

January-March 2025



Social Sciences & Humanity Research Review



Heart Disease Risk PredictionStatistical Analysis and Classification of Heart Diseases Risk Using Clinical Parameter

Bushra Shehzadi¹, Hafiz Abdul Sami², Dr. Shahbaz Nawaz³, Dr. Anam Javaid⁴, Tooba Nihal⁵, Fatima Bibi⁶

¹Department of Statistics ,The Women University Multan, Pakistan ²Bureau Of Statistics , Planning & Development Board, Govt of the Punjab, Pakistan ³Bureau Of Statistics, Planning & Development Board, Govt of the Punjab, Pakistan ⁴ Assistant Professor , Department of Statistics , The Women University Multan, Pakistan

⁵Lecturer, Department of Statistics, The Women University Multan, Pakistan. ⁶Department of Statistics, The Women University Multan, Pakistan

ARTICLE INFO

Keywords:Heart Disease Risk, Statistical Modeling, Logistic Regression, Machine Learning, Random Forest, Neural Networks

Corresponding Author: Bushra Shehzadi,

Department of Statistics, The Women University Multan, Pakistan

ABSTRACT

Heart disease remains a leading cause of morbidity and mortality worldwide, necessitating robust statistical methods for early detection and risk prediction. This study applies multiple statistical and machine learning techniques, including Logistic Regression, Random Forest, and Neural Networks, to analyze a clinical dataset of 1,025 observations with 14 variables related to demographic, physiological, biochemical parameters. The study evaluates model performance using accuracy, sensitivity, specificity, and AUC metrics, while also assessing multicollinearity, correlation structures, and variable significance through standardized coefficients and information criteria (AIC/BIC). The Logistic Regression model achieved an AUC of 0.83, indicating strong predictive capability, whereas ensemble methods required parameter tuning to improve specificity. The results highlight key risk factors including chest pain type, ST depression, maximum heart rate, and number of major vessels, offering data-driven insights for preventive healthcare strategies.

1. Introduction

Cardiovascular diseases (CVDs), particularly heart disease, are among the leading causes of morbidity and mortality worldwide. According to the World Health Organization (WHO), an estimated 17.9 million people die each year from CVDs, accounting for nearly one-third of global deaths. Early detection and accurate classification of individuals at risk are critical to reducing the burden of these diseases and improving population health outcomes.

Statistical modeling and machine learning techniques provide powerful tools to analyze complex relationships between clinical parameters and disease outcomes. By examining demographic, physiological, and biochemical variables simultaneously, researchers can identify key risk factors and develop predictive models for early diagnosis.

In recent years, techniques such as logistic regression, Random Forest, and Neural Networks have been increasingly applied in medical research to improve diagnostic accuracy and support clinical decision-making. However, there is a need for comparative studies evaluating the performance and interpretability of these methods in heart disease risk prediction.

This study aims to bridge this gap by performing a comprehensive statistical analysis and classification of heart disease risk using multiple clinical parameters. The findings offer insights into the most significant predictors and provide recommendations for integrating statistical and machine learning approaches into preventive healthcare strategies.

2. Literature Review

Previous studies have explored risk factors such as hypertension, cholesterol levels, diabetes, and behavioral aspects using regression and survival models (Poulter, 1999; Haffner, 2000). Machine learning approaches, including Random Forests and Neural Networks, have shown promise in capturing non-linear relationships between variables (Carney & Freedland, 2017). However, comparisons between classical statistical models and modern machine learning algorithms remain limited in heart disease prediction contexts.

3. Methodology

The dataset, sourced from the Kaggle Heart Disease repository, comprises 1,025 observations across 14 variables including age, sex, chest pain type, resting blood pressure, cholesterol, fasting blood sugar, resting ECG, maximum heart rate, exercise-induced angina, ST depression, slope of ST segment, number of major vessels, and thalassemia type. Descriptive statistics summarize variable distributions, followed by correlation analysis and multicollinearity assessment using Variance Inflation Factor (VIF). Predictive models include Logistic Regression, Random Forest, and Neural Networks. Model selection relies on accuracy, AUC, sensitivity, specificity, and information criteria (AIC/BIC). The best-performing model was selected based on higher AUC and balanced sensitivity-specificity performance.

