

A Lightweight Predictive Model for Heart Disease Based on Routine Clinical Data

Hanyu Wang

1 Introduction

Heart disease sometimes goes unnoticed, especially in asymptomatic individuals. However, occult heart disease can still cause shock and fatal consequences. Currently, diagnosing the disease requires special tests that are expensive and time-consuming. The aim of this study is to use low-cost, readily available physical data to predict the presence of heart disease.

Since liver disease, certain brain diseases and heart disease have similar characteristics, understanding statistical methods for predicting them can help predict heart disease. In previous studies, logistic regression was used in the prediction of liver disease (Abdalrada et al., 2019) and Alzheimer's disease (Johnson et al., 2014), and stepwise variable selection based on AIC was used in that of Alzheimer's disease (Johnson et al., 2014). In addition, Ture et al. (2008) used five models to predict the presence of coronary heart disease and suggested logistic regression to be a better performing technique. These suggest that it might be effective to use these methods to predict heart disease. Although Ture's study was similar to this one, the variables they ultimately considered included family history of CAD, diabetes mellitus, hypercholesterolemia and other difficult-to-obtain data. To enhance the applicability of the model, I utilized more accessible data and expanded the range of predicted diseases to the full spectrum of heart disease.

2 Methods

2.1 Choice of Method

Since the response variable is a binary, indicating whether the person has heart disease, we can use a binary logistic regression model. A generalized linear model (GLM) is proposed:

$$\log\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_q x_q = X\beta$$

2.2 Variable Selection

Akaike information criterion (AIC) is an estimator that can evaluate the fitness of the models, where $AIC = -2L(\hat{\theta}) + 2p$; $L(\hat{\theta})$ is the maximized log-likelihood value, and p is the number of parameters. Stepwise selection based on AIC starts with a full model and drops a variable if it decreases AIC until all predictors are statistically significant. This method is used not only for it had been found efficient in previous studies, but also considering that it suggests a model that is more favorable for prediction instead of one with fewer predictors.

2.3 Model Violations and Diagnostics

a. Assumptions of Binary Logistic Regression:

- (i) Binary Response: The response variable is dichotomous (two possible responses).
- (ii) Independence: The observations are independent of each other.
- (iii) Linearity: The log of the odds ratio, $\log\left(\frac{\pi}{1-\pi}\right)$, is a linear function of x .

(iv) Sufficient Sample Size: Each independent variable requires at minimum of 10 cases with the lowest frequency of outcomes.

b. Diagnostics: Dfbeta plots and deviance residual plots are used to check model diagnostics. Plots of an ideal model should display a random scatter of points symmetrically distributed around the identity line without discernible patterns.

c. Validation: Calibration and ROC plots are used to evaluate model performance. The bias-corrected line in the calibration plot of an ideal model should be a straight line starting from (0, 0) to (1, 1), and AUC of the ROC curve should be large and close to 1.

3 Results

3.1 Data Description

The data set contained information on 918 individuals, including eleven physical data and their heart disease status. After cleaning up the missing data, there were 746 observations remaining. To enhance the robustness of the model, all 746 observations were used.

The EDA plots are shown in Figure 4 in the Appendix.

Table 1: Dataset Summary

Numerical Variables						
Variables	Min.	1 st Qu.	Median	Mean	3 rd Qu.	Max
Age	28.00	46.00	54.00	52.88	59.00	77.00
RestingBP	92	120	130	133	140	200
Cholesterol	85.0	207.2	237.0	244.6	275.0	603.0
MaxHR	69.0	122.0	140.0	140.2	160.0	202.0
Oldpeak	-0.1000	0.0000	0.5000	0.9016	1.5000	6.2000
Nominal Variables						
Sex	Female: 193			Male: 725		
ChestPainType	ASY: 496		ATA: 173		NAP: 203	
					TA: 46	
FastingBS	No: 704			Yes: 214		
RestingECG	LVH: 188		Normal: 552		ST: 178	
ExerciseAngina	No: 547			Yes: 371		
ST_Slope	Down: 63		Flat: 460		Up: 395	
HeartDisease	No: 410			Yes: 508		

