



A Recommendation System Based on AI for Storing Block Data in the Electronic Health Repository

Vinodhini Mani^{1*}, C. Kavitha¹, Shahab S. Band^{2*}, Amir Mosavi^{3,4,5*}, Paul Hollins⁶ and Selvashankar Palanisamy⁷

¹ Department of Computer Science and Engineering, School of Computing, Sathyabama Institute of Science and Technology, Chennai, India, ² Future Technology Research Center, College of Future, National Yunlin University of Science and Technology, Yunlin, Taiwan, ³ Faculty of Civil Engineering, TU-Dresden, Dresden, Germany, ⁴ Institute of Information Society, University of Public Service, Budapest, Hungary, ⁵ John von Neumann Faculty of Informatics, Obuda University, Budapest, Hungary, ⁶ Cultural Research Development School of Arts, Institute of Management, University of Bolton, Bolton, United Kingdom, ⁷ Intelligent Automation, Ford Motors Pvt. Ltd., Chennai, India

OPEN ACCESS

Edited by:

Thippa Reddy Gadekallu,
VIT University, India

Reviewed by:

Praveen Kumar,
VIT University, India
Abdul Rehman Javed,
Air University, Pakistan

*Correspondence:

Vinodhini Mani
vinodhini.cse@sathyabama.ac.in
Shahab S. Band
shamshirbands@yuntech.edu.tw
Amir Mosavi
amir.mosavi@mailbox.tu-dresden.de

Specialty section:

This article was submitted to
Digital Public Health,
a section of the journal
Frontiers in Public Health

Received: 08 December 2021

Accepted: 20 December 2021

Published: 21 January 2022

Citation:

Mani V, Kavitha C, Band SS,
Mosavi A, Hollins P and Palanisamy S
(2022) A Recommendation System
Based on AI for Storing Block Data in
the Electronic Health Repository.
Front. Public Health 9:831404.
doi: 10.3389/fpubh.2021.831404

The proliferation of wearable sensors that record physiological signals has resulted in an exponential growth of data on digital health. To select the appropriate repository for the increasing amount of collected data, intelligent procedures are becoming increasingly necessary. However, allocating storage space is a nuanced process. Generally, patients have some input in choosing which repository to use, although they are not always responsible for this decision. Patients are likely to have idiosyncratic storage preferences based on their unique circumstances. The purpose of the current study is to develop a new predictive model of health data storage to meet the needs of patients while ensuring rapid storage decisions, even when data is streaming from wearable devices. To create the machine learning classifier, we used a training set synthesized from small samples of experts who exhibited correlations between health data and storage features. The results confirm the validity of the machine learning methodology.

Keywords: artificial intelligence, machine learning, health repository, patients, health data, storage, deep learning

INTRODUCTION

In the modern era, clinicians no longer manage health data exclusively, but are increasingly responsible for obtaining consent from patients (1). The rights of patient's access to, analysis of, and exchange of their health information have evolved dramatically (2). The majority of patients are dissatisfied with their health care providers after sharing self-tracking data (3). It is still possible to enhance patient health care by incorporating patient health data into the current health data systems. Literature has identified various categories of patient health information (4). These categories include information about medications, biometrics, behavioral information, data about social interactions, genetics, psychological data, data about symptoms, and reports. Blockchain-based interplanetary file system secondary storage of health data has been implemented to safeguard the privacy and security of patient health information (5). Yet very few studies have evaluated how patients' health data is stored. A key component of the proper management of health data is protecting the privacy and confidentiality of the patient while maintaining data accessibility for relevant stakeholders. Studies indicate that health data security poses a massive threat. This is evidenced by the proliferation of medical devices with limited memory and power (6, 7) and substantial medical data repositories (8). Many types of organizations are responsible for managing the massive amount of health data.

Health data is often portrayed as being sensitive to all patients with the same level of privacy and confidentiality; however, this is not true in practice because it is not equally sensitive to everyone at the same time. When a patient reaches a high level of public prominence, she may surrender the ECG data she generated on her own and to her cardiologist. This data can be accessed by other healthcare providers through an electronic health record. A patient who wishes to keep her pregnancy test results private may be forced to allow her provider to store her pregnancy test results. The dissemination of health data between multiple providers who manage data repositories now enables the storage medium to be customized based on patient needs. This includes the cost, size, security, confidentiality, and privacy of each chunk of data. Hybrid execution models, such as those described by the author (9), allow sensitive data to be stored in private clouds while no sensitive data is maintained in public clouds. Nevertheless, it does not specifically address health data processing. Communication between the two cloud platforms also takes time, and computations that rely on bandwidth use a lot of resources. A hybrid cloud platform was developed by (10) for solving this problem. Medical sensors, apps, and devices provide data to artificial intelligence, which enables the automatic diagnosis of health conditions. Health data, including ECG, blood pressure, and pulse rate, can be classified as normal or abnormal by algorithms based on a range of conditions and thresholds set by healthcare professionals. Clinical research and clinical care are usually aided by abnormal data. Using the Body Area Sensor Network, (8) developed an agent-based system developed for elderly people to preserve abnormal data. Health information is generated in enormous quantities nowadays, so a diverse storage solution is needed (11). Several researchers have examined the performance and cost parameters of various Cloud Service Providers (CSPs) to design methods for selecting suitable CSPs for storing consumers' data (12–14). High-performance cloud services minimize the time spent in operations but incur high costs. Additionally, researchers are investigating blockchain technology for its promise of security and privacy for health data management. Combining blockchain-based eHealth with traditional health databases is possible, which can be arranged based on users' preferences and the possibility of utilizing the data in the future. However, due to the design of blockchains, they are not suitable for hosting large amounts of health data. A software agent that knows the patient's preferences is inserted inside the application in (15). Nonetheless, they never described a way to make this decision. To assist in choosing storage repositories, we developed a model that incorporated not only (8)'s criteria, but also aspects like data confidentiality, privacy, and quality of performance.

Motivation

Every Blockchain miner owns a local ledger, so this technology allows transactions to be verified and processed without the need for third parties. Verifying transactions does not require a centralized server. Document alterations cannot be guaranteed through conventional database storage and blockchain-based hash management. Data is only detectable in a blockchain if a hash pointer holds a pointer to it. Depending on the

patient, personal preferences, and other factors, the sensitivity and significance of the health information are also different from repository to repository. Choosing the right repository is extremely crucial. As wearable sensors continuously stream health data, the challenges are exacerbated. In (16), the author has surveyed the importance of artificial intelligence in healthcare. The prediction of COVID-19 infected patients using artificial intelligence has been implemented in (17), but there is a need for an appropriate repository to store the data.

