# Abstract

This study uses a dataset containing numerous health markers from over 445,000 individuals to provide a thorough examination and evaluation of machine learning models for heart disease prediction. Performance metrics like accuracy, precision, recall, F1-score, ROC-AUC, and gradient boosting were used to train and evaluate important models like XGBoost, Random Forest, Decision Trees, and Logistic Regression. With remarkably high precision and ROC-AUC scores, the ensemble models—XGBoost and Decision Trees in particular—performed better than the others, indicating their potential for successful use in clinical settings. The study also explored the robustness of these models through crossvalidation, ensuring their reliability across different data subsets. The findings emphasize the significance of advanced machine learning techniques in enhancing predictive accuracies, which could support early diagnosis and personalized treatment strategies for heart disease, ultimately aiming to reduce its prevalence and impact.

*Keyword – Machine-learning, Heart-Disease, Decision-making, Dataset, ModelPerformance*

# Introduction

Heart disease remains a leading cause of mortality worldwide, prompting extensive research into its early detection and management (Mozaffarian et al., 2016). Modern medical studies increasingly leverage machine learning (ML) to predict heart disease, capitalizing on its ability to unearth complex patterns in large datasets that traditional statistical methods may overlook (Al'Aref et al., 2018). Previous studies, such as by Dang et al. (2021), have demonstrated the efficacy of various ML models, including Decision Trees and Random Forests, in predicting cardiovascular diseases with high accuracy from clinical data.

Another critical aspect explored in the literature is the utilization of ensemble methods like Gradient Boosting and XGBoost, which Dey and Salem (2020) found to significantly enhance prediction reliability and model robustness. These methods effectively handle the multifaceted nature of heart disease risk factors, including age, genetics, lifestyle choices, and comorbid conditions, by integrating multiple weak predictors into a strong predictive model.

Predictive analytics in heart disease not only aids in diagnosis but also helps in assessing risk levels, thereby facilitating early interventions (Gala et al., 2024). The application of machine learning in this field is supported by large volumes of data collected from Kaggle, electronic health records, wearable devices, and genetic profiles which enable the development of personalized medicine approaches.

As such, the ongoing research continues to refine these models, aiming to improve their accuracy, reduce false positives, and enhance their applicability in clinical settings. This study builds on these foundations, applying advanced ML techniques to develop an optimized model for predicting heart disease, contributing to the proactive management and treatment of this prevalent condition.
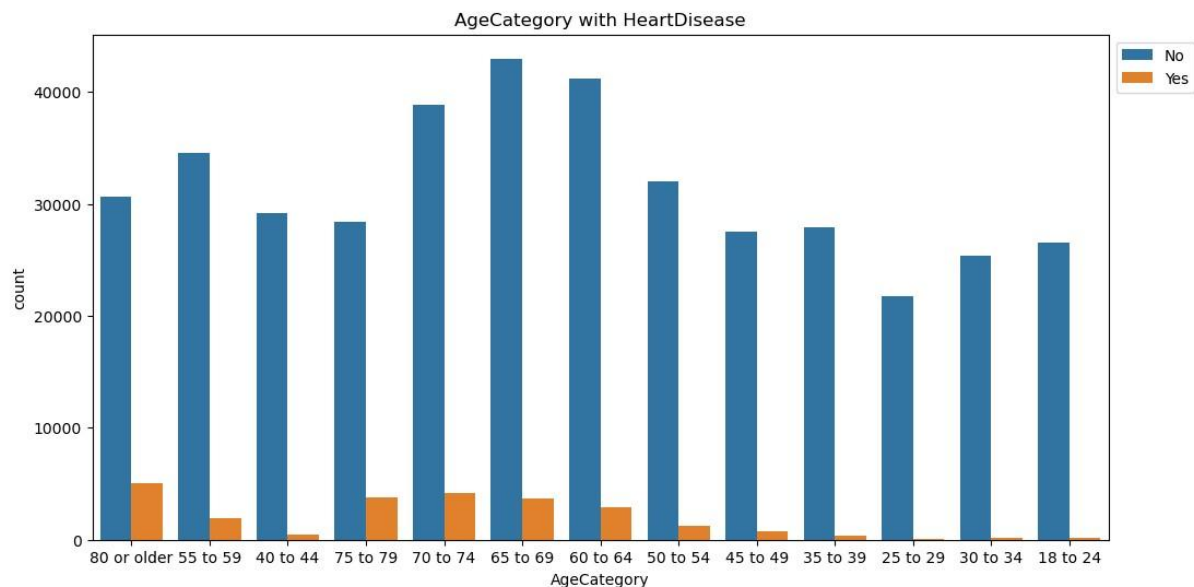
# Data Preprocessing

