

Hadoop Map Reduce Framework for Keyword-Aware Recommendation System with Sentiment Analysis

Shireen Khan

M. Tech. Scholar, All Saints' College of Technology, Bhopal, India

Sarwesh Site

Professor, Dept. of CSE, All Saints' College of Technology, Bhopal, India

ABSTRACT: Recommender systems apply knowledge discovery techniques to the problem of making personalized product recommendations using customers usage pattern. Service recommender systems provide an important tool for information gathering, information filtering, information processing, and recommending services to the users with appropriate results. The regularly and rapidly growing number of users, services, and data are presenting new challenges. New recommender systems technologies are needed to quickly produce high quality recommendations, even for very large-scale problems. There exist a lot of recommendation

methods currently but these existing recommender systems needs to be improved, efficient, and trustworthy to meet the present requirements. To address these challenges, we propose a Keyword-Aware Service Recommendation method with Sentiment Analysis to recommend trustworthy services to the user. This research is based on the work done in a article named Keyword-Aware Service Recommendation method (KASR) on MapReduce for Big Data Applications. KASR is based on Preference-Aware Service Recommendation method (PASR) and it uses a user-based collaborative filtering algorithm that searches for specific keywords. KASR is implemented on Hadoop environment to handle the big data generated by recommendation systems to improve the scalability and efficiency. KASR considers user's preferences but lacks the sentiment analysis that may result in an effective number of wrong recommendations. We proposed the keyword aware recommendation system with sentiment analysis algorithm. We address performance issues by implementing the algorithm using Hadoop Map-Reduce framework combined with similarity based collaborative filtering.

KEYWORDS: Sentiment Analysis, Hadoop, KASR, PASR, Machine Learning.

Date of Submission: 12-11-2021

Date of acceptance: 28-11-2021

I. INTRODUCTION

Recommendation systems have turned out to be greatly useful for number of practical applications in recent years, and these applications have wide domain of variety. Films, music, news, research articles, e-marketing, social networking, and consumer based online productions are very common applications where a recommendation system is a need. Generally, recommendation systems utilizes the response of previous users who purchased or used and reviewed products to provide suggestions for current or future users. These systems take the reviews of past customers similar to current customer and usually current preferences of the customer as keywords, then aggregates them on the basis of certain rules to provide recommendations as output.

Due to heavily increasing amount of data, these algorithms need to process intensive large information sets. Now days applications like social networking websites and search engines generate several gigabytes to several terabytes and even up to several petabytes of data [1][2]. Several years ago Google published a paper presenting the MapReduce framework to deal with its ever growing enormous amount of data in distributed environment. MapReduce framework simplifies the complexity of processing large amount of information stored on distributed node clusters in parallel manner, because it facilitates programmers so that they can map a instruction among different nodes in the cluster to execute on different set of data. MapReduce deals with the distributed environment and parallelism omits own, reducing the complexity for programmers, and results in minimalistic required data [3].

Self data replication in MapReduce model offers great fault tolerance also [4]. MapReduce deals with data distribution, parallel processing, and big data in quite efficient way. Many major companies have implemented their own MapReduce models with their own set of requirements. Information filtering algorithms are well known for prediction of preferences or ratings that user would give to a product. Recommender algorithms are subset of information filtering algorithms. Most of these filtering algorithms are static in nature

providing a generic experience to every user, meaning users just search and buy products. But recommendation systems provide greater user interaction to seller, and richer experience to users. To provide effective recommendations, recommender systems take the advantage of old data such as past records of purchases, searches, reviews, and users' behavior. Correct implementation of recommender systems can be extremely effective at increasing user engagement and purchasing. Today, many of the world's most heavily trafficked websites, such as Google, Facebook, LinkedIn, Amazon, and Twitter employ recommendation systems to engage their users with appropriate and personalized content [5][6]. Figure 1 shows the Recommendation System Process



Figure 1. Framework for Keyword-Aware Recommendation System

To solve the recommendation problem the concept of knowledge data discovery is used by taking in users' previous responses in the form of purchased items, user's behavior on websites, user feedbacks, and ratings given by users. Recommendation systems enhance the user experience by providing suggestions to users such as what products to buy, which article to read, and which movie to watch. In research field many different approaches have been made to address the problem of making efficient and accurate personalized recommendations. Some of these approaches are Data mining, machine learning, Collaborative Filtering, User based, Item based [7][8].

