

A Novel Information retrieval system for distributed cloud using Hybrid Deep Fuzzy Hashing Algorithm

Dr. V. Suma,
Professor,
Department of Information Science & Engineering,
Dayananda Sagar College of Engineering,
Bangalore, India.
E-mail id: suma-ise@dayanandasagar.edu

Abstract: The recent technology development fascinates the people towards information and its services. Managing the personal and public data is a perennial research topic among researchers. In particular retrieval of information gains more attention as it is important similar to data storing. Clustering based, similarity based, graph based information retrieval systems are evolved to reduce the issues in conventional information retrieval systems. Learning based information retrieval is the present trend and in particular deep neural network is widely adopted due to its retrieval performance. However, the similarity between the information has uncertainties due to its measuring procedures. Considering these issues also to improve the retrieval performance, a hybrid deep fuzzy hashing algorithm is introduced in this research work. Hashing efficiently retrieves the information based on mapping the similar information as correlated binary codes and this underlying information is trained using deep neural network and fuzzy logic to retrieve the necessary information from distributed cloud. Experimental results prove that the proposed model attains better retrieval accuracy and accuracy compared to conventional models such as support vector machine and deep neural network.

Keywords: - Information Retrieval, Hashing, Deep Neural network, Fuzzy Logic, Cloud Computing

1. Introduction

Information retrieval is an essential process in cloud computing in order to store and retrieve the necessary information from the environment to cloud and vice versa. The technology development and resource availability drags the information management system into a drastic shift with in few years. Moreover, technology started to expunge the trace of conventional information management process through its innovative web based information management systems. Similarly, internet based services, network mediums, electronic libraries and recent advanced search engines makes the information management systems always in demand. For those systems, information management is not only to store the data, also it requires to manage the unstructured and structured data in a large scale manner. Information retrieval system is a core support to internet based services and search engines. This makes a demand to the researchers to develop and fine-tuned information retrieval system as a sophisticated application.

The importance of information retrieval presents in its extracting nature of most suitable information for a query from a database. But the issue is, on what basis the best relevant information could be retrieved for the query. Since the user gives a common representation as a query and the system must analyse and need to produce information which ensures the retrieved items are most relevant to the user query. Extracting information based on keywords will improve the precision and recall, but due to technology development user can able to include query in natural language format and the system must search based on that. For this purpose, from a large set of documents, based on the query it is subdivided into small sets to retrieve the relevant information. Generally, information retrieval system uses techniques to predict the documents and once it is retrieved, those documents are ordered and ranked in a decreasing order. Several retrieval models are evolved based on structures, similarity and weightage measures.

The ultimate goal of data mining is to collect the information based on the extracted patterns to obtain essential knowledge over the vast collected data. Since, data mining application is not limited into science and engineering, it provides wide range of supports in the field of games, medical, business analysis, etc., However, conventional mining applications are hardware based and it is considered as hurdles to the organizations to adopt those mining applications. Also, the cost of data management for storage and retrieval is huge due to physical attributes, so that the organizations particularly small scale organizations doesn't show much interest to move on into data mining applications. When web based information management systems are evolved, these small scale organizations show more interest due to its cost effective features. Cloud computing is a best example for information management and now it is considered as an ultimate platform for data mining. Since the service provider takes responsibility of technologies and its expenditures, it is widely adopted in large scale to small scale organizations. Cloud offers service based on the needs and the organization doesn't spend money over hardware and environment setup, it reduces the major expenditure for organizations. The adoptability rate of cloud computing is high so that the user could increase or decrease the storage based on their necessity which further reduces the organizational expenditures. Simple process of information retrieval system is depicted in figure 1.