3.2 Analysis Procees

- (i) A logistic regression model including all predictors was fitted.
- (ii) Stepwise selections based on AIC and BIC, and LASSO method were used to identify the relevant predictors.
- (iii) Models were fitted based on the selected predictors from the previous step.
- (iv) Compare the modeling results (the model with smaller AIC and has predictors with larger estimators, smaller p-value and VIF has better performance), refer to previous similar studies, and select the most appropriate model.

The summary of the final model is shown in Table 2.

Table 2: Summary of the Final Model (Selected by Stepwise Selection Based on AIC)

(AIC: 542.33)	Estimate	Std. Error	p-value	VIF
Intercept	-2.764717	1.373911	0.044189	/
Age	0.034637	0.012895	0.007230	1.087973
SexMale	1.721240	0.296365	6.33e-09	1.101012
ChestPainTypeATA	-0.131288	0.516333	0.799287	2.787430
ChestPainTypeNAP	0.070530	0.480299	0.883254	3.374990
ChestPainTypeASY	1.600773	0.460741	0.000512	4.243226
RestingBP	0.011170	0.006929	0.106961	1.084407
Cholesterol	0.003334	0.001908	0.080556	1.051047
ExerciseAngina	1.076158	0.248973	1.54e-05	1.146023
Oldpeak	0.343885	0.138275	0.012884	1.282333
ST_Slope	-1.805422	0.236034	2.03e-14	1.272491

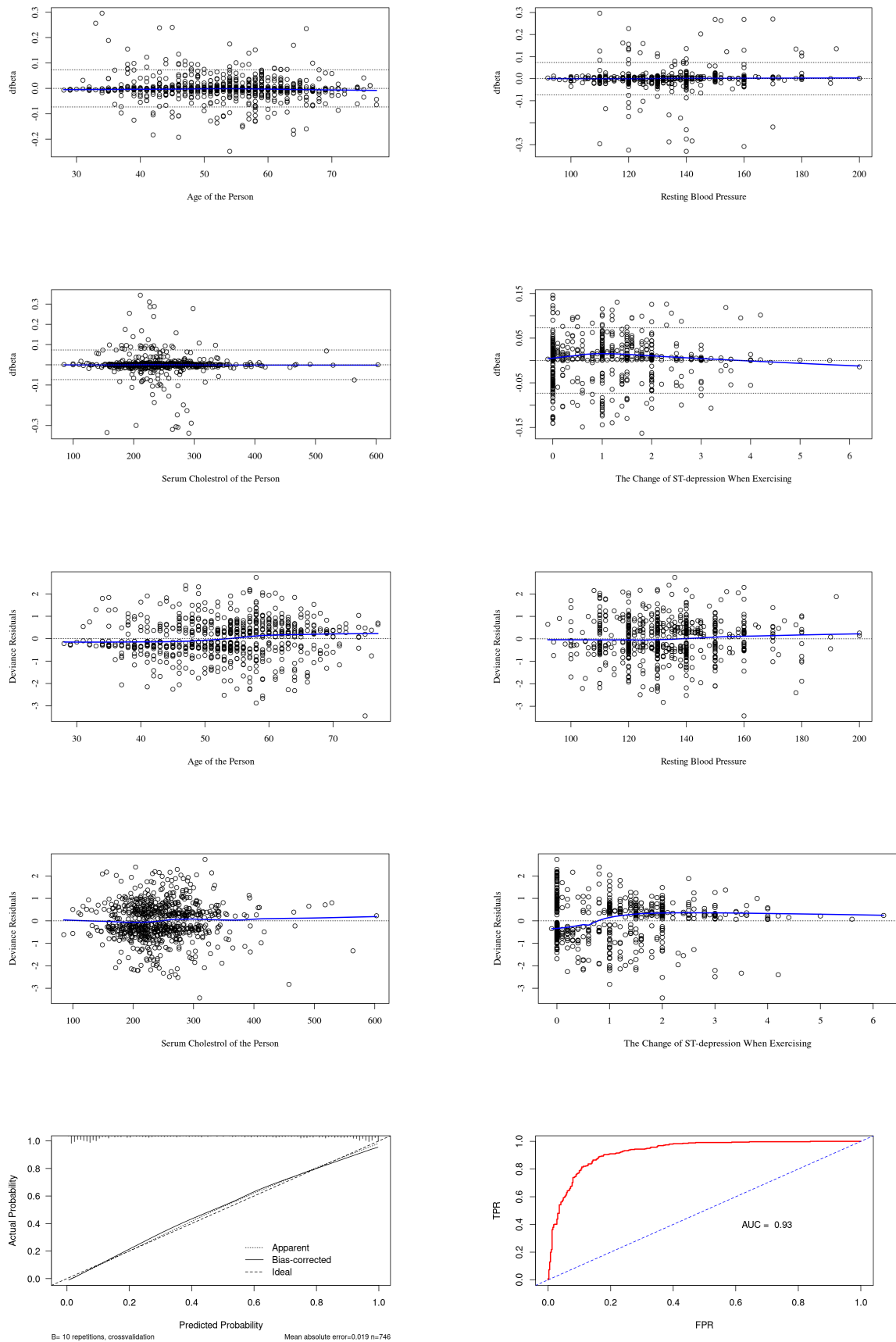
3.3 Diagnostics and Accuracy of the Final Model

