

CLASSIFICATION OF POSTS RELATED TO FACEBOOK BASED ON SPECIFIC CATEGORY

¹G.Sravani, ²A.Vamsi Krishna, ³Dr..B.GeethaVani

¹Assistant Professor, ² Assistant Professor, ³Professor
Dept. of CSE,
Narayana Engineering College, Nellore, A.P, India

ABSTRACT- The increase of interest in using social media as a source for research has motivated tackling the challenge of automatically geolocating messages , given the lack of explicit location information in the majority of messages. In contrast to much previous work that has focused on classification of messages(posts) restricted to a specific category, here we undertake the task in a broader context by classifying global messages, which is so far unexplored in a real-time scenario. To develop a system that classifies Facebook posts , Extract texts from posts, images, videos of social media as Face book to know how people feel about different posts, Classify the post gives probable effect of the original post, so its easier to understand social affect of any post on Facebook or any such social media.So we focused on Facebook statuses, which we can view as opinions of users or their reaction on concern we want to analyze. We develop tool status puller that automatically collects random Facebook statuses. Then we make classifier that performs classifications on that corpus collected from Facebook.The dramatic and exponential growth of content available on web and its classification has now become an efficient methodology to make the contents of large repository in an organized manner. Social networking websites are the new era of expressing views. Today every fifth person put their opinions, views, comments on these micro-blogging and social sites like, facebook and many more. Now people are using internet as a communication tool among their social network including friends, family, friends of friends to these micro-blogging and social network sites. In this we gradually put and share their opinions among their friends on these sites which finally becomes huge and relevant repository for any of particular entity or organization

Index terms- geolocating messages, micro-blogging

1. INTRODUCTION

Social media are increasingly being used in the scientific community as a key source of data to help understand diverse natural and social phenomena, and this has prompted the development of a wide range of computational data mining tools that can extract knowledge from social media . Thanks to the availability of a public API that enables the cost-free collection of a significant amount of data. Having Facebook as a new kind of data source, researchers have looked into the development of tools for real-time trend analytics or early detection of newsworthy events , as well as into analytical approaches for understanding the sentiment expressed by users towards a target , or public opinion on a specific topic. They use ground truth labels to pre-filter posts originating from other regions and/or written in languages other than English. The classifier built on this pre-filtered dataset may not be applicable to a stream where every post needs to be classified. An ability to classify,in real-time is crucial for applications exploiting social media updates as social sensors that enable tracking topics and learning about location-specific trending topics, emerging events and breaking news.To the best of our knowledge, our work is the first to deal with global posts in any language, using only those features present within the content of a post and its associated metadata. We also complement previous work by investigating the extent to which a classifier trained on historical posts can be used effectively on newly harvested posts.To develop a system that classifies Facebook posts ,extract texts from posts, images, videos of social media as Facebook to know how people feel about different posts and classifying the post gives probable effect of the original post, so its easier to understand social affect of any post on Facebook or any such social media.So we focused on Facebook statuses, which we can view as opinions of users or their reaction on concern we want to analyze. We develop tool status puller that automatically collects random Facebook statuses. Then we make classifier that performs classifications on that corpus collected from Facebook.The dramatic and exponential growth of content available on web and its classification has now become an efficient methodology to make the contents of large repository in an organized manner. Social networking websites are the new era of expressing views. Today every fifth person put their opinions, views, comments on these micro-blogging and social sites like, FACEBOOK and many more. Now people are using internet as a communication tool among their social network including friends, family, friends of friends to these micro-blogging and social network sites. In this we gradually put and share their opinions among their friends on these sites which finally becomes huge and relevant repository for any of particular entity or organization

II. PROBLEM DEFINITION

The motivation of this project is to implement design for the functional minimum-storage regenerating (FMSR) codes. The system FMSR code implementation maintains double-fault tolerance and has the same storage cost as in traditional erasure coding schemes based on RAID-6 codes, but uses less repair traffic when recovering a single-cloud failure. The aim of our project is to minimize the cost of storage repair (due to the migration of data over the clouds) for a permanent single-cloud failure. In this work, we focus on comparing two codes: traditional RAID-6 codes and our FMSR codes with double-fault tolerance. The main objective of our project is to minimize the cost of storage repair (due to migration of data over the clouds) for a permanent single-cloud failure. We focus on comparing: traditional RAID-6 codes and FMSR codes with double-fault tolerance.