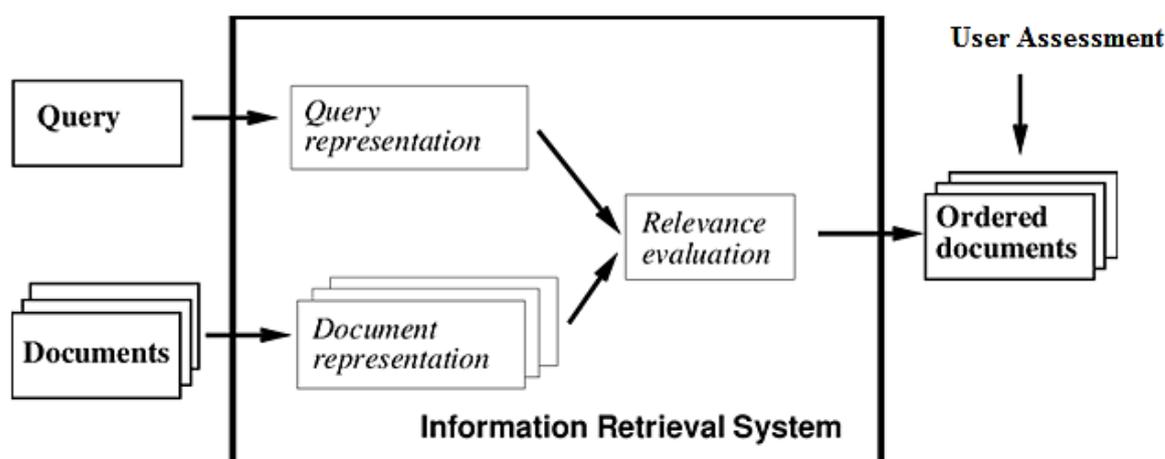


Figure 1 Information Retrieval System

Cloud computing has impact over every field and almost 75% of business environment are moved into cloud due to the ever increasing demand and better service. In the beginning cloud uses mainframe to provide multiple user support, later implementation of virtual environment supports multiple platform which attracts more organizations to get more benefits to utilize the infrastructure. basic elements in cloud infrastructure is a server, database and device. In which, the servers establish links between the host and user so that the system will intact and keep the data flow for the service request. Data manipulation is performed in the database as storage and manipulation process. Technology development in cloud has invented different types of cloud such as private cloud for organization and single person, public cloud for general use with almost cost free. Multiple organizations share the cloud resources based on their needs as community cloud and recently hybrid cloud are evolved which offers service with the features of public and private cloud altogether as a single cloud.

IT industries uses cloud as an essential strategical service to obtain competitive advantage over the customers. In order to improve organizational performance, cloud computing is used a core technology now a day. The cloud service delivery model in cloud computing is used to analyse the cloud adoption and service evaluation which helps the organizations. Not only the IT industry, recently government and community organizations, business, educational institutions and individuals adopted into cloud environment. This helps them to concentrate over the development of individual to organization and not on the technologies. Since the cloud has high performance computing service as combined infrastructure which has clusters and grid which supports the user in terms of discovery and exploration. The multitenancy option in cloud offers shared resource service among large number of users. So that, it reduces the resource cost for small scale organizations and increases load capacity and resource utilization effectively. Improved scalability, elasticity and reliability are the other advantages of cloud. However,

cloud has few limitations such as communication failure, data loss, data traffic surveillance and other potential risks in terms of privacy and security.

The growth of information service such as storage, computing and retrieving the data through cloud computing is a prevalent method in this era. Even the user considered sensitive information is also transferred to cloud. However, there is some uncertainty arises in data privacy as the cloud service provider has full rights to manage the user data. This lack of confidentiality between user and service provider is overcome by encryption process. Prior to outsourcing those confidential files are encrypted and attached to the cloud, so that data owners do not need to worry about data privacy. In some cases, it is essential to share the data with multiple users of different domain, the user must perform information retrieval to obtain the necessary information from the encrypted data. From this, it could be visible, encryption is essential to cloud environment to maintain security and privacy among confidential data. Rest of the research work is organized as a in depth research analysis of existing information retrieval models in section 2, followed by proposed hybrid information retrieval system in section 3, experimental analysis and its discussion is presented in section 4. Conclusion is given in the last section with limitation of proposed work along with future scope.

2. Related works

The research towards data mining and its applications is a still in progress and researchers working on it to obtain better model. This section provides summary of such existing research works in information retrieval to obtain the issues while implementing retrieval system. Data mining and information retrieval system is analysed by Jiaying Liu.*et.al.* [1] research work. The entire development of information retrieval system for this 21st century is analysed in the survey works which provides details about applications of information retrieval systems in detail. Research work describes the merits of various algorithms based on text, graph and map based retrieval models.

