

Novel Mechanism to Improve Hadith Classifier Performance

Kawther A. Aldhlan
Kulliyyah of Information and
Communication Technology,
International Islamic
University Malaysia
k_aldhlan@hotmail.com

Akram M. Zeki
Kulliyyah of Information and
Communication Technology,
International Islamic
University Malaysia
akramzeki@iiu.edu.my

Ahmad M. Zeki
Department of Information
System College of Information
Technology, University of
Bahrain, Kingdom of Bahrain
amzeki@uob.edu.bh

Hamad A. Alreshidi
Department of Instruction
technology, College of
Education, University Of Hail,
Hail, Saudi Arabia
mr_hamad15@hotmail.com

Abstract— Muslims believe that the Sunnah of the Prophet Muhammad (SAAW) is the second of the two revealed fundamental sources of Islam, after the Holy Qur'an. Hadith provides a Gold Standard "ground truth" for Artificial Intelligent (AI) knowledge extraction and knowledge representation experiments. In the present study, the extracted Islamic knowledge represented the focal point of the research; three famous books in Hadith science framed the corpus of the study. This study attempted to explore new approach to classify Hadith using a combination of the expert system and data mining techniques to classify Hadith according to its validity degree (Sahih, Hasan, Da'eef and Maudu'), the proposed Hadith Classifier model was built through learning process, Decision Tree (DT) classifier modeling had been represented by the tree structure model, and the attributes of the instances originally were obtained from the source books. Whilst some attributes were indicated as null values, or missing values. A novel mechanism called missing data detector (MDD) was employed to handle these missing data. This mechanism simulated the Isnad verification methods in Hadith science. The results of the research were compared with the resource books, concurrently with the point of view of the experts in the Hadith science. The findings of the research showed that the performance of DT Hadith classifier had significant effect with MDD, the CCR was sharply increased from (50.1502 %) to (97.597%) Furthermore, the favorable obtained results indicated that the DT Modeling is a viable approach to classify Hadith due to the ease of rules induction and results interpretation.

Keywords- Data Mining; Decision Tree; Hadith Classifier; Missing Data; Supervised Learning Algorithm.

I. INTRODUCTION

Data mining (DM) is the process of finding patterns that lie within large collections of data. Data mining is the analysis of observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data owner [1]. DM has become a widely used tool in a number of fields, including business, finance [2], security [3], medicine [4] and human science. DM methods include neural networks [5], decision trees (DT) [6], cluster analysis, market basket analysis, and regression analysis, among others.

The tree structured modeling is a data mining technique used to recursively partition a dataset into relatively homogeneous subgroups in order to make more accurate

predictions on the future instances. Moreover, decision tree algorithms have the ability to deal with missing values, while this ability is considered to be advantage, the extreme effort which is required to achieve it is considered a drawback. The algorithm must employed enhanced mechanisms to handle missing values. However, ignoring these missing data may cause critical decision. In the research case, the ignoring of missing values may cause incorrect Hadith classification that misleads to reject or accept Hadith. Thus, current study is conducted to propose approach to classify Hadith according to the validity of its Isnad (Sahih, Hasan, Dae'f and Maudu'). The target approach using a novel mechanism to deal with missing data in the Isnad attributes. The sample of the study is collected from three books in Hadith includes Sahih Al-Bukhari, Jami'u Al-Termithi and Silsilat Al-Ahadith Al-Dae'ifah w' Al-Mawdhu'ah. The evaluation of the proposed algorithm is carried out by comparing the results of classification with the point of view of the expert in Hadith science

