

Predictive Models for Cardiovascular Health Based on Social Determinants of Health

Noreen Mayat¹

Abstract—How may we improve existing predictive metrics for cardiovascular health outcomes using machine learning?

I. INTRODUCTION

One of the questions people have been grappling with is whether or not to incorporate race (or other genetic markers), or SDOH correlated with race to predict cardiovascular health outcomes. Sometimes, incorporating race into these predictive equations for different health outcomes has adverse effects; for example, existing risk scores for kidney disease, such as eGFR, may disproportionately disadvantage Black communities, as Black patients with similar creatinine levels to white counterparts can be scored as having “healthier” kidneys, which may lead to under-treatment by health professionals. Such issues point to why using social determinants of health alongside race in these predictive models may allow for a more nuanced understanding of an individual’s health risks by considering the broader picture of their environment, socioeconomic status, BMI, mental health, and access to care, which may all also be critical factors in determining an individual’s health risks in addition to their race.

Our research goal is to build a machine learning model incorporating other variables, such as social determinants of health, not currently being leveraged in computing existing risk scores for cardiovascular health outcomes, and to analyze how such a machine learning model compares to existing risk scores in predicting cardiovascular health outcomes.

II. RELATED WORK

A. ASCVD Risk Scores

There are various existing risk scores used for predicting cardiovascular health outcomes: the first is a risk score known as the Framingham score, developed by the Framingham Heart Study [1]. Other risk scores for predicting cardiovascular health outcomes include the Systematic Coronary Risk Evaluation (SCORE) algorithm in Europe, the QRISK3 in England and Wales, and the risk score for atherosclerotic cardiovascular disease (ASCVD), developed by the American College of Cardiology/American Heart Association using pooled cohort equations (PCE): the ACC/AHA 2013 pooled cohort risk equation [1]. This is typically the ASCVD risk score referenced in most studies.

¹Noreen Mayat is a senior majoring in data science at Barnard College, Columbia University New York, NY. nm3224@barnard.edu

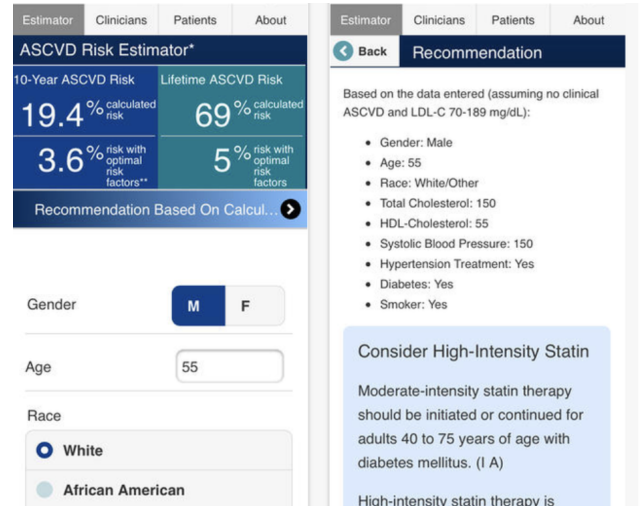


Fig. 1. Score demonstrating the variables used in the current ASCVD Risk Estimator Calculator using the ACC/AHA PCE to compute ASCVD risk scores for a 55 year-old white male [4].

PCEs leverage a combination of cohort studies in public health where they recruited patients from various demographics and followed their cardiovascular health for varying amounts of time, and “pool” the studies together to increase the diversity of the sample used to develop the metric. The current health variables leveraged to estimate ASCVD risk using the ACC/AHA PCE includes age, sex, race, total cholesterol, HDL cholesterol, systolic blood pressure, blood pressure lowering medication use, diabetes status, and smoking status [2]. According to a 2016 study, researchers found that in a large, multi-ethnic population, the ACC/AHA Pooled Cohort Risk Equation for ASCVD substantially overestimated actual 5-year risk in adults without diabetes, overall and across sociodemographic subgroups [3].

B. Social Determinants of Health

The World Health Organization defines social determinants of health to be the non-medical factors that influence health outcomes; they are the “conditions in which people are born, grow, work, live, and age, and the wider set of forces and systems shaping the conditions of daily life” [5]. Past research has demonstrated that integrating individual-level social determinants of health into electronic health records can assist in overall risk assessment models and in predicting holistic healthcare utilization and health outcomes, motivating efforts to collect and standardize patient-level social determinants of health information to benefit all types

of risk scores, not just cardiovascular risk scores [6].