Contribution

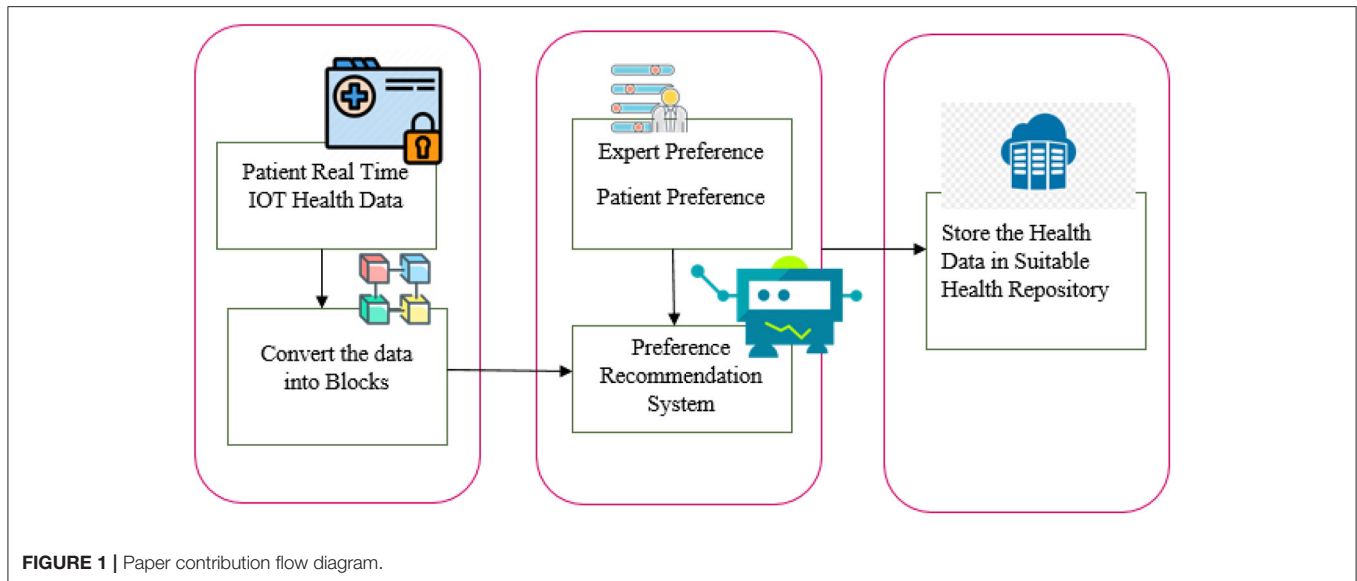
In our research, we considered the variation in data sensitivity, volume, and other factors to locate the appropriate system to manage health records. The flow diagram of the paper contribution is shown in **Figure 1**. Collect the health data and health repository parameters. Evaluations of both health information and health repository parameters are given a score. The machine learning-based recommendation model for health data storage proposes a way to distribute health data among multiple repositories. A model for automated health data storage recommendation is being developed to determine appropriate storage repositories. Through correlation analysis, user preferences, and clinical heuristics, a machine learning-based classifier is used to map health data characteristics to each repository. Patients' security and privacy preferences are taken into account as well as the sensitivity of health data.

Organization

Following are the sections of the paper: section Background addresses related work. In section Model for Recommendation of Health Repositories, we present the proposal for a recommendation model for a health repository. Section Implementation describes how the system will be implemented. The results and evaluation of performance will be discussed in section Results and Discussion. Conclusions and future work will be discussed in section Conclusion.

BACKGROUND

Big Data cannot be stored, accessed, or analyzed with a single health record system. Patients can lose medical information when their electronic health records are malfunctioning (18). Due to the manual uploading of data generated by wearable sensors to personal health records, caregiver responses were delayed. For this reason, (19) developed methods for storing patient-generated health information on commercial blood glucose monitors. The electronic health record system could be made to fit the streamed data if it is filtered or compressed (20). In (21–24), a number of action plans and standards were advocated for the adoption of an electronic health record system. A selection of an electronic health record should take into account functional requirements, troubleshooting, and optimization features (22). The author provides a list of steps to follow before buying an electronic health record system. Checklists mostly cover client meetings on site, site visits, and maintaining live workflows. Health data sources such as hospitals, clinics, insurers, and patients should be integrated into



centralized databases, according to the author (25). In particular, patient-centered health data with high degrees of structural heterogeneity must be stored and processed quickly because of their high volume and rate. For health data, to provide useful insights, precision is essential, but some sources produce vague and inaccurate information. Distributed data storage systems do offer some relief to these issues (26). Various cloud storage mediums have been examined. A machine learning and deep learning model is used to predict the thermal sensation vote system (27). Utilization of a compression algorithm to retrieve the health repository data as fast as possible using blockchain and interplanetary file systems (IPFS) without data loss (28). Diabetic Retinopathy is efficiently classified using a deep learning and machine learning algorithm (29). Genetic algorithm with fuzzy logic is a tool to help medical practitioners diagnose heart disease at an early stage using adaptive genetic algorithm with fuzzy logic (AGAFI) (30). Health data storage systems and data properties were not considered in the selection of repositories. Furthermore, no machine learning mechanisms were developed to cater to user preferences.

In the next section, we describe how we facilitate distributed health data management.

MODEL FOR RECOMMENDATION OF HEALTH REPOSITORIES

As data streams increase, the need for storage decisions becomes more frequent, making manual consultation with patients an inefficient process that requires an automated solution. It is, however, impossible to prespecify the data storage requirements for each patient that will apply to all possible future contexts. The learning classifier may generalize to a broader range of mappings based on a manual mapping specification by an expert.

The following sections explain in detail the overall approach described in **Figures 2, 3**. Data storage requirements - an

illustration of which is displayed in layer 1 of **Figure 2**, consists of a set of variables or features that characterize the requirements for storing a chunk of data. Some of the attributes' values have been shown to be numerical [1–10] and others to be qualitative. Secondly, each instance of the dataset contains the specifications required to store each chunk of data as shown in **Figure 2**.

Health Repository Evaluation Criteria are calculated in layer 3 by adding a rating provided by an expert group. These criteria reflect the characteristics of storage repositories as shown in **Figure 2**. Three standards apply to rank five storage repositories. Medical professionals and patients themselves may create clinical heuristic rules in layer-3 of **Figure 2** and each instance in the dataset is categorized according to the preferences of the users. A storage repository can be assigned to an instance based on heuristic rules in a real-world situation. The correlation coefficient offers an inference of a class label when preferences and heuristics do not match well. The health repository requirements can be mapped to layer-4 (user and expert expectations) by a machine learning classifier, as shown in **Figure 2**. In **Figure 3**, a recommendation framework for health repositories is illustrated. There are two parts to the framework: determining which standards should be used for the storage and assessment of data and implementing machine learning.

IMPLEMENTATION

This recommendation system assumes that a patient is in full control of his or her decision regarding storage. It is impossible to make decisions manually in many cases because they are made so frequently. Hence, automated processes are essential. In the mapping process, the characteristics of a repository managed by an agent group are matched with the characteristics of data about the storage requirements of patients. Because patients' storage requirements vary so much, it is impossible to predetermine every possible scenario. By utilizing a set

