# Project2 :Breast Cancer Prediction Model

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

March 31, 2020

## Introduction

Breast cancer is the most common invasive cancer and the second leading cause of cancer death 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 a predictive model based on binary logistic regression that classifies between malignant and benign cases. Using the Breast Cancer Diagnosis dataset, logistics model and logistic-LASSO model will be implemented to predict the diagnosis. As two important methods for numerical optimization, Newton-Raphson algorithm and Pathwise Coordinate optimization will be developed to estimate the logistic model and the lasso model respectively.

### Breast Cancer Diagnosis Dataset

The Breast Cancer Diagnosis dataset contains the diagnosis and a set of 30 features capturing the characteristics of the cell nuclei present in the digitized image of breast mass. Ten features are collected for each cell nucleus:

- 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 (perimeter^2 / 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 error (SE) and largest values of these features are computed for each image, resulting in 30 features, which will be analyzed to understand the predictive value for diagnosis of cancer(malignant) or benign cases.

## Methods

### Model define

For $i_{th}$ observation, the response variable $Y_i$ follows binary distribution:

$$Y_i \sim Bin(\pi_i)$$

the log-liklihood function:

$$l(\mathbf{Y}, \boldsymbol{\pi}) = \sum_{i=1}^{n} l(y_i, \pi_i) = \sum_{i=1}^{n} (y_i log \frac{\pi_i}{1-\pi_i} + log(1-\pi_i)) = \sum_{i=1}^{n} (y_i log \frac{\pi_i}{1-\pi_i} + log(1-\pi_i))$$

Where $\pi_i$ denotes the probablity of the $i_{th}$ observation to be maglinant.

To build relationship between the response and predictors, the GLM is defined as:

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

### Full model

Firstly, Newton-Rapson method is used to fit the full model. To find the maximum likelihood estimation of coefficients, iteration process is set as follows:

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

where $\delta$ is the step coefficient to ensure the increasing of likelihood funciton, and $\gamma$ is the modification coefficient to ensure the aescent direction of the iteration vector.

$$\nabla l(\theta|\boldsymbol{X}) = \mathbf{X}^T (\mathbf{Y} - \boldsymbol{\pi})$$

$$\nabla^2 l(\theta|\boldsymbol{X}) = -\mathbf{X}^T diag(\pi_i(1-\pi_i)) \mathbf{X}$$

### logit-lasso pathwise coordinate-wise update algorithm

To select variables and increase the prediction efficiency, lasso was integrated into the coordinate-wise logit regression.

The target function:

$$min\{-l(\boldsymbol{\beta}) + \lambda \sum_{j=1}^{p} |\beta_j|\}$$

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

$$\pi_i = \frac{exp(\boldsymbol{X_i} \boldsymbol{\beta})}{1 + exp(\boldsymbol{X_i} \boldsymbol{\beta})}$$

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

$$z_i = \boldsymbol{X_i}\boldsymbol{\beta} + \frac{y_i - \pi_i}{\pi_i(1 - \pi_i)}$$

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

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

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

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

**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 MSE and person chi-square statistics, which are defined as:

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

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

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

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

## Results

**Model define**

For $i_{th}$ observation, the response variable $Y_i$ follows binary distribution:

$$Y_i \sim Bin(\pi_i)$$

the log-liklihood function:

$$l(\mathbf{Y}, \boldsymbol{\pi}) = \sum_{i=1}^n l(y_i, \pi_i) = \sum_{i=1}^n (y_i log \frac{\pi_i}{1 - \pi_i} + log(1 - \pi_i)) = \sum_{i=1}^n (y_i log \frac{\pi_i}{1 - \pi_i} + log(1 - \pi_i))$$

Where $\pi_i$ denotes the probablity of the $i_{th}$ observation to be maglinant.

To build relationship between the response and predictors, the GLM is defined as:

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

**Full model**

Firstly, Newton-Rapson method is used to fit the full model. To find the maximum likelihood estimation of coefficients, iteration process is set as follows:

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

where $\delta$ is the step coefficient to ensure the increasing of likelihood funciton, and $\gamma$ is the modification coefficient to ensure the aescent direction of the iteration vector.

$$\nabla l(\theta|\boldsymbol{X}) = \mathbf{X}^T(\mathbf{Y} - \boldsymbol{\pi})$$

$$\nabla^2 l(\theta|\boldsymbol{X}) = -\mathbf{X}^T diag(\pi_i(1 - \pi_i))\mathbf{X}$$

**logit-lasso pathwise coordinate-wise update algorithm**

To select variables and increase the prediction efficiency, lasso was integrated into the coordinate-wise logit regression.

The target function:

$$min\{-l(\boldsymbol{\beta}) + \lambda \sum_{j=1}^{p} |\beta_j|\}$$

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

$$\pi_i = \frac{exp(\boldsymbol{X_i}\boldsymbol{\beta})}{1 + exp(\boldsymbol{X_i}\boldsymbol{\beta})}$$

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

$$z_i = \boldsymbol{X_i}\boldsymbol{\beta} + \frac{y_i - \pi_i}{\pi_i(1 - \pi_i)}$$

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

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

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

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

## 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 MSE and person chi-square statistics, which are defined as:

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

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

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

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

## Discussion