

Linguistic-Based and user-Based Recommending Posts using Two-Level Clustering Methods

Sayali Joag, Rupali Dalvi

Abstract: Online social networks have produced bunches of online social groups where individuals can collaborate and shuffle their thoughts. In spite of, the real problems that conflicts with the user security and convenience are confidentiality break, groups without inception, confusion created from various groups of which user is a member of and difficulty in moderating groups. This can be moderated to an extent by an automated filtering method required to categorizing group members based on their response patterns. This paper proposes clustering of group posts on stylistic, thematic, emotional, sentimental and psycholinguistic methods and members of the group are categorized based on their responses to the posts belonging to different clustering methods. The categorization affords security just like the conflict associated with irrelevant notifications received from more than one groups, via recommending the users, posts which might be probable to be of interest to them. It also helps to identify the group members meant closer to spreading posts that violate group policies. The categorization post shows increased performance where there are large numbers of members in a social group by performing linguistic clustering. The contribution work is to implement location-aware personalized posts recommendation using users' behavioral patterns and their geographic location. Another, important work is to implement text-to-speech system converting English text into speech using speech synthesis technique. The system gives rating to the users who shares posts depending on clustering. Also system provides read later functionality to the user side. The system has been tested on Twitter API group data where a significant solution to an unaddressed problem associated with social networking groups is offered.

Keywords: Emotion analysis, Multi-level clustering, Psycholinguistics, Sentiment analysis, Stylistics Clustering.

I. INTRODUCTION

The most popular social networking site Twitter has groups with over 100K members in it. Thus, it becomes difficult for the admin to track the members violating group policies. This shows that there exists a need for a measure to categorize the posts made by members in it based on acceptability and group behavior. A community is a fraternity that seeks a platform to discuss subjects that relate to the common cause that motivated the establishment of the group. The members

Revised Manuscript Received on September 05, 2019.

* Correspondence Author

Sayali Joag, ME Student, Department of Computer Engineering, Marathwada Mitramandal's College of Engineering, Pune, India
Email: joagsayali@gmail.com

Ms. Rupali Dalvi, Professor, Department of Computer Engineering, Marathwada Mitramandal's College of Engineering, Pune, India,
Email: rupalidalvi@mmcoe.edu.in

within it enjoy discussion with regard to the purpose of the group. This is especially applicable in the case of academic and interest groups. Thus posting articles irrelevant to these groups can create unnecessary clutter that affects the comfort of the members within the group. Unnecessary advertisements and marketing matters targeting a different audience should be avoided from a group. There is no existing method to deal with this.

A method to control the influx of irrelevant messages within a community is very much essential for the smooth functioning of a social networking group. The existing policies of popular social networking site Twitter enable the administrator of a group to delete and monitor posts made by the members of a community. To screen all the messages posted by the members of a heavily populated group is very difficult. The existing settings provide no option for automated notification for group admin with regard to the members involved in posting articles that do not match the general interest of the group. Our proposed method aims at providing an automated notification to the group moderators regarding suspicious members who frequently post articles that appear to gather negative response from the members within the group. This novel approach has been experimented to account for the hypothesis that, though members may be socially connected; there might be disparity in their likings, thoughts, sentimental orientation and thematic inclination.

Motivation

To provide an automated notification to the group moderators regarding suspicious members who frequently post articles that appear to gather negative response from the members within the group.

II. REVIEW OF LITERATURE

The paper [1] proposes a new content based method for personalized tweet recommendation, based on conceptual relations between users' topics of interest. The Concept Graph is a way to exploit logical relations between topics of interest in order to provide interesting and efficient tweet recommendations. Advantages are: Provides a social media user with a new timeline that contains messages that strongly match ones interests and that are not necessarily posted by ones followings. This model is effective and efficient to recommend interesting tweets to users. Disadvantages are: The recommender still recommends some tweets that were not retreated by the user.

The paper [2] proposes a TWIMER framework the use of language models as a basis for analyzing strategies and techniques for tweet advice based on person's interest profiles. TWIMER consists of

several components, including tweet retrieval (query formation and relevance model), tweet relevance verification, and final relevance ranking. Advantages are: Automatic query expansion. Higher performance in the language model retrieval entities.