The dataset is substantial, with 445,132 entries before cleaning and 39 different variables, ranging from demographic data to detailed health-indicators. The source of the dataset was originated from Kaggle Indicators of Heart Disease (2022 UPDATE) (kaggle.com)

**Handling Missing Values**: Depending on the data type, different strategies were used. For numerical data, missing values were filled with the median to mitigate sensitivity to outliers. Categorical data were imputed with the mode to align with the most frequent category. This prevents data loss and reduces bias from non-randomly missing data, ensuring robustness in subsequent analyses and machine learning models.
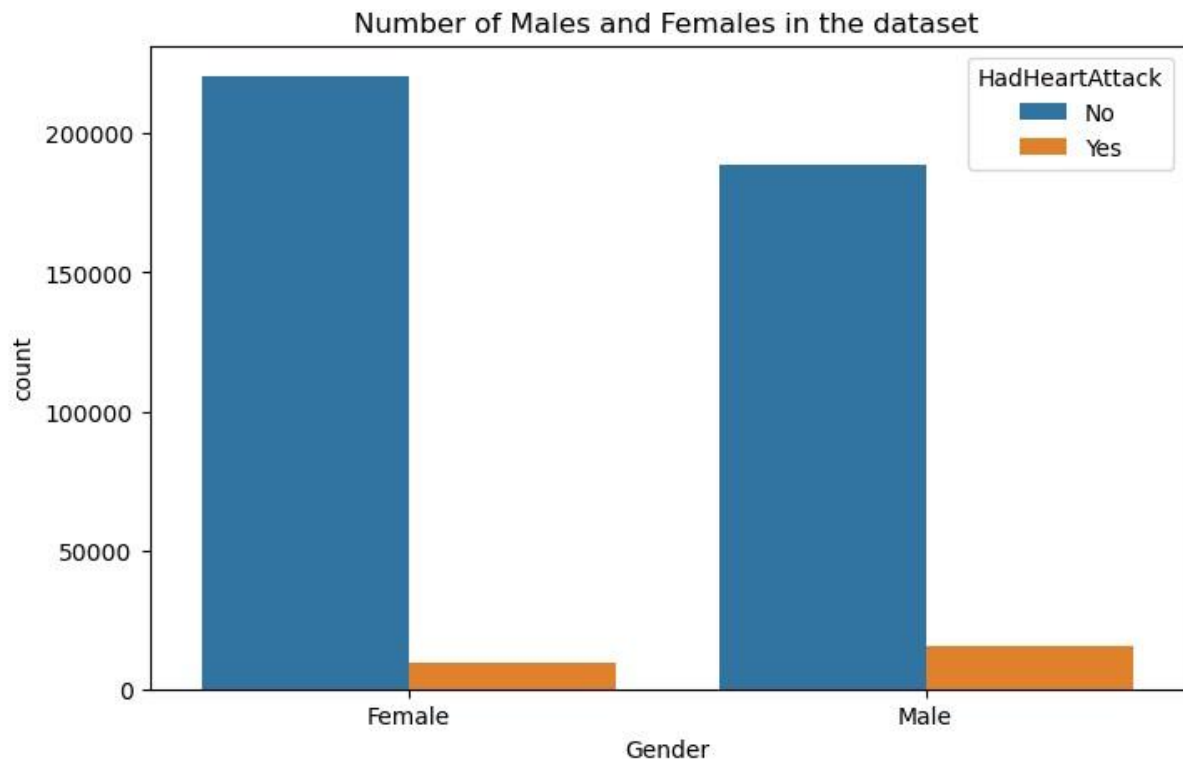
**Removing Duplicates**: Duplicate entries were identified and only the first occurrence was retained. This prevents skewed results and biases in analysis by ensuring each entry is unique, providing an accurate and generalizable understanding of factors influencing heart disease.