The results of Dfbeta plots and deviance residual plots are mostly favorable, with most values very close to 0 and no clear pattern or trend. Although in the deviance residual plot of “The Change of ST-segment decline during exercise”, there are some fluctuations in the residuals, there is no clear systematic pattern overall, meaning that the model does not show significant heteroskedasticity or nonlinear trends.

The bias-corrected line in the calibration plot is close to the ideal line. The AUC of the ROC curve is 0.93, indicating that the model can correctly distinguish 93% of the data.

These show the usability and good predictability of the model.

Figure 3: Final Model Diagnostics and Accuracy



4 Discussion

4.1 Final Model Interpretation and Importance

Since the estimated coefficient of Age is 0.034637, when it increases by 1, the odds ratio is $\exp(0.034637) = 1.035$. Other numeric predictors can be interpreted in the same way.

For nominal variables, the predictor is 1 if that condition is met or 0 otherwise. For example, since the estimated coefficient of SexMale is 1.721240, holding other predictors constant, the odds ratio of male versus female is $\exp(1.721240) = 5.591$, so males have larger probabilities to have heart disease than females.

To predict the presence of heart disease, if the predicted probability is larger than 0.5, the person is likely to have heart disease; otherwise, the person is more likely to be healthy.

The model has good performance, and all physical indicators needed are readily available, so this model solves the research question.

4.2 Limitations of the Analysis

The estimated coefficients of several predictors are small, with p-values larger than 0.05, indicating the predictors have small contributions to model prediction. However, since a model with accurate predictability is expected, this model was ultimately chosen considering that a model with a smaller AIC is more effective for prediction and previous similar studies have used AIC for model selection.