Researchers have been exploring if incorporating social determinants of health into existing risk score equations for cardiovascular health outcomes improves the accuracy of these risk score predictions. A 2020 systematic review analyzing articles reporting on the use of machine learning models for cardiovascular disease prediction, which incorporated social determinants of health, found that most studies that compared performance with or without social determinants of health showed increased performance with them [7]. The most commonly included social determinants of health variables in these studies were gender, race/ethnicity, marital status, occupation, and income. The researchers note that there were a limited variety of sources and data in the reviewed studies, and thus, there is not as much research on how other social determinants of health variables, such as environmental ones, are known to impact cardiovascular disease risk, would impact model performance. Recording such data in electronic databases, as previous studies have also recommended, would enable their use. Further, a 2022 study found that adding social determinant of health risk factors alongside the existing variables used in ACC/AHA PCE to train a model in fact improved ASCVD risk prediction in specifically an African American cohort, a historically disadvantaged group in the healthcare system [8]. In this study, social determinants of health such as BMI, depression, weekly stress, insurance status, family income, and neighborhood violence were determined as the most important for prediction in this demographic, and were independently associated with 10-year ASCVD risk [8]. Other studies have found similar results in ASCVD risk prediction models leveraging social determinants of health such as education, income, and employment in addition to the existing variables used in PCEs for ASCVD risk prediction in both Black populations and non-Black female populations [9].

Other studies have investigated the effects of social determinants of health on ASCVD risk scores by establishing a baseline Social Disadvantage Score (SDS) and examining its relationship with atherosclerotic cardiovascular disease (ASCVD) and overall mortality, as well as its influence on the prediction of ASCVD risk scores. The SDS ranged from 0 to 4, and was calculated by tallying the following social factors: (1) household income less than the federal poverty level; (2) educational attainment less than a high school diploma; (3) single-living status; and (4) experience of lifetime discrimination. However, this study found that Although SDS is independently associated with incident ASCVD and all-cause mortality, it does not improve 10-year ASCVD risk prediction beyond pooled cohort equations [10]. This may mean some social determinants of health may improve ASCVD risk score prediction, but not all, so we must be careful which ones to include in future models based on this existing research.

We must note that previous studies have found that removing race from machine learning models predicting ASCVD risk scores did not improve model performance in any subgroup, while various studies have found including race

Figure 1

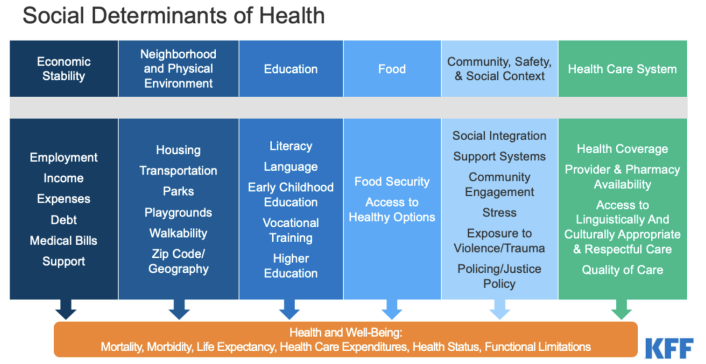


Fig. 2. Visualization of various social determinants of health and how they contribute to overall health outcomes, developed by KFF [11].

alongside other social determinants of health have improved model performance [12]. This information points to the idea that race is definitely still a significant variable in computing ASCVD risk, and should continue to be included in any future models predicting ASCVD risk, alongside other social determinants of health.

III. BACKGROUND

A. NHANES

The National Health and Nutrition Examination Survey (NHANES) is a downloadable public use data set used to document health care utilization, health status of various age groups, and related personal and lifestyle characteristics [13]. The data files are prepared and disseminated through the Centers for Disease Control and Prevention (CDC) to provide full access to data. More specifically, NHANES is a population-based survey designed to collect information on the health and nutrition of the US household population.

B. Machine Learning Models

Predictive machine learning models analyze datasets to predict a specific target variable: such datasets are composed of multiple data points, or samples, where each data point represents an entity we want to analyze [14]. Each of these entities has a list of various features associated with it: these features can be categorical (predefined values of no particular order like male and female), ordinal (predefined values that have an intrinsic order to them like a disease stage), or numerical (e.g., real values), and they are used to train the model and predict the target variable [14].