II. LITERATURE REVIEW

As the rest of the Islamic science, researchers have increased intension to process Hadith and to simulate its methods in detecting and validation Hadith which is called Takhreej Al-Hadith. With respected to the efforts that are provided in computerized Hadith, even for the software that are produced by commercial companies, a few researches are conducted to implement Takhreej Al-Hadith [7]. In this regards, Alraza[8] adopted DM techniques to extract Islamic knowledge from the tradition books. He presented a practical experiment to explain the mechanism of these techniques. He used "Al-Resalh" book for AL-Shafe'i and "Al-Mu'amlat" book by Al-Shatebi as samples for this experiment. The findings of his study indicated that DM techniques can determine the main features of the Hadith methodology in Al-Shafei's book, as well as, the basic characteristics for speech in Al-Shatebi's book. While [9] explored the Implementation of a text classification method to classify Prophet Mohammed's traditions. The corpus of this study contained eight books separated into eight files, the study used testing set contained eighty Hadith from the same collection. The average accuracy of this sample is approximately 83.2%. In most relevant work to Hadith validation Alraza [10],[11] Established theoretical frames to represent Hadith literature, also he has adapted the use of

expert systems to acquire principles of information that Compatible with Hadith scientists methods in Criticism Maten and Tracing transmitters. Furthermore, he has developed rules based on tradition books to formulate the rules of the Knowledge system. Ghazizadeh [12] agreed with [10],[11] to use expert system software to implement the fuzzy system where the data knowledge base has been designed and the essential rules have been extracted to determine the rate of validity of Hadith, The deduced results from designed expert system were compared with their expert. The comparison showed that the system was correct in 94% of the cases. Meanwhile, Hyder & Ghazanfer [13] defined a graph theoretic representation of the chain of narrators and an aligned database structure suitable for storing the biographical data of the narrators and other historical events. Their study aimed to use computer science concepts for algorithmic research, database queries, and data-warehouses besides using of advanced data-mining techniques to assist Hadith research and research in Islamic history and literature. Their way to represent Hadith was amenable for cross verification and analysis in a computationally feasible manner, they found the nodes and arcs with various kinds of weights and then evaluating the aggregate averages over different paths and over the entire graph to yield numerical grades of evaluations. According to their findings the classifications of Hadith are qualitative, and these kinds of aggregate functions would enable quantitative grading of these classifications. Such quantitative grades would make it easier to compare and contrast criteria for evaluations.

III. RESEARCH METHOD

The current study attempts to reach the same goal of classification using supervised learning algorithms, 999 Hadiths from Sahih Al-Bukhari, Jami'u Al-Termithi and Silsilat Al-Ahadith Al-Dae'ifah w' Al-Mawdu'ah are framed the sample of the study, the attributes of the Hadith database are calculated according to the validate methods of Hadith science. The sample is divided into two parts (75%) as training dataset to build the classification model, while the rest of the sample (25%) is used to evaluate the performance of the Hadith classifier model. Moreover, the experiment applied C4.5 algorithm to extract the rules of classification. Figure1 illustrates the research framework using Missing Data Detector method (MDD).

The summary of the process in Figure1 are as follows: There is a training data set including four classes. Different shapes denote different classes. The whole training data set is portioned in to four classes A1, A2, A3 and A4. Some objects from A1, A2 and A3 have missing attributes that may classify them into incorrect class.

Step1: Applying the proposed mechanism into the training dataset to detect the missing attributes.

Step2: Applying DT algorithm to classify Hadith.

Step3: Some objects are correctly classified, while other objects are still in the incorrect class.

Step4: Building the tree and inducing the rules.

A. Hadith database

According to Tahan [14] there are five conditions must be satisfied to validate the Isnad of Al- Hadith:

- (1)All narrators in Isnad were renowned for their honesty.(2) All narrators in Isnad were renowned for their accuracy
- (3)There is no interrupting in the Isnad. (4)There is no irregular statement in the Hadith Maten (5)There is no defective in the Hadith Maten. Therefore, the experiment corpus consists of five basic features (link, defective, irregular, grade of reliability, grade of preservation). Table 1 shows the attributes with the possible values.

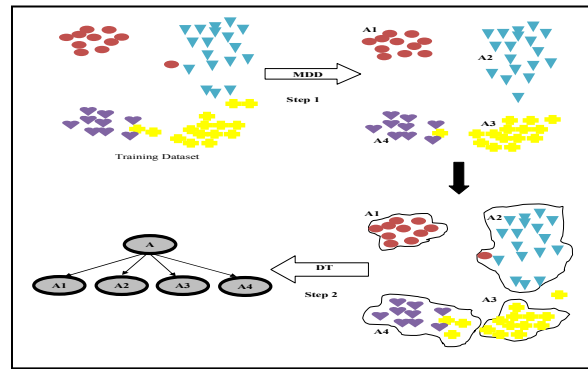


Figure1: Research frame work

IV. THE PROPOSED CLASSIFICATION APPROACH

The proposed approach consists of four phases; first one is the data pre-processing. Followed by the training phase, the input of this phase is a set of pre-classified documents, while the output is the Hadith classifier model. Whilst, the third phase is the classification (testing) phase which is responsible to test the prediction ability of the proposed classifier. Finally, evaluation phase.