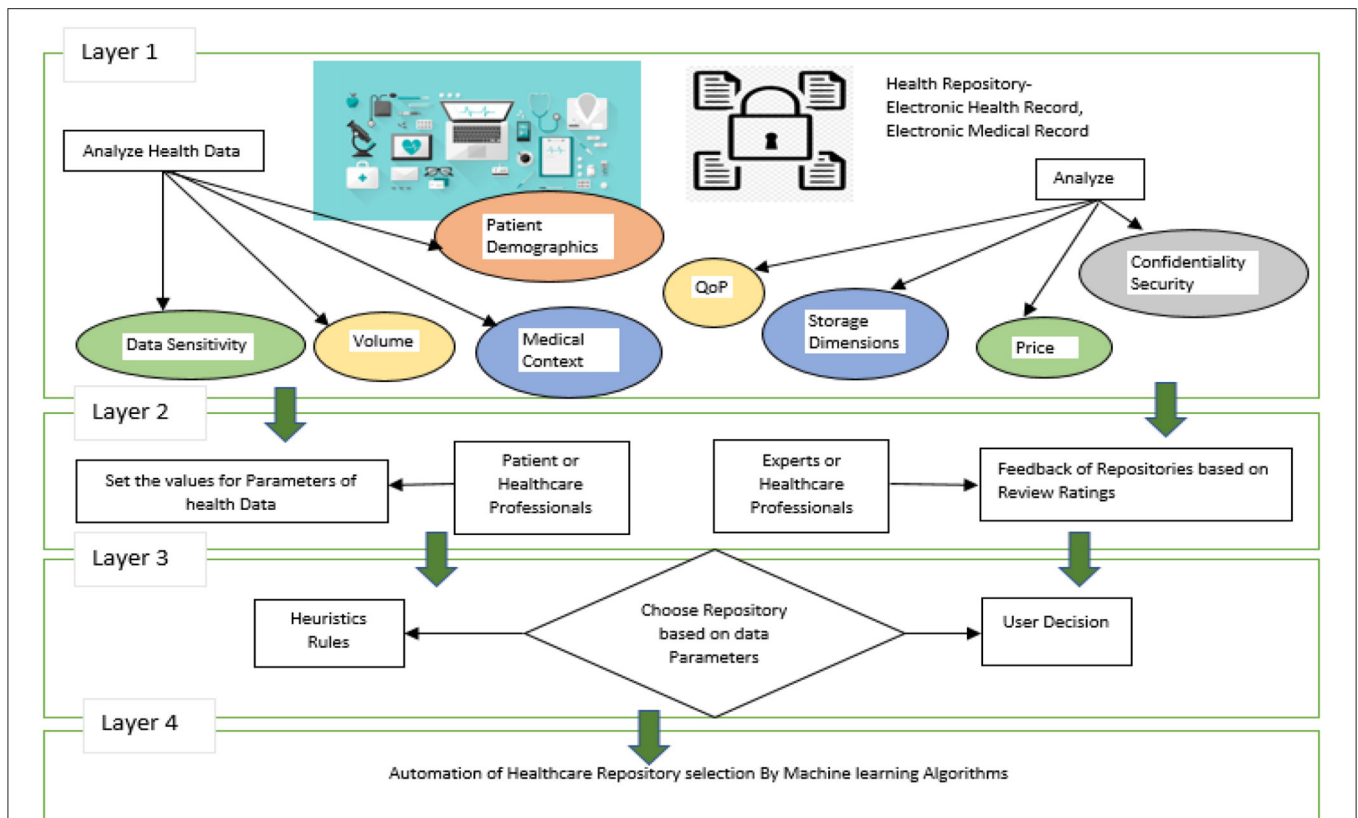


FIGURE 2 | Proposed system architecture.

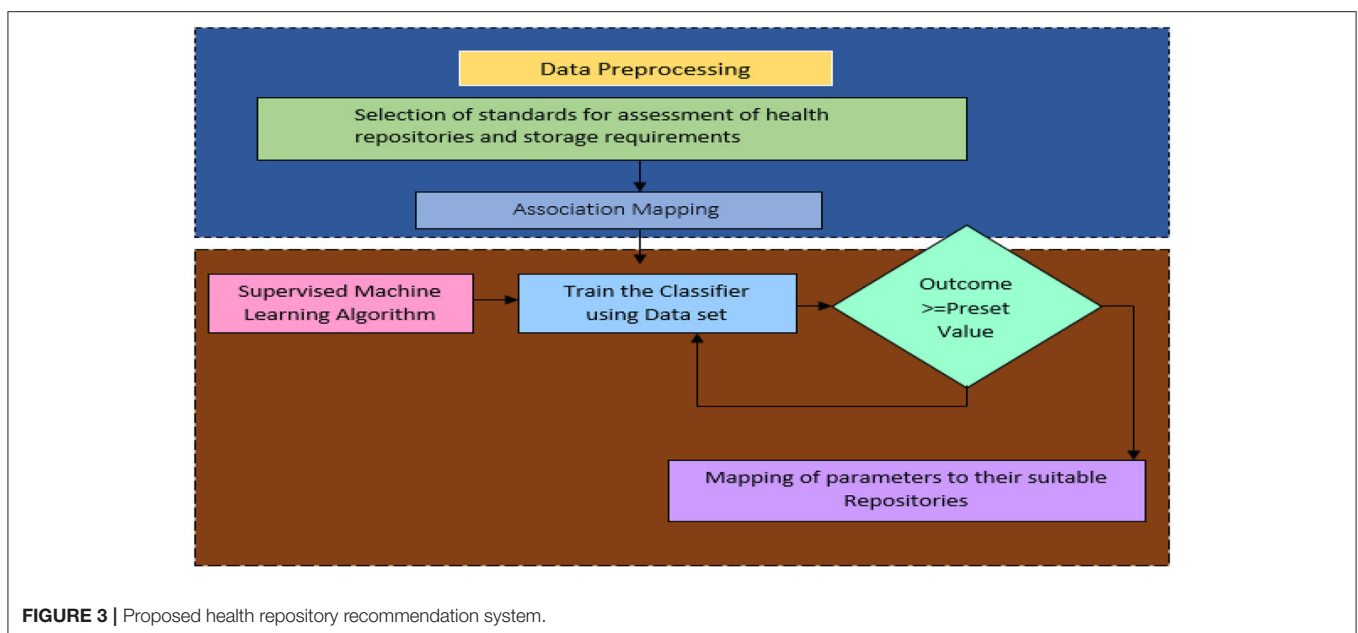


FIGURE 3 | Proposed health repository recommendation system.

of mappings that is specified manually by experts, machine learning is used to generalize a mapping over a wide range of patient contexts. This methodology involves defining a set of attributes that describe what chunk of data needs to be stored.

There are numerical values and categorical values assigned to those attributes. Thus, a dataset containing these attributes will be created, with each instance representing a different set of storage requirements. A group of experts' ratings are then

used to determine the characteristics of the available storage mediums. To determine what class each instance falls into, statistical correlation and heuristic rules are employed. Based on the training datasets, the supervised machine learning classifier maps the data into a storage repository. **Figure 3** illustrates two components of the recommendation system: Data Pre-processing and Supervised Machine Learning. According to **Figure 3**, the upper portion of the framework contains the characteristics of the data storage requirements. There are a number of features that demonstrate the characteristics of health repositories. A number of associations were found between the two groups of features.

Data Preprocessing

The data collected from hospitals and patients undergoes a preprocessing process, which includes analyzing data storage requirements, identifying sensitive data areas, analyzing the volume of each record, analyzing the patient health profile, determining the demographics of patients, and analyzing health repository parameters as well as storage, cost, security, privacy, and performance.

Characteristics of Data Storage Requirements

To determine which repository is the best option, consideration is given to the sensitivity of the data, the volume of the data, medical care context, and demographics of the patient.

Sensitivity of the Data

It is imperative to prevent unauthorized access to all health-related data. Depending on the data type, some breaches are more likely than others. Depending on the individual's preferences and context, the level of data sensitivity may vary.

