

**APPLICATION OF MACHINE LEARNING FOR DETECTION OF HYPERTROPHIC
CARDIOMYOPATHY PATIENTS FROM ECHOCARDIOGRAM MEASUREMENTS**

Nasibeh Zanjirani Farahani, Ph.D.
Mayo Clinic
Rochester, MN

Moein Enayati, Ph.D.
Mayo Clinic
Rochester, MN

**Divaaakar Siva Baala
Sundaram, M.S.**
University of Minnesota
Rochester, MN

**Devanshi Damani,
MB.BS**
Mayo Clinic
Rochester, MN

Vinod C. Kaggal, M.S.
Mayo Clinic
Rochester, MN

April L. Zacher, M.S.
Mayo Clinic
Rochester, MN

Jeffrey B. Geske, M.D.
Mayo Clinic
Rochester, MN

Garvan Kane, M.D., Ph.D.
Mayo Clinic
Rochester, MN

**Shivaram Poigai Arunachalam,
Ph.D.**
Mayo Clinic
Rochester, MN

Kalyan Pasupathy, Ph.D.
Mayo Clinic
Rochester, MN

**Adelaide M. Arruda-Olson,
MD., Ph.D.**
Mayo Clinic
Rochester, MN

ABSTRACT

Hypertrophic cardiomyopathy (HCM) is a heritable, phenotypically diverse, and often asymptomatic heart muscle disease which is a major cause of sudden cardiac death (SCD) in young adults. The gold-standard for the diagnosis of HCM is echocardiography (echo), which is an ultrasound-based cardiac imaging modality. Across all sites of the Mayo Clinic enterprise, echo images and measurement data are reviewed, interpreted, and reported via the Echo Information Management System (EIMS). The objective of this paper is to develop a machine learning model for the identification of HCM from cardiac measurements obtained by the echo. We developed a novel machine learning model on patient demographic information and echo measurements that were retrieved from the EIMS digital data registry and selected by cardiologists.

Random forest (RF) was utilized to investigate the predictive performance of these features on the identification of HCM patients. The HCM cohort consists of 3,548 patients with at least one HCM diagnostic billing code (ICD-9 or ICD-10), from 2014 to 2019. The class labels “HCM yes” and “HCM no” were assigned by manual review of medical records as well as the outcomes of the gold standard imaging tests for HCM diagnosis. The developed model performed well in finding HCM patients with an accuracy of 95%, recall of 99%, and precision of 97%. The F1 score was 98 %, while 4% of patients

were misclassified. This model will be translated into clinical practice for a clinical decision support system in EIMS to assist providers in the accurate diagnosis of HCM from echo data automatically while ensuring high-quality echo interpretation.

Keywords: Machine Learning, Echocardiography, Hypertrophic Cardiomyopathy, Random Forest, Decision Support System

1. INTRODUCTION

Hypertrophic cardiomyopathy is the most common heritable cardiomyopathy in the United States with an estimated prevalence of 1 in 500 [1-4] and the most frequent cause of sudden cardiac death (SCD) in the young. [5, 6] It is a genetic disease and it is extremely heterogeneous, and its phenotypes have been described extensively in the literature. Some experts suggest that HCM should be defined genetically and not morphologically and some others recommend a morphological classification for the identification of such patients.[7] HCM manifests phenotypically as left ventricular hypertrophy and is transmitted in an autosomal dominant manner. The relatively high prevalence of HCM in the general population (estimated to affect >700,000 Americans) contrasts sharply with less frequent recognition in clinical practice, inferring that many individuals remain undiagnosed throughout life.[5, 8]

Identification of HCM cases by manual review of electronic health record (EHR) data is time-consuming and

laborious. A prior study has shown that one-third of HCM patients are incorrectly classified by billing codes [4] underscoring the need for machine learning models for the identification of HCM utilizing echo data. The gold standard for the diagnosis of HCM is based on echo, which is an ultrasound-based cardiac imaging modality. To establish the diagnosis of HCM by echo an enormous amount of data is acquired and interpreted by providers in Echocardiography laboratories. Machine learning on echo data would enable efficient and accurate identification of HCM cohorts.

This paper proposes the use of historical echo measurements to train a machine learning model for the identification of HCM patients. We will first describe the dataset, data preparation, feature selection, and the manual labeling process. Then we describe how we adapted random forests in training and test to perform the prediction. And finally, we present the results and discussions.