A model would analyze these feature variables and learn which variables are generally correlated/significant for predicting ASCVD using a training set, which is usually 75 percent of the data, and then makes predictions on a test set it hasn't seen before, which is usually 25 percent of the data. Different metrics are then used to evaluate model performance and accuracy, such as root mean square error, mean absolute percentage error, and r-squared value [14]. Popular machine learning models include KNN, Random Forest, and Decision Trees [15]. Past research in Taiwan

has used machine learning models with appropriate transfer learning as a tool for the development of cardiovascular risk prediction (ASCVD) models for Asian populations [16]. Transfer learning is the improvement of learning in a new task through the transfer of knowledge from a related task that has already been learned [17].

IV. DESIGN / IMPLEMENTATION / ALGORITHM

A. High-Level Overview

I am interested in leveraging the data in NHANES to incorporate different social determinant of health variables into existing risk score metrics, such as the ACC/AHA 2013 pooled cohort risk equation, to predict cardiovascular health outcomes. I want to pull all relevant data to social determinants health and data used in existing ASCVD risk score calculators, such as general health and demographic data, as well as data surrounding cardiovascular health outcomes.

For our research purposes, our target variable would be negative cardiovascular health outcomes, and our feature variables used to train the model would be various numerical and categorical values related to demographic data (age, gender, sex), general health data used to compute current ASCVD scores (total cholesterol, HDL cholesterol, systolic blood pressure, blood pressure lowering medication use, diabetes status, and smoking status), and data on various social determinants health, such as potential environmental data on air quality and neighborhood violence, income data, nutrition data and BMI, mental health data, drug data, and insurance data. Below is a visualization of different types of feature variables, both those used in to predict current ASCVD scores, and new features I plan implement, that will be used to train a machine learning model to successfully predict our target variable, CVD health outcomes, in adult patients.

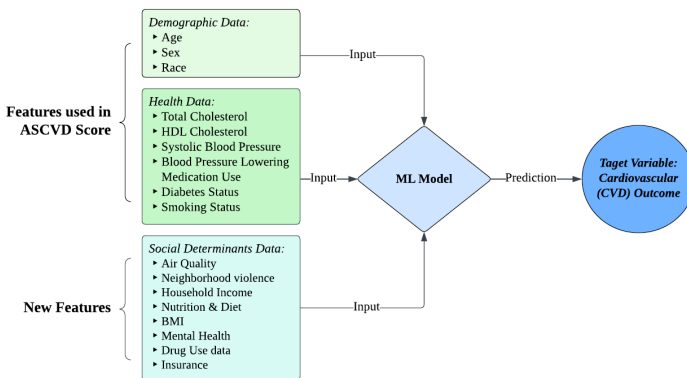


Fig. 3. Machine Learning Pipeline representing the feature variables to be used to train the model to predict CVD outcomes in adult patients.

After building a model and outputting predictions, I would want to compare the model's predictions of CVD health outcomes to what the ASCVD risk score predictions would be for adult individuals in my dataset, in order to observe which is more accurate (my model, or the risk score) in predicting cardiovascular health outcomes.

B. Machine Learning Pipeline

1) *Data Pre-Processing*: The first step in building a machine learning model is data pre-processing and preparation [18]. This usually begins with data cleaning, which can either remove or impute missing values, correcting errors, and removing outliers [18]. Later steps may also include data integration, which may require merging and joining various datasets; all of the data I need from NHANES cannot be downloaded in one file, it will require downloading various files. Each of these files will need to be cleaned and filtered to only retain the relevant data for my project, and joined into one, cleaned comprehensive dataset for training.

2) *Data Transformation*: Data transformation involves converting data into a format that is more appropriate for modeling by normalizing the data (scaling all numeric attributes in the dataset to a specific range, ex. converting them into a proportion or percentage) and transformation: for categorical values, this usually involves transforming and encoding categorical data into numerical data (0s and 1s, or 0s, 1s, 2s, 3s, etc.) using techniques like one-hot encoding [19]. Larger datasets may also be reduced to include less features using techniques such as dimensionality reduction (ex. PCA Principal Component Analysis) [19].

Finally, the dataset is split into training, validation, and test sets. Typically, about 70-75 percent of the data is used to train the model, and the remaining 25-20 percent is used to validate and test [19]. The training set is used to train the model, the validation set is used to tune the hyperparameters, and the test set is used to evaluate the model's performance [19].

3) *Model Training and Tuning*: Next, the training set is used to train the model. This involves feeding the the model our data, and allowing it to learn from the data by adjusting its parameters [19]. We evaluate the model's performance using the validation datasets to pick an optimal hyperparameter [19]. Evaluation metrics include accuracy, precision, recall, F1 score, and root mean squared error [19]. These metrics allow us to compare various hyper-parameters and adjust the model's hyperparameters to improve performance based on which hyper-parameters result in lower error scores and higher accuracy scores. This can be done manually or through automated processes like grid search or random search [19]. We will explain what these hyper-parameters are in more depth in our section outlining the chosen model for our project: Random Forest.