III. LITERATURE REVIEW

In this paper, we present Google, a prototype of a large-scale search engine which makes heavy use of the structure present in hypertext. Google is designed to crawl and index the Web efficiently and produce much more satisfying search results than existing systems. The prototype with a full text and hyperlink database of at least 24 million pages is available at <http://google.stanford.edu/>. To engineer a search engine is a challenging task. Search engines index tens to hundreds of millions of web pages involving a comparable number of distinct terms. They answer tens of millions of queries every day. Despite the importance of large-scale search engines on the web, very little academic research has been done on them. Furthermore, due to rapid advance in technology and web proliferation, creating a web search engine today is very different from three years ago. This paper provides an in-depth description of our large-scale web search engine -- the first such detailed public description we know of to date. Apart from the problems of scaling traditional search techniques to data of this magnitude, there are new technical challenges involved with using the additional information present in hypertext to produce better search results. This paper addresses this question of how to build a practical large-scale system which can exploit the additional information present in hypertext. Also we look at the problem of how to effectively deal with uncontrolled hypertext collections where anyone can publish anything they want.

Expertise identification using email communications.

A common method for finding information in an organization is to use social networks---ask people, following referrals until someone with the right information is found. Another way is to automatically mine documents to determine who knows what. Email documents seem particularly well suited to this task of "expertise location", as people routinely communicate what they know. Moreover, because people explicitly direct email to one another, social networks are likely to be contained in the patterns of communication. Can these patterns be used to discover experts on particular topics? Is this approach better than mining message content alone? To find answers to these questions, two algorithms for determining expertise from email were compared: a content-based approach that takes account only of email text, and a graph-based ranking algorithm (HITS) that takes account both of text and communication patterns. An evaluation was done using email and explicit expertise ratings from two different organizations. The rankings given by each algorithm were compared to the explicit rankings with the precision and recall measures commonly used in information retrieval, as well as the d' measure commonly used in signal-detection theory. Results show that the graph-based algorithm performs better than the content-based algorithm at identifying experts in both cases, demonstrating that the graph-based algorithm effectively extracts more information than is found in content alone.

Mining the network value of customers.

One of the major applications of data mining is in helping companies determine which potential customers to market to. If the expected profit from a customer is greater than the cost of marketing to her, the marketing action for that customer is executed. So far, work in this area has considered only the intrinsic value of the customer (i.e., the expected profit from sales to her). We propose to model also the customer's network value: the expected profit from sales to other customers she may influence to buy, the customers those may influence, and so on recursively. Instead of viewing a market as a set of independent entities, we view it as a social network and model it as a Markov random field. We show the advantages of this approach using a social network mined from a collaborative filtering database. Marketing that exploits the network value of customers -- also known as viral marketing -- can be extremely effective, but is still a black art. Our work can be viewed as a step towards providing a more solid foundation for it, taking advantage of the availability of large relevant databases.

Measuring user influence on twitter using modified k-shell decomposition.

We survey the several measures that exists in literature to rank influential users in Twitter network. We propose a classification of these measures according to different criteria, such as the kind of metrics, the use of Page Rank algorithm, the use of content analysis, among others. Besides the influential users, we also study measures of activity and popularity. We finish by mentioning some aspects of this topic related with computational complexity and correlation measures. Centrality is one of the most studied concepts in social network analysis. There is a huge literature regarding centrality measures, as ways to identify the most relevant users in a social network. The challenge is to find measures that can be computed efficiently, and that can be able to classify the users according to relevance criteria as close as possible to reality. We address this problem in the context of the Twitter network, an online social networking service with millions of users and an impressive flow of messages that are published and spread daily by interactions between users. Twitter has different types of users, but the greatest utility lies in finding the most influential ones. The purpose of this article is to collect and classify

the different Twitter influence measures that exist so far in literature. These measures are very diverse. Some are based on simple metrics provided by the Twitter API, while others are based on complex mathematical models. Several measures are based on the Page Rank algorithm, traditionally used to rank the websites on the Internet. Some others consider the timeline of publication, others the content of the messages, some are focused on specific topics, and others try to make predictions. We consider all these aspects, and some additional ones. Furthermore, we include measures of activity and popularity, the traditional mechanisms to correlate measures, and some important aspects of computational complexity for this particular context.