Yongjun.*et.al.* [2] reported the issues in natural language interface to bibliographic information retrieval system. Since the information retrieval using natural language is difficult to process as the database management system faces difficulty in organizing natural language data. To reduce the issues an interface is proposed in the research work which helps the user to search bibliographic data using natural language. Graph based information retrieval system Sidali Hocine Farhi.*et.al.* [3] is familiar and it is widely adopted in many applications, based on those graph based system, the proposed bibliographic information system is developed which process the queries as text and retrieves information through the interface. Similarly, Joby.*et.al.* [4] reported the issues in information retrieval from large data set through natural language model. Considering the limitations in probabilistic, space vector and other conventional retrieval models, the proposed research work emphasis the natural language based retrieval system which extracts relevant information from large dataset.

Ranking based information retrieval model is reported in andrei.*et.al.* [5] research work. Based on document description and term frequency model ranks are allocated considering the user request. The differences in the documents and relativity to the user request are considered to assess the quality of the model. Proposed research used modified genetic algorithm and provides relevant information with minimum stagnation. The structural complexity in conventional information retrieval is reduced in the proposed approach using ranking models and genetic based map criterion. However, the proposed system fails in processing natural language requests which is considered as the limitation of the model.

Deepanwita.*et.al.* [6] reported the issues in multimodal retrieval system while retrieving document images. In general text caption is used to describe an image in a document, but the process is bit complex compared to other retrieval process. since the system needs to analyse the text and relevant images in the database which consumes more time and generates false results which affects the efficiency of the retrieval system. In order to reduce the complexity, key phrase extraction techniques are used in the proposed model which yields better retrieval efficiency. Similar multi criteria model is reported in Stefania Marrara.*et.al.* [7] research work to reduce the issues in decision making in information retrieval. Since it is essential for a system to decide whether the retrieved information is relevant to the user query. In order to define the decision making process various dimensions like novelty, topic relevance and user needs are considered in the proposed information retrieval

model. The limitation of these model is present in its dimensionality based decision. System doesn't able to retrieve relevant information if the given query didn't fall on the predefined dimensions.

Hamid Khalifi.*et.al.* [8] reported the issues in information retrieval systems while using large database. In case of text based retrieval in a huge database, the similarity probability will be high and the retrieval process becomes complex. In order to restructure the user query without deviating the request, the semantic relationships are identified using support vector machine in the proposed work to obtain the necessary result. Machine learning based models mostly performs well and provides better classification results which helps the retrieval system to extract the information from the large database. The issues in conventional text based information retrieval system is reported in youssef Chouni.*et.al.* [9] research model. Using graphs, the words in a document is represented which helps to measure the similarity. Further the similarity measure is enhanced by synonymy and semantic index to retrieve the necessary information for the user query which provides better performance.

Private information retrieval model is reported in Jianchang Lai.*et.al.* [10] research work which allows the user to retrieve the information from the database based on user preference. In case of conventional private information model, each data is need to be published with description which leads into information leakage. To reduce such information leakage, attribute based information retrieval model is proposed in the research work. The proposed work attains better data privacy which doesn't reveal any information about the data. Complexity of this model is its data description process as it is difficult to describe the data which is present in large data set. Similar model is reported in Razane Tajeddine.*et.al.* [11] and Heecheol Yang.*et.al.* [12] research work which used distributed database for information retrieval. Retrieving information from distributed database provides better data security and enhanced retrieval performance. Compared to conventional models the performance of distributed database dependent retrieval system performs better in terms of data segregation, data classification and retrieval efficiency.

Youcef Djenouri.*et.al.* [13] proposed a cluster based information retrieval model. The proposed approach identifies the frequency of the information based on the user query and provides the most frequent items as results. The proposed model has advantages, as the user gets most predominant terms as results which is widely adopted by others. But the limitation of this model is its irrelevant information. In case the user need to search unfamiliar items, then the system produces random results which includes the user query as a part of the document which affects user preferences. Structural equation modelling based information retrieval model is presented in Massimo Melucci.*et.al.* [14] research work. Proposed model classifies the experimental which is collected using testing system across the datasets. It is essential to obtain the relationship between latent variables and other variables, so that the essential system affecting parameters could be identified and removed. This helps to evaluate and improve the information retrieval performance. Marco Angelini.*et.al.* [15] proposed an information retrieval system based on combinatorial visual analytics. The proposed research work explores and increases the performance of retrieval system through case based test collections. Using combinatorial composition and consolidated deep statistical analysis, the proposed approach attains better retrieval performance than conventional models.

