

Project2: Breast Cancer Prediction Model

Group 7: Melanie Mayer, Yuqi Miao, Sibe Liu, Xue Jin

March 31, 2020

Introduction

Breast cancer is the most common invasive cancer and the second leading cause of death from cancer in women worldwide, marked by the uncontrolled growth of breast cells. Non-cancerous breast tumors do not metastasize and are usually not life-threatening, while malignant tumors are cancerous, aggressive and deadly. Therefore, it's important to have breast lumps accurately diagnosed so that decision with regard to medical treatment, rehabilitation and personal matters can be made appropriately.

Objectives

The main objective of this 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, logistics model and logistic-LASSO model will be implemented to predict the diagnosis. A Newton-Raphson algorithm and Pathwise Coordinate optimization will be developed to estimate the logistic model and the lasso-logistic model respectively. We aim to find the model with the best performance in terms of predicting when breast tissue is malignant.

Breast Cancer Diagnosis Dataset

The Breast Cancer Diagnosis dataset contains the diagnosis and a set of 10 features capturing the characteristics of the cell nuclei present in the digitized image of breast mass. The ten features collected for each cell nucleus are:

- radius (mean of distances from center to points on the perimeter)
- texture (standard deviation of gray-scale values)
- perimeter
- area
- smoothness (local variation in radius lengths)
- compactness ($\text{perimeter}^2 / \text{area} - 1.0$)
- concavity (severity of concave portions of the contour)
- concave points (number of concave portions of the contour)
- symmetry
- fractal dimension ("coastline approximation" - 1)

The mean, standard deviation (SD) and largest values of these features are computed for each image, resulting in 30 possible predictive variables. We will analyze the predictive ability for diagnosis of malignant or benign cases of these covariates.

Methods

Model Parameters

For logistic regression we assume the response variable Y_i for the i_{th} observation follows a binary distribution:

$$Y_i \sim \text{Bin}(\pi_i)$$

where π_i denotes the probability that the i_{th} observation's tissue is malignant. We can assume all observations are independent from one another, hence the likelihood function for the vector $\boldsymbol{\pi}$ can be written as:

$$L(\boldsymbol{\pi}) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$$

For logistic regression, the logit link function is used:

$$\log\left(\frac{\pi_i}{1 - \pi_i}\right) = \boldsymbol{\beta}^T \mathbf{x}_i = \theta_i$$

where $\mathbf{x}_i^T = [1 \ x_{1i} \ x_{2i} \ \dots \ x_{pi}]$ and $\boldsymbol{\beta}^T = [\beta_0 \ \beta_1 \ \beta_2 \ \dots \ \beta_p]$. One can solve for $\pi_i = \frac{e^{\theta_i}}{1 + e^{\theta_i}}$. We aim to find the best estimate of the vector of coefficients, $\boldsymbol{\beta}$. The log-likelihood function for this vector can be written as:

$$l(\boldsymbol{\theta}) = \log L(\boldsymbol{\pi}) = \sum_{i=1}^n (Y_i \log \frac{\pi_i}{1 - \pi_i} + \log(1 - \pi_i)) = \sum_{i=1}^n (Y_i \theta_i - \log(1 + e^{\theta_i}))$$

The maximum likelihood is thus achieved when the gradient is equal to zero and the Hessian is negative definite. The gradient can be found to be:

$$\nabla l(\boldsymbol{\theta}|\mathbf{X}) = \sum_{i=1}^n (Y_i - \pi_i) \mathbf{x}_i = \mathbf{X}^T (\mathbf{Y} - \boldsymbol{\pi})$$

The Hessian matrix is thus:

$$\nabla^2 l(\boldsymbol{\theta}|\mathbf{X}) = - \sum_{i=1}^n \pi_i (1 - \pi_i) \mathbf{x}_i \mathbf{x}_i^T = -\mathbf{X}^T \text{diag}(\pi_i (1 - \pi_i)) \mathbf{X}$$

Full model

In order to estimate $\boldsymbol{\beta}$ we need to maximize the loglikelihood function. There is no closed form, hence we turn to numerical methods. The Newton-Raphson method is used to fit the full logistic model. To find the maximum likelihood estimate of each element of $\boldsymbol{\beta}$, an iterative process is set as follows:

$$\theta_{i+1} = \theta_i - \delta (\nabla^2 l(\theta_i|\mathbf{X}) - \gamma I)^{-1} \nabla l(\theta_i|\mathbf{X})$$

This is the Newton-Raphson method with two modifications. δ is included in the process to accomplish the step-halving modification, the step coefficient ensures the likelihood is always increasing in order to achieve quicker convergence. γ is the modification coefficient to ensure the ascent direction of the iteration vector at θ_i .