A Huber man Influence and passivity in social media

The ever-increasing amount of information flowing through Social Media forces the members of these networks to compete for attention and influence by relying on other people to spread their message. A large study of information propagation within Twitter reveals that the majority of users act as passive information consumers and do not forward the content to the network. Therefore, in order for individuals to become influential they must not only obtain attention and thus be popular, but also overcome user passivity. We propose an algorithm that determines the influence and passivity of users based on their information forwarding activity. An evaluation performed with a 2.5 million user dataset shows that our influence measure is a good predictor of URL clicks, outperforming several other measures that do not explicitly take user passivity into account. We demonstrate that high popularity does not necessarily imply high influence and vice-versa.

Online health communities (OHCs) have become a major source of support for people with health problems. This research tries to improve our understanding of social influence and to identify influential users in OHCs. The outcome can facilitate OHC management, improve community sustainability, and eventually benefit OHC users.

IV. SYSTEM ANALYSIS

Twitter user ranking algorithm was proposed to identify authoritative users who often submit useful information. Page Rank algorithm named Twitter Rank was developed to rank Twitter users based on their influence. K-shell decomposition algorithm is developed to measure the user influence in Twitter. Most of these measurements quantify the influence in an isolated way, rather than in a collective way. The proposed algorithm mainly works based on the user-tweet graph, rather than the user-user social graph. The focus of these methods is on influence, which is still different from the vitality. Many interactions often keep going on within online social networks over time. Examples of interaction include but are not limited to the retweeting, mention, and sending message. Our goal is to rank user vitality based on all interactions in a time period. Suppose that we have a social network S that contains N users (nodes) denoted as $fUjg1_j_N$ and L links among users denoted as $fEjkg1_j;k_N$, where j and k are indices. We have recorded all interactions between them within M consecutive time periods Ti ($1 \leq i \leq M$). Our goal is to rank all users from high vitality to low vitality for a time period Ti based on all previously observed interactions. Such a vitality-based ranking list of users may provide a good guidance for the social networking service providers to understand the dynamics of systems. They may directly find the relatively most active users and make better operation and business decisions upon the findings. The accurate results of both user vitality ranking and prediction could benefit many parties in different social networking services, a user vitality ranking list could help ads providers to better display their ads to active users and reach more audiences.

V.IMPLEMENTATION

OSN MODULE

In the first module, we develop the Online Social Networking (OSN) system module. We build up the system with the feature of Online Social Networking. Where, this module is used for new user registrations and after registrations the users can login with their authentication. Where after the existing users can send messages to privately and publicly, options are built. Users can also share post with others. The user can able to search the other user profiles and public posts. In this module users can also accept and send friend requests. With all the basic feature of Online Social Networking System modules is build up in the initial module, to prove and evaluate our system features.

VITALITY RANKING MODULE

The accumulated number of interactions SA_i of a node j ($1 \leq j \leq N$) in time period i ($1 \leq i \leq M$) within a social network I is defined as:

$$SA_j^i = \sum_{k \in \{Nei(j)\}} \theta_{kj},$$

Calculating Average Interaction

The relative increase of interaction of a node j ($1 \leq j \leq N$) in time period i ($1 \leq i \leq M$) within a social network I is defined as:

$$IA_j^i = \frac{SA_j^i}{SA_j^{i-1}}.$$

Calculating Vitality Ranking

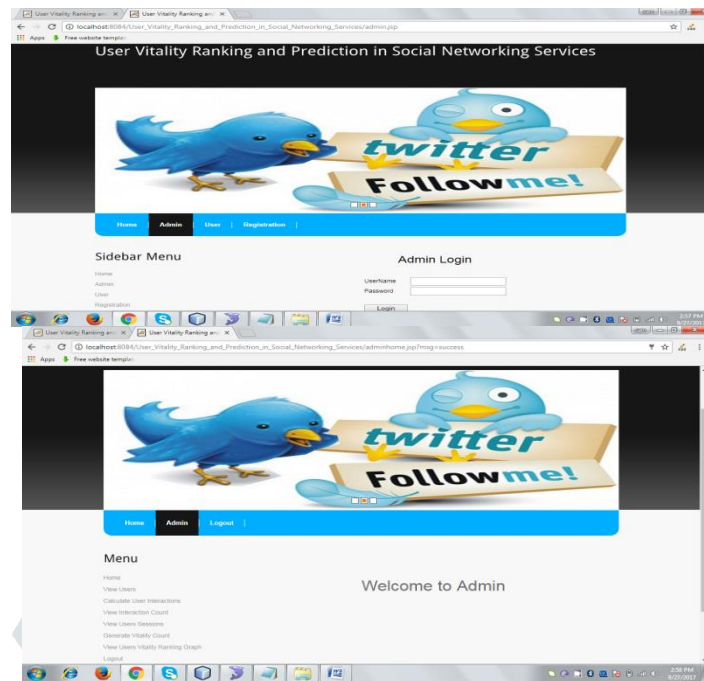
Given the number of interactions between all pairs of users, we may count the number of all interactions for each user and rank them based on the count using Initial Ranking Algorithm. However, given the number of interactions between two nodes (users), The unified vitality score $_{ij}$ of a node j ($1 \leq j \leq N$) in time period i ($1 \leq i \leq M$) within a social network I is defined as:

ALGORITHM

- The accumulated number of interactions SA_{ij} of a node j ($1 \leq j \leq N$) in time period i ($1 \leq i \leq M$) within a social network I is defined as $SA_{ij} = \sum_{k \in fNei(j)} g_{kj}$, where $fNei(j)$ denotes the set of users that are connected to user j .
 - First, if the accumulated number of interaction of node (i.e., SA_{ij}) increases a lot over that in the previous time period.
 - The relative increase of interaction of a node j ($1 \leq j \leq N$) in time period i ($1 \leq i \leq M$) within a social network I is defined as: $IA_{ij} = SA_{ij} / SA_{ij}^{i-1}$.
 - we can get that the relative increase of interaction for node A and node C in time period n .
 - The average interaction for user S_i is defined as: $Average_{ij} = SA_{ij} / degree_{ij}$ where $degree_{ij}$ denotes the number of connected friends for user j .
 - The term $Average_{ij}$ represents the average number of interactions of User j in period i .
 - The unified vitality score $_{ij}$ of a node j ($1 \leq j \leq N$) in time period i ($1 \leq i \leq M$) within a social network I is defined by average of SA_{ij} and IA_{ij} .
 - we define the user's vitality score $_{ij}$ with two terms and combine them in a linear way.
 - The first part, SA_{ij} , indicates the dynamic vitality level in period i .
 - The second part denotes the static vitality level of user in one period.
 - By tuning the parameter, we may balance the impact of the relative increase of interaction and the average interaction.
 - In the experiment, we will empirically examine the impact of $_{ij}$ on the performance of our algorithm. Given a social network, we will compute the $_{ij}$ for all nodes (users) for a specified time period, and then rank all nodes according to the value.

VI. EXPERIMENTAL RESULTS





VII. CONCLUSION AND FUTURE SCOPE

In this project, we presented a study on user vitality ranking and prediction in social networking services such as micro blog application. Specifically, we first introduced a user vitality ranking problem, which is based on dynamic interactions between users on social networks. To solve this problem, we developed two algorithms to rank users based on vitality. While the first algorithm works based on the developed two user vitality measurements, the second algorithm further takes into account the mutual influence among users while computing the vitality measurements. Then we presented a user vitality prediction problem and introduced a regression based method for the prediction task. Intensive experiments on two real-world data sets that are collected from different domains clearly demonstrate the effectiveness of our ranking and prediction methods. The accurate results of both user vitality ranking and prediction could benefit many parties in different social networking services, e.g., a user vitality ranking list could help ads providers to better display their ads to active users and reach more audiences.