V. THE EXPERIMENT PROCEDURES

A. Data Pre-processing

As mentioned earlier, the dataset was collected from different books, therefore, data pre-processing is conducted on each Hadith in the training and testing sets to reduce redundancy and to uniform the style of Hadith.

This phase includes:

- 1) *Attaching Isnad:* Some Hadith were separated from their Isnad either for suspicion in the narrator chain or redundancy. This process aimed to attach the Isnad at the beginning of the Maten to facilitate the narrators' chain scanning.

2) *Removing punctuation and diacritical marks:*

Removing diacritical and punctuation marks is important since these marks are prevalent in AL-Hadith and have no effect on determining the class of Hadith.

3) *Adding special character:* Adding special character to distinguish between the narrators while scanning Isnad. Table 2 shows the results of the pre-processing stage.

B. Experiments Specifications

The target approach is supervised classification. The training dataset is used to be applied by learning algorithm, in purpose to build Hadith classifier model .In the experiments author uses (75%) of AL-Hadith database as training set , while the rest (25%) of the sample is used as testing set. Two algorithms of learning are chosen to run using the same corpus after and before applying the detector method these are C4.5 and naïvebayes.

TABLE 1
The Attributes of the sample

ID	link	Irregular	Defective	Grdae Of Reliability	Grade Of Preservation	Class
1	True	False	False	True	True	Sahih
2	True	False	False	True	False	Hasan
3	False	False	False	True	True	Hasan
4	False	False	False	True	False	Hasan
5	True	False	True	True	True	Hasan
6	False	False	False	True	False	Daeef
7	True	False	False	True	Poor	Daeef
8	True	False	False	Daeef	True	Daeef
9	True	False	False	Daeef	poor	Daeef
10	True	False	False	False	True	Daeef
11	True	False	False	Any	Poor	Daeef
12	False	True	False	Null	Any	Maudof
13	False	Any	Any	Matrook	Any	Maudof
14	False	Any	Any	Monker	Any	Maudof
15	False	Any	Any	Liar	Any	Maudof
:	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:

TABLE 2
Results of Preprocessing Phase

Step	Result of the step
Attaching Isnad	عن عبد الله بن سعد الرقي حدثني والدتي مروة بنت مروان قالت حدثني والدتي عاتكة بنت بكر عن أبيها قالت: سمعت الزهري يحدث عن سالم بن عبد الله عن ابن عمر أن رسول الله صلى الله عليه وسلم قال: (ما ترك عبد شيئا لله إلا يتركه إلا عوضه منه ما هو خير له في دينه ودنياه)
Removing punctuation and diacritical marks	عن عبد الله بن سعد الرقي حدثني والدتي مروة بنت مروان قالت حدثني والدتي عاتكة بنت بكر عن أبيها قالت سمعت الزهري يحدث عن سالم بن عبد الله عن ابن عمر أن رسول الله صلى الله عليه وسلم قال ما ترك عبد شيئا لله لا يتركه إلا عوضه منه ما هو خير له في دينه ودنياه
Adding special character	عن .عبد الله بن سعد الرقي. قال حدثني والدتي .مروة بنت مروان. قالت حدثني والدتي .عاتكة بنت بكر. عن أبيها قالت سمعت .الزهري. يحدث عن .سالم بن عبد الله. عن .ابن عمر. أن رسول الله صلى الله عليه وسلم قال ما ترك عبد شيئا لله لا يتركه إلا عوضه منه ما هو خير له في دينه ودنياه

C. Attributes selection

The attributes are selected according to the information gained criteria. Table 3 illustrates the ranking of the features according to this criterion. See Figure2 for the resulted tree according to these attributes.