Data is classified into structured and unstructured types based on the characteristics and source. Web based data mining and its issues are reported in Saravana Kumar.*et.al.* [16] research work. since internet has enormous flow of data and by using ontological and semantic structures, the issues in web mining is addressed in the research work. combining the feature extraction and selection process data mapping is performed and the necessary information is retrieved from the web. The advantages of proposed approach are its reduced dimensionality and complexity in information retrieval process. Scarcity theory based information retrieval model is presented in Ruixiang Ou.*et.al.*[17] research work which uses matching process to interconnect the user cognition and system cognition. In the present situation information retrieval as cognitive view is considered as important due to its strong theoretical foundation. Scarcity theory helps to define the user cognitive nature and based on that information retrieval system has constructed in the proposed work. however, the proposed approach is not convenient for different types of applications.

Jennifer.*et.al.* [18] reported the issues in multi cloud environment information processing. Multi cloud provides better consistency and reduces the severity. But information maintenance in multi cloud environment is a complex process due to its interface, service renders and technologies. Information maintenance such as store, secure and retrieve process needs highly reliable and flexible system. Proposed model overcome the limitation in

traditional cloud based information processing system through cryptography integrated computational intelligence which is widely adoptable for multi cloud models. From the above survey it is observed that most of the research works are performed based on the text or graph based retrieval. Security is a major concern in multi model data processing system. Though the performance of multimodal information system is better, it is essential to improve the features of the system. Cloud based data mining is widely used in various applications and this research work considers the above issues and proposed a distributed cloud based information retrieval system. machine learning is implemented in few research models but the performance could be improved if it based on deep neural networks. In order to define the security and other privacy features of data, a hybrid data processing model is required. Based on this research gaps, the proposed hybrid model is developed in the next section as hybrid deep fuzzy hashing algorithm.

3. Proposed work

A hybrid deep fuzzy hashing algorithm is presented in this section. Mathematical formulation for the hashing approach, deep fuzzy for information retrieval are systematically obtained to improve the retrieval efficiency and security of the data in distributed cloud environment. The intuition for implementing fuzzy with deep neural network is due to fuzzy transform knowledge which is expressed as fuzzy rules. This combination enhances the prediction accuracy through fuzzy logic and retrieval accuracy through learning ability of neural network as a simultaneous process. Hashing function is used as initial process which helps to establish relationship between the query and database which reduces the memory consumption. Figure 2 depicts the hybrid deep fuzzy hashing algorithm model.