In [3] paper, collaborative topic Poisson factorization (CTPF) can be used to build recommender systems through gaining knowledge of from reader histories and content material to advocate personalized articles of hobby. CTPF models both reader behavior and article texts with Poisson distributions, connecting to the latent subjects that represent the texts with the latent choices that show to the readers. Advantages are: CTPF performs well in the face of massive, sparse, and long-tailed data. CTPF offers a natural mechanism to resolve the “cold start” hassle. CTPF scales more easily and provides significantly better recommendations than CTR.

In [4] paper, represents the trouble of concurrently predicting consumer decisions and modeling customers’ interests in social media with the aid of reading rich statistics collected from Twitter. Proposes Co-Factorization Machines (CoFM), which deal with two (multiple) aspects of the dataset where each aspect is a separate FM. This type of model can easily predict user choices even as modeling user interests via content material at the equal time. Advantages are: CoFM can easily predict user decisions while modeling user interests through content at the same time. Factorization Machines to text data with constraints can mimic state-of-the art topic models and yet benefit from the efficiency of a simpler form of modeling. Disadvantages are: The services are only interacting with the user’s interest not on user’s behaviors.

The paper [5] focuses on recommending useful tweets that users are really interested in personally to reduce the users’ effort to find useful information. The topic level latent factors of tweets to capture users’ common interests over tweet content, which helps us to solve the problem of information sparsity in users’ retweet actions. This allows us to adjust the collaborative filtering technique to solve the recommendation problem. Advantages are: A collaborative ranking method is better than collaborative filtering for different optimization criterion. The proposed CTR method greatly improves the recommendation performance. The CTR method is generic; it is easy to incorporate more information by adding extra features. Disadvantages are: It only works on user’s interests over time not on user’s history and tags of the tweet.

In [6] paper, fill the distance among current assessment-summarization and evaluation selection strategies with the aid of selecting a small subset of reviews that together hold the statistical properties of the whole assessment corpus. The proposed three heuristic algorithms for selecting a characteristic review set, which evaluate via experiments on a wide range of datasets of real reviews from different domains. Advantages are: To accurately emulate the opinion distribution in the underlying corpus. Improvement by previous work. Disadvantages are: Positive and negative comment can’t generalize to arbitrary domain.

In [7] paper, consider the review set selection problem where given a set of reviews for a specific item, and want to select a comprehensive subset of small size. Provides authentic review using TOPQLTY algorithm sorting technique problem is based on limited review set. Advantages are: Performance is statically significant. High quality of the review.

The paper [8] presents proposed methods to compute quality, diversity and coverage properties using multidimensional content and context data. The proposed metrics so as to evaluate the picture summaries based totally on their illustration of the bigger corpus and the capacity to meet user’s information needs. Advantages are: The greedy algorithm for summarization performs better than the baselines. Summaries help in effective sharing and browsing of the personal photos. Disadvantages are: Computation is expensive.

In paper [9], there are two methods for incorporating social context in the quality prediction: either as features, or as regularization constraints, based on a set of hypotheses. The method proposes quite generalizable and applicable for quality (or attribute) estimation of other types of user-generated content. Advantages are: Improve the accuracy of review quality prediction. The resultant predictor is utilizable even when social context is unavailable. Disadvantages are: A portal may lack an explicit trust network.

In multi-document summarization [10], redundancy is a particularly important issue since textual units from different documents might convey the same information. A high quality (small and meaningful) summary should not only be informative about the remainder but also be compact (non-redundant). Advantages are: The best performance is achieved. Submodular summarization achieves better ROUGE-1 scores. Disadvantages are: The proposed system very expensive to solve.

III. OPEN ISSUES

A method to control the influx of irrelevant messages within a community is very much essential for the smooth functioning of a social networking group. The existing policies of popular social networking site Twitter enable the administrator of a group to delete and monitor posts made by the members of a community. To screen all the messages posted by the members of a heavily populated group is very difficult. The existing settings provide no option for automated notification for group admin with regard to the members involved in posting articles that do not match the general interest of the group.

Disadvantages are:

1. Does not solve the clutter issues and group policy management hassles associated with social networking groups.
2. None of the methods deal with the issue of providing customized suggestions to posts in a community to reduce clutter associated with notifications generated by large groups.
3. Does not provide automated notification for group admin with regard to the members involved in posting articles.
4. Need to categorize the posts made by members in it based on acceptability and group behavior.

