. 293 - 302, February 2014.

### International Journal of Advance Research in Science and Engineering Volume No.07, Special Issue No.(02), March 2018 www.ijarse.com

- [2] Y. Zhang, J. Tang, J. Sun, Y. Chen, and J. Rao, "MoodCast: Emotion prediction via dynamic continuous factor graph model," IEEE Int. Conf. on Data Mining, pp. 1193–1198, 2010.
- [3] Y. Wu, F. Wei, S. Liu, N. Au, W. Cui, H. Zhou, and H. Qu, "OpinionSeer: Interactive visualization of hotel customer feedback," IEEE Transactions on Visualization Comput. Graph., vol. 16, No. 6, pp. 1109–1118, Nov 2010.
- [4] Chih-Yi Chiu, Po-Chih Lin, Sheng-Yang Li, Tsung-Han Tsai, and Yu-Lung Tsai, "Tagging Webcast Text in Baseball Videos by Video Segmentation and Text Alignment," IEEE Transactions on Circuits and Systems for Video Technology, Vol. 22, No. 7, pp. 999 - 1013, July 2012.
- [5] Hassan Ghasemzadeh, and Roozbeh Jafari, "A Signal Processing Model to Evaluate Baseball Swings," IEEE Sensors Journal., Vol. 11, No. 3, pp. 603 - 610, March 2011.
- [6] Pradeep D. Prasad, Harsha N. Halahalli, John P. John and Kaushik K. Majumdar, "Single-Trial EEG Classification Using Logistic Regression Based on Ensemble Synchronization," IEEE Journal of Biomedical and Health Informatics, Vol. 18, No. 3, pp. 1074 - 1080, May 2014.