2. MATERIALS AND METHODS

The objective of this paper is to develop a machine learning model applied to echo measurement data, to automatically identify HCM patients. The process of developing this model had four steps, from data gathering to practical prediction which is explained in the following sections. (Figure 1)

Study setting: Mayo Clinic is a tertiary referral center for patients with HCM with over 1,600 patients evaluated each year for this diagnosis across the enterprise. The Mayo Cardiovascular Ultrasound Imaging and Hemodynamic Laboratory (Echo Lab) is a large, clinical, educational, and research echocardiographic laboratory facility. The Echo Information Management System (EIMS) was developed by Mayo and was originally deployed in 2001. EIMS enables custom Mayo Clinic workflows and prompt generation of echo reports which include hundreds of cardiac measurements (numeric values) and standardized sentences (termed impressions which are in semi-structured format) describing various cardiac structures and specific diagnosis. Echo data is reviewed, interpreted, and reported using EIMS. Echo reports documenting HCM diagnosis are generated in EIMS and send to the electronic health record (EHR) in HL7 format. These data are also stored in the EIMS database which is replicated in a near real-time fashion in the institutional data warehouse generating an echo digital data registry. In this study, echo data was extracted from this digital registry.

Data gathering: Patients with at least one diagnostic billing code for HCM from 2014 to 2019 were identified in the Mayo Clinic data warehouse. Subsequently, echo report data of these patients were retrieved from the EIMS digital registry. These data consist of three groups: demographics, echo measurements, and impressions. There were 15 thousand echo reports of the adults with age over 18 years old with 22 demographic variables, 1,047 echo measurements which are numeric variables, and 3,655 impressions. 16% of these data

originated from the stress echo, which was excluded to focus on resting echo data.

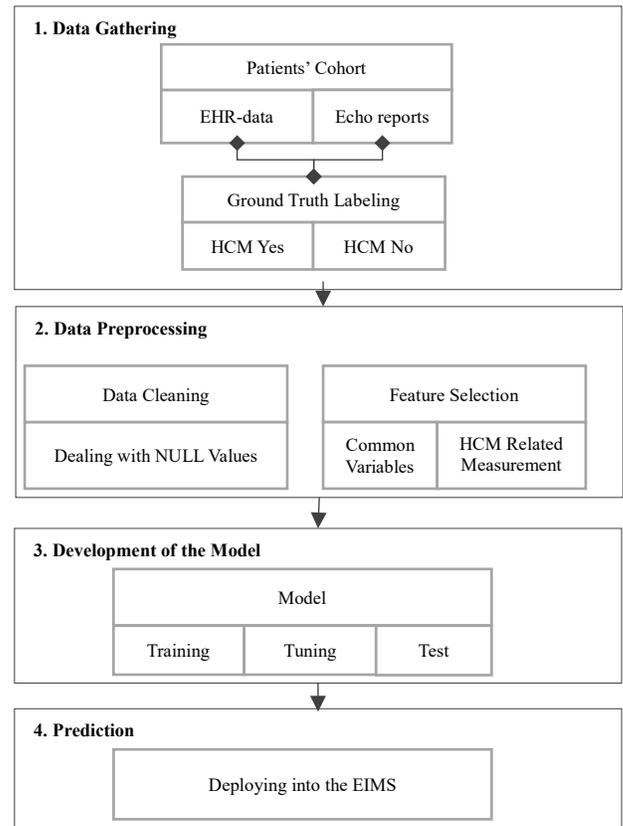


FIGURE 1: STEPS OF DEVELOPMENT AND DEPLOYING OF ECHO MACHINE LEARNING MODEL. ECHO RECORDS ARE LABELLED AS YES AND NO.

In the next step, each echo report was labeled as “HCM Yes” and “HCM No”. The ground truth for labeling “HCM Yes” was impressions indicating the diagnosis of hypertrophic cardiomyopathy, non-obstructive, or obstructive. The impression “no evidence of hypertrophic cardiomyopathy” was used for labeling “HCM No”. Of 3,548 patients in this cohort, 93% of patients were identified as “HCM Yes” and 7% were “HCM No”. The median patient age was 61 years with a mean BMI of 29 kg/m² and BSA of 2 m²; 43% were women and 57% men. The median age, BMI, and BSA in women were respectively 64, 28.89 kg/m² and 1.8 m²; and for men were 58, 29.64 kg/m², and 2.12 m².