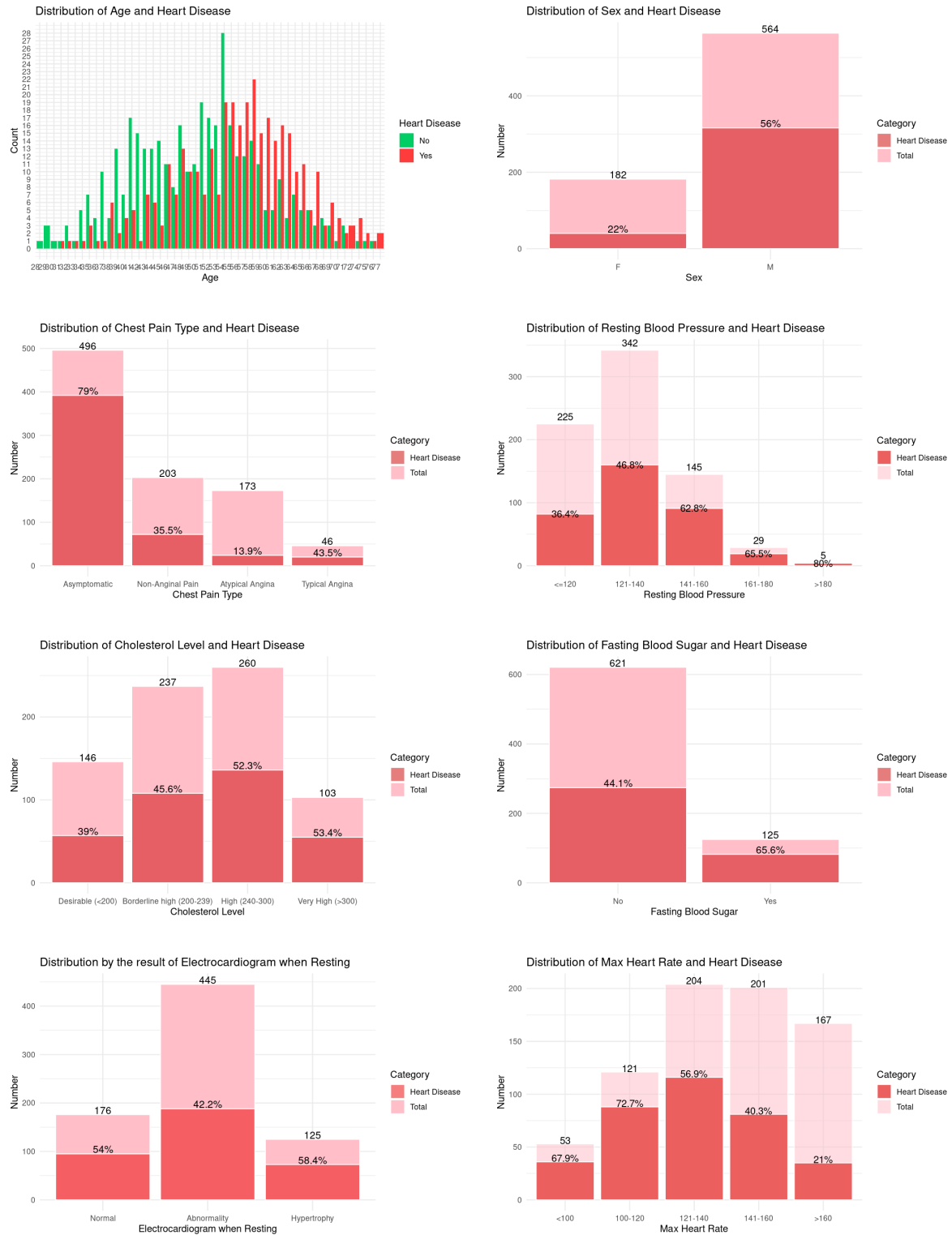
However, a model with most effective predictors was also fitted based on the results of stepwise selections based on AIC and BIC, and LASSO method. The details are shown in Table 5 in the Appendix for reference.