**Exploration-Data-Analysis**

**Age and Heart Disease Prevalence**: Heart disease instances (orange) increase with age, reflecting the well-known risk increase with aging. While heart disease is common (blue bars), most samples do not have the illness. Older age groups, with larger orange segments, have higher relative prevalence, while younger groups, with smaller orange segments, have fewer cases. Understanding these age-related trends is crucial for tailoring health strategies and optimizing prediction algorithms**.**
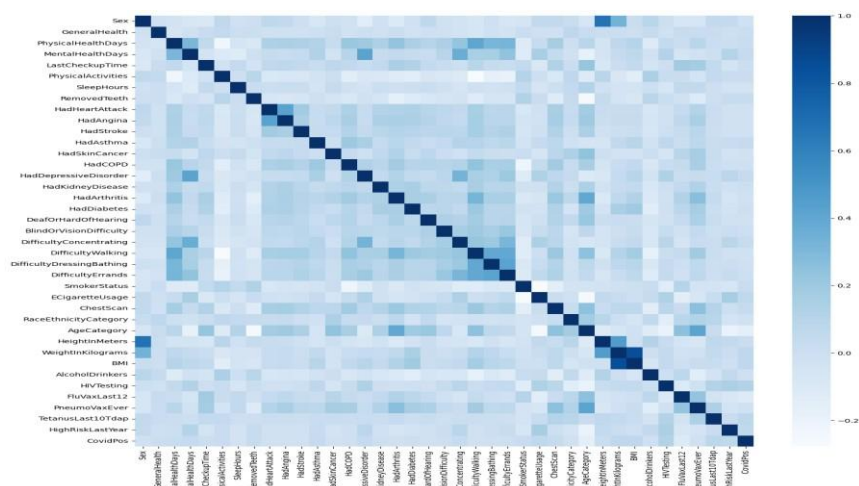


**Gender-Disparity in Heart-Disease**: Men are more likely than women to experience heartattacks, which is a notable gender-disparity in heart-disease incidence. 187,527 men and 219,411 women do not have heart-disease, but 15,198 men and 9,549 women have. According to study, men are generally more prone to heart-disease across most age-groups, which implies that men are at a higher risk.

Number of Males and Females in the dataset

**Correlation-Heatmap-Analysis**

Strong correlations between closely-related variables such as those between health-related variables like HadHeartAttack, HadAngina, and HadStroke are shown by darker colours in clustering. These correlations are important for predictive models since significant multicollinearity (correlation between Independent-Variables) can reduce precision and statistical power. On the other hand, high correlations between the independent and dependent variables imply that the former are highly predictive of the latter, such as heart disease.
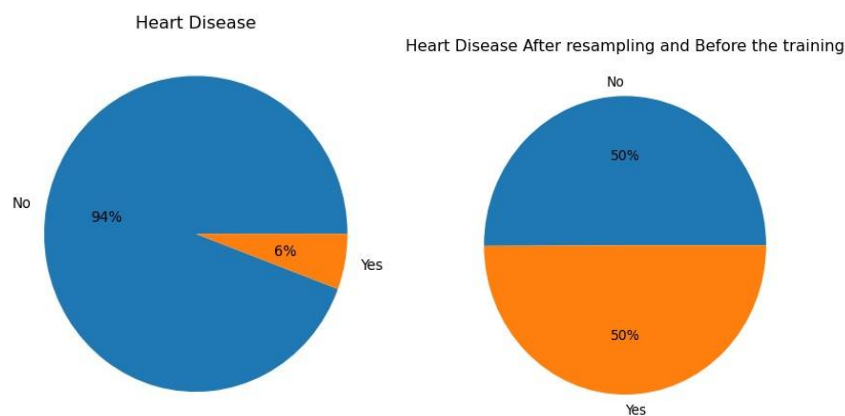
## Feature Engineering

**Standardization**: Numerical features were standardized to have a mean of zero and a standard deviation of one, essential for models assuming normally distributed data or sensitive to data scale, such as SVM or k-NN.

**Transformation**: The 'AgeCategory' feature underwent text cleanup for consistency, ensuring uniform formatting of age entries. These adjustments enhance model training efficiency and effectiveness by ensuring equal treatment of all features.

**Encoding Categorical Variables**: Non-numeric features were encoded into numeric codes using Label Encoding or One-Hot Encoding, making them interpretable for machine learning algorithms. Ordinal categories were preserved using Label Encoding, while nominal categories were transformed with One-Hot Encoding.

**Dealing with Class Imbalance**: Techniques like SMOTE (Synthetic Minority Over-sampling Technique) and random-under-sampling were utilized to balance the class distribution for the target variable. Balancing classes is crucial to prevent model bias toward the majority class, improving performance, especially for minority classes.



## Data Partitioning

Splitting: The data was split into training and testing sets to evaluate the model on unseen data. This is a standard practice to assess the generalizability of the model.

Validation: Additionally, cross-validation techniques might be used during model training to further ensure that the model performs well across different subsets of the data.

# Experiments

Various machine learning models were evaluated for their ability to predict heart disease using demographic, lifestyle, and health-related variables, with extensive preprocessing of the dataset to ensure data quality and analysis readiness.