TABLE3
The Information Gained Of the Hadith Features

Feature	Information gain after splitting
Link	0.8711
Irregular	0.7927
Defective	0.704
Reliability_Grade	1.0201
Preservation_Grade	1.10296

D. Detection Of Missing Attributes

The present study proposed enhanced mechanism to handle the missing attributes in the Hadith database. This mechanism is based on the validate methods of the Isnad [14]:

1) *The status of reliability attribute in the Isnad chain:* Each narrator must be reliable and well known in the narration of Hadith. There are a lot of terms that indicate the reliability status of the narrator. Table 4 summarized these terms and the definitions regarding to the research goals.

2) *The status of the narrators' retention or preservation in the Isnad chain:* In this process the approach determines the value of the preservation for each narrator in the Isnad chain. Table 5 illustrates the terms of narrator's retention.

3) *The status of the link attribute in Isnad chain :* There are three methods to evaluate the status of the Isnad link (a) Tracing the student and the teachers for each narrator. (b) Check the time period between two consecutive narrators. (c) Check the place of each narrator and his journey.

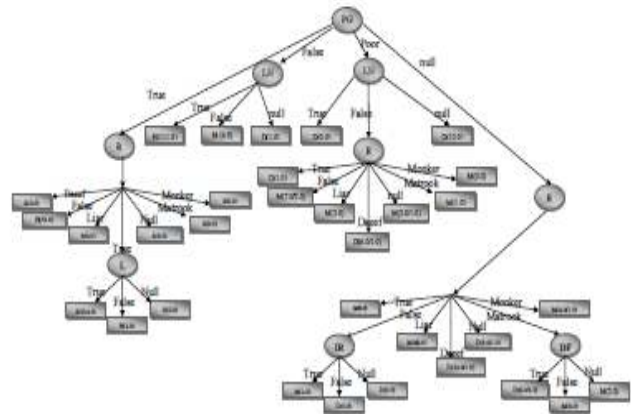


Figure2: The decision tree of the target Hadith Classifier

PG: preservation_Grade; RG: Reliability_Grade; LN: link status; IR: Irregularity; DF: defective; S: Sahih; H: Hasan; D:Da'eef; M:Maudoo'

TABLE 4
Hadith Terms Used To Indicate The Narrator's Reliability

Hadith Term	The attribute value
صحابي، أو ثق الناس ، ثقة ثقة ، ثقة حافظ ، إمام، ثبت ، عدل، ثقة	True
صدوق، لا بأس به، ليس به بأس، مقبول	False
صدوق سيء الحفظ، صدوق بهم، أو له أوهام، أو يخطئ، تغيير بأخرة	False
رمي ببدعة ، رمي بالتشيع، رمي بالقدر، لين الحديث، مستور ، مجهول، ضعيف	Daeef
متروك، متروك الحديث واهي الحديث، ساقط	Matrook
منكر الحديث	Monker
متهم بالكذب ، متهم بالوضع ، كذاب ، وضاع	Liar

TABLE 5
Hadith Terms Used To Indicate The Narrator's Retention

Hadith Term	The attribute value
الضبط	True
خفيف الضبط	False
سوء الضبط	Poor

4) *The status of the defective attribute in the Isnad chain:* This process aims to evaluate the value of the defective attribute of the narrators' chain.

E. Evaluation Strategy

It is important to measure the performance of classification model to determine how well the model will perform with new cases. The model performance evaluated after and before applying the detector in the testing phase. Four important measurements are used:

1) Correct Classification Rate (CCR):

CCR is the number of correctly predicted scores by the classifier. It is also known as the accuracy of the classifier. This measurement is represented by (1).

$$CCR=(NCP/NOP)*100 \quad (1)$$

Where CCR, NCP,NOP are the Correct Classification Rate, Number of Correct Prediction and total Number of Predictions, respectively.

2) Error Rate(ER):

Equation (2) represents the mathematical form of the number of incorrect prediction.

$$ER=(NWP/NOP)*100 \quad (2)$$

Where ER, NWP and NOP are the Error Rate, Number of wrong Prediction and total Number of Predictions, respectively.