5 References

- Johnson, P., Vandewater, L., Wilson, W., Maruff, P., Savage, G., Graham, P., Macaulay, L. S., Ellis, K. A., Szoek, C., Martins, R. N., Rowe, C. C., Masters, C. L., Ames, D., & Zhang, P. (2014). Genetic algorithm with logistic regression for prediction of progression to Alzheimer's disease. *BMC Bioinformatics*, 15(Suppl 16), S11.
- Kurt, I., Ture, M., & Kurum, A. T. (2008). Comparing performances of logistic regression, classification and regression tree, and neural networks for predicting coronary artery disease. *Expert Systems with Applications*, 34(1), 366–374.
- Abdalrada, A., Yahya, O., Alaidi, A., Hussein, N., Alrikabi, H., & Al-Quraishi, T. (2019). A predictive model for liver disease progression based on logistic regression algorithm. *Periodicals of Engineering and Natural Sciences (PEN)*, 7(3), 1255.

6 Appendix

Figure 4: EDA plots



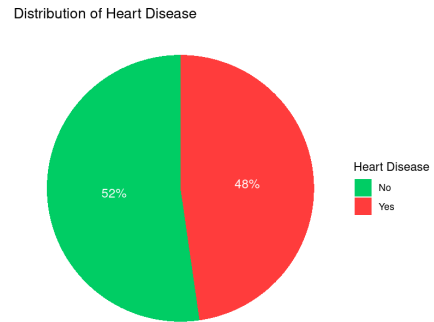
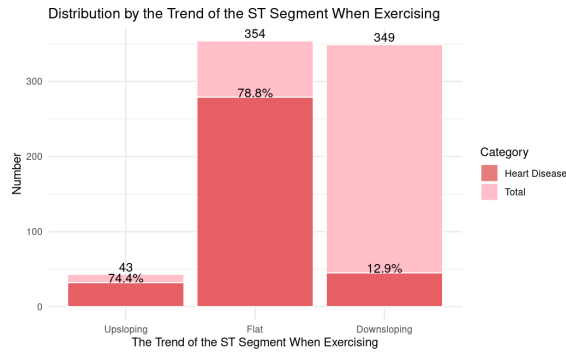
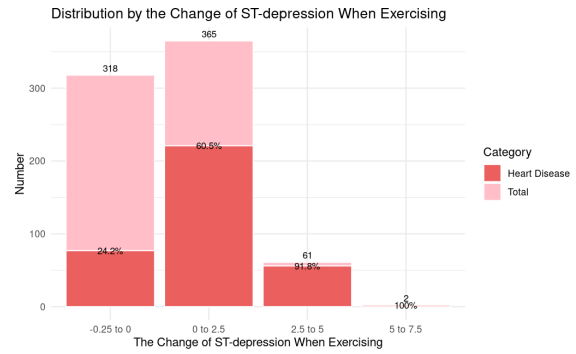
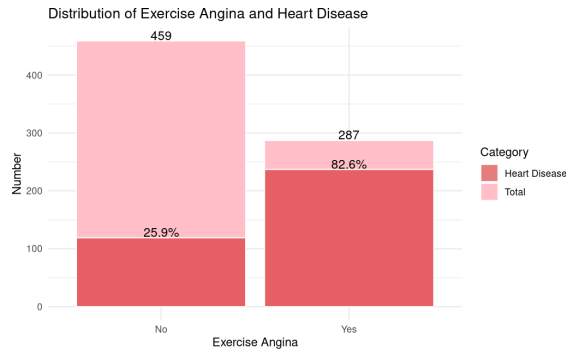


Table 5: Summary of the Fitted Model

(AIC: 559.15)	Estimate	Std. Error	p-value	VIF
Intercept	2.6749	0.6378	2.74e-05	/
SexMale	1.5528	0.2818	3.59e-08	1.053973
ChestPainTypeATA	-0.4343	0.4981	0.38323	2.709254
ChestPainTypeNAP	-0.1488	0.4642	0.74857	3.315481
ChestPainTypeASY	1.3918	0.4425	0.00166	4.111303
ExerciseAngina	1.3364	0.2361	1.51e-08	1.080341
ST_Slope	-2.1167	0.2151	< 2e-16	1.111114