3. MACHINE LEARNING MODEL

After data gathering and preparation, based on the literature review, the Random Forest model was used for the classification of patients with HCM. This model is an ensemble of multiple decision trees with a voting aggregation mechanism to produce the final classification prediction. According to the literature, the random forest has shown high accuracy in a variety of clinical use cases and practical applications. [9] Our

model uses demographic information of patients and echo measurements to train a binary classification of patients into “HCM yes” and “HCM no”. The model was trained on a training dataset and it was tuned for the best performing hyperparameter values (e.g. number of trees, depth of each tree, and the number of variables to be sampled).

Data preprocessing: From demographic variables, the most relevant ones, age at echo and sex, were selected. For the selection of Echo measurements, a cardiologist reviewed a list of 1,047 measurements and categorized them into three groups of “green”, “yellow”, and “red” which indicated if that measurement is highly related, partially related, and not related to HCM diagnosis. The red category measurements were excluded from the list of features whereas green and yellow were included. Correlated measurements were recognized from the correlation matrix and measurements with less dependency were selected. Finally, 53 measurements and 2 demographic variables were selected as input variables for the model.

Data preparation was subsequently performed. The major component of data preparation is breaking down the data sets into test and train data and also refining the data sets. This could include dealing with unbalanced data, removing duplicates, and dealing with Null values.

Model training: Data was divided into training (75%) and test (25%) sets. The null values were substituted by the median of each variable. The training data was balanced by oversampling before running the model. To run the random forest model, two parameters including the number of trees and the number of branches at each tree split were tuned. Then, the values of 100 trees and 7 splits for trees’ depth were chosen.

4. RESULTS

Random forest classification was performed on 53 measurements and 2 demographic variables (55 input variables) from the echo reports of patients in the study cohort. The confusion matrix of this model is shown in Table 1, for the classification of patients into “HCM yes” and “HCM no” employing echo report features.

TABLE 1: CONFUSION MATRIX OF THE MODEL

Actual \ Predicted	Predicted	
	HCM No	HCM Yes
HCM No	66	121
HCM Yes	132	3432

TABLE 2: MODEL PERFORMANCE METRICS

Metric	percentage
Accuracy	95 %
Sensitivity	97 %
Specificity	33 %
Detection Rate	0.9 %
Precision	97 %
Recall	99 %
F1 score	98 %

The model accuracy was 95% with a 95% confidence interval - from 94.75 % to 96.11 % (Table 2). Four percent of HCM patients in the test data were misclassified patients using this model. No Information Rate was 94.72 % and the p-value was 0.2×10^{-16} .

The overview of the most important variables in the model, mean decrease accuracy and mean decrease Gini index for the top thirty variables were calculated. These variables are demonstrated in Figure 2. In the following section, these variables will be discussed in more detail.

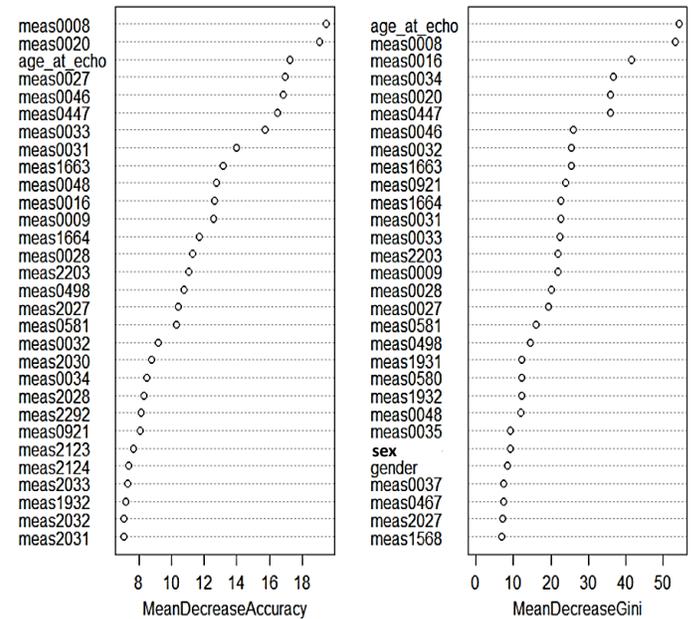


FIGURE 2: MEAN DECREASE OF ACCURACY AND MEAN DECREASE OF GINI INDEX OF THE TOP THIRTY MODEL VARIABLES

Most Important Variables: In our model, among demographic variables, age at echo was the most important demographic predictor and it had the highest increase in accuracy compared to all other variables while sex was not among the most important variables.