The Volume of the Data

Reports, medical diagnoses, and medication summaries are not frequently created, which means that their storage needs are less than those of health data sets.

Context of Medical Care

The context may be palliative care, critical care, chronic illness, or no chronic illness. The context may also differ based on the country.

Demographics of Patients

Several factors can play a significant role in determining which storage medium to use, such as socioeconomic status, occupation, education, and nationality.

Health Repository Evaluation Parameters

Evaluation parameters for health repository such as security, privacy, cost, storage capacity, and performance. **Table 1** shows the parameters and criteria of the health repository evaluation.

The Relationship Between Repository Evaluation Standards and Data Features

Medical records, in particular those generated by patients, are to be transferred to a health record system that reflects the preferences of the user and the data requirements. Health data

TABLE 1 | Health repository evaluation.

Assessment parameters	Survey questions for health repository ratings
Storage	Can the repository be used to store Big Data? Regarding processing Big Data, what is the repository's role? Are there any benefits to storing continuously streamed data in the repository?
Cost	Does deployment cost a lot? Does maintenance cost much? What is the service cost?
Security	Is the storage repository capable of maintaining data integrity? Does the storage repository have 24/7 accessibility? Are storage repositories resistant to cyberattacks?
Privacy	Is data accessible to third parties? Is the access control right given to the owner of the health records?
Performance	How fast can you upload files? Is it possible to retrieve data quickly? Is it possible to process data quickly?

Algorithm 1: Association mapping ().

```

Step 1:  Begin
Step 2:  Let Data Source as DS;
Step 3:  Let Storage Requirements as SR;
Step 4:  Let Health Repository Parameters as HRP;
Step 5:  For each data  $\in$  DS do
Step 6:    For each Storage Requirement  $\in$  SR do
Step 7:      Collect the data;
Step 8:      Identify the SR;
Step 9:      Collect the HRP;
Step 10:   For each SR and HRP do
Step 11:     Analyze the parameters using Evaluation
              Criteria;
Step 12:     If (SR  $\in$  HRP)
Step 13:       SR (SR1...n)  $\rightarrow$  HRP (HRP1...n);
Step 14:       Create Association Dataset as AD;
Step 15:     Else
Step 16:       Print Not Associated;
Step 17:     End; End; End; End; End;

```

requirements and criteria for evaluating storage are correlated in a one-to-many fashion as implemented in **Algorithm 1**. Some associations are strong, and some are weak. To facilitate the rapid processing of highly confidential data, a health record system may accept data blocks in plaintext format. Data with relatively low confidentiality can be highly sensitive due to the demographic characteristics of patients. Data about a patient's demographics, such as their educational background and professional experience, may affect their privacy concerns. Users can then choose from a variety of storage repositories that

Table 2 | Association mapping.

S. No	Characteristics of data storage requirements	S. No	Health repository evaluation parameters	Association mapping
1	Sensitivity of the data	A	Storage	1→(B,C,D,E)
2	The volume of the data	B	Cost	2→(A)
3	Context of Medical Care	C	Security	3→(E)
4	Demographics of patients	D	Privacy	4→(B,C,D,E)
		E	Performance	

Algorithm 2: Health repository recommendation system ().

- Step 1: Begin
- Step 2: data collected from various data sources;
- Step 3: Call Association Mapping ();
- Step 4: For each Health Data Block \in HB do
- Step 5: Select the Supervised Machine learning algorithm;
- Step 6: Train the Data block HB;
- Step 7: Apply Heuristic Rule;
- Step 8: If (Accuracy \geq Threshold)
- Step 9: Test data;
- Step 10: Allocate the Health Data Block
 HB→Health Repository HR;
- Step 11: Send (Recommend Repository to Patients);
- Step 12: Break;
- Step 13: Else
- Step 14: Continue;
- Step 15: End; End; End;

protect their confidentiality. The sample association mapping as shown in **Table 2**.

Supervised Machine Learning Algorithm

Dynamically suggest health repositories based on supervised learning for particular data blocks, which is implemented using **Algorithm 2**. A training dataset must be generated for every instance of the dataset in addition to the labeled training datasets. Health repositories will be assigned data blocks that have a number of attributes. Among the attributes are some that are directly linked to the data block and others that are directly linked to the patient. Attributes include data sensitivity, volume, context of care, and demographics of the patients. The health repository should consider for evaluation such as electronic health records, cloud based electronic health records, blockchain based electronic health records, patient health record, and Electronic Medical Records. We considered the following health repository parameters in this study: security, privacy, cost, storage capacity, and performance. Each repository has been assigned a rating value ranging from 1 to 10. Whenever other attributes are not significant in determining the health repository, a linear regression Y (1) is calculated to label the instance as

shown in Equation 1.

$$Y = A + RX \tag{1}$$

$$R = n \left(\sum_{i=1}^n xiyi - \left(\sum_{i=1}^n xi \right) \left(\sum_{i=1}^n yi \right) \right) \tag{2}$$

$$A = \frac{\left(\sum_{i=1}^n yi \right) - R \left(\sum_{i=1}^n xi \right)}{n} \tag{3}$$

Where R is the Coefficient which contains $R1, R2, R3, R4, R5, R6, R7, R8, R9, R10$ are calculated between the set of data storage requirements(DR) as shown in equation 2. Here are the evaluation criteria for Electronic health record (D1), Patient health record (D2), Cloud-based electronic health record (D3), Blockchain-based electronic health records (D4), and Electronic Medical records (D5). The calculation of health repository recommendation Di is estimated using the equation:

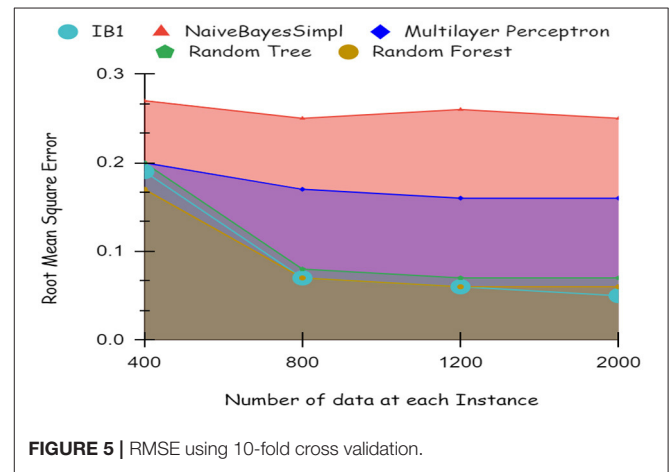
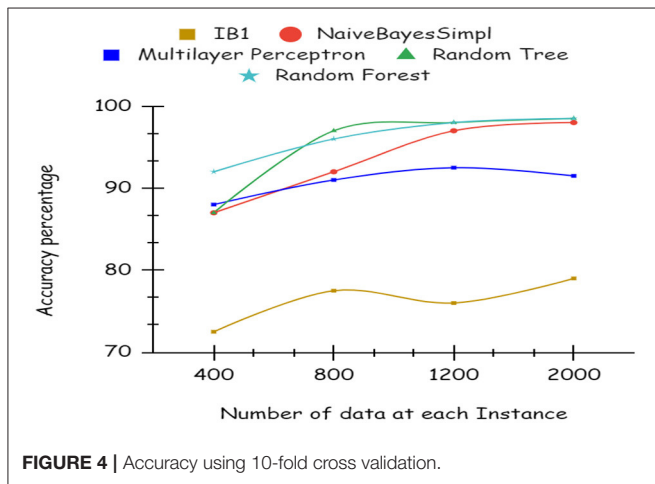
$$Di = High (R1, R2, \dots, Rm) \tag{4}$$