IV. PROPOSED METHODOLOGY

1) Sentiment Clustering

The posts in a network group vary based on the sentimental status associated with the contents. The overall sentiment associated with a post can be



positive, negative or neutral [3]. The assessment of the sentiment is made at the keyword level. This enables to classify the posts based on the attitude towards the trends. The response of a member to a post shows his/her attitude towards the entities. This can be illustrated by an example below:

- User A says “I love ice cream. I Love to have it every time”
 - While User B says “I hate ice cream. I don’t want it at all”
- Here the keyword is ice-cream and the sentiment of the keyword with respect to user A is positive while that of user B is negative. Thus, analyze the sentiments associated with prominent keywords associated with a group.

2) Theme Clustering

The next type of clustering performed on the posts is the theme based clustering. Here the posts are grouped based on the thematic similarities. After consider concepts, entities and topics for determining the theme. Use the semantically rich common sense knowledge base Concept Net for extracting contextual and conceptual information of a post. The Concept Net toolkit provides numerous assertions related to a word [1]. Utilize them to arrive at a general concept of a post. Initially the posts are chunked to obtain keywords and key phrases. These phrases are fed to the Concept Net to obtain generalized concepts. After exploit the analogy making and topic gusting features of Concept Net to arrive at a conclusion regarding the important concepts of a post.

3) Emotional Clustering

Emotion based aspects can be investigated to gain insight into a person’s emotional attitude. The similarity in these aspects can be exploited to predict the liking and sharing probability of members within a community. Thus, the users within a community are grouped together based on the emotional aspects. The clustering based on emotion is performed by taking into account the score of the entire post [5] with respect to the emotional categories namely anger, sadness, fear, disgust and joy.

4) Stylistic Clustering

Stylistics refers to the writing-style followed in a document. Each person has a unique writing style of his or her own. The writing-style factor has been considered because the writing style followed determines the popularity of an article. There are phrases and usages peculiar to authors that can be of interest to the audience. This aspect has been exploited to find the like-minded audience of a particular style of writing. The stylistic features such as words and character n-grams can capture the writing style effectively. These features are clustered by using K-means clustering.

5) Psycholinguistic Clustering

The posts have been subjected to the analysis of psycholinguistic orientation of the posts. The psycholinguistic differences contribute to the nature of the posts and hence there will be difference in the audience based on the psycholinguistic perspective of members within a group. Psycholinguistic aspects throw light on various factors that vary based on a person’s personality, hobbies, passions, intellects, perception and context of references [2]. Thus it reflects a person’s way of responding to a post.

The contribution of the proposed method can be summed up as follows:

1. The method proposes a two-level clustering mechanism that aims at categorizing the members of a group based on interest.

2. The method is extended to provide features for group policy management.
3. The contribution work is to implement location-aware personalized posts recommendation using both the users’ personal interests and their geographical contexts.
4. Another, important work is to implement text-to-speech system converting English text into speech using speech synthesis technique.
5. The system provides the read later posts functionality to the users.
6. The system gives rating to the users who shares posts depending on clustering.
7. The experiments have been conducted on limited content available for processing which makes the experiment standout from the state-of-the-art techniques that rely on lengthy text for processing.
8. The method is scalable in huge data-set.

Advantages are:

- The approach is independent of the pre-existing social relations for members within a group and the approach considers behavioral similarity between the users in terms of their response pattern towards different aspects.
- The approach proves to be beneficial in tackling cold start problem, i.e. the cases where a new post has not been rated yet.
- Like minded members within a group are obtained by considering similarities in various dimensions.
- The first work is based on recommending posts relevant to a user based on interests expressed by them to the content of the posts, within a social network group and also one that helps in identifying members aimed at spreading unpopular posts within a group.
- The system is beneficial for users to listen the posts.
- The system also provides read later posts functionality to the users.
- The system also provides rating to the user based on the posts clustering.

A. Architecture

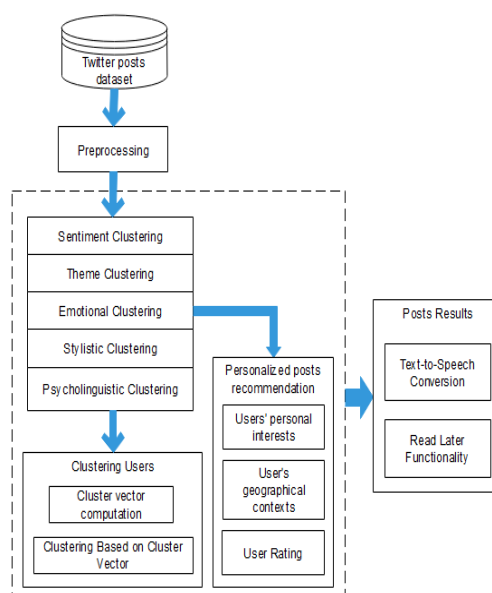


Fig. 1 Proposed System Architecture



B. Mathematical Model

1. Clustering Based On Cluster Vector:

To check the validity of the similar users formed from initial clusters, the members within a cluster are made to undergo similarity check. This is to check the confidence values of the users. It is measured using Pearson’s correlation coefficient, which can be computed as follows:

$$sim_{m,n} = \frac{\sum_{p \in P} (r_{mp} - \bar{r}_m)(r_{np} - \bar{r}_n)}{\sqrt{\sum_{p \in P} (r_{mp} - \bar{r}_m)^2} \sqrt{\sum_{p \in P} (r_{np} - \bar{r}_n)^2}} \quad (1)$$

Where m and n are users whose similarity is to be identified. P refers to the set of posts which are liked or shared by at least one of them. r_{mp} Refers to the posts liked by m and \bar{r}_m refers to the average rating of user m.

2. Group Policy Management:

Acceptability:

This term is defined in association with the posts. A post is checked for its acceptability with respect to how the different aspects (sentimental, psycholinguistic, stylistic, emotional and thematic) of the post are liked by the members of the community. It shows an aspect based score.

$$Acceptability(p) = \sum_{k} \sum_{i \in clues} val(C_{mi}) \quad (2)$$

Where, clues is a list that stores the cluster of the post p with respect to different aspects. $val(C_{mi})$ Gives the vector values (0 or 1) associated with the cluster vector of member m with respect to the aspect i.

Popularity: The term popularity is defined for both posts and members within a community.

$$popularity(p) = Acceptability(p) + share(p) + (com_{pos} + like_{s_{pos}}) - (com_{neg}(p) + like_{s_{neg}}) \quad (3)$$

Where $share(p)$ is the number of shares for the post p, com_{pos} is the number of positive comments for p, $like_{s_{pos}}$ is the like obtained for positive comments, com_{neg} is the number of negative comments for p and $like_{s_{neg}}$ is the negative comments obtained for p.

C. Algorithms

1. Sentiment Analysis using Sentiwordnet Dictionary

```

polarizedTokensList ← newList()
while tokenizedTicket.hasNext() do
token ← tokenizedTicket.next()
lemma ← token.lemma
polarityScore ← null
if DomainDictionary.contains(lemma,pos) then
if SentiWordNet.contains(lemma,pos) and
SentiWordNet.getPolarity(lemma,pos) != 0 then
polarityScore ← SentiWordNet.getPolarity(lemma, pos)
else
domainDicToken ← DomainDictionary.getToken(lemma, pos)
if domainDicToken.PolarityOrientation == "POSITIVE" then
polarityScore ← DefaultPolarity.positive
else
polarityScore ← DefaultPolarity.negative
    
```

end if

end if

polarizedTokensList.add(token, polarityScore)

end if

end while

return polarizedTokensList

2. Latent Dirichlet Allocation (LDA) Algorithm:

First and major, LDA offers a generative model that describes how the files in a dataset were created. In [9] context, a dataset is a group of D files. Document is a set of phrases. So our generative version describes how each document obtains its phrases. Initially, permits anticipate that recognize K subject matter distributions for our dataset, meaning K multinomial containing V elements each, where V is the wide variety of terms in our corpus. Let β_i represent the multinomial for the ith topic, where the size of β_i is V: $|\beta_i|=V$. Given these distributions, the LDA generative method is as follows:

Steps:

1. for every document:
 - (a) Randomly choose a distribution over subjects (a multinomial of length K)
 - (b) For each word within the document:
 - (i) Probabilistically draw one of the K subjects from the distribution over topics obtained in (a), say topic β_j
 - (ii) Probabilistically draw one of the V words from β_j

3. K-means Clustering Algorithm

Step 1. Let $X = \{x_1, x_2, \dots, x_n\}$ be the set of data points and $V = \{v_1, v_2, \dots, v_c\}$ be the set of centers.

Step 2. Randomly select ' c ' cluster centers.

Step 3. Calculate the distance between each data point and cluster centers.

Step 4. Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers.

