

# P8160 Group Project 2: Breast Cancer Diagnosis and Optimizations

Wenhan Bao | Tianchuan Gao | Jialiang Hua | Qihang Wu | Paula Wu

3/25/2022

## Objective

The main objective of our project is to build an accurate predictive model based on logistic regression that classifies between malignant and benign images of breast tissue. Using the Breast Cancer Diagnosis dataset, we will implement logistic model and logistic-LASSO model to predict the diagnosis. A Newton-Raphson algorithm and Pathwise Coordinate optimization will be developed to estimate the logistic model and the lasso model respectively. Our aim is to find the model with the best performance in predicting the diagnosis of breast tissue image.

## Background & Methods

### Background

Breast cancer, which affects 1 in 7 women worldwide, is the most common invasive cancer in women around the world. It can start from different parts of the breast and is marked by the uncontrolled growth of breast cells. Benign breast tumors do not metastasize and are usually not life-threatening, while malignant tumors are aggressive and deadly. Nowadays, substantial support for breast cancer awareness and research funding has created great advances in the diagnosis and treatment of breast cancer. The prognosis of the disease has been greatly improved once it is detected and treated early. Therefore, it's important to have breast lumps accurately diagnosed so that timely clinical intervention can be conducted.

### Methods

#### Exploratory Data Analysis (EDA)

Firstly, we imported and cleaned the data and conducted exploratory data analysis. We plotted the correlation plot of all the predictors by R. From the plot, we can find that there is a high correlation between many of the predictors. This is because our predictors include mean, standard deviation and the largest values of the distributions of 10 features, which means that some predictors can be calculated by other predictors. Then, we made the feature plots to explore the relationship between binary diagnosis outcome and all our predictors to determine the potential effective predictors.

#### Modified Newton-Raphson Algorithm

#### Logistic-LASSO

***Least Absolute Shrinkage and Selection Operator (LASSO)*** To estimate coefficients through Newton-Raphson method, it is necessary to compute the corresponding inverse of Hessian Matrix  $[\nabla^2 f(\theta_{i-1})]^{-1}$ . However, the computational burden of calculation will increase as the dimension of predictors increases and the collinearity will also be a problem. Therefore, we use a regularization method, LASSO, to shrink

coefficients and perform variable selections. For regression lasso, the objective function is:

$$f(\beta) = \frac{1}{2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{i,j} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|,$$

where the first term is residual sum of squares (RSS) and the second term is the lasso l1 penalty. Noted that the  $x_{i,j}$  needs to be standardized before LASSO so that the penalty will be equally applied to all predictors. For each single predictor, the LASSO solution is like:

$$\hat{\beta}^{lasso}(\lambda) = S(\hat{\beta}, \lambda) = \begin{cases} \hat{\beta} - \lambda, & \text{if } \hat{\beta} > 0 \text{ and } \lambda < |\hat{\beta}| \\ \hat{\beta} + \lambda, & \text{if } \hat{\beta} < 0 \text{ and } \lambda < |\hat{\beta}| \\ 0, & \text{if } \lambda > |\hat{\beta}|, \end{cases}$$

where  $S(\hat{\beta}, \lambda)$  is called soft threshold. The basic idea of this function is to shrink all  $\beta$  coefficients between  $-\lambda$  and  $\lambda$ .

**Coordinate-wise Descent Algorithm** Another approach to solve the complex computation of inverse Hessian Matrix is coordinate descent approach, which starts with an initial guess of parameters  $\theta$ , optimize one parameter at one time based on the best knowledge of other parameters, and use the results as the start values for the next iteration. We then repeat the above steps until convergence. Finally, this approach tries to minimize the following objective function when considering a lasso penalty:

$$f(\beta_j) = \frac{1}{2} \sum_{i=1}^n (y_i - \sum_{k \neq j} x_{i,j} \tilde{\beta}_k - x_{i,j} \beta_j)^2 + \lambda \sum_{k \neq j} |\tilde{\beta}_k| + \lambda |\beta_j|,$$

where  $\tilde{\beta}$  represents the current estimates of  $\beta$  and therefore constants. If we also consider a weight  $\omega_i$  is associated with each observation, then the updated  $\beta_j$  in this case is:

$$\tilde{\beta}_j(\lambda) \leftarrow \frac{S(\sum_i \omega_i x_{i,j} (y_i - \tilde{y}_i^{(-j)}), \lambda)}{\sum_i \omega_i x_{i,j}^2},$$

where  $\tilde{y}_i^{(-j)} = \sum_{k \neq j} x_{i,k} \tilde{\beta}_k$ . From here, we use the Taylor expansion. So the log-likelihood around “current estimate” is:

$$l(\beta) = -\frac{1}{2n} \sum_{i=1}^n \omega_i (z_i - X_i \beta)^2,$$

where working weights  $\omega_i = \tilde{\pi}_i(1 - \tilde{\pi}_i)$ , working response  $z_i = X_i \tilde{\beta} + \frac{y_i - \tilde{\pi}_i}{\tilde{\pi}_i(1 - \tilde{\pi}_i)}$ , and  $\tilde{\pi}_i = \frac{\exp(X_i \tilde{\beta})}{1 + \exp(X_i \tilde{\beta})}$ .

Finally, similar to regression lasso, the logistic lasso can be written as a penalized weighted least-squares problem like this:

$$\min_{\beta} L(\beta) = -l(\beta) + \lambda \sum_{j=0}^p |\beta_j|$$

**Pathwise Coordinate-wise Algorithm** The difference between pathwise coordinate-wise method and coordinate-wise method is that a sequence value of lambda is required to input. There are some steps taken as followed:

- Select a  $\lambda_{max}$  for all the estimate  $\beta = 0$  which is the inner product ( $\max_l \langle X_l, y \rangle$ ).
- Compute the solution with a sequence of descending  $\lambda$  from maximum to zero ( $\lambda_{max} > \lambda_k > \lambda_{k-1} \dots > 0$ ).
- Initialize coordinate descent algorithms for  $\lambda_k$  by the calculated estimate  $\beta$  from previous  $\lambda_{k+1}$  as a warm start
- By repeating the two steps above, a sequence of optimal coefficients  $\beta$  for each descending  $\lambda$ . When objective function is taken the minimum value ( $\min_{\beta} L(\beta)$ ), the best  $\lambda$  could be identified and optimal coefficients  $\beta$  could be selected under this best  $\lambda$

## **5-fold Cross-validation**

In a nutshell, a 5-fold cross-validation first splits the shuffled training data (80% of the whole dataset) into 5 parts and takes one group as the hold-out set (validation set) while fitting the model using the remaining 4 training sets. After a model is fitted, we evaluate and retain the model performance, namely AUC and RSS, and move on to the next iteration. After 5 iterations, we calculated the mean RSS and the mean AUC with standard deviations. Our goal is to find an optimal  $\lambda$  that maximizes the mean AUC.

## **Conclusions**

### **Logistic Regression Model**

### **Logistic-LASSO Model**

## **Discussion**

## **Contributions**

We contributed to this project evenly.

## **Appendix**

## **Reference**