M is the number of health repositories and n is the rating criteria. Secondly, the choice of a health data repository can be influenced by the decision of the healthcare professional, the preferences of the user, and a variety of factors such as normal or abnormal behavior patterns and patient health status, as well as other demographic factors. Patients with unusual health patterns should store their health records in a repository that health care professionals can access quickly. A less secured and less expensive repository can be used to store data which is hardly ever accessed by health care professionals. Different users may have different privacy preferences, and those preferences may change over time based on different contexts (31). The health record system for a patient should take into account a variety of factors. There are several factors involved, such as medical conditions, personal characteristics, socioeconomic status, as well as the type and significance of data. The level of privacy and security preferences of individuals may change over time as well. In contrast to patients with terminal illnesses, young individuals may be more concerned with privacy and security. By considering author preference, some of the sample user preference and health professional preference heuristic rules were implemented, as shown below:

1. If (Data= standard && volume=large)
 - Then
 - Storage Repository=Cloud based Health Record Management System
2. If (Data= standard && volume=low)
 - Then
 - Storage Repository=Blockchain enabled Personal Health Record System
3. If (Data=Unusual patterns && volume=low)

Table 3 | Mapped sample training data set.

Information block	Sensitivity data	Volume	Context of medical care	Social status	Profile visibility	Patient status	Health repository
Data Block 1	1	2	3	3	high	Typical	Blockchain based electronic health record
Data Block 2	2	5	3	5	Low	Typical	Cloud electronic health record
.....
Data Block n	3	2	3	2	1	Abnormal	Electronic medical record



Then
 Storage Repository=Blockchain based Electronic Medical Record

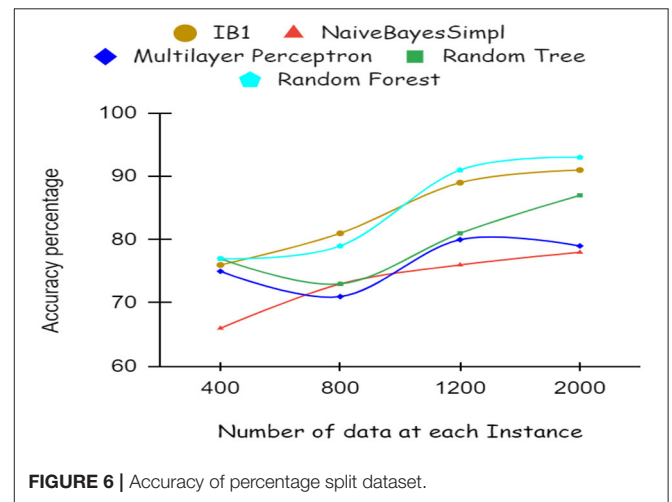
4. If (Patient= Famous Personality && health condition = Good)
 Then
 Storage Repository=Blockchain based Electronic Health Record

5. If (Patient= Famous Personality && health condition = Serious)
 Then
 Storage Repository=Blockchain based Electronic Medical Record

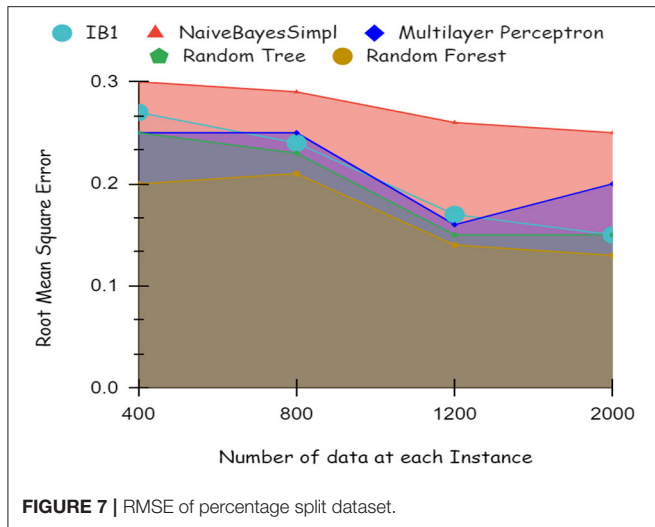
6. If (Data of type Disease)
 Then
 Store data in Disease Registry

RESULTS AND DISCUSSION

Research was conducted on supervised machine learning classification techniques. Using the WEKA tool, different classification algorithms were tested. The study used an Intel



Core i7 6700H processor with up to 3.5 GHz and 16 GB of RAM. The dataset was divided into training and test sets. Data preprocessing is performed prior to analysis. To train the data in the recommended health repository, linear regression data blocks and user and health professional preference rules have been used. During this experiment, we determine whether the classifiers can learn how to classify data distributions. The training datasets each



contain 400, 800, 1200, and 2000 instances. **Table 3** shows the mapped sample training dataset.

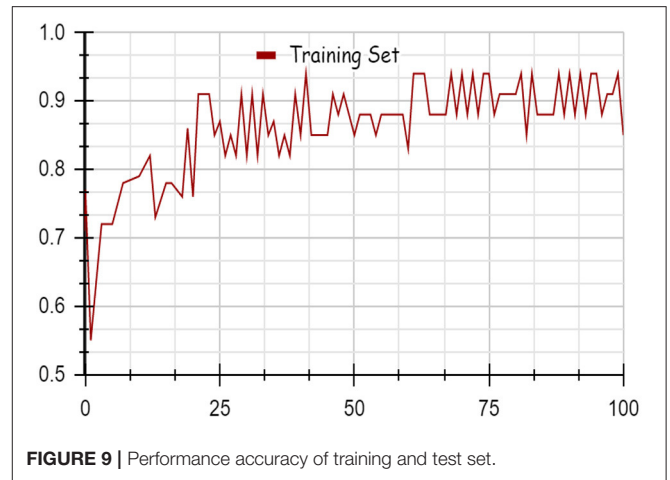
Four different classifiers were run on four datasets to test whether a machine learning algorithm could choose an appropriate storage medium, NaïveBayesSimple, Multilayer Perceptron, Random Forest Classifier, Random Tree and the IB1 algorithm are four different types of classifiers trained here. Several classification techniques were compared using Python to determine their accuracy scores (32).

Classification Model Accuracy

1. Confusion matrix.
2. Classification measure.

Confusion Matrix

In the confusion matrix, N is the number of target classes, and N is the number of rows. It is used to evaluate the performance of a classification model. Machine learning is used to predict target values from the actual values in the matrix. True Positive (TP) and True Negative (TN) rates should be high and False Positive (FP) and False Negative (FN) rates are low for a successful model. A confusion matrix as is always more appropriate as a



machine learning model evaluation criterion when working with an imbalanced dataset.