4. Results and Discussion

Descriptive analysis revealed a mean age of 54.4 years with 51.3% of participants diagnosed with heart disease. Correlation analysis identified chest pain type, maximum heart rate, and slope as positively associated with heart disease, while exercise-induced angina and ST depression exhibited negative correlations. Logistic Regression yielded significant predictors (p<0.05) including chest pain, ST depression, number of vessels, and maximum heart rate, achieving an AUC of 0.83. Random Forest models demonstrated high sensitivity but low specificity, suggesting overfitting risks without parameter tuning. Neural Networks underperformed due to limited training data and parameter optimization needs.

Descriptive statistics of Various Factor

Variabl e	Observation s	with missin	t	Maximu m	Mean	Std. deviatio n
		g data	missin			

			g data				
target	1025	0	1025	0.000	1.000	0.513	0.500
age	1025	0	1025	29.000	77.000	54.434	9.072
sex	1025	0	1025	0.000	1.000	0.696	0.460
ср	1025	0	1025	0.000	3.000	0.942	1.030
trestbps	1025	0	1025	94.000	200.000	131.61	17.517
chol	1025	0	1025	126.000	564.000	246.00 0	51.593
fbs	1025	0	1025	0.000	1.000	0.149	0.357
restecg	1025	0	1025	0.000	2.000	0.530	0.528
thalach	1025	0	1025	71.000	202.000	149.11 4	23.006
exang	1025	0	1025	0.000	1.000	0.337	0.473
oldpeak	1025	0	1025	0.000	6.200	1.072	1.175
slope	1025	0	1025	0.000	2.000	1.385	0.618
ca	1025	0	1025	0.000	4.000	0.754	1.031

Correlation matrix

	age	sex	ср	trest bps	cho l	fbs	rest ecg	thal ach	exa ng	oldp eak	slo pe	ca	thal	targ et
age	1	0.1 03	- 0.0 72	0.27	0.2	0.1 21	0.13	- 0.39 0	0.0	0.20	- 0.1 69	0.2 72	0.0 72	- 0.2 29
sex	0.1 03	1	- 0.0 41	- 0.07 9	- 0.1 98	0.0 27	- 0.05 5	- 0.04 9	0.1 39	0.08	- 0.0 27	0.1 12	0.1 98	- 0.2 80
ср	- 0.0 72	- 0.0 41	1	0.03	- 0.0 82	0.0 79	0.04	0.30	- 0.4 02	- 0.17 5	0.1 32	- 0.1 76	0.1	0.4