Among echo measurements, Interventricular Septum Diastolic Thickness by 2-D was the most important variable. The most important variables of the random forest model are shown in Table 3. Among the top thirty variables, approximately half of them were chosen from left ventricular variables. The next group of measurements was Doppler velocities across the mitral valve. Also, nine strain measurements were on this list.

TABLE 3: THE MOST IMPORTANT VARIABLES BASED ON MEAN INCREASE IN ACCURACY

#	Variable	Description
1	meas0008	Interventricular Septum Diastolic Thickness by 2-D
2	meas0020	LV Outflow Tract Systolic TVI by Pulsed Wave Doppler
3	age	Age at the time of echo
4	meas0027	Mitral Valve E-Wave Peak Velocity by Pulsed Wave Doppler
5	meas0046	Tricuspid Valve Regurgitant Systolic Peak Velocity by Continuous Wave Doppler
6	meas0447	LV EF MOD by 2-D Biplane Apical Views
7	meas0033	Mitral Valve Medial Annulus e' Velocity by Tissue Doppler Imaging
8	meas0031	Mitral Valve E to A Ratio at Baseline by Pulsed Wave Doppler
9	meas1663	LA Maximum Volume by 2-D Method of Disks Four Chamber
10	meas0048	Estimated Right Atrial Pressure
11	meas0016	LV Mass Index by 2-D
12	meas0009	LV Posterior Wall Diastolic Thickness by 2-D
13	meas1664	LA Maximum Volume by 2-D Method of Disks Two Chamber
14	meas0028	Mitral Valve A-Wave Peak Velocity by Pulsed Wave Doppler
15	meas2203	LA Maximum Volume by 2-D Method of Disks Biplane
16	meas0498	LV Outflow Tract Obstruction Systolic Peak Velocity by Continuous Wave Doppler
17	meas2027	LV Basal Averaged Peak Systolic Strain for SRI Analysis
18	meas0581	Mitral Valve Lateral Annulus e' Velocity by Tissue Doppler Imaging
19	meas0032	Mitral Valve Deceleration Time by Pulsed Wave Doppler
20	meas2030	Left Ventricular Mid Inferior Septum Peak Systolic Strain for SRI Analysis
21	meas0034	Mitral Valve Medial E to e' Ratio by Pulsed Wave Doppler
22	meas2028	LV Mid Anterior Peak Systolic Strain for SRI Analysis
23	meas2292	LV Apical Inferior Peak Systolic Strain for SRI Analysis
24	meas0921	Mitral Valve Lateral Annulus E to e' Ratio Diastolic Pulse Wave Doppler
25	meas2123	LV Apical Septum Peak Systolic Strain for SRI Analysis
26	meas2124	LV Apical Lateral Peak Systolic Strain for SRI Analysis
27	meas2033	LV Mid Infero-Lateral Peak Systolic Strain for SRI Analysis
28	meas1932	Mitral Valve Lateral Annulus Systolic Velocity by Tissue Doppler Imaging
29	meas2032	LV Mid Anterior Septum Peak Systolic Strain for SRI Analysis
30	meas2031	LV Mid Inferior Peak Systolic Strain for SRI Analysis

5. DISCUSSION

In this paper, the random forest model described herein accurately identified HCM cases from demographics and cardiac measurements obtained by the echo. This model could efficiently identify HCM cases for large EHR-based cohort studies and quality improvement projects. While the application of machine learning in cardiology is still in its infancy, future integration of these models into everyday clinical systems could greatly increase the efficiency of workflow processes. This integration would ensure high-quality echo interpretation and empowers practitioners to diagnose HCM.