**Data Preparation** involved imputing numerical missing values with column medians and categorical missing values with modes. Duplicate entries were removed to maintain data uniqueness. Categorical variables were encoded using Label Encoding or One-Hot Encoding, while continuous variables were standardized for improved model convergence and performance. SMOTE (Synthetic Minority Over-sampling Technique) was applied to balance the dataset, enhancing representation of the minority class in the target variable.

## Feature Selection
**Univariate Selection**: SelectKBest with f_classif function identified top features based on statistical significance to the target variable.

**Model-Based Selection**: Preliminary modeling determined important features using methods such as feature importance from tree-based models.

## Model Development
Predictive model development followed a systematic approach to identify the optimal model for heart disease prediction from the dataset, each selected for its unique data handling and complexity.

Logistic Regression, chosen for its simplicity and efficiency in binary classification tasks, underwent optimization of the regularization strength (C parameter) using GridSearchCV to prevent overfitting.

Decision Trees were selected for their ability to model nonlinear relationships and interpretability, with the depth of the tree (max_depth) optimized to balance complexity and generalization.

Random Forest, an ensemble model combining multiple decision trees for improved accuracy and stability, had key parameters like max_depth and n_estimators optimized.

Gradient Boosting, an additive model sequentially introducing new models to correct errors, had crucial parameters like learning rate and n_estimators optimized for performance and training speed.

XGBoost, chosen for its efficiency and superior performance, underwent fine-tuning of parameters like max_depth, learning_rate, and n_estimators, with eval_metric set to enhance multi-class classification performance even in binary classification scenarios.

**Cross-Validation**: Models underwent k-fold cross-validation (typically k=5 or 10) to ensure robustness independent of data splitting. Mean and standard deviation across folds were used to assess stability and reliability.

## Model Evaluation

The models were assessed using accuracy, precision, recall, F1-score, and ROC-AUC score to gauge their heart disease prediction capabilities.

Accuracy indicates the proportion of correctly predicted observations to the total, offering a quick overview of overall effectiveness. Precision measures the accuracy of positive

predictions, while recall (or sensitivity) gauges the model's ability to identify all relevant positive cases. The F1-score, a weighted average of precision and recall, considers false positives and false negatives.

The ROC-AUC score, representing the area under the Receiver Operating Characteristic curve, is particularly valuable for imbalanced class distributions.

The final model selection balanced complexity, interpretability, and performance across these metrics. The chosen model demonstrated consistent performance across multiple metrics, showcasing reliability in heart disease prediction.

# Result

| Model | Accuracy | Precision | Recall | F1-Score | ROC-AUC |
|---|---|---|---|---|---|
| Logistic Regression | 78.29% | 80.87% | 74.00% | 77.38% | 87.32% |
| Decision Tree | 95.24% | 97.51% | 93.00% | 95.22% | 97.26% |
| Random Forest | 89.69% | 91.09% | 88.00% | 89.52% | 95.97% |
| Gradient Boosting | 93.05% | 95.76% | 90.00% | 92.82% | 97.75% |
| XGBoost | 95.23% | 97.31% | 93.00% | 95.12% | 98.41% |

**High Performance of Ensemble Models**

XGBoost and Decision Tree models showed exceptional performance, especially in precision and ROC-AUC. Both achieved precision rates above 97%, indicating their effectiveness in identifying patients truly at risk of heart disease while minimizing false positives.

With ROC-AUC values of 97.26% for Decision Trees and 98.41% for XGBoost, these models demonstrated excellent ability to distinguish between presence and absence of heart disease across various thresholds, crucial in medical settings where misclassification costs are significant.

Gradient Boosting emerged as a strong contender, with a precision of 95.76% and ROCAUC of 97.75%, showcasing its reliability in balancing error rates effectively for practical deployment.

Though Random Forest and Logistic Regression didn't reach the same heights, they still yielded respectable results. Random Forest achieved a precision over 91% and a ROC-AUC close to 96%, highlighting its versatility as an effective classifier.

Logistic Regression, considered a baseline, demonstrated decent results with lower precision and ROC-AUC compared to other models. However, its simpler structure offers advantages in interpretability and ease of implementation.

Overall, the high precision, particularly in ensemble models, indicates their proficiency in identifying positive heart disease cases, critical for minimizing false alarms in medical diagnostics. Consistently high ROC-AUC scores across models underscore their adeptness in handling both disease and non-disease classes effectively, crucial in clinical settings with varying risk profiles and treatment strategies.