Classification Measure

As an evaluation measure, the classification measure is used in addition to the confusion matrix. They are:

1. Accuracy.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad 0.0 < \text{Accuracy} < 1.0 \quad (5)$$

2. Precision.

$$\text{Precision} = \frac{TP}{TP+FP} \quad (6)$$

3. Recall.

$$\text{Recall} = \frac{TP}{TP+FN} \quad (7)$$

4. F1-Score.

$$\text{F1-Score} = 2 \frac{\text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}} \quad (8)$$

5. Sensitivity and specificity.

$$\text{Sensitivity} = \frac{TP}{TP+FN} \quad (9)$$

$$\text{Specificity} = \frac{TN}{TN+FP} \quad (10)$$

6. Root mean square error.

Modified Mean Square Error (MSE) is a variation of Root Mean Square Error (RMSE). Measuring the mean square error squared is equivalent to this metric. The RMSE of an ideal model is zero, just as the MSE and MAE are zero.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\text{Actual Values} - \text{Predicted Values})^2} \quad (11)$$

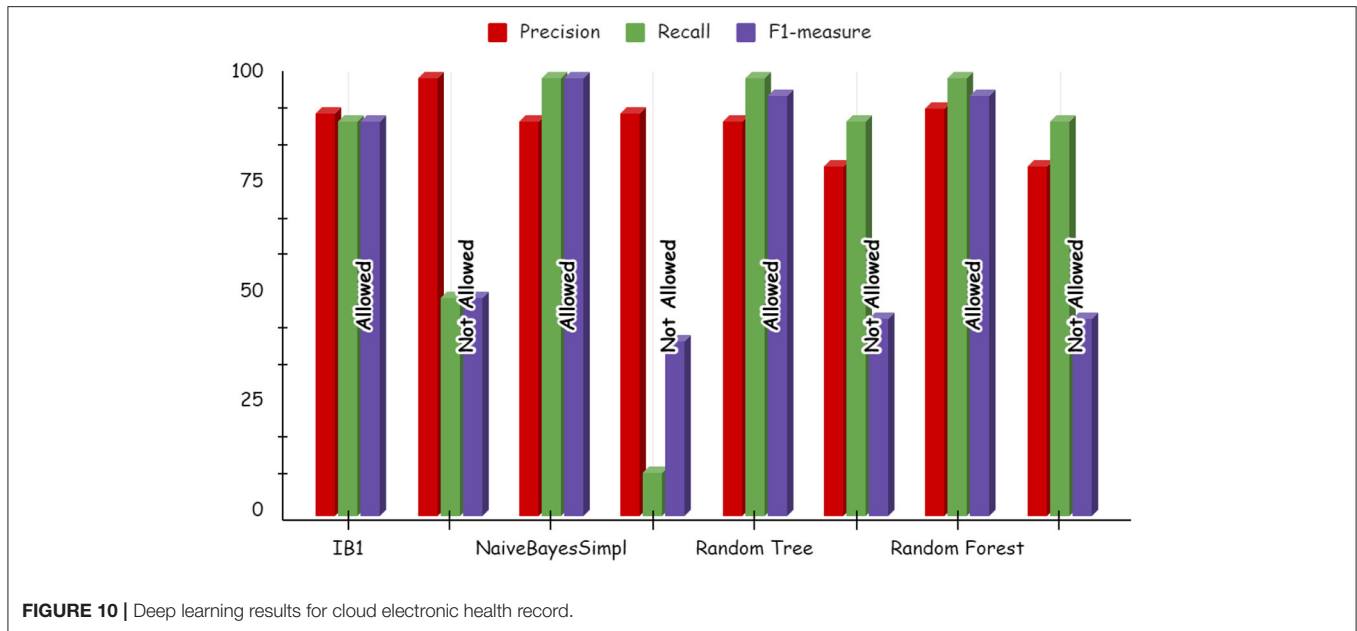


FIGURE 10 | Deep learning results for cloud electronic health record.