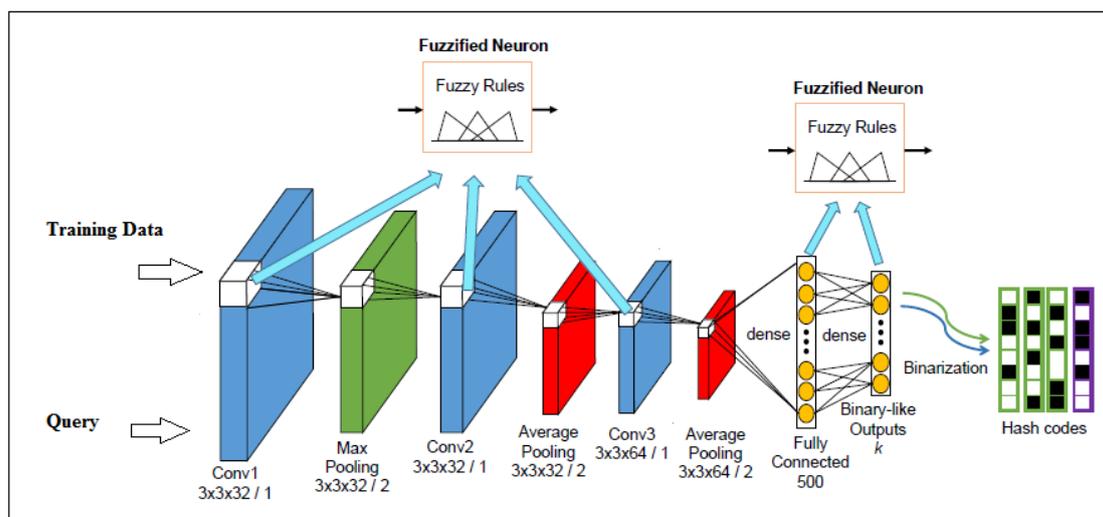


Fig.2 Hybrid Deep Fuzzy Hashing Model

Deep supervised hashing function incorporates convolution neural network to reduce the insufficiency issues and avoid linear predictions. In the proposed work, the learning capacity of the system must satisfy the requirement as (i) codes for the given query must be of same class and it must similar to hamming space (ii) the training process must have high efficiency. Let us consider the given query has ρ initial details and the necessary hash codes are generated and mapped for the query in the range of n-bit binary codes as

$$\delta: \rho \rightarrow \{+1, -1\}^n \quad (1)$$

From the codes, the information of same class is mapped as much as possible. Based on the similarity the mapping process is performed and the huge queries are classified and segregated as large distance codes. In order to reduce the distance codes, a loss function is introduced in the system which defines the differences.

Let p_1, p_2, p_3 be the pairs for the given query ρ and the respective hash codes are $h_1, h_2, h_3 \in \{+1, -1, +1\}^n$ and the loss function is defined as

$$l_f(h_1, h_2, h_3, t) = \frac{1}{2}(1 - t)d_h(h_1, h_2, h_3) + \frac{1}{2}t_{max}(n - d_h(h_1, h_2, h_3), 0) \quad (2)$$

where d_h is the hamming distance and t is the class which is 0 for same class 1 for others. The queries of same class are mapped with similar codes, and others are identified as different binary codes. To obtain better similarity, the other classes are mapped together based on a marginal threshold function. Loss function is defined based on the threshold function, if the hamming distance is near to margin, then loss function will be high. This inverse proportional condition allows only dissimilar pairs so that collision in the loss function is avoided. The overall loss function for the trained data sets is given as

$$L_f = \sum_{n=1}^N l(h_{n1}, h_{n2}, h_{n3}, t_n) \quad (3)$$

Combining the fuzzy and neural network is performed based on fuzzy logic splicing and parallelize the slices into neural network, so that each neuron in the network attempts to operate the fuzzy rules and it is given as

$$g^l(w \oplus x) = \frac{1}{L_f} \sum_{n=1}^l w_l + \frac{1}{L} \sum_{n=1}^l x_l - \frac{2}{L} wx \quad (4)$$

where L is the neuron weight length, w is the neuron weight and x is the output function. The neuron operations are interrupted based on the similarity measurements of output and neuron weight functions. The activation function like ReLU is used to set the threshold based on the hamming distance which is measured earlier. The hamming distance is generalized for positive and negative values so that ReLU will interrupt either the positive side information or negative side information and block the information in the other half which reduces the presence of irrelevant information in the retrieval process. The process flow of proposed hybrid model is depicted in figure 3.

