

THE THEORETICAL STUDY OF RARE EVENT WEIGHTED LOGISTIC REGRESSION FOR CLASSIFICATION OF IMBALANCED DATA

Dian Eka Apriana Sulasih, Santi Wulan Purnami, Santi Puteri Rahayu

Institut Teknologi Sepuluh Nopember, Department of Statistics,
Jl. Arif Rahman Hakim, Surabaya 60111, Indonesia

dian.eka14@mhs.statistika.its.ac.id (Dian Eka A S)

Abstract

One of the problems in data classification is imbalanced data. In two-class classification, imbalance problem occurs where one of the two classes has more samples than another class. In such situation, most of the classifier will be biased towards the major class, while the minor class will be subordinated eventually which leads to inaccurate classification. Therefore, a method to classify the imbalanced data is required. Rare Event Weighted Logistic Regression (RE-WLR) which is developed by Maalouf and Siddiqi is a method of classification applied to large imbalanced data and rare event. This study showed the review of RE-WLR for the classification of imbalanced data. It explicated the steps to obtain the estimator specifically, particularly for IRLS. RE-WLR is a combination of Logistic Regression (LR) rare events corrections and Truncated Regularized Iteratively Re-weighted Least Squares (TR-IRLS). Rare event correction in LR is applied to Weighted Logistic Regression (WLR). Regularization was added to reduce over-fitting. The estimation of β is performed by using the method of *maximum likelihood* (ML), while WLR *maximum likelihood estimates* (MLE) were obtained by using IRLS method of Newton-Raphson algorithm. In order to solve large optimization problems, Truncated-Newton method is applied.

Keywords: two-class classification, imbalanced data, logistic regression, rare event, RE-WLR

Presenting Author's biography



Dian Eka Apriana Sulasih earned her bachelor's degree from Sekolah Tinggi Ilmu Statistik Jakarta in 2007. She is an government employee at Badan Pusat Statistik. Currently, she is also a post-graduate student of Statistics in Institut Teknologi Sepuluh Nopember.

1. Introduction

Data classification is an important process in the field of data mining. Classification is the processing of finding a set of models (or functions) which describe and distinguish data classes or concepts, for the purposes of being able to use the model to predict the class of objects whose class label is unknown [1]. The methods often used for classification are discriminant analysis and logistic regression. In addition, many classification methods are used by approach to programming, such as Artificial Neural Network (ANN), Naive Bayes, Adaptive Classification Regression Tree (CART) and Support Vector Machine (SVM).

Linear Discriminant Analysis and Logistic Regression methods search for linear functions which are subsequently used for classification purposes. The use of linear functions enables better interpretation of the results by analyzing the value of the coefficients obtained. Not every classification method permits this type of analysis and, in fact, some are classified as “black box” models. Hence, Classic Discriminant Analysis and Logistic Regression continue to be interesting methodologies [2].

Linear classification is an extremely important machine-learning and data-mining tool. Compared to other classification techniques, such as the Kernel methods, which transform data into higher dimensional space, linear classifiers are implemented directly on data in their original space. The main advantage of linear classifiers is their efficient training and testing procedures, especially when implemented on large and high-dimensional data sets [3].

Besides discriminant analysis, logistic regression is a method of classification with linear classifiers that is often used. Logistic regression, which is a linear classifier, has been proven to be a powerful classifier by providing probabilities and by extending to multi-class classification problems [4,5]. Logistic regression has been extensively studied. Moreover, LR requires solving only unconstrained optimization problems. Hence, with the right algorithms, the computation time can be much less than that of other methods, such as Support Vector Machines (SVM), which require solving a constrained quadratic optimization problem [6].

One of the problems in data classification is the composition of the data that is out of balance (imbalanced data). In the binary or two-class classification, one class has a greater number of samples than the other class. The majority is negative class while the minority is positive class. The problem that occurs is a good prediction accuracy of the negative class and poor prediction accuracy for positive class. In other words, classifier tends to predict class which has more data composition.

Problems of imbalanced data occur in various fields such as information retrieval and filtering [7], detection of oil spills from satellite imagery [8], medical diagnosis [9], text classification [10], credit card fraud detection [11], telecommunications [12], and others.

Most of algorithm is more preoccupied in classifying major sample and ignoring or misclassifying minor sample. The minor samples are those that rarely occur but very important. There are various methods available for classification of imbalance data set which is divided into three main categories, the algorithmic approach, data preprocessing approach, and feature selection approach. Each of this technique has their own advantages and disadvantages [13].

Some Logistic Regression method development has been done to improve the accuracy of classification in imbalanced data. Maalouf and Trafalis [14] developed a method of Rare Event Weighted Kernel Logistic Regression (RE-WKLR) that is suitable for small to medium-sized data. Rahayu [15] developed a method of AdaBoost Newton Truncated Regularized Weighted Kernel Logistic Regression (AB-WKLR) and AdaBoost NTR Weighted Regularized Logistic Regression (AB-WLR) that significantly increase the performance of the accuracy and stability of general classifiers at NTR-KLR and NTR-LR, Furthermore, Maalouf and Siddiqi [6] developed a method of Rare Event Weighted Logistic Regression (RE-WLR) for the classification of imbalanced data on large-scale data. The study concluded that the RE-WLR has better performance than Truncated-Regularized Iteratively Re-weighted Least Squares (TR-IRLS). Wang, Xu and Zhou [16] used Lasso-Logistic Regression to evaluate unbalanced credit scoring.