4) *Model Evaluation*: Lastly, these same metrics are used to evaluate the final model's performance and accuracy when ran on new data, the test set [19]. We will compare the model's predictions and accuracy to the computed ASCVD risk scores to see which predictive metric is better; for example, if our model successfully predicts a negative CVD outcome in a patient that was said to have a low ASCVD risk score, we can conclude our model is a better predictor for

CVD outcomes in this patient than the ASCVD risk score. This would have to be the case for more than half, or the majority, of our dataset for this conclusion to be true. We can evaluate how well our model compares to ASCVD score predictions for each demographic group as well to see if it performs better on some groups but worse than others and gain more insights.

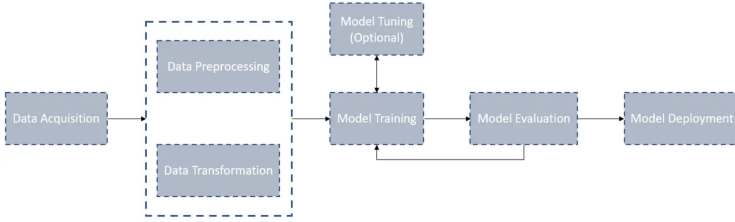


Fig. 4. Machine Learning Pipeline representing the various steps required to build a predictive machine learning model [18].

C. Random Forest

Random Forest is an ensemble machine learning method used for both classification and regression tasks [20]. It operates by constructing a random “forest” of decision trees during training, and outputting the class that is the average of the classes of the individual trees [20]. A decision tree resembles a flowchart, where the root node represents a sample row in the dataset containing feature variables, and each node in the tree leads to a different path based on the features to predict the final class [20]. Each of the child nodes of the root tree considers a different subset of the training data—this technique is known as bagging. Bagging results in a wide diversity that generally results in a better performing model.

After this initial bagging of the training data, the following levels of nodes consider different subsets of the training data’s *features* as well to predict a class [20]. This means that at each split in each tree, the algorithm is randomly selecting from the set of all features, from its respective subset of training data, and continuously splitting the data based on the “best” feature to predict a class [20]. For example, some trees may only consider age, BMI, and income data, while other trees might include an entirely different set of 3 features—if the tree finds that “age” is the most significant feature out of the subset of features it chose for predicting CVD, it will split on the age variable and choose another subset of 3 features (ex. smoking status, sex and cholesterol) to again choose which feature is most helpful in predicting the target variable, based on values such as information gain, and repeat the process [20].

In this example, the size of the feature subset considered at each split is fixed; this is where hyper-parameter tuning may come into play. We may find the model performs better when more than 3 features are considered at each split based on the evaluation metrics—in that case, our final model may consider not just 3, but 5 features at each split. Hyper-parameter tuning may also include altering the number of decision

trees used to train the model, and the minimum number of samples (subset of our training data) to be considered at each split. Each of these hyper-parameters have the potential to impact our final model’s performance, which is what makes validating our dataset and choosing optimal hyper-parameters for our final model so important.

Finally, each child node from the root outputs a final target outcome (in our case, a CVD outcome) for that row of data, and the model outputs “average” outcome as the final target variable for the row. Below is a visual representation of random forest and decision trees.

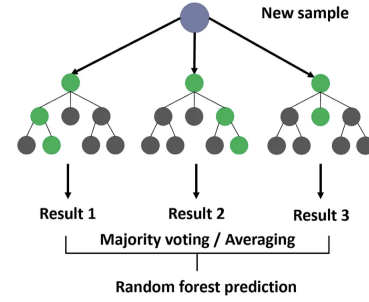


Fig. 5. Random Forest diagram representing how the model uses and averages various decision trees to output a prediction [18].

1) *Advantages of Random Forest:* The random forest model has various advantages: for example, it can handle both classification and regression tasks, as well as deal with large datasets with higher dimensionality and successfully estimate which variables are important in the classification by nature of the algorithm’s splitting process [20]. This will be useful for our dataset, which is relatively large and contains a lot of features.

2) *Disadvantages of Random Forest:* Some disadvantages to random forest are that it can be complex and computationally intensive, and may overfit to datasets, making it difficult to generalize on a new dataset [20]. This is particularly harmful when working with health data, where datasets may drastically differ among different patient types and demographics. This is important to consider in the scope of our project.