trest bps	0.2 71	- 0.0 79	0.0 38	1	0.1 28	0.1 82	- 0.12 4	- 0.03 9	0.0 61	0.18 7	- 0.1 20	0.1 05	0.0 59	- 0.1 39
	0.2	- 0.1	- 0.0	0.12		0.0	- 0.14	- 0.02	0.0	0.06	- 0.0	0.0	0.1	- 0.1
chol	20	98	82	8	1	27	7	2	67	5	14	74	00	00
fbs	0.1 21	0.0 27	0.0 79	0.18 2	0.0 27	1	- 0.10 4	- 0.00 9	0.0 49	0.01 1	- 0.0 62	0.1 37	- 0.0 42	- 0.0 41
reste cg	- 0.1 33	- 0.0 55	0.0 44	- 0.12 4	- 0.1 47	- 0.1 04	1	0.04	- 0.0 66	- 0.05 0	0.0 86	- 0.0 78	- 0.0 21	0.1 34
thala	0.3	- 0.0	0.3	- 0.03	- 0.0	- 0.0	0.04		0.3	- 0.35	0.3	- 0.2	- 0.0	0.4
ch	90	49	07	9	22	09	8	1	80	0	95	08	98	23
exan g	0.0 88	0.1 39	- 0.4 02	0.06 1	0.0 67	0.0 49	- 0.06 6	- 0.38 0	1	0.31 1	- 0.2 67	0.1 08	0.1 97	- 0.4 38
oldp eak	0.2 08	0.0 85	- 0.1 75	0.18 7	0.0 65	0.0 11	- 0.05 0	- 0.35 0	0.3 11	1	- 0.5 75	0.2 22	0.2 03	- 0.4 38
slop	- 0.1	- 0.0	0.1	- 0.12	- 0.0	- 0.0	0.08	0.39	0.2	- 0.57		- 0.0	- 0.0	0.3
e	69	27	32	0	14	62	6	5	67	5	1	73	94	46
ca	0.2 72	0.1 12	- 0.1 76	0.10 5	0.0 74	0.1 37	- 0.07 8	- 0.20 8	0.1 08	0.22	0.0 73	1	0.1 49	0.3 82
thal	0.0 72	0.1 98	- 0.1 63	0.05 9	0.1 00	- 0.0 42	- 0.02 1	- 0.09 8	0.1 97	0.20	- 0.0 94	0.1 49	1	- 0.3 38

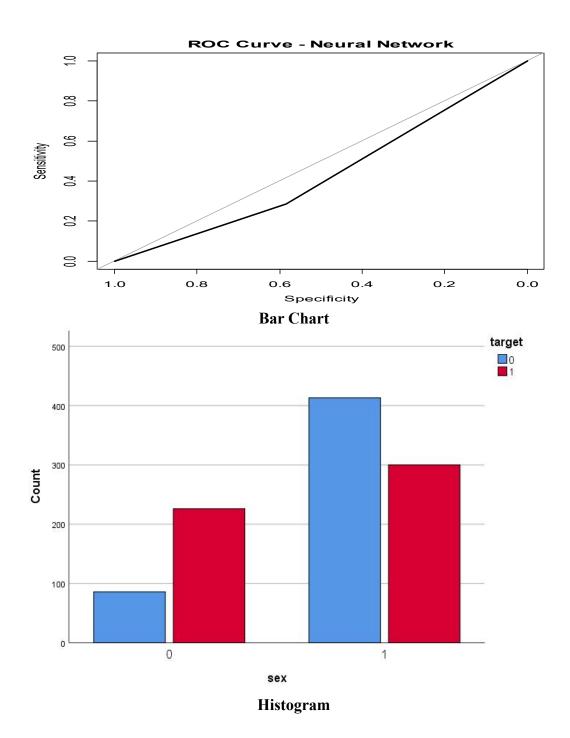
targe t	- 0.2 29	- 0.2 80	0.4 35	- 0.13 9	- 0.1 00	- 0.0 41	0.13	0.42	- 0.4 38	- 0.43 8	0.3 46	- 0.3 82	- 0.3 38	1