In most traditional recommender systems, users are either presented with a preferred recommendations or a user can give some textual or selection inputs (user's preferences) [9] to sort the list out. In some modern recommender systems, users are recommended using previous user's selections and feedback along with the user's preferences [10]. So in the first case they are not considering the user's preferences. Users are not same, so a good system should consider the user's requirement. While in the second case they are not considering the sentiments of the feedback provided by the previous users. They are only taking it as positive feedback and not considering negative feedback, which may provide wrong recommendations as in their method they are filtering previous users' feedbacks with the keyword entered by the current user.

This may not work in case of sentences with negative sense means negative feedback. In this research, keyword aware service recommendation systems are researched broadly and a new approach for keyword-based recommendation system is proposed to provide better recommendations.

Motivated by these observations, in this paper, we address these problems by introducing some improvement in one pre-existing recommendation approach, named as "Keyword-Aware Service Recommendation Method on MapReduce for Big Data Applications".

The Thesis report is organized into 5 sections. Section 2 gives the literature review by explaining the major research and papers about keyword-aware recommendation systems. In Section 3, we present the design and implementation of the proposed algorithm. In Section 4, we provide the experiments along with the results of this algorithmic approach. We discuss the future work that can be done in this research and conclude our work in Section 5.

II. RELATED WORK

The First recommender system was developed by Goldberg, Nichols, Oki Terry [11] in 1992. Tapestry was an electronic messaging system that allowed users to either rate messages ("good" or "bad"). Recommender system as defined by M. Deshpande and G. Karypis: A personalized information filtering technology used to either predict whether a particular user will like a particular item (prediction problem) or to identify a set of N items that will be of interest to a certain User. Recommender systems form or work from a specific type of information filtering system technique that attempts to recommend information items (movies, TV program/show/episode, video on demand, music, books, news, images, web pages, scientific literature etc.) or social elements (e.g. people, events or groups) that are likely to be of interest to the user. Typically, a recommender system compares a user profile to some reference

characteristics, and seeks to predict the 'rating' or 'preference' that a user would give to an item they had not yet considered.

Shunmei Meng in 2014 proposed a KASR [12] method for personalized recommendation. In this user based collaborative filtering is used. For more efficiency the method is implemented on Hadoop. For evaluation Jaccard coefficient and Cosine similarity measure is used. User's positive and negative reviews are not

considered separately. Sentiments in the reviews are not considered. To make the method more efficient and scalable Hadoop MapReduce model is used.

Xi Wang Yang in 2013 proposed a Bayesian inference based recommendation in online social networks [13]. In this users share their content ratings with friends. Rating similarity is measured using conditional probability. Based on similarity score ranking and recommendation is done.

There is a Cold start and rating sparseness problem. In [14], the authors propose a Bayesian inference based recommendation system for online social networks. They show that the proposed Bayesian inference based recommendation is better than the existing trust based recommendations and is comparable to Collaborative Filtering recommendation.

In [15], Adomavicius and Tuzhilin give an overview of the field of recommender systems and describe the current generation of recommendation methods. They also describe various limitations of current service recommendation methods, and discuss possible extensions that can improve recommendation capabilities and make recommender systems applicable to an even broader range of applications. Most existing service recommender systems are only based on a single numerical rating to represent a service's utility as a whole. In fact, evaluating a service through multiple criteria and taking into account of user feedback can help to make more effective recommendations for the users.

Guosheng Kang [16] proposed an active web service recommendation. Web usage history and QoS are the main criteria for recommendation. Using this recommendation top k services are generated for users. Passive users reviews are not considered. Using this approach top k services are generated for users. Passive users reviews about the website is not considered.

Yan Ying Chen [17] proposed a probabilistic personalized travel recommendation model. People attributes and photos are used which are effective for mining demographics for travel landmarks and paths, and thus greatly benefiting personalized travel recommendation.