Step 5. Recalculate the new cluster center using:

$$v_i = \left(\frac{1}{c_i}\right) \sum_{j=1}^{c_i} x_j$$

Where, ' c_i ' represents the number of data points in i^{th} cluster.

Step 6. Recalculate the distance between each data point and new obtained cluster centers.

Step 7. If no data point was reassigned then stop, otherwise repeat from step 3.

4. Cluster Vector Formulation of Users Algorithm

1. $M \leftarrow$ {members in a group}
2. $A \leftarrow$ {sentiment, theme, emotion, stylistics, psycholinguistic}
3. $M_like[m] \leftarrow$ posts liked by member m
4. $Cluster[p][a] \leftarrow$ cluster of post p for aspect a
5. for $p \in M$ do
6. for $a \in A$ do
7. $x_i = 0$
8. for $p \in M_like[m]$ do
9. $c \leftarrow cluster[p][a]$



10. $m_cluster_a[x_i] += 1$
11. $cv[a] \leftarrow \text{argmax}_i(m_cluster_a[x_i])$
12. for $i := 1$ to $n(m_cluster_a)$ do
13. if $i = cv[a]$ then
14. $x_i = 1$
15. Else
16. $x_i = 0$
17. $merge(m_cluster_a)$

V. RESULT AND DISCUSSIONS

Experiments are done by a personal computer with a configuration: Intel (R) Core (TM) i3-2120 CPU @ 3.30GHz, 4GB memory, Windows 7, MySQL 5.1 backend database and jdk 1.8. The application is web application used tool for design code in Eclipse and execute on Tomcat server. The real time tweet posts collection for dataset of this application using Twitter API with the help of Twitter4j-core and Twitter4j-stream jars. Some functions used in the algorithm are provided by list of jars like stanfordcore-nlp jar for POS tagging etc.

Proposed work is expected to implement posts recommendation system which collects input dataset of list of posts from Twitter API. Apply all the clustering methods like, sentiment clustering, theme clustering, emotional clustering, stylistic clustering and psycholinguistic clustering posts to provide clutter free group environment.