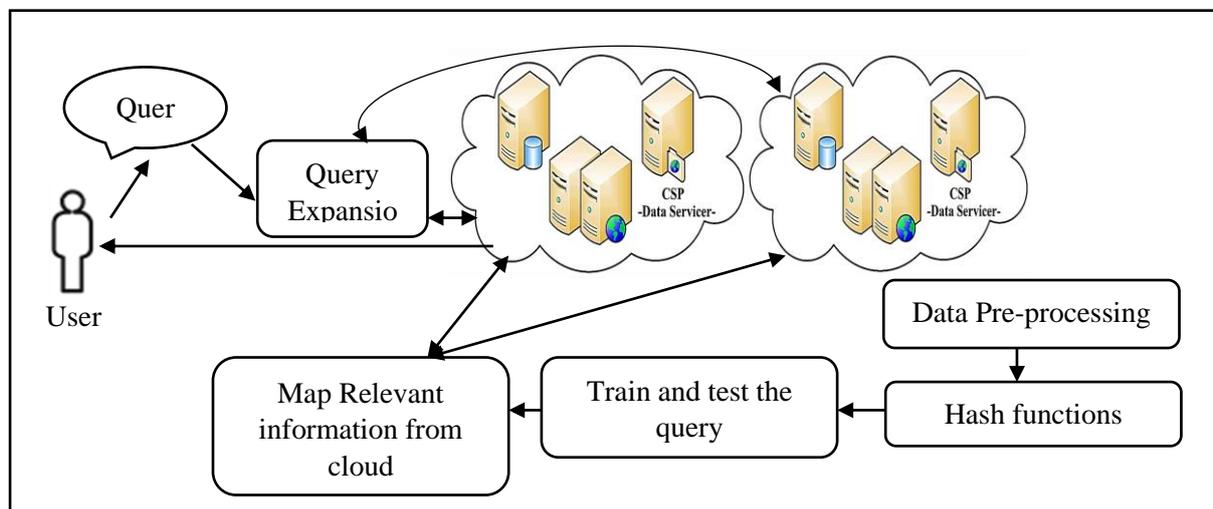


Figure 3 Proposed Hybrid Information Retrieval System

The process starts from user query, once the query is received it is expanded based on the query type and then it is fed into cloud service provider. Proposed model is designed with distributed cloud environment which has predefined data. The data in pre-processed initially to remove the unnecessary information, then using deep hash function, a hash table is created for the data in the database for convenient retrieval operation. Using deep fuzzy neural network, the data are trained and validated so that when the query arrives the relevant information are mapped and produced as result to the user. Summary of learning procedure is given as follows

Input: Query features, Class, code length, learning rate and coefficient function

Output: Model parameter and binary codes

Initialize $h_1 h_2 h_3$ as $\{+1, -1, +1\}^n$

for $n \leq l_{epochs}$ or loss function in eqn. (2) converge do

sample map $\{(h_1, t_1)(h_2, t_2)(h_3, t_3) \dots (h_N, t_N)$

compute hash codes

compute loss $L_f = \sum_{n=1}^N l(h_{n1}, h_{n2}, h_{n3}, t_n)$

update the parameter

$$g^l(w \oplus x) = \frac{1}{L_f} \sum_{n=1}^l w_l + \frac{1}{L} \sum_{n=1}^l x_l - \frac{2}{L} wx$$

end for

4. Result and Discussion

The proposed hybrid deep fuzzy hashing algorithm for distributed cloud environment is experimentally verified in CloudSim and the results are observed. KDD Cup 2004 Database is used in the experimentation which has 50,000 training examples, 1,00,000 test examples and 78 numerical attributes. The classification and retrieval efficiency of the proposed model is measured in terms of specificity, sensitivity. To validate the efficiency of proposed model conventional support vector based information retrieval system and deep neural network models are compared. Table 1 gives the performance measures of proposed model.

Table 1 Performance Comparison of Proposed Model

Parameter	Support Vector Machine Based Information Retrieval System (%)	Deep Neural Network Based Information Retrieval System (%)	Proposed Hybrid Deep Fuzzy Hashing Algorithm (%)
Specificity	88.42	94.26	96.41
Sensitivity	87.38	95.44	97.58
f-measure	89.24	94.61	96.52

Figure 4 depicts the performance comparison of proposed work and conventional models as a collective illustration. It is observed that proposed model attains better sensitivity, specificity and f-measure compared to other models due to the hashing and learning functions. Compared to support vector machine, deep neural network model performs better which shows 2% lesser in an average compared to proposed model.

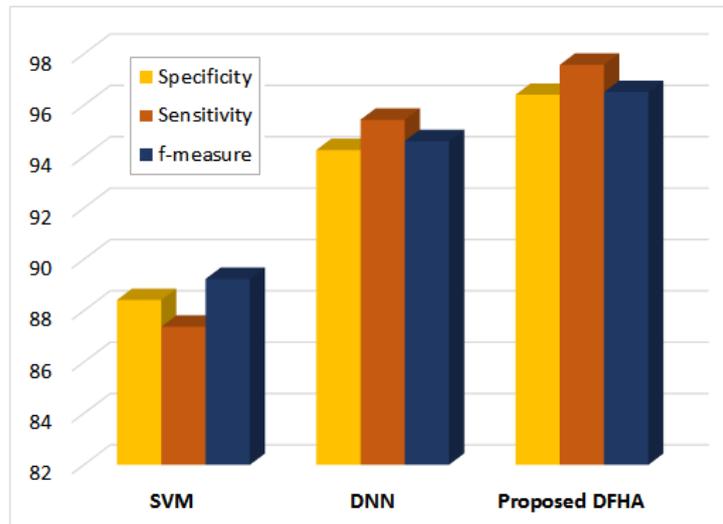


Fig.4 Performance comparison of proposed model

Retrieval efficiency of the proposed mode is compared and depicted in figure 5. It is observed that proposed model attains better performance even the number of query increases while the other system lags in performance for large number of queries.