This study showed the review of RE-WLR developed by Maalouf and Siddiqi [6] for the classification of imbalanced data. It mainly explicated the steps to obtain estimator specifically, particularly for IRLS.

2. Logistic Regression

Let $\mathbf{X} \in \mathbb{R}^{N \times d}$ be a data matrix where N is the number of instances (examples) and d is the number of features (parameters or attributes), and \mathbf{y} be a binary outcomes vector. For every instance $\mathbf{x}_i \in \mathbb{R}^d$ (a row vector in \mathbf{X}), where $i = 1 \dots N$, the outcome is either $y_i = 1$ or $y_i = 0$. Let the instances with outcomes of $y_i = 1$ belong to the positive class occurrence of an event), and the instances with outcomes $y_i = 0$ belong to the negative class (non occurrence of an event). The goal is to classify the instance \mathbf{x}_i as positive or negative. An instance can be treated as a Bernoulli trial with an expected value $E(y_i)$ or probability p_i . The logistic function commonly used to model each instance \mathbf{x}_i with its expected outcome is given by the following formula [17]:

$$E(y_i | \mathbf{x}_i, \boldsymbol{\beta}) = p_i = \frac{e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \quad (1)$$

where $\boldsymbol{\beta}$ is the vector of parameters with the assumption that $x_{i0} = 1$ so that the intercept β_0 is a constant term. From then on, the assumption is that the intercept is included in the vector $\boldsymbol{\beta}$.

The logistic (logit) transformation is the logarithm of the odds of the positive response and is defined as

$$\boldsymbol{\eta}_i = \ln\left(\frac{p_i}{1 - p_i}\right) = \mathbf{x}_i \boldsymbol{\beta} \quad (2)$$

In matrix form, the logit function is expressed as

$$\boldsymbol{\eta} = \mathbf{X} \boldsymbol{\beta} \quad (3)$$

3. Regularized Logistic Regression

There are two general methods of parameter estimation. They are least-squares estimation (LSE) and maximum likelihood estimation (MLE). This study uses Maximum Likelihood (ML). The ML method is based on the joint probability density of the observed data, and acts as a function of the unknown parameters in the model [18]. Recall that the outcome y is a Bernoulli random variable with mean p_i in the LR model. Therefore we may interpret the expectation function as the probability that $y = 1$, or equivalently that \mathbf{x}_i belongs to the positive class. Thus we may compute the probability of the i -th experiment and outcome in the dataset \mathbf{X}, \mathbf{y} as

$$P(\mathbf{x}_i, y_i | \boldsymbol{\beta}) = \begin{cases} p_i & , \text{if } y = 1 \\ 1 - p_i & , \text{if } y = 0 \end{cases} \quad (4)$$

$$= (p_i)^{y_i} (1 - p_i)^{(1 - y_i)} \quad (5)$$

The likelihood function is

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n (p_i)^{y_i} (1 - p_i)^{(1 - y_i)} = \prod_{i=1}^n \left(\frac{e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{y_i} \left(\frac{1}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{(1 - y_i)} \quad (6)$$

and hence, the log-likelihood is then

$$\ln L(\boldsymbol{\beta}) = \sum_{i=1}^n \left(y_i \ln \left(\frac{e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right) + (1 - y_i) \ln \left(\frac{1}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right) \right) \quad (7)$$

There is no closed form solution to maximize $\ln L(\beta)$ with respect to β . Therefore, LR *maximum likelihood estimates* (MLE) are obtained using numerical optimization methods. One of the most commonly used numerical methods is the Newton-Raphson method, for which, both the gradient vector and the Hessian matrix are needed. The gradient vector is obtained from the first derivatives of the log likelihood and the Hessian matrix is the second derivatives.

The first derivatives with respect to β are given by

$$\frac{\partial}{\partial \beta_j} \ln L(\beta) = \sum_{i=1}^n \left(y_i \left(\frac{x_{ij}}{1 + e^{x_i \beta}} \right) + (1 - y_i) \left(\frac{-x_{ij} e^{x_i \beta}}{1 + e^{x_i \beta}} \right) \right) \quad (8)$$

$$= \sum_{i=1}^n \left(y_i x_{ij} \left(\frac{1}{1 + e^{x_i \beta}} \right) - (1 - y_i) x_{ij} \left(\frac{e^{x_i \beta}}{1 + e^{x_i \beta}} \right) \right) \quad (9)$$

$$= \sum_{i=1}^n (y_i x_{ij} (1 - p_i) - (1 - y_i) x_{ij} p_i) \quad (10)$$