3) Sensitivity :

The True Positive Rate (TPR) -called also recall- given that the actual value is positive. As represented in (3).

$$TPR=TP/(TP+FN) *100 \quad (3)$$

Sensitivity measures the proportion of actual positives which are correctly identified.

4) Specificity:

The True Negative Rate (TNR) of the classification model given that the actual value is negative, the fraction value classified as true negative [15].

$$TNR= TN/(TN+FP) \quad (4)$$

$$Sp = 1- FP \quad (5)$$

Specificity measures the proportion of negatives which are correctly identified.

5) Receiver Operating Characteristic (ROC) Curve:

ROC curves provide a visual model that displays the trade-off between sensitivity and specificity. The ROC curve is produced by graphing the false positive rate (FPR) which is the same as "1-Specificity" against the true positive rate (TPR) [16]. Figures 3 and 4 illustrate the ROC of the Hadith classifier before and after using MDD.

VI. RESULTS AND DISCUSSION

This section presents the main results of the experiment, then capped with a brief discussion. Table 6 illustrates the detailed accuracy by class. It can be seen from this table that the average of sensitivity of the case (2) has sharply increased with score (97.6%). Furthermore, the average of specificity of the same trial recorded better results (99.4%) than case (2) which indicates that the proposed model performance improved by MDD. And an ROC value result is (0.996) which indicates that the classifier with MDD is performed well with sharp increase of CCR (97.597%).

VII. CONCLUSION

All of all, the researchers can use any book as training data for knowledge extraction research. The holy Qur'an, Hadith and Islamic books are special case. They stand out as the source of a large collection of analysis and interpretation texts, which could provide a gold standard "ground truth" for AI (artificial intelligent) knowledge extraction and knowledge representation experiments. In addition researchers must cross-check for compatibility and consistency with knowledge extraction results from the Islamic corpus. Some computational results may be incompatible with specific inferences, which will shed new light on traditional interpretations. On the other hand, new outcomes may result from these experiments, thus adding to the canon of Islamic wisdom. The system that would

implement an Islamic knowledge must be reliable because it will be used by billions of Muslims, and non-Muslims.

In the present study, the extracted of Islamic knowledge represent the focal point of the research, three famous books in Hadith science represent the corpus of the study. The proposed Hadith classifier model was built through learning process, DT modeling had represented the structure model of the classifier, and the attributes of the instances originally were obtained from the source books. Whilst some attributes were indicated as null values, or missing data. A novel mechanism was employed to handle these missing data. This mechanism was generated based on the Isnad validity methods in Hadith science. As mentioned earlier, the implementation of the Islamic knowledge is very critical step due to its effects on the Muslim's life. Thus, the results of the research were compared with the resource books, concurrently with the point of view of the expert in Hadith science. The extracted knowledge represented the methods of Al-Imam Al-Bukhari, Al-Termithi and Al-Albani in Takhreej Al-Hadith, their approaches are slightly different. Therefore, it is difficult to claim that the proposed model represent all the Mohadeethen methods. The findings of the research showed that the performance of DT Hadith classifier had significant effect with the MDD. Whilst, the CCR was sharply increased from (50.1502 %) to (97.597%) Furthermore, the favorable results of the present research indicated that the DT Modeling is a viable approach to classify Hadith due to the ease of rules induction and results interpretation.

TABLE 6
Hadith Terms Used To Indicate the Narrator's Retention

Measurement Class	Case(1)Before MDD			Case(2) After MDD		
	SEN.	SEP.	ROC	SEN.	SEP.	ROC
Sahih	1	0	0.5	1	0.9994	0.997
Hasan	0	1	0.5	0.988	1	0.994
Da'eef	0	1	0.5	0.971	0.98	0.994
Maudu'	0	1	0.5	0.875	0.996	0.996
Weighted average	0.502	0.498	0.5	0.976	0.994	0.996
CCR	50.1502 %			97.597%		
ER	49.8498 %			2.4024%		