Faustino Sanchez [18] proposed a recommender system for sport videos, transmitted over the Internet and/or broadcast, in the context of large-scale events, which has been tested for the Olympic Games.

Zhibin Zheng [19] proposed quality of service ranking prediction for cloud services. This paper investigate the combination of rating based approaches and ranking based approaches, so that the users can obtain QoS ranking prediction as well as detailed QoS value prediction. How to detect and exclude malicious QoS values provided by users is not proposed here. With the development of cloud computing software tools such as Apache Hadoop, MapReduce, and Mahout, it becomes possible to design and implement scalable recommender systems in "Big Data" environment. [20] The authors of implement a CF algorithm on Hadoop. They solve the scalability problem by dividing dataset. But their method doesn't have favorable scalability and efficiency if the amount of data grows.

Jin et al. [21] propose a large-scale video recommendation system based on an item-based CF algorithm. They implement their proposed approach in Quist, which is a .Net MapReduce framework, thus their system can work for large scale video sites.

III. PROPOSED METHODOLOGY

This chapter explains the proposed design used in the implementing keyword aware service recommendation system with sentiment analysis on MapReduce. Figure 2 shows the basic framework for key-aware recommendation system with sentiment analysis on MapReduce.

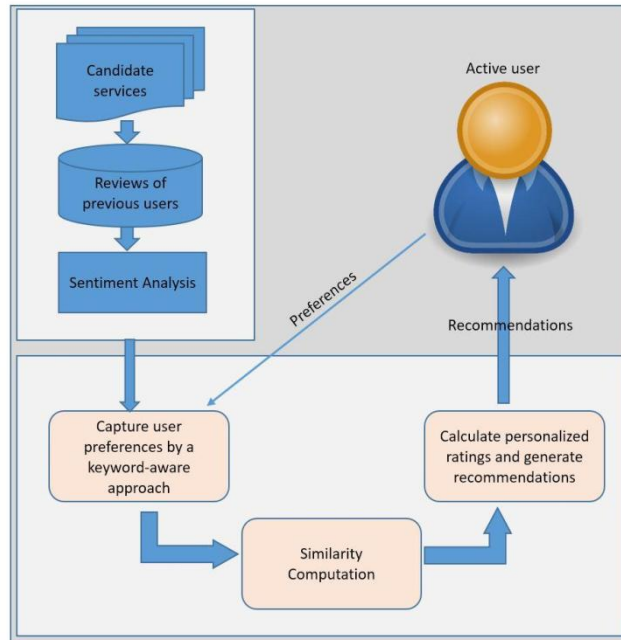


Figure 2: Framework of Keyword-Aware Recommendation System with Sentiment Analysis

Capturing User’s Preferences

There are two types of users: Active user and Previous users. Active user is the user which is currently using the services by providing set of keywords according to his preferences from the Keyword List. Previous users are the users which has already used the services and reviewed for the services. Based on the reviews from the previous user, the keywords are extracted from the reviews according to Keyword List and Domain Thesaurus.

Keyword Extraction

In this phase, each review will be transformed into a corresponding keyword set according to the keyword-candidate list and domain thesaurus. If the review contains a word in the domain thesaurus, then the corresponding keyword should be extracted into the preference keyword set of the user. For example, if a review of a previous user for a hotel has the word “spa”, which is corresponding to the keyword “Fitness” in the domain thesaurus, then the keyword “Fitness” should be contained in the preference keyword set of the previous user. If a keyword appears more than once in a review, the times of repetitions will be recorded. In this paper, it is regarded that keywords appearing multiple times are more important. The times of repetitions will be used to calculate the weight of the keyword in preference keyword set in the next step.

Similarity Computation

The second step is to identify the reviews of previous users who have similar tastes to an active user by finding neighborhoods of the active user based on the similarity of their preferences. Before similarity computation, the reviews unrelated to the active user’s preferences will be filtered out by the intersection concept in set theory. If the intersection of the preference keyword sets of the active user and a previous user is an empty set, then the preference keyword set of the previous user will be filtered out.