V. RESULTS

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi,

congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

VI. CONCLUSION AND FUTURE WORK

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

REFERENCES

- [1] P. L. Temporelli, "Risk scores, atherosclerotic cardiovascular disease and the crystal ball," in *European Journal of Preventive Cardiology*, vol. 28, no. 14, 2020, pp. 14–15.
- [2] "Ascvd risk calculator," in *American College of Cardiology and American Health Association*.
- [3] J. Rana, G. Tabada, and M. Solomon, "Accuracy of the atherosclerotic cardiovascular risk equation in a large contemporary, multiethnic population," in *J Am Coll Cardiol*, 2016, p. 2118–2130.
- [4] "Ascvd risk estimator," in *MDedge Federal Practitioner*, vol. 31, no. 5, 2014.
- [5] "Social determinants of health," in *World Health Organization*.
- [6] M. Chen, X. Tan, and R. Padman, "Social determinants of health in electronic health records and their impact on analysis and risk prediction: A systematic review," in *Journal of the American Medical Informatics Association*, vol. 27, no. 11, 2020, p. 1764–1773.
- [7] Y. Zhao, E. Wood, N. Mirin, S. Cook, and R. Chunara, "Social determinants in machine learning cardiovascular disease prediction models: A systematic review," in *American Journal of Preventive Medicine*, vol. 61, no. 4, 2021, pp. 596–605.
- [8] J. Hammond, R. Waken, M. Sims, K. Henderson, and K. J. Maddox, "The addition of social determinants of health improves the predictive accuracy of the pooled cohort equations for 10-year ascvd events in african americans," in *Epidemiology, Big Data and Precision Medicine: Social, Structural and Systemic Determinants of Cardiovascular Disease Risk*, 2022.
- [9] M. Xia, J. An, M. M. Safford, L. Colantonio, P. Muntner, K. Reynolds, A. E. Moran, and Y. Zhang, "Atherosclerotic cardiovascular disease risk associated with social determinants of health at individual and area levels," in *Epidemiology, Big Data and Precision Medicine: Neighborhood and Multi-Level Social Determinants of Health in Association with Cardiovascular Health*.
- [10] A. Hammoud, H. Chen, A. Ivanov, J. Yeboah, K. Nasir, M. Cainzos-Achirica, S. U. K. Alain Bertoni, M. Blaha, D. Herrington, and M. D. Shapiro, "Implications of social disadvantage score in cardiovascular outcomes and risk assessment: Findings from the multi-ethnic study of atherosclerosis," in *HomeCirculation: Cardiovascular Quality and Outcomes*, vol. 16, 2023.
- [11] S. Artiga and E. Hinton, "Beyond health care: The role of social determinants in promoting health and health equity," in *KFF*, 2018.
- [12] A. Ghosh, S. Venkatraman, and M. G. Nanna, "Risk prediction for atherosclerotic cardiovascular disease with and without race stratification," in *JAMA Cardiol*, vol. 9, no. 1, 2024, pp. 55–62.
- [13] J. A. Fain, "Nhanes: Use of a free public data set," in *European Journal of Preventive Cardiology*, vol. 43, no. 2, 2017.
- [14] S. Badillo, B. Banfai, F. Birzele, I. I. Davydov, L. Hutchinson, T. Kam-Thong, J. Siebourg-Polster, B. Steiert, and J. D. Zhang, "An introduction to machine learning," in *Clinical Pharmacology Therapeutics*, 2020.
- [15] D. Chicco, M. J. Warrens, and G. Jurman, "The coefficient of determination r-squared is more informative than smape, mae, mape, mse and rmse in regression analysis evaluation," in *PeerJ Computer Science*, 2021.
- [16] Y.-C. Hsiao, C.-Y. Kuo, F.-J. Lin, Y.-W. Wu, T.-H. Lin, H.-I. Yeh, J.-W. Chen, and C.-C. Wu, "Machine learning models for ascvd risk prediction in an asian population — how to validate the model is important," in *Acta Cardiol Sin*, vol. 39, no. 6, 2023.
- [17] L. Torrey and J. Shavlik, "Transfer learning," in *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*, 2010.
- [18] D. Hurley, "Simplify machine learning workflows," in *Medium, Towards Data Science*, 2020.
- [19] T. L. Wiemken and R. R. Kelley, "Machine learning in epidemiology and health outcomes research," in *Annu Rev Public Health*, vol. 41, no. 1, 2020, pp. 21–36.
- [20] A. Jehad, "Random forests and decision trees," in *International Journal of Computer Science Issues (IJCSI)*, vol. 9, no. 5, 2012, pp. 272–278.