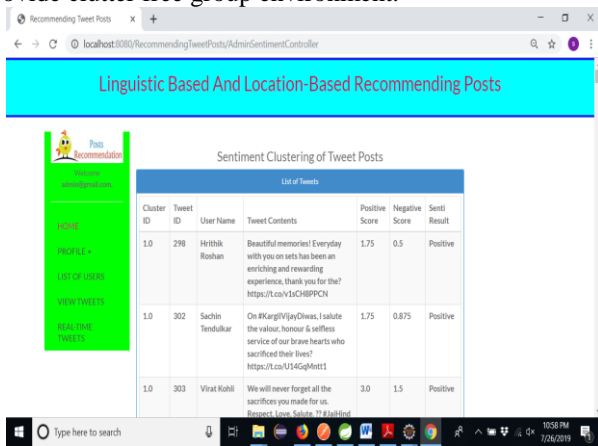


Fig. 2 Sentiment clustering of tweet posts

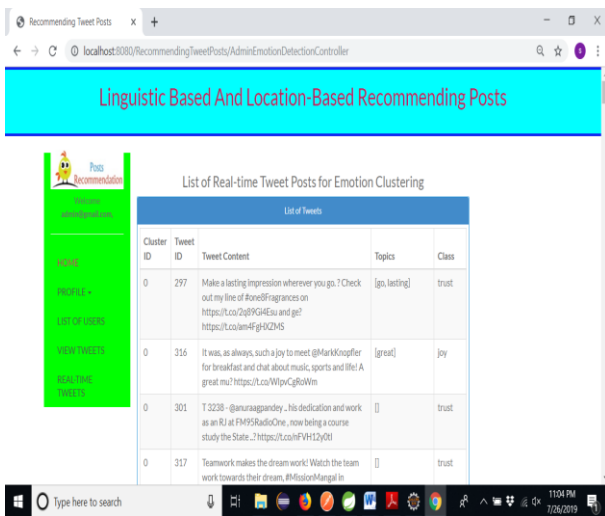


Fig. 3 Emotion clustering of tweet posts

7	518	Was deeply saddened to hear about Sheila ji's demise. A remarkable woman, former CM of New Delhi and a keen admirer? https://t.co/YuZuBq9As9	[Person: Sheila, Location: New Delhi]
7	520	Deeply saddened by the demise of Sheila Dikshit Ji. Blessed with a warm and affable personality, she made a notewor? https://t.co/PBGVwo0FEp	[Person: Sheila Dikshit Ji]
21	569	RT @aspalod: Another proud moment for India ??? athlete #HimaDas got 4th gold for India ??? congratulations @SrBachchan ji @EF_MahekShuk?]	[Location: India, Location: India]
21	508	T 3233 - Hima Das .. the pride of India .. to the Moon and beyond .. indeed but we need to add another Moon for she? https://t.co/gBK80Px5e2	[Location: India, Location: Moon, Location: Moon]
22	478	Congratulations @BorisJohnson on assuming office as Prime Minister of the United Kingdom. I wish you success and? https://t.co/eBTSJR5Vj	[Location: United Kingdom]

Fig. 4 Theme clustering of tweet posts

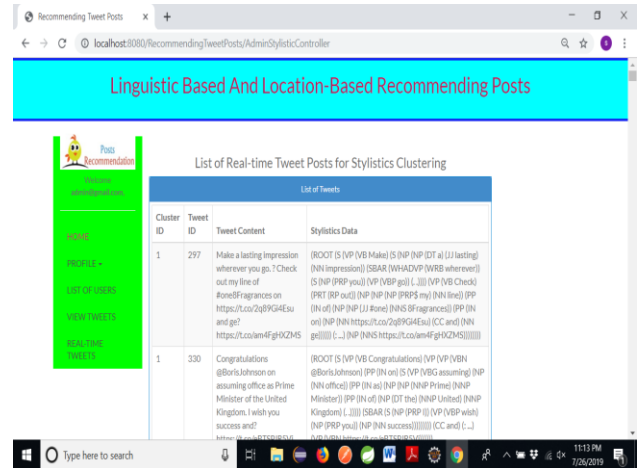


Fig. 5 Stylistics clustering of tweet posts

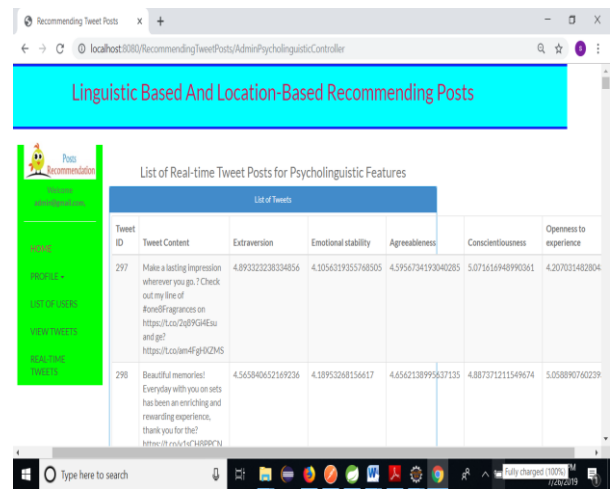


Fig. 6 Psycholinguistic clustering of tweet posts

Expected outcome of this project is providing clutter free group environment to Twitter users. With the help of various aspects based clustering and the location-based clustering [2] the posts with the help of personalized notifications. Finally, formulating groups of individuals within a group sharing similar interests. Fig. 7 represents number of clusters generated for each aspect. The Fig. 8 shows performance on combination of features for number of clusters effectively gives notifications to the users.



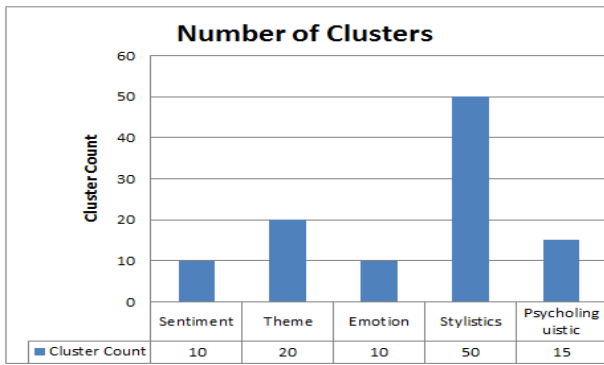


Fig. 7 Performance of Twitter posts of number of clusters for each aspect

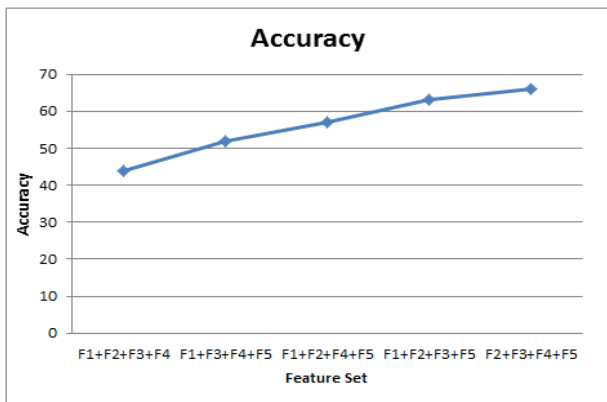


Fig. 8. Accuracy associated with clustering users by combination of four features