Approximate Similarity Computation

A frequently used method for comparing the similarity and diversity of sample sets, Jaccard coefficient, is applied in the approximate similarity computation. Jaccard coefficient is measurement of asymmetric information on binary (and non-binary) variables, and it is useful when negative values give no information. The similarity between the preferences of the active user and a previous user based on Jaccard coefficient is described as follows: The formula used by Jaccard Coefficient is

$$simASC(APK, PPKj) = \frac{|APK \cap PPKj|}{|APK \cup PPKj|}$$

Where APK is the preference keyword set of the active user, PPK is the preference keyword set of a previous user. And the weight of the keywords is not considered in this approach.

Exact Similarity Computation

In Exact Similarity Computation, the similarity between preference of active user and previous users is calculated based on Cosine-Based approach. The formula used by Cosine-Based approach is

$$simESC(APK, PPK) = \cos(WAP, WPP) = \frac{WAP \cdot WPP}{\|WAP\|_2 \times \|WPP\|_2}$$

Preference weight vector: In this cosine based approach, the preference keyword sets of active user and previous users is transformed into n-dimensional weight vectors respectively, namely preference weight vector. The preference weight vectors of the active user and a previous user are noted as WAP and WPP, respectively. In our research, we use the Analytic Hierarchy Process (AHP) model to decide the weight of the keywords in the preference keyword set of the active user. AHP method is provided by Saaty in 1970s to choose the best satisfied business role for its hierarchy nature. The weight computing based on the AHP model is decided as follows:

Firstly, we construct the pair-wise comparison matrix in terms of the relative importance between each two keywords. The pair-wise comparison matrix $A_m = (a_{ij})_m$ must satisfy the following properties, a_{ij} represents the relative importance of two keywords:

$$a_{ij} = 1, i = j = 1, 2, 3, \dots, m.$$

$$a_{ij} = 1/a_{ji}, i, j = 1, 2, 3, \dots, m \text{ and } i \neq j.$$

$$a_{ij} = a_{ik} / a_{jk}, i, j, k = 1, 2, 3, \dots, m \text{ and } i \neq j.$$

After checking the consistence of the matrix, then we calculate the weight by the following function:

$$w_i = 1/m \sum_{j=1}^m \frac{a_{ij}}{\sum_{k=1}^m a_{kj}}$$

Here, a_{ij} is the relative importance between two keywords; m is the number of keywords in the preference keyword set of active users.

IV .RESULT ANALYSIS

The algorithm is designed to derive recommendations to the users based on their preference keywords and reviews of previous users. Experiments are done on both single node and multi-node clusters with varying data sizes and performance metrics are recorded. All the map-reduce jobs are run in sequence as the output of the previous jobs are given as input for the next job

Provided experiment results are based on the comparison between execution times taken by KASR_ESC algorithm and proposed approach KASRwithSentiAnalysis.

Figure 3 shows the comparative analysis of both the algorithms on 250MB data set in single node cluster, 2-node cluster, 4-node cluster, and 6-node cluster. From the graph it is clear that KASRwithSentiAnalysis is performing better in single node cluster and 2-node cluster. But in 4-node cluster and 6-node cluster very slight improvement in proposed algorithm.

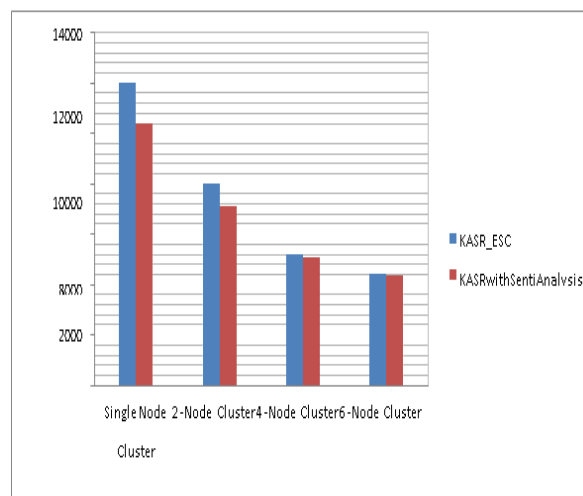


Figure 3:- Comparative Analysis of KASR_ESC and KASRwithSentiAnalysis on 250MB data set