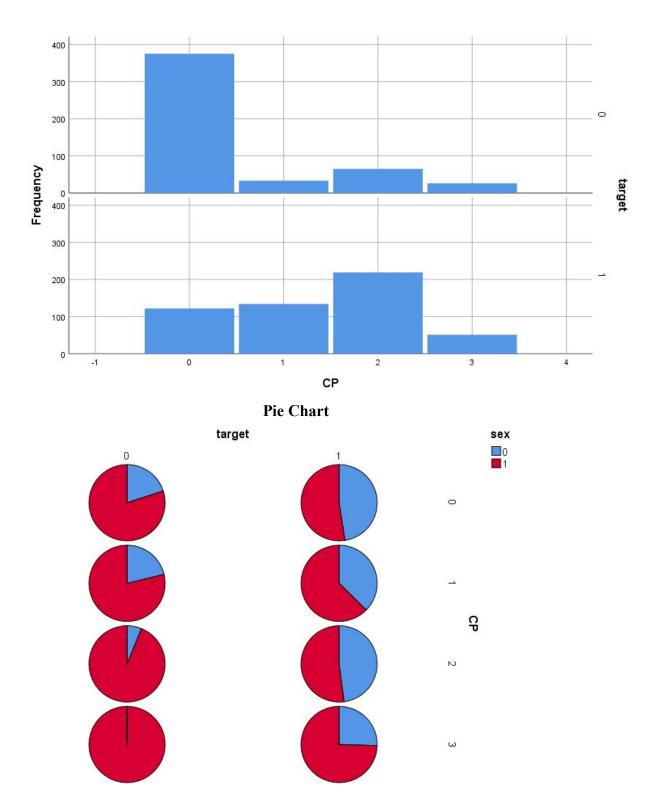
Multicollinearity statistics

	ag e	sex	ср	trest bps	ch ol	fbs	rest ecg			oldp eak		ca	th al
Toler ance	0.7 00	0.8 65	0.7 73	0.85	0.8 73	0.9 17	0.93	0.61 9	0.7 05	0.58 5	0.6 09	0.8 35	0.8 79
VIF	1.4 29	1.1 56	1.2 93	1.16 8	1.1 46	1.0 90	1.06 4	1.61 5	1.4 19	1.70 9	1.6 43	1.1 97	1.1 38

Standardised coefficient

Source	Value	Standard error	Chi- $Pr > Chi^2$		Wald Lower bound (95%)	Wald Upper bound (95%)
age	-0.041	0.063	0.423	0.516	-0.165	0.083
sex	-0.468	0.065	51.797	< 0.0001	-0.596	-0.341
ср	0.485	0.057	72.521	< 0.0001	0.373	0.596
trestbps	-0.176	0.054	10.529	0.001	-0.282	-0.070
chol	-0.161	0.058	7.603	0.006	-0.276	-0.047
fbs	-0.020	0.056	0.126	0.723	-0.130	0.090
restecg	0.120	0.055	4.781	0.029	0.012	0.228
thalach	0.300	0.072	17.292	< 0.0001	0.158	0.441
exang	-0.258	0.058	19.514	< 0.0001	-0.373	-0.144
oldpeak	-0.370	0.075	24.206	< 0.0001	-0.517	-0.222
slope	0.182	0.064	8.012	0.005	0.056	0.308
ca	-0.429	0.059	53.603	< 0.0001	-0.543	-0.314
thal	-0.303	0.053	32.415	< 0.0001	-0.407	-0.199





5. Conclusion and Recommendations

Statistical and machine learning techniques offer complementary strengths in heart disease risk prediction. Logistic Regression provides interpretable results with strong discriminatory power, while Random Forest and Neural Networks require further optimization for clinical deployment. Future research should focus on larger datasets, feature engineering, and ensemble learning to enhance specificity and overall predictive accuracy.