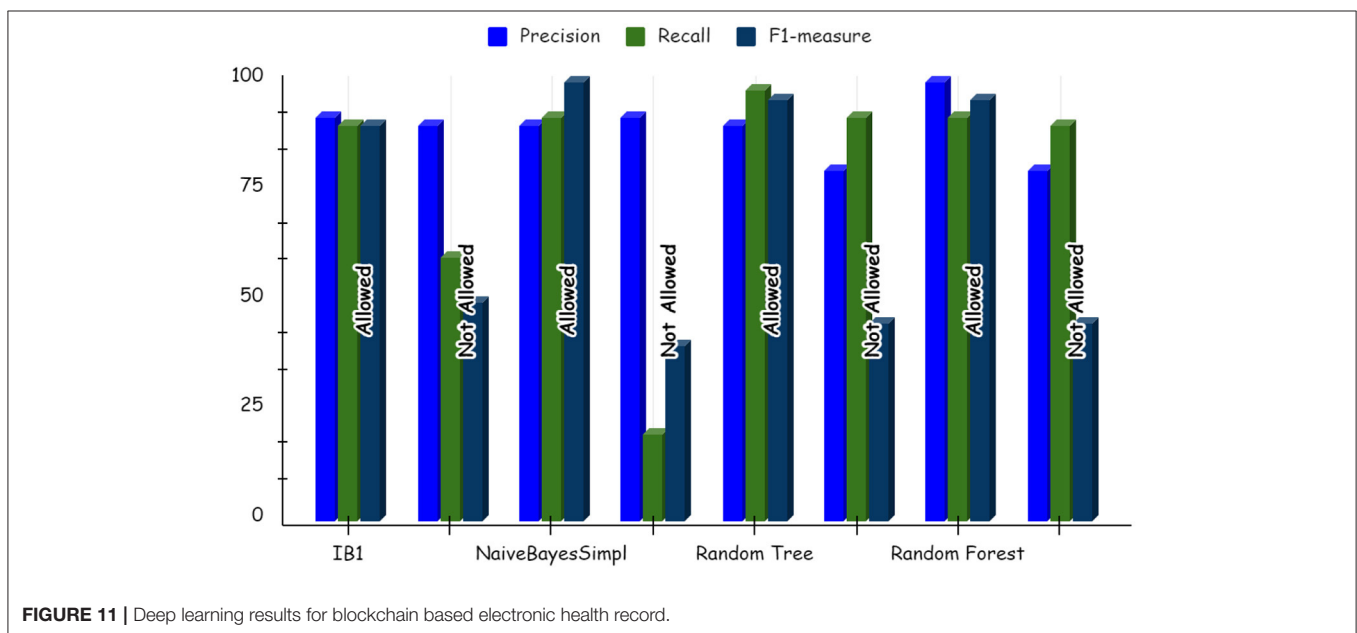


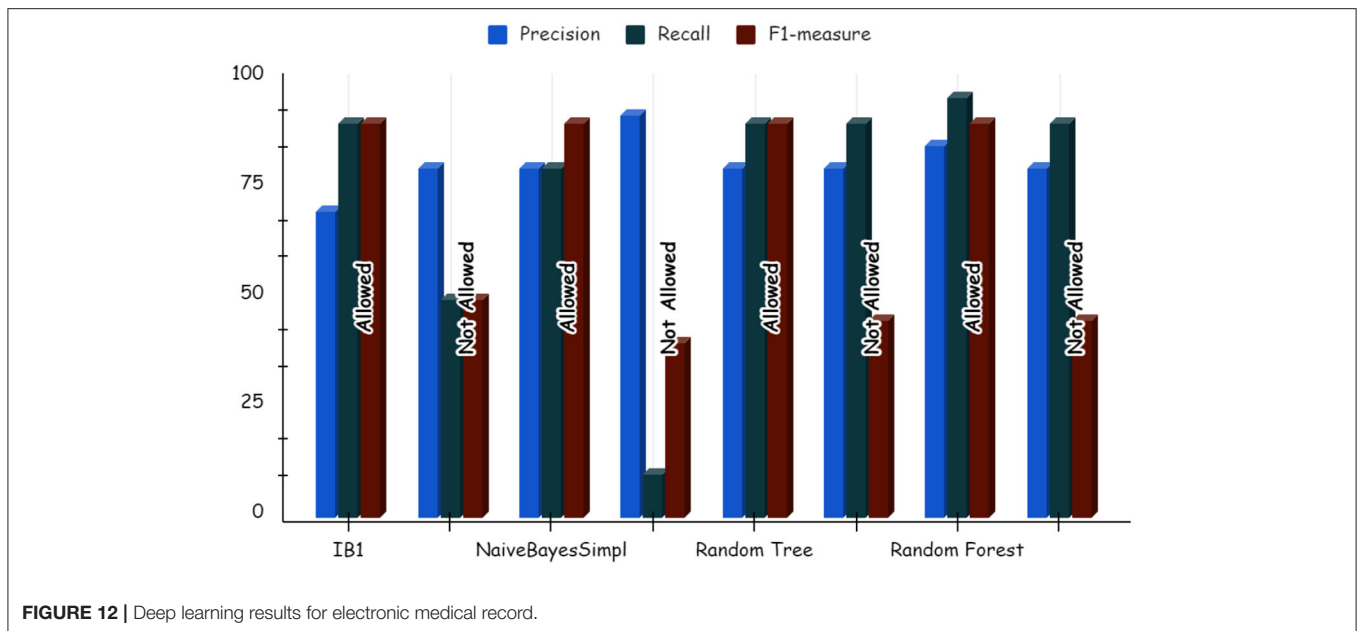
FIGURE 11 | Deep learning results for blockchain based electronic health record.

Result Analysis

As illustrated by the graph in **Figure 4**, Random Forest classifiers become more accurate as the number of instances increases, as shown by a 10-fold cross-validation analysis. A balanced ratio of each class was found in the dataset of 1,200 records, thus all classifiers performed better. The Random Forest performed best, with 98.21% accuracy. On the 2,000-record dataset, however, all classifiers had lower accuracy, largely because the dataset was skewed. Compared to other classifiers, Random Forest exhibits lower root mean square error in **Figure 5**. **Figure 6** illustrates the percentage split results, which are less accurate than the cross-validation results presented in 10-fold cross-validation.

By using a percentage split, 80% of the data were used for training and 20% for testing. The classifier is trained only once, as seen in **Figure 7**, which demonstrates low accuracy and large RMSE. Artificial intelligence is a technique for deep learning.

Using deep learning networks, unstructured or unlabeled data can be learned unsupervised. Real-world health repositories are usually recommended based on unstructured and unlabeled datasets. For our synthetic dataset, we analyzed the accuracy using a deep learning algorithm. A deep learning model is run on the synthetic dataset, and it shows 88.70 percent accuracy. It is implemented in Python. There are three hidden layers in the model; the first of these layers has 100 output nodes,



while the second and third have five output nodes each. Training is done with 100 iterations and eight batches are used. The training dataset is shown in **Figures 8, 9**, with a Y-axis showing the loss and X-axis showing the number of iterations. A deep learning classifier and a machine learning classifier are displayed in **Figures 10–12** for the classification. With reference to recall, F1-measure, and precision, the Random Forest classifier outperformed the other tested classifiers. Classes that were allowed and those that were not included in the experiment. In terms of recall, precision, and F1-measure, the Random classifier scored 93, 100, and 96% for cloud electronic health records, 100, 92, and 96 for blockchain-based electronic health records, and 85, 96, and 90 for electronic medical records. In terms of the allowed class, the rest of the experimented models perform well. In terms of the disallowed class, they did not perform well.

The accuracy of the classifier supports the use of machine learning to map the health storage mediums to health data blocks. Given the growing volume of health data that will need to be stored and accessed globally, this machine learning model may play a crucial role in improving storage and access arrangements in the future. This will make health data storage easy and straightforward for consumers. In addition, they would be able to ensure that the size of the data store is manageable. It can help to determine which storage solution best fits the requirements of different data assets using a machine learning model.

Mapping of Health Data Parameters to Repositories

Medical technology is expected to develop health record systems in the future. Health records are taking on novel forms as a result of the expansion of medical data. As described below, the proposed system will support various data variations and

health records. First, the system requests the ratings for the latest health record on the basis of health parameters from the IT staff and healthcare professionals. Second, the system relabels instances from the entire training dataset. As soon as a new instance is created, the old instances' labels do not change.

CONCLUSION

Health data will increasingly be preserved in a variety of repositories, so patients can select the repository that best meets their needs. Patients are realistically expected to avoid using a single repository for all their health data because the context of treatment, patterns of data, and legal constraints may change. To automate the storage decision, a selection algorithm must be developed. This is especially relevant in the case of constantly streaming health data. The process of choosing the right repository is complicated. In addition to knowledge of storage features used for interoperability, data security, and privacy, regulatory concerns must also be considered. To preserve confidentiality, we propose distributing health data among various vendors. By keeping medical records together, confidentiality will also be preserved. Based on factors like data type, sensitivity level, significance, patient safety, and privacy requirements, this model can recommend which health data blocks should be stored on which storage medium. When applied to the dataset generated, random forest yielded the highest accuracy of 96.4%. Accuracy of algorithms depends on the dimension, origin, and nature of the data. As a result, we intend to evaluate these various algorithms with different characteristic datasets in the near future. In the future, we will implement a role-based access control system to store medical record information by

integrating the health repository recommendation system to allow access to the health records based on the permission of patients.

DATA AVAILABILITY STATEMENT

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author/s.