Figure 4 shows the comparative analysis on 1GB data set. We can see that as the data size increases the performance of KASRwithSentiAnalysis increases or atleast same as of KASR in single node and 2-node cluster. is performing better in single node cluster and 2-node cluster

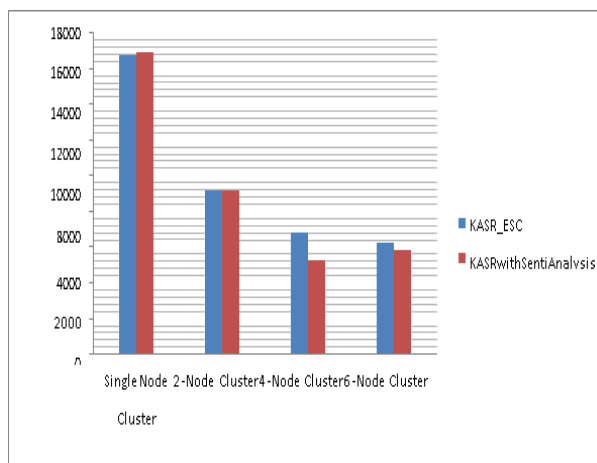


Figure 4 Comparative Analysis of KASR_ESC and KASRwithSentiAnalysis on 1GB data set

Next we experimented both the algorithms on 10GB data set. As the data size increases more, the performance of proposed algorithm degrades. This is due to the overhead of the computation for sentiment analysis, but also gives better result than KASR.

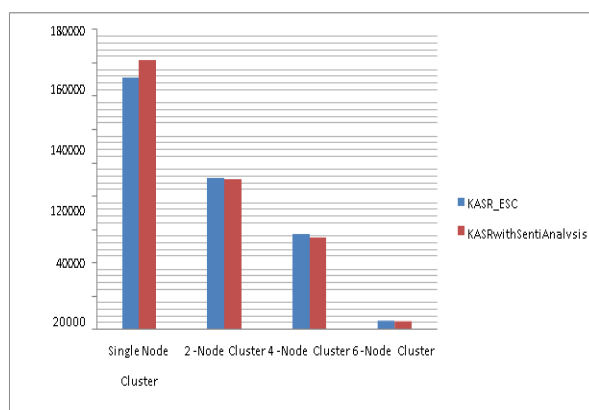


Figure 5. Comparative Analysis of KASR_ESC and KASRwithSentiAnalysis on 10GB data set

V. CONCLUSION

In this paper we discussed about recommendation system algorithms and MapReduce programming model by Hadoop. A sequential algorithm to compute similarity is explained. Then we proposed keyword-aware recommendation system with sentiment analysis to provide better recommendations to users. Our approach provides recommendations by taking the user preferences and utilizing the reviews given by previous users. The major improvement was to analyze the sentiments of users by processing their reviews. Quality wise we achieved better results than the original algorithm.

Different experiments are conducted to show the performance improvement of using a multi-node cluster. Test data for experiments is taken from data bases of a hotel review system.

The Hadoop map-reduce approach saves a lot of resources in computing similarities and generates recommendations in short time.

For future work, lot of improvements can still be done. Sentiment analysis algorithm can be more accurate and optimized. Another better approaches can be applied to identify the sentimental behavior more accurately.

REFERENCES

- [1]. Cheikh Kacfeh Emani, Nadine Cullot and Christophe Nicolle "Understandable Big Data: A survey" in Computer Science Review Volume 17, August 2015, Pages 70–81.
- [2]. H. Hu, Y. Wen, T.-S. Chua, and X. Li, "Towards scalable systems for big data analytics: A technology tutorial," IEEE Access, vol. 2, pp. 652–687, 2014.
- [3]. Rodrigo Agerri, Xabier Artola, Zuhaitz Beloki, German Rigau, Aitor Soroa "Big data for Natural Language Processing: A streaming approach" in Knowledge-Based Systems Volume 79, May 2015, Pages 36–42.
- [4]. A. Rabkin and R. H. Katz, "How Hadoop Clusters Break," IEEE Software, vol. 30, pp. 88–94, 2013.
- [5]. T. Jiang, Q. Zhang, R. Hou, L. Chai, S. A. McKee, Z. Jia, and N. Sun, "Understanding the behavior of in-memory computing workloads," in Workload Characterization (IISWC), IEEE International Symposium on, 2014, pp. 22–30.