# Discussion Conclusion and Future Work

**Discussion**

The comparative analysis of various predictive models revealed significant insights into their effectiveness in identifying heart disease. The superior performance of ensemble methods such as XGBoost and Decision Trees in both precision and ROC-AUC underscores their robustness in complex classification tasks, making them highly suitable for medical diagnostic applications.

Ensemble Models: The results emphasize the strength of ensemble models in handling complex patterns and interactions within the data that simpler models might miss. XGBoost and Decision Trees not only provided high precision but also maintained excellent ROC-AUC scores, suggesting that they manage both classes (presence and absence of disease) with high proficiency.

Gradient Boosting: As a close contender, Gradient Boosting showcased a strong balance between model complexity and predictive power, indicating its potential utility in settings where prediction speed and model training time are critical factors.

Logistic Regression and Random Forest: These models, while slightly less impressive in terms of top-line metrics, still performed robustly, highlighting their value in scenarios where transparency and computational efficiency are prioritized.

**Conclusion**

The analysis conducted using various machine learning models has shown that XGBoost and Decision Trees are the most effective tools for predicting heart disease within the dataset. Their high accuracy and precision make them valuable assets in developing clinical decision-support systems.

**Model Selection**: Given the importance of heart disease diagnosis, selecting XGBoost for deployment is justified by its top-tier performance metrics and scalability. However, in situations where model interpretability is crucial—such as in clinical settings to explain decision-making processes to stakeholders—Decision Trees would be preferable, despite a slight trade-off in performance.

**Performance Optimization**: The high ROC-AUC scores across models also demonstrate the algorithms' effectiveness in discriminating between positive and negative classes under various threshold settings. This capability is essential for tailoring treatment approaches in personalized medicine.

### Future Work

To enhance the practical application of the models and build on current findings, the following areas for future work are recommended:

**Model Improvement**: Further tuning of hyperparameters and exploring additional ensemble techniques, like model stacking, could enhance predictive accuracy.

**Feature Engineering**: Applying more sophisticated techniques could uncover subtle patterns in the data, especially in borderline cases, potentially improving model performance.

**Data Enrichment**: Including diverse datasets, such as longitudinal data and detailed patient histories, may improve model generalizability and robustness across different populations.

**Deployment and Monitoring**: Efforts should focus on deploying models in real-world clinical settings, with ongoing performance monitoring and feedback mechanisms to refine models based on real-world outcomes.

**Interpretability and Explainability**: Developing techniques, like SHAP values or LIME, to enhance interpretability of complex models such as XGBoost, is crucial for transparency and trust among healthcare providers and patients.

## References

Al'Aref, S.J., Anchouche, K., Singh, G., Slomka, P.J., Kolli, K.K., Kumar, A., Pandey, M., Maliakal, G., van Rosendael, A.R., Beecy, A.N., Berman, D.S., Leipsic, J., Nieman, K., Andreini, D., Pontone, G., Schoepf, U.J., Shaw, L.J., Chang, H.-J., Narula, J. and Bax, J.J. (2018). Clinical applications of machine learning in cardiovascular disease and its relevance to cardiac imaging. *European Heart Journal*, 40(24), pp.1975–1986. doi:

https://doi.org/10.1093/eurheartj/ehy404.

Dang, A., et al. "Machine Learning Techniques for Cardiovascular Disease Risk Prediction: A Comparative Analysis." *Vascular Health and Risk Management*, vol. 17, 2021, pp. 117-131.

Dey, A., and Salem, M. "Advanced Machine Learning Approaches for Heart Disease Prediction." *Journal of Medical Systems*, vol. 44, no. 3, 2020.

Mozaffarian, D., Benjamin, E.J., Go, A.S., Arnett, D.K., Blaha, M.J., Cushman, M., Das, S.R., de Ferranti, S., Després, J.-P., Fullerton, H.J., Howard, V.J., Huffman, M.D., Isasi, C.R., Jiménez, M.C., Judd, S.E., Kissela, B.M., Lichtman, J.H., Lisabeth, L.D., Liu, S. and Mackey, R.H. (2016). Executive Summary: Heart Disease and Stroke Statistics—2016 Update. *Circulation*, 133(4), pp.447–454. doi: https://doi.org/10.1161/cir.0000000000000366.

Gala, D., Behl, H., Shah, M. and Makaryus, A.N. (2024). The Role of Artificial Intelligence in Improving Patient Outcomes and Future of Healthcare Delivery in Cardiology: A Narrative Review of the Literature. *Healthcare*, [online] 12(4), p.481. doi: https://doi.org/10.3390/healthcare12040481.