REFERENCES

1. Eytan Bakshy, Jake M Hofman, Winter A Mason, and Duncan J Watts. Everyone's an influencer: quantifying influence on twitter. In Proceedings of the fourth ACM international conference on Web search and data mining, pages 65–74. ACM, 2011.
2. Sergey Brin and Lawrence Page. Reprint of: The anatomy of a largescalehypertextual web search engine. Computer networks, 56(18):3825–3833, 2012.
3. Robert Goodell Brown. Smoothing, forecasting and prediction of discrete time series. Courier Corporation, 2004.
4. Christopher S Campbell, Paul P Maglio, Alex Cozzi, and Byron Dom. Expertise identification using email communications. In Proceedings of the twelfth international conference on Information and knowledge management, pages 528–531. ACM, 2003.
5. Meeyoung Cha, Hamed Haddadi, Fabricio Benevenuto, and P Krishna Gummadi. Measuring user influence in twitter: The million follower fallacy. ICWSM, 10(10-17):30, 2010.
6. Pedro Domingos and Matt Richardson. Mining the network value of customers. In Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining, pages 57–66. ACM, 2001.
7. Philip E Brown Junlan Feng. Measuring user influence on twitter using modified k-shell decomposition. 2011.
8. Jian Jiao, Jun Yan, Haibei Zhao, and Weiguo Fan. ExpertRank: An expert user ranking algorithm in online communities. In New Trends in Information and Service Science, 2009. NISS'09. International Conference on, pages 674–679. IEEE, 2009.

9. Jon M Kleinberg. Authoritative sources in a hyperlinked environment. *Journal of the ACM (JACM)*, 46(5):604–632, 1999.
10. Shamanth Kumar, Fred Morstatter, and Huan Liu. *Twitter data analytics*. Springer, 2014.
11. Wenting Liu, Guangxia Li, and James Cheng. Fast pagerank approximation by adaptive sampling. *Knowledge and Information Systems*, 42(1):127–146, 2015.
12. Rada Mihalcea. Graph-based ranking algorithms for sentence extraction, applied to text summarization. In *Proceedings of the ACL 2004 on Interactive poster and demonstration sessions*, page 20. Association for Computational Linguistics, 2004.
13. Amin Omidvar, Mehdi Garakani, and Hamid R Safarpour. Context based user ranking in forums for expert finding using wordnet dictionary and social network analysis. *Information Technology and Management*, 15(1):51–63, 2014.
14. Tore Opsahl, Filip Agneessens, and John Skvoretz. Node centrality in weighted networks: Generalizing degree and shortest paths. *Social networks*, 32(3):245–251, 2010.
15. Daniel M Romero, Wojciech Galuba, Sitaram Asur, and Bernardo A Huberman. Influence and passivity in social media. In *Machine learning and knowledge discovery in databases*, pages 18–33. Springer, 2011.
16. Yuanfeng Song, Wilfred Ng, Kenneth Wai-Ting Leung, and Qiong Fang. Sfp-rank: significant frequent pattern analysis for effective ranking. *Knowledge and Information Systems*, 43(3):529–553, 2015.
17. Stanley Wasserman and Philippa Pattison. Logit models and logistic regressions for social networks: I. an introduction to markov graphs and p. *Psychometrika*, 61(3):401–425, 1996.
18. Jianshu Weng, Ee-Peng Lim, Jing Jiang, and Qi He. Twitter rank: finding topic-sensitive influential twitterers. In *Proceedings of the third ACM international conference on Web search and data mining*, pages 261–270. ACM, 2010.
19. Yuto Yamaguchi, Tsubasa Takahashi, Toshiyuki Amagasa, and Hiroyuki Kitagawa. Turank: Twitter user ranking based on user-tweet graph analysis. In *Web Information Systems Engineering–WISE 2010*, pages 240–253. Springer, 2010.
20. Yang Yang, Ryan N Lichtenwalter, and Nitesh V Chawla. Evaluating link prediction methods. *Knowledge and Information Systems*, 45(3):751–782, 2015.
21. Kang Zhao, John Yen, Greta Greer, Baojun Qiu, Prasenjit Mitra, and Kenneth Portier. Finding influential users of online health communities: a new metric based on sentiment influence. *Journal of the American Medical Informatics Association*, 21(e2):e212–e218, 2014.