

Bayesian Logistic Regression for Heart Disease

By Daniel Vitale





Introduction

According to the CDC National Center for Health Statistics, heart disease represents the number one cause of death for both men and women in the United States.

The objective of this project is to produce a predictive model and examine the potential relevance of a handful of metrics related to heart disease in male and female patients



UCI Heart Disease Dataset

<https://archive.ics.uci.edu/ml/datasets/Heart+Disease>

The dataset contains metrics from 1025 individual patients including the following:

1. age
2. sex
3. chest pain type (4 values)
4. resting blood pressure
5. serum cholesterol in mg/dl
6. fasting blood sugar > 120 mg/dl
7. resting electrocardiographic results (values 0,1,2)
8. maximum heart rate achieved
9. exercise induced angina
10. oldpeak = ST depression induced by exercise relative to rest
11. the slope of the peak exercise ST segment
12. number of major vessels (0-3) colored by fluoroscopy
13. thal: 0 = normal; 1 = fixed defect; 2 = reversible defect



Bayesian Logistic Regression Model

The model was implemented as follows:

$$\mu = \text{logistic}(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)$$

$$y \sim \text{Bernoulli}(\mu)$$

Where:

$$\text{logistic}(x) = 1/(1 + \exp(-x))$$

Standardization:

$$z = (x - \bar{x})/s_x$$

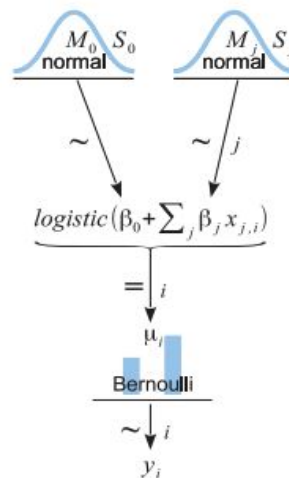


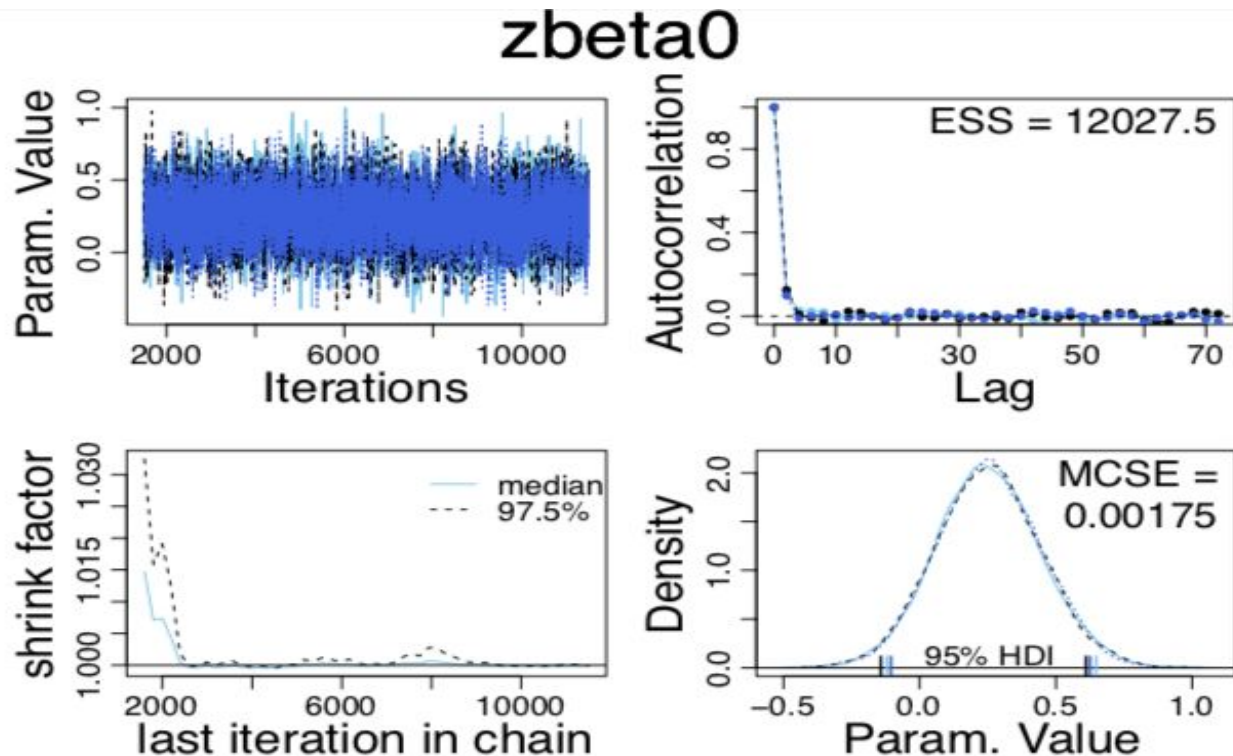
Figure 21.2: Dependency diagram for multiple logistic regression. Compare with the diagram for robust multiple linear regression in Figure 18.4 (p. 498). Copyright © Kruschke, J. K. (2014). *Doing Bayesian Data Analysis: A Tutorial with R, JAGS, and Stan. 2nd Edition*. Academic Press / Elsevier.