- [6]. Martin Zwilling. What Can Big Data Ever Tell Us About Human Behavior [online]. Available: <http://www.forbes.com/sites/martinzwilling/2015/03/24/what-can-big-data-ever-tell-us-about-human-behavior/#1580346f1bed>. Date accessed: (March 24, 2015).
- [7]. Saeed Shahrivari, "Beyond Batch Processing: Towards Real-Time and Streaming Big Data", *Computers*, Vol. 3, pp. 117-129, 2014.
- [8]. Sachin Agarwal. Monitoring and Troubleshooting Apache Storm [online]. Available: <https://dzone.com/articles/monitoring-and-troubleshooting-apache-storm-with-o>. Date accessed: (April 7, 2016).
- [9]. M. R. Evans, D. Oliver, K. Yang, X. Zhou, S. Shekhar, "Enabling Spatial Big Data via CyberGIS: Challenges and Opportunities," Ed. S. Wang, M. F. Goodchild, *CyberGIS: Fostering a New Wave of Geospatial Innovation and Discovery*. Springer, 2014.
- [10]. Wikipedia. "Spatialdatabase"[online]. Available https://en.wikipedia.org/wiki/Spatial_database.
- [11]. Shunmei Meng, Wanchun Dou, Xuyun Zhang, Jinjun Chen, "KASR: A Keyword-Aware Service Recommendation Method on MapReduce for Big Data Applications" *IEEE Transactions On Parallel And Distributed Systems*, TPDS-2013-12-11.
- [12]. X. Yang, Y. Guo, Y. Liu, "Bayesian-inference based recommendation in online social networks," *IEEE Transactions on Parallel and Distributed Systems*, Vol. 24, No. 4, pp. 642-651
- [13]. A. Tuzhilin and G. Adomavicius, "Toward the Next Generation of Recommender Systems: A Survey of the State of the Art and Possible Extensions," *IEEE Transactions on Knowledge and Data Engineering*, vol. 17, no. 6, pp. 734-749, 2005.
- [14]. G.Kang, J. Liu, M. Tang, X. Liu and B. cao, "AWSR: Active Web Service Recommendation Based on Usage History," 2012 IEEE 19th International Conference on Web Services (ICWS), pp. 186-193, 2012.
- [15]. Yan-Ying Chen, An-Jung Cheng, "Travel Recommendation by Mining People Attributes and Travel Group Types From Community-Contributed Photos" *IEEE Transactions on Multimedia*, Vol. 15, No. 6, October 2013.
- [16]. M. Alduan, F. Alvarez, J. Menendez, and O. Baez, "Recommender System for Sport Videos Based on User Audiovisual Consumption," *IEEE Transactions on Multimedia*, Vol. 14, No.6, pp. 1546-1557, 2013.
- [17]. Zibin Zheng, Xinmiao Wu, YileiZhang, Michael R. Lyu, Fellow, and Jianmin Wang, "QoS Ranking Prediction for Cloud Services" *IEEE Transactions On Parallel And Distributed Systems*, Vol. 24, No. 6, June 2013.
- [18]. D Zhao and M. S. Shang, "User Based Collaborative Filtering Recommendation Algorithms on Hadoop," In the third International Workshop on Knowledge Discovery and Data Mining, pp. 478-481, 2010.
- [19]. Konstan, J., Miller, B., Maltz, D., Herlocker, J., Gordon, L., and Riedl, J. (2012). *GroupLens: Applying Collaborative Filtering to Usenet News*. *Communications of the ACM*, 40(3), pp. 77-87.. [Online; accessed June 2014].
- [20]. Adomavicius G., Tuzhilin A. (2005). *Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions..* [Online; accessed December 2014].
- [21]. A. Rabkin and R. H. Katz, "How Hadoop Clusters Break," *IEEE Software*, vol. 30, pp. 88- 94, 2013.