REFERENCES

- Plastiras P, O'Sullivan D. Exchanging personal health data with electronic health records: A standardized information model for patient generated health data and observations of daily living. *Int J Med Inform.* (2018) 120:116–25. doi: 10.1016/j.ijmedinf.2018.10.006
- Cortez A, Hsui P, Mitchell E, Riehl V, Smith P. *Conceptualizing a data infrastructure for the capture, use, and sharing of patient-generated health data in care delivery and research through 2024 (white paper)*. (2018).
- Chung CF, Dew K, Cole A, Zia J, Fogarty J, Kientz JA, et al. Boundary negotiating artifacts in personal informatics: Patient-provider collaboration with patient generated data. In: *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing*, (2016). p. 770–786. doi: 10.1145/2818048.2819926
- Lordon RJ, Mikles SP, Kneale L, Evans HL, Munson SA, Backonja U, et al. How patient-generated health data and patient-reported outcomes affect patient-clinician relationships: A systematic review. *Health Inform J.* (2020) 26:2689–706. doi: 10.1177/1460458220928184
- Mani V, Manickam P, Alotaibi Y, Alghamdi S. Hyperledger Healthchain: patient-centric IPFS-based storage of health records. *Electronics.* (2021) 10:3003. doi: 10.3390/electronics10233003
- Albahri A, Zaidan A, Albahri O, Zaidan B, Alsalem M. Real-time fault-tolerant mhealth system: Comprehensive review of healthcare services, opens issues, challenges and methodological aspects. *J Med Syst.* (2018) 42:137. doi: 10.1007/s10916-018-0983-9
- Isern D, Moreno A. A systematic literature review of agents applied in healthcare. *J Med Syst.* (2016) 40:43. doi: 10.1007/s10916-015-0376-2
- Vaidehi V, Vardhini M, Yogeshwaran H, Inbasagar G, Bhargavi R, Hemalatha CS. Agent based health monitoring of elderly people in indoor environments using wireless sensor networks. *Procedia Comput Sci.* (2013) 19:64–71. doi: 10.1016/j.procs.2013.06.014
- Ko SY, Jeon K, Morales R. The hybex model for confidentiality and privacy in cloud computing. *HotCloud.* (2011) 11:8. doi: 10.5555/2170444.2170452
- Zhang H, Ye L, Du X, Guizani M. Protecting private cloud located within public cloud. In *Global Communications Conference (GLOBECOM)*. IEEE, (2013). p. 677–681.
- Stranieri A, Balasubramanian V. Remote patient monitoring for healthcare: a big challenge for big data. In *Managerial Perspectives on Intelligent Big Data Analytics*. IGI Global, (2019). p. 163–179. doi: 10.4018/978-1-5225-7277-0.ch009
- Ruiz-Alvarez A, Humphrey M. A model and decision procedure for data storage in cloud computing. In *Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012)*. IEEE Computer Society, (2012). p. 572–579. doi: 10.1109/CCGrid.2012.100
- Ruiz-Alvarez A, Humphrey M. Toward optimal resource provisioning for cloud mapreduce and hybrid cloud applications. In: *Proceedings of the 2014 IEEE/ACM International Symposium on Big Data Computing*. IEEE Computer Society, (2014). p. 74–82. doi: 10.1109/BDC.2014.14
- Yoon MS, Kamal AE. Optimal dataset allocation in distributed heterogeneous clouds. In: *2014 IEEE Globecom Workshops (GC Wkshps)* IEEE, (2014). p. 75–80. doi: 10.1109/GLOCOMW.2014.7063389
- Zhang Q, Lu J. Artificial intelligence in recommender systems. *Complex Intell Syst.* (2021) 7:439–57. doi: 10.1007/s40747-020-00212-w

AUTHOR CONTRIBUTIONS

VM and CK: conceptualization, methodology, investigation, data curation, and writing—original draft preparation. SB, AM, and PH: software, validation and visualization, and resources. VM, AM, and PH: formal analysis. VM, SB, and AM: writing—review and editing and supervision. VM, CK, SB, AM, and PH: project administration. All authors have read and agreed to the published version of the manuscript.

- Yang Y, Chen T. Analysis and visualization implementation of medical big data resource sharing mechanism based on deep learning. *IEEE Access.* (2019) 7:156077–88. doi: 10.1109/ACCESS.2019.2949879
- Stock C, Dias S, Dietrich T, Frahsa A, Keygnaert I. Editorial: How can We Co-Create Solutions in Health Promotion with Users and Stakeholders? *Front. Public Health.* (2021) 9:773907. doi: 10.3389/fpubh.2021.773907
- Andy YY, Shen CP, Lin YS, Chen HJ, Chen AC, Cheng LC, et al. Continuous, personalized healthcare integrated platform. In *TENCON 2012 IEEE Region 10 Conference*. IEEE, (2012). p. 1–6. doi: 10.1109/TENCON.2012.6412226
- Peleg M, Shahar Y, Quaglini S, Fux A, García-Sáez G, Goldstein A, et al. Mobiguide: a personalized and patient-centric decision-support system and its evaluation in the atrial fibrillation and gestationaldiabetes domains. *User Model User-Adapt Interact.* (2017) 27:159–213. doi: 10.1007/s11257-017-9190-5
- Hohemberger R, da Rosa CE, Pfeifer FR, da Rosa RM, de Souza PS, Lorenzon AF, et al. An approach to mitigate challenges to the electronic health records storage. *Measurement.* (2020) 154:107424. doi: 10.1016/j.measurement.2019.107424
- Busis NA. How can i choose the best electronic health record system for my practice? *Neurology.* (2010) 75:S60–4. doi: 10.1212/WNL.0b013e3181fc9888
- Weathers AL, Esper GJ. How to select and implement an electronic health record in a neurology practice. *Neurol Clin Pract.* (2013) 3:141–8. doi: 10.1212/CPJ.0b013e31828d9fb7
- Hart EM, Barmby P, LeBauer D, Michonneau F, Mount S, Mulrooney P, et al. Ten simple rules for digital data storage. *PLoS Comput. Biol.* (2016) 12:10. doi: 10.1371/journal.pcbi.1005097
- Wilson G, Bryan J, Cranston K, Kitzes J, Nederbragt L, Teal TK. Good enough practices in scientific computing. *PLoS Comput Biol.* (2017) 13:e1005510. doi: 10.1371/journal.pcbi.1005510
- Khan SI, Hoque ASML. Towards development of health data warehouse: Bangladesh perspective. In *2015 International Conference on Electrical Engineering and Information Communication Technology (ICEEICT)*. IEEE (2015). p. 1–6. doi: 10.1109/ICEEICT.2015.7307514
- Mackey TK, Kuo TT, Gummadi B, Clauson KA, Church G, Grishin D, et al. 'Fit-for-purpose?'—challenges and opportunities for applications of blockchain technology in the future of healthcare. *BMC Med.* (2019) 17:68. doi: 10.1186/s12916-019-1296-7
- Rehman SU, Javed AR, Khan MU, Nazar Awan M, Farukh A, Hussien A. PersonalisedComfort: a personalised thermal comfort model to predict thermal sensation votes for smart building residents. *Enterpr Inf Syst.* (2020). 1852316. doi: 10.1080/17517575.2020.1852316
- Mubashar A, Asghar K, Javed AR, Rizwan M, Srivastava G, Gadekallu TR. Storage and proximity management for centralized personal health records using an ipfs-based optimization algorithm. *J Circ Syst Comput.* (2021) 2250010. doi: 10.1142/S0218126622500104
- Gadekallu TR, Khare N, Bhattacharya S, Singh S, Maddikunta PKR, Srivastava G. (2020). Deep neural networks to predict diabetic retinopathy. *J Ambient Intell Human Comput.* (2020) 1–14. doi: 10.1007/s12652-020-01963-7

30. Reddy GT, Reddy MPK, Lakshmana K, Rajput DS, Kaluri R, Srivastava G. Hybrid genetic algorithm and a fuzzy logic classifier for heart disease diagnosis. *Evol Intell.* (2020). 13:185–96.
31. Trojer T, Katt B, Schabetsberger T, Mair R, Breu R. The process of policy authoring of patient-controlled privacy preferences. In: *International Conference on Electronic Healthcare*. Springer, (2011). p. 97–104. doi: 10.1007/978-3-642-29262-0_14
32. Analytics V. What is confusion matrix (2020). Available online at: <https://medium.com/analytics-vidhya/what-is-a-confusion-matrix-d1c0f8feda5> (accessed November 17, 2020).

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Copyright © 2022 Mani, Kavitha, Band, Mosavi, Hollins and Palanisamy. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.