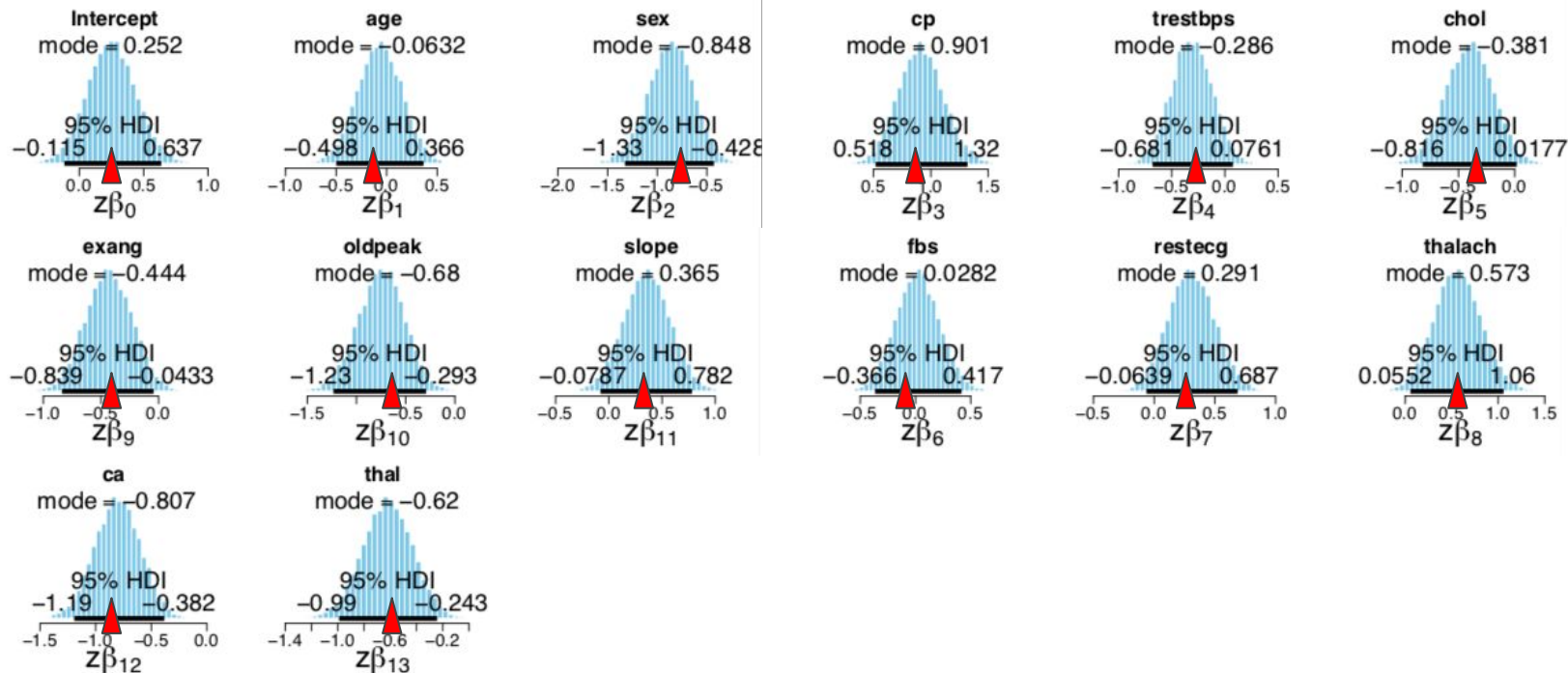
Bayesian Logistic Regression Model (continued)

- Due to lack of priors, β_0 and β_n were assumed to have been generated by normal distributions with mean = 0 and standard deviation = $1/2^2$ (presets of DBDA logistic regression script) though this is likely not the case in real life
- This will have effects on the posteriors that we generate!