VI. CONCLUSION

The paper discusses a novel solution for the least addressed problem of clutter created out of group messages in online social networking sites. The method follows a different approach by considering linguistic features of data that are readily available from a social networking group. It does not depend on any of the pre-existing social relations between users unlike customary methods. The method shows a considerable degree of accuracy in predicting the response of a member to a post. The methodology proves to be an efficient means for managing group policies by providing a trusty environment. It offers a means to provide notification to group admin regarding activities of members within a community whose posting patterns do not suit the group principle.

REFERENCES

1. D. P. Karidi, Y. Stavarakas, Y. Vassiliou, "A Personalized Tweet Recommendation Approach Based on Concept Graphs", In Ubiquitous Intelligence Computing, Advanced and Trusted Computing, Scalable Computing and Communications, Cloud and Big Data Computing, Internet of People, and Smart World Congress (UIC/ATC/ScalCom/CBDCOM/IoP/SmartWorld), 2016, pp. 253–260
2. R. Makki, A. J. Soto, S. Brooks, E. E. Miliotis, "Twitter Message Recommendation Based on User Interest Profiles", In Advances in Social Networks Analysis and Mining (ASONAM), IEEE/ACM International Conference, 2016, pp. 406–410
3. P. K. Gopalan, L. Charlin, D. Blei, "Content-based recommendations with Poisson factorization", In Advances in Neural Information Processing Systems, 2014, pp. 3176–3184
4. L. Hong, A. S. Doumith, B. D. Davison, "Co-factorization machines: modeling user interests and predicting individual decisions in twitter"

5. In Proceedings of the sixth ACM international conference on Web search and data mining, 2013, pp. 557–566
6. K. Chen, T. Chen, G. Zheng, O. Jin, E. Yao, Y. Yu, "Collaborative Personalized Tweet Recommendation", In Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval, ACM, 2012, pp. 661–670.
7. T. Lappas, M. Crovella, and E. Terzi, "Selecting a characteristic set of reviews," in Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining, 2012, pp. 832–840.
8. P. Tsaparas, A. Ntoulas, and E. Terzi, "Selecting a comprehensive set of reviews," in Proc. 17th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining, 2011, pp. 168–176.
9. P. Sinha, S. Mehrotra, and R. Jain, "Summarization of personal photologs using multidimensional content and context," in Proc. 1st ACM Int. Conf. Multimedia Retrieval, 2011, p. 4.
10. Y. Lu, P. Tsaparas, A. Ntoulas, and L. Polanyi, "Exploiting social context for review quality prediction," in Proc. 19th Int. Conf. World Wide Web, 2010, pp. 691–700.
11. H. Lin and J. Bilmes, "Multi-document summarization via budgeted maximization of sub modular functions," in Proc. Human Lang. Technol.: Annu. Conf. North Amer. Chapter Assoc. Comput. Linguistics, 2010, pp. 912–920.

ACKNOWLEDGMENT

It is optional. The preferred spelling of the word "acknowledgment" in American English is without an "e" after the "g." Use the singular heading even if you have many acknowledgments. Avoid expressions such as "One of us (S.B.A.) would like to thank." Instead, write "F. A. Author thanks" *Sponsor and financial support acknowledgments are placed in the unnumbered footnote on the first page.*

AUTHORS PROFILE



Ms Sayali Joag has completed BE degree in Information Technology and pursuing ME in Computer Engineering from Marathwada Mitra Mandal's college of engineering Pune. Her area of interest is Data Mining



Mrs Rupali Dalvi is an Assistance Professor and ME coordinator of Computer Engineering in Marathwada Mitra Mandals college of engineering Pune. She has academic experience of 13 years. Areas of interest include data mining, information security and IoT. She is a recognized post graduate teacher of computer engineering at SPPU.