Logit-lasso pathwise coordinate-wise update algorithm

In the case of large dimensionality of predictors or multicollinearity it can be beneficial to perform a regularization method which shrinks coefficients and performs variable selection. Here we implement the Least Absolute Shrinkage and Selection Operator (LASSO) method for logistic regression with a path-wise coordinate-wise optimization algorithm to select variables from the full model and increase the prediction

efficiency and avoid overfitting. Seeing how we have 30 predictors describing 10 features, we are likely to experience multicollinearity and expect LASSO to outperform the classical logistic regression.

In LASSO we add an L1 penalization term to the squared loss function, such that we try to find the coefficients to minimize:

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

Where we find the likelihood to be:

$$l(\beta) = -\frac{1}{2n} \sum_{i=1}^n \omega_i (z_i - \mathbf{X}_i \beta)^2$$

$$\pi_i = \frac{\exp(\mathbf{X}_i \beta)}{1 + \exp(\mathbf{X}_i \beta)}$$

$$\omega_i = \pi_i (1 - \pi_i)$$

$$z_i = \mathbf{X}_i \beta + \frac{y_i - \pi_i}{\pi_i (1 - \pi_i)}$$

Pre-define the tuning parameter sequence $\{\lambda_1, \dots, \lambda_s\}$, starting point $\beta_{\text{start}} = \{\beta_0^{(0)}, \dots, \beta_p^{(0)}\}$. Here we elaborately explain the optimal process for λ_u : Using the optimal $\beta_{\mathbf{u}-1}$ from last iteration as the warm start. within every iteration, find the optimal β coordinate-wise. For β_j in $t_i h$ iteration

$$\beta_j^{(t)} = \begin{cases} \frac{\sum_{i=1}^n \omega_i (z_i - \sum_{j=1}^p \mathbf{X}_i \beta_j)}{\sum_{i=1}^n \omega_i}, & j = 0 \\ \frac{s(\beta_j^{(t^*)}, \lambda_u n)}{\sum_{i=1}^n \omega_i x_{ij}^2}, & j = 1, 2, \dots, p \end{cases}$$

$$\beta_j^{(t^*)} = \sum_{i=1}^n \omega_i x_{ij} z_{ij}^*$$

$$z_{ij}^* = z_i - \sum_{\substack{k=0 \\ \beta_k \neq 0}}^{j-1} \beta_k^{(i)} x_{ik} - \sum_{\substack{k=j+1 \\ \beta_k \neq 0}}^p \beta_k^{(i-1)} x_{ik}$$

Cross Validation

In order to check the model performance, we use 5-fold cross validation. Using training dataset to choose the best tuning parameter and fit model, and using test dataset to evaluate the final prediction. The statistics we use to compare the validation is SSE and pearson chi-square statistics, which are defined as:

$$SSE = \sum_{i=1}^n (y_i - \hat{\pi}_i)^2$$

$$\hat{\pi}_i = \log \frac{\exp(\mathbf{X}_i \beta)}{1 + \exp(\mathbf{X}_i \beta)}$$

$$G = \sum_{i=1}^n \frac{y_i - \hat{\pi}_i}{\hat{\pi}_i (1 - \hat{\pi}_i)}$$

By taking average of above 2 statistics of 5 fold, we get the final index to evaluate the model fitting.

Results

Newton-Raphson

After doing the Newton-Raphson modified with step and direction, the estimation of coefficients are:

intercept	0.4870168
radius_mean	-7.2218505
texture_mean	1.6547562
perimeter_mean	-1.7376303
area_mean	14.0048456
smoothness_mean	1.0749533
compactness_mean	-0.0772346
concavity_mean	0.6751231
concave points_mean	2.5928743
symmetry_mean	0.4462563
fractal_dimension_mean	-0.4824842

Table 1. Estimated coefficients under Newton-Raphson method

Coordinate-Wise

The range of λ we tried is (3,0) with length 100. The initial guess of all β including the intercept is 0.02. The g.statistics, SSE, AUC are introduced to select best λ in Cross Validation. The optimal λ would be selected based on the minimum SSE, maximum AUC and minimum g-statistics in test data. Below is the results in 5-Fold Cross Validation:

	Enter	Fold1	Fold2	Fold3	Fold4	Fold5
k	0.00	1.0000000	2.0000000	3.0000000	4.0000000	5.0000000
best_lambda	0.00	0.0000000	0.0000000	0.0000000	0.0000000	0.0000000
beta_vec1	0.02	-0.5851898	-0.5921021	-0.6782398	-0.6305423	-0.5250743
beta_vec2	0.02	1.9370600	1.8598428	1.8706748	1.0162875	1.7805325
beta_vec3	0.02	0.8620696	0.9438123	0.9124855	0.8276538	0.9191474
beta_vec4	0.02	-0.0027197	0.0110750	0.0517034	1.0505263	-0.0067178
beta_vec5	0.02	-0.0045932	-0.0077140	-0.0160055	-0.0010837	-0.0105529
beta_vec6	0.02	0.4154545	0.4039963	0.3839284	0.3662432	0.5432791
beta_vec7	0.02	-0.0406285	0.0146759	-0.0267084	0.0795153	-0.0452543
beta_vec8	0.02	0.1820403	0.1956254	0.3122696	0.2686984	0.1299537
beta_vec9	0.02	2.0666740	1.9595822	1.9458032	1.8233049	2.3142211
beta_vec10	0.02	0.0725989	0.1110710	0.1140780	0.1066402	0.1031025
beta_vec11	0.02	-0.1413081	-0.1589079	-0.1699873	-0.1369779	-0.1387135
g.stat_tr	Inf	137.3082360	135.9975379	116.9462145	135.4900541	114.5090024
auc_te	0.00	0.9913435	0.9923154	0.9849530	0.9826870	0.9743770
g.stat_te	Inf	22.6159703	22.1514224	48.1975784	34.5067570	48.0076306
SSE_test	Inf	4.2783563	4.1138120	4.9735803	5.9718400	6.8409613

Table 2. Cross validation results

In all 5 folds, all three critria indicates the same optimal λ , whcih is 0. In the test data of each fold, the AUC ranges from 0.99 to 0.97, g-statistics ranges from 22.6 to 48.0, while SSE changes from 4.1 to 6.8.

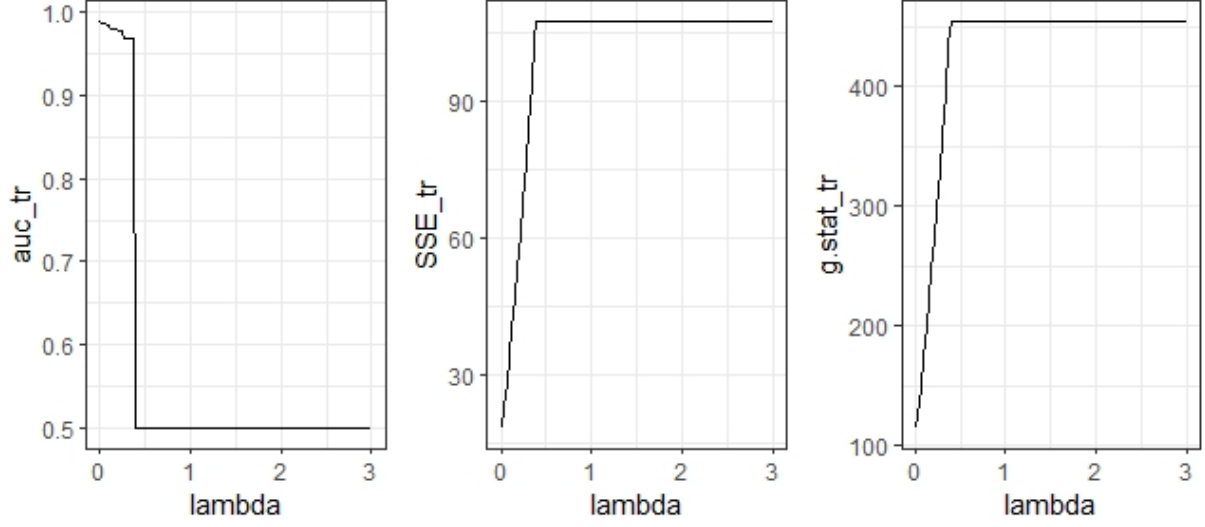


Fig 1. Changes in three criteria vs λ

The Figure 1 shows the trend of AUC SSE and g-statistics in train data. With the increse of λ , both SSE and g-statistics have a dramatic soar. While AUC has a great drop from 1 to 0.5. All of them indicate that the bigger the λ is, the worse the model would perform.

Discussion