Machine learning models for HCM diagnosis would be especially beneficial to junior cardiologists, sonographers, and cardiologists with limited experience in the interpretation of echo for the diagnosis of HCM. Real-time and on-demand machine learning models would also serve as a second set of eyes to remind sonographers and cardiologists to consider the possibility of HCM diagnosis. The model reported herein is the first step to realize this vision. [3, 10]

The datasets used for developing this model were unbalanced and oversampling was used to resolve this issue. Although this method was very effective in increasing the sensitivity and specificity of the model, it might imbalance the data in favor of one specific variable.

Another concern is about data imputation for null values. Unfortunately, the number of null values is high in echo reports (70% for the chosen echo variables) this is mostly related to clinical workflow processes. In clinical practice, echo measurements are driven by the presence of specific conditions detected by echo, which largely differ across individual HCM patients. For example, if a patient does not have mitral regurgitation by echo, mitral jet measurements are consequently not performed and mitral regurgitation measurements will have null values. In this model, null values were replaced with the median of variables in training datasets. Choosing other methods of imputation may change model performance.

Prior studies have shown BSA and sex as top predictors for HCM in comparison with other demographic variables [8]. Marian and Braunwald [6] also discussed the importance of ejection fraction on HCM patient identification. In our model, these variables also were on the list of the most important variables while age, BSA, BMI, and ejection fraction were more important than gender. Among the top 30 variables, there were measurements which are utilized for evaluation of cardiac wall thickness, left ventricular diastolic and systolic functions, right ventricular systolic pressure, left ventricular outflow tract obstruction, and mid-ventricular obstruction. These phenotypic characteristics are used in the thought process of cardiologists diagnosing HCM by echo. Hence these findings indicate the potential for translation of this model for clinical practice to assist providers diagnosing HCM by the echo.

6. CONCLUSIONS

HCM has been considered a silent killer among heart diseases. EHR-based cohorts of HCM patients are identified manually while machine learning models may assist providers in accurate and efficient identification of HCM patients. Mayo Clinic EIMS system and robust digital infrastructure are capable of running such models in real-time.

In this paper, different steps of developing a machine learning model based on historical data of echocardiogram reports are explained. Future deployment of this model on EIMS may benefit patients, sonographers, cardiologists, and referring clinicians by expediting the interpretation process while maintaining high-quality and accurate echo reports and promoting practice standardization.

REFERENCES

- [1] Centers for Disease Control and Prevention, "Underlying Cause of Death, 1999–2018," 2018.
- [2] M. Devitt, "ACCF and AHA Update Guidelines on the Diagnosis and Treatment of Hypertrophic Cardiomyopathy," *American Family Physician*, vol. 86, no. 7, pp. 694-697, 2012.
- [3] A. Ghorbani *et al.*, "Deep learning interpretation of echocardiograms," *NPJ digital medicine*, vol. 3, no. 1, pp. 1-10, 2020.
- [4] P. Magnusson, A. Palm, E. Branden, and S. Morner, "Misclassification of hypertrophic cardiomyopathy: validation of diagnostic codes," *Clin Epidemiol*, vol. 9, pp. 403-410, 2017, doi: 10.2147/CLEP.S139300.
- [5] Ashley EA and N. J., "Understanding the echocardiogram," in *Cardiology Explained*. London: Remedica, 2014, ch. 4.
- [6] A. J. Marian and E. Braunwald, "Hypertrophic cardiomyopathy: genetics, pathogenesis, clinical manifestations, diagnosis, and therapy," *Circulation research*, vol. 121, no. 7, pp. 749-770, 2017.
- [7] V. M. Parato *et al.*, "Echocardiographic diagnosis of the different phenotypes of hypertrophic cardiomyopathy," *Cardiovascular ultrasound*, vol. 14, no. 1, pp. 1-12, 2015.
- [8] R. Huurman *et al.*, "Effect of body surface area and gender on wall thickness thresholds in hypertrophic cardiomyopathy," *Netherlands Heart Journal*, vol. 28, no. 1, pp. 37-43, 2020.
- [9] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [10] N. Zanjirani Farahani, D. S. B. Sundaram, M. Enayati, S. P. Arunachalam, K. Pasupathy, and A. M. Arruda-Olson, "Explanatory Analysis of a Machine Learning Model to Identify Hypertrophic Cardiomyopathy Patients from EHR Using Diagnostic Codes," in *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2020: IEEE, pp. 1932-1937, doi: 10.1109/BIBM49941.2020.9313231.