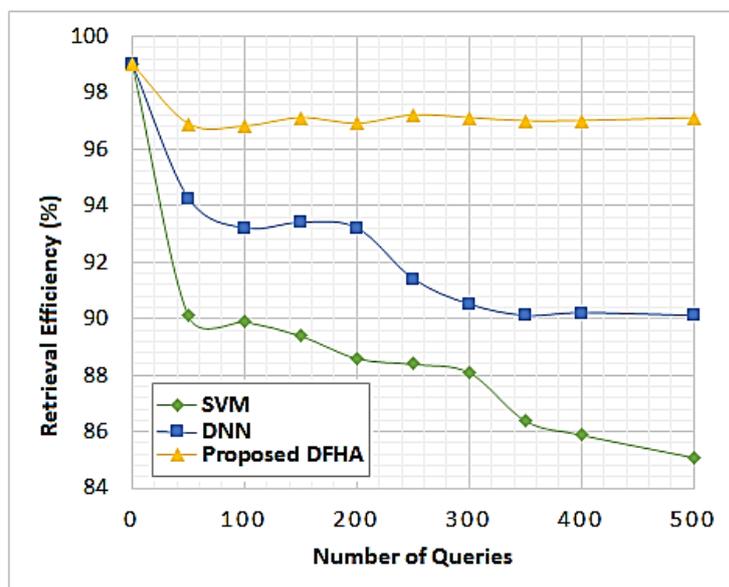


Fig.5 Retrieval efficiency comparison

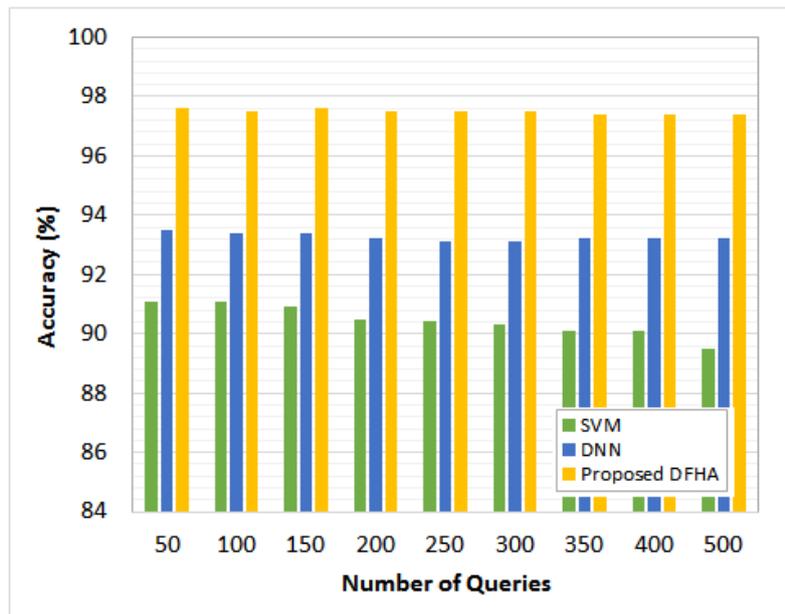


Fig.6 Accuracy comparison

Figure 6 depicts the accuracy comparison of proposed model. Since the proposed system is designed for distributed cloud environment, the accuracy is an essential factor to measure the performance of proposed work. it is observed that proposed model has better accuracy range of 97.56% which is 4% greater than deep neural network and 6% greater than support vector machine based information retrieval system.

5. Conclusion

A Hybrid Deep fuzzy hashing algorithm for information retrieval from distributed cloud environment is presented in the paper. Considering the issues in conventional information retrieval systems the proposed model is designed to achieve high retrieval accuracy and efficiency. Hybrid combination of deep learning and fuzzy along with hashing algorithm improves the data management in cloud environment. Proposed work is experimentally verified and compared with conventional support vector machine based information retrieval system and deep neural network. Proposed model achieves 97.6% retrieval efficiency which is considered as a remarkable improvement in information retrieval systems. Handling wide range of features is the limitation of proposed work. further the research work could be improved using optimization models to reduce the number of features in feature selection process.

References

1. Jiaying Liu, Xiangjie Kong, Xinyu Zhou, Lei Wang, Da Zhang, Ivan Lee, Bo Xu, Feng Xia (2019). Data Mining and Information Retrieval in the 21st century: A bibliographic review. *Computer Science Review*, 34:1-15.
2. Yongjun Zhu, Erjia Yan, Il-Yeol Song (2017). A natural language interface to a graph-based bibliographic information retrieval system. *Data & Knowledge Engineering*, 111:73-89.
3. Sidali Hocine Farhi, Dalila Boughaci. Graph based model for information retrieval using a stochastic local search. *Pattern Recognition Letters*, 105:234-239.
4. Joby, P. P (2020). Expedient Information Retrieval System for Web Pages Using the Natural Language Modeling. *Journal of Artificial Intelligence* 2(02):100-110.
5. Andrei S. Kulunchakov, Vadim Strijov (2017). Generation of simple structured information retrieval functions by genetic algorithm without stagnation. *Expert Systems with Applications*, 85:221-230.

6. Deepanwita Datta, Shubham Varma, Ravindranath Chowdary, Sanjay K. Singh (2017). Multimodal Retrieval using Mutual Information based Textual Query Reformulation. *Expert Systems with Applications*,68: 81-92.
7. Stefania Marrara, Gabriella Pasi, Marco Viviani (2017). Aggregation operators in Information Retrieval. *Fuzzy Sets and Systems*, 324:3-19.
8. Hamid Khalifi, Abderrahim Elqadi, Youssef Ghanou (2018). Support Vector Machines for a new Hybrid Information Retrieval System. *Procedia Computer Science*, 127:139-145.
9. Youssef Chouni, Mohamed Erritali, Youness Madani, Hanane Ezzikouri (2019). Information retrieval system based semantique and big data. *Procedia Computer Science*, 151:1108-1113.
10. Jianchang Lai, Yi Mu, Fuchun Guo, Peng Jiang, Willy Susilo (2018). Privacy-enhanced attribute-based private information retrieval. *Information Sciences*, 454–455:275-291.
11. Razane Tajeddine, Oliver W. Gnilke, Salim El Rouayheb (2018). Private Information Retrieval from MDS Coded Data in Distributed Storage Systems. *IEEE Transactions on Information Theory*, 64(11): 7081-7093.
12. Heecheol Yang, Wonjae Shin, Jungwoo Lee (2018). Private Information Retrieval for Secure Distributed Storage Systems. *IEEE Transactions on Information Forensics and Security*, 13(12):2953-2964.
13. Youcef Djenouri, Asma Belhadi, Philippe Fournier-Viger, Jerry Chun-Wei Lin (2018). Fast and effective cluster-based information retrieval using frequent closed itemsets. *Information Sciences*, 453:154-167.
14. Massimo Melucci, Adriano Paggiaro (2019). Evaluation of information retrieval systems using structural equation modelling. *Computer Science Review*, 31:1-18.
15. Marco Angelini, Vanessa Fazzini, Nicola Ferro, Giuseppe Santucci, Gianmaria Silvello (2018). CLAIRE: A combinatorial visual analytics system for information retrieval evaluation. *Information Processing & Management*, 54(6):1077-1100.
16. Saravana Kumar.C.S, Santhosh.R (2020). Effective information retrieval and feature minimization technique for semantic web data. *Computers & Electrical Engineering*, 81:1-13.
17. Ruixiang Ou, Yao Huang, Feng Pan, Hui Pan (2019). Research on information retrieval model under scarcity theory and user cognition. *Computers & Electrical Engineering*, 76:353-363.
18. Raj, Jennifer S (2019). Efficient Information Maintenance Using Computational Intelligence in The Multi-Cloud Architecture. *Journal of Soft Computing Paradigm (JSCP)*, 1(02):113-124.