Results: Metropolis Algorithm Diagnostics

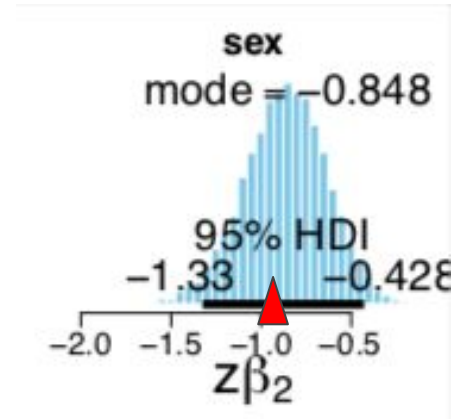


Results



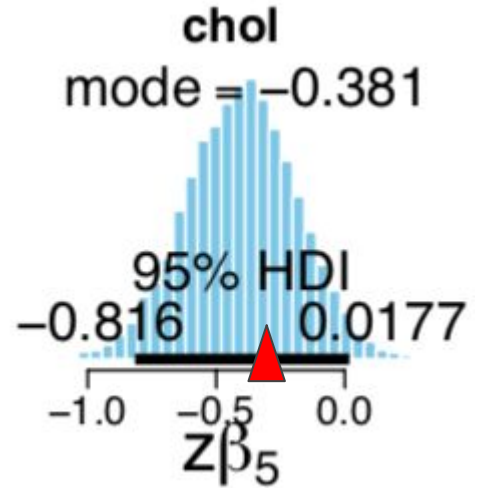
Discussion- Biological Sex and Heart Disease

- HDI for $z\beta$ falls entirely in the negative
- Log-odds of heart disease in men (sex=1) versus women (sex=0) is negative
- This is accounted for bias in the dataset:
 - 45% of men have heart disease
 - 75% of women have heart disease



Discussion- Cholesterol and Heart Disease

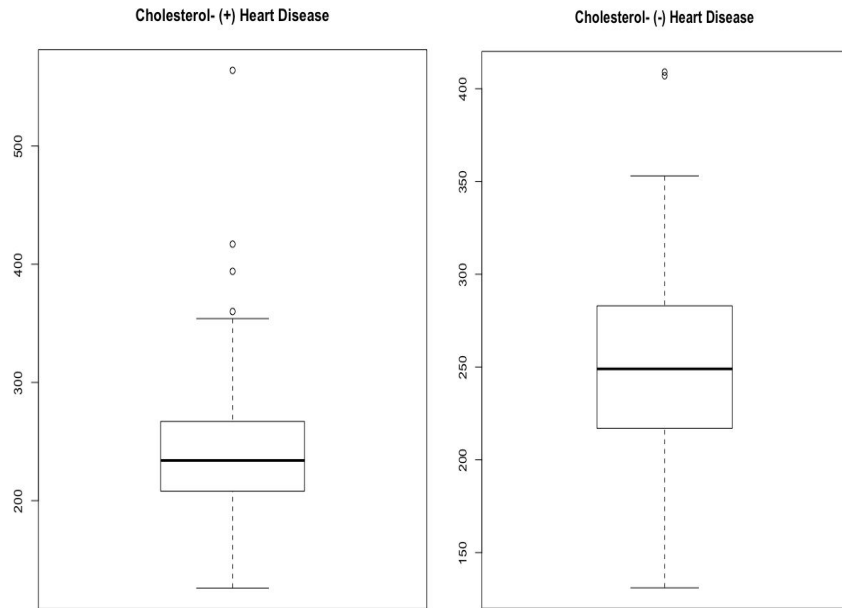
- HDI for $z\beta$ falls almost entirely in the negative
- Log-odds of heart disease and total cholesterol have negative relationship





Discussion- Cholesterol and Heart Disease (continued)

- Dataset is biased for high cholesterol
- Cholesterol means:
 - + for Heart Disease: 242.2303
 - - for Heart Disease: 251.087





Conclusions

- Bias in data matters!
- Priors matter!
 - Without including the prior distributions for specific features (cholesterol and sex) I have very likely produced a very overfit model!
 - I would expect the results to be much different for cholesterol and sex if priors were included (but I ran out of time!)



Going Forward

- Examination of sex individually
- Examination of this dataset with the inclusion of prior distributions of features



Sources

Kruschke, J. K. (2014). Doing Bayesian Data Analysis: A Tutorial with R, JAGS, and Stan. 2nd Edition. Academic Press / Elsevier

CDC - NCHS - National Center for Health Statistics. (n.d.). Retrieved from <https://www.cdc.gov/nchs/>