REFERENCES

[1] Hand, D., Mannila, H., and Smyth, P. *Principles of Data Mining*, Cambridge, 2001, MA: The MIT Press.
[2] Brachman, R. J., Khabaza, T., & Kloesgen, W. (1996). Mining business databases. *Communications of the ACM* 39, 42-48.
[3] Lee, w., Stolfo, S. J., & Mok., K. W. (1999). A Data Mining Framework for Buildin Gintrusion Detection Models. In *Proceedings of the 1999 IEEE Symposium on Security and Privacy*, Oakland, CA: IEEE, 120-132.
[4] Lavrac, N. (1999). Selected Techniques for Data Mining in Medicine. *Artificial intelligence in medicin*16 (1), 3-23 .
[5] Solomon, S., Nguyen, H., Liebowitz, J., & Agresti, W. Using data mining to improve traffic safety programs. *Industrial Management and Data Systems* , 5, 2006, pp. 621-643.
[6] Kotsiantis, S. B., Supervised Machine Learning: A Review of Classification Techniques. *Informatca* , 31, 2007,PP: 249-268

[7] Aldhlan, K. A., Zeki, A., & Zeki, A. *Encyclopedias of Hadeeth: The current status and future direction* . Seminar Warisan Nabawi Kali kedua ,2010 (p. 91). KL, Malaysia: universiti Sains Islam Malaysia.
[8] Alrazo, H. (2003). Al-'Utur al-ma'umatiah le tadawel al-ma'arafah al-islamiah fi zaman al-a'wlamh: Information frame works to deal with Islamic Knowledge in globalization era. *Journal of Islamic knowledge* 4, pp. 33-34.
[9] Al-Kabi, M. N., Kanaan, G., & Al-Shalabi, R. (2005). Al- Hadith Text Classifier. *Journal of Applied Sciences* 5(3), 584-587
[10]H.M. Alrazo, "الأنموذج المحوسب للسنة النبوية" Computerized frame of the Prophetic tradition', 17th National conferences for computer ,pp. 597-611.Madenh: scientific publishing center,2004.
[11] H.Alrazo,"تطبيقات التنقيب المعلوماتي على موارد المعرفة الإسلامية" Data mining application on the Islamic knowledge resource", 2008 . Retrieved JAN 13, 2010, from Alukah : <http://www.alukah.net/Culture/0/3123/>
[12]M.Ghazizadeh, M.H. Zahedi, M.Kahani,andB.M. Bidgoli, "Fuzzy Expert system in determining Hadith validity", *advances in computer and information sciences and engineering* ,2008, PP.354-359.
[13] S.I.Hyder and S.Ghazanfer, " Towards a database Oriented Hadith Research Using Relational, Algorithmic and Data-warehousing Techniques", *The Islamic Culture, Quarterly Journal of Shaikh Zayed Islamic Center for Islamic and Arabic Studies*, Vol. 19, University of Karachi, 2008,PP. 14.
[14]M.Tahan,"أصول التخریج ودراسة الأسانید", Riyadh: Al-Maref publishing ,1996.
[15] Kelly, H., Bull, A., Russo, P., & McBryde, E., Estimating sensitivity and specificity from positive predictive value, negative predictive value and prevalence: application to surveillance systems for hospital-acquired infections. *Journal of Hospital,Elsevier* , 2008, pp. 164-168.
[16] Fawcett, T. ,An Introduction to ROC Analysis. *Pattern Recognition Letters, Elsevier*, 2006, pp. 861-87.