References

- Poulter, N. R. (1999). Coronary heart disease is a multifactorial disease. *European Heart Journal*, 20(9), 685-692.
- Haffner, S. M. (2000). Coronary heart disease in patients with diabetes. *New England Journal of Medicine*, 342(14), 1040-1042.
- Carney, R. M., & Freedland, K. E. (2017). Depression and coronary heart disease. *Nature Reviews Cardiology*, 14(3), 145-155.
- Barker, d. J. 1995. Fetal origins of coronary heart disease. Bmj, 311, 171-174.
- Carney, r. M. & freedland, k. E. 2017. Depression and coronary heart disease. Nature reviews cardiology, 14, 145-155.
- CHANDOLA, T., BRITTON, A., BRUNNER, E., HEMINGWAY, H., MALIK, M., KUMARI, M., BADRICK, E., KIVIMAKI, M. & MARMOT, M. 2008. Work stress and coronary heart disease: what are the mechanisms? European heart journal, 29, 640-648.
- CHEN, Y., WANG, L., MA, D., CUI, Z., LIU, Y., PANG, Q., JIANG, Z. & GAO, Z. 2024. Research on rheumatic heart disease from 2013 to early 2024: a bibliometric analysis. Journal of Cardiothoracic Surgery, 19, 659.
- DANESH, J., COLLINS, R. & PETO, R. 1997. Chronic infections and coronary heart disease: is there a link? The lancet, 350, 430-436.
- DAWBER, T. R., MOORE, F. E. & MANN, G. V. 2015. II. Coronary heart disease in the Framingham study. International journal of epidemiology, 44, 1767-1780.
- DRAZNER, M. H. 2011. The progression of hypertensive heart disease. Circulation, 123, 327-334.
- GRUNDY, S. M. 1986. Cholesterol and coronary heart disease: a new era. Jama, 256, 2849-2858.
- HAFFNER, S. M. 2000. Coronary heart disease in patients with diabetes. Mass Medical Soc.
- HAINES, A., IMESON, J. & MEADE, T. 1987. Phobic anxiety and ischaemic heart disease. Br Med J (Clin Res Ed), 295, 297-299.
- HANCOCK, S. L., TUCKER, M. A. & HOPPE, R. T. 1993. Factors affecting late mortality from heart disease after treatment of Hodgkin's disease. Jama, 270, 1949-1955.
- HOFFMAN, J. I. & KAPLAN, S. 2002. The incidence of congenital heart disease. Journal of the American college of cardiology, 39, 1890-1900.
- HOFFMAN, J. I., KAPLAN, S. & LIBERTHSON, R. R. 2004. Prevalence of congenital heart disease. American heart journal, 147, 425-439.
- KEYS, A. 1957. Diet and the epidemiology of coronary heart disease. Journal of the American Medical Association, 164, 1912-1919.
- KOENIG, W. 2001. Inflammation and coronary heart disease: an overview. Cardiology in review, 9, 31-35.
- LIU, J. L., MANIADAKIS, N., GRAY, A. & RAYNER, M. 2002. The economic burden of coronary heart disease in the UK. Heart, 88, 597-603.
- LLOYD-JONES, D. M., LARSON, M. G., BEISER, A. & LEVY, D. 1999. Lifetime risk of developing coronary heart disease. The Lancet, 353, 89-92.
- MCCARTHY, M., LAY, M. & ADDINGTON-HALL, J. 1996. Dying from heart disease. Journal of the Royal College of Physicians of London, 30, 325.
- NASCHITZ, J. E., SLOBODIN, G., LEWIS, R. J., ZUCKERMAN, E. & YESHURUN, D. 2000. Heart diseases affecting the liver and liver diseases affecting the heart. American heart journal, 140, 111-120.

PAUL, O., LEPPER, M. H., PHELAN, W. H., DUPERTUIS, G. W., MACMILLAN, A., MCKEAN, H. & PARK, H. 1963. A longitudinal study of coronary heart disease. Circulation, 28, 20-31.

PEARSON, T. A. 1996. Alcohol and heart disease. Circulation, 94, 3023-3025.

PELL, J. & COBBE, S. 1999. Seasonal variations in coronary heart disease. Qjm, 92, 689-696.

POULTER, N. 1999. Coronary heart disease is a multifactorial disease. American Journal of Hypertension, 12, 92S-95S.

UEBING, A., STEER, P. J., YENTIS, S. M. & GATZOULIS, M. A. 2006. Pregnancy and congenital heart disease. Bmj, 332, 401-406.

VIRANI, S. S., ALONSO, A., APARICIO, H. J., BENJAMIN, E. J., BITTENCOURT, M. S., CALLAWAY, C. W., CARSON, A. P., CHAMBERLAIN, A. M., CHENG, S. & DELLING, F. N. 2021. Heart disease and stroke statistics—2021 update. Circulation.