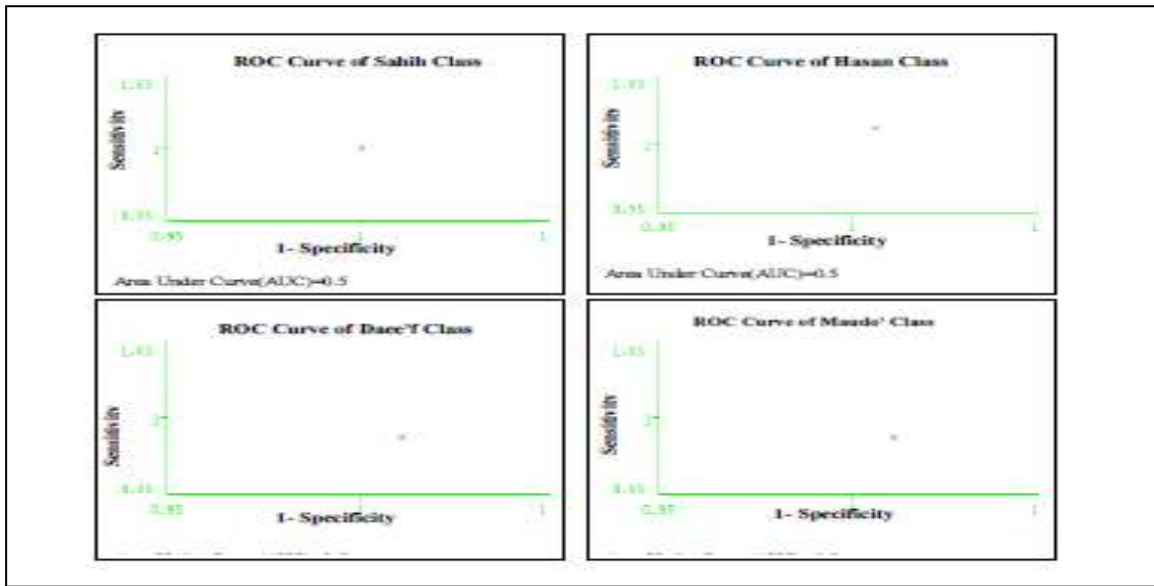


Figure.3:ROC Curves of the classes in Hadith classifier before using MDD

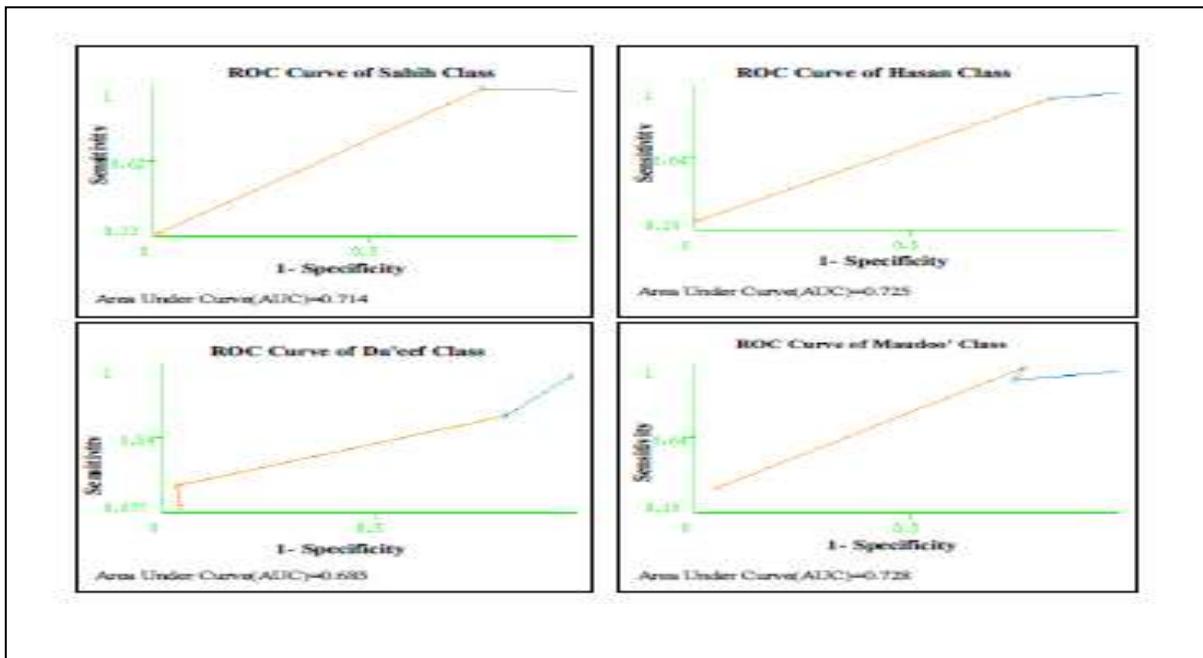


Figure 4: ROC Curves of the classes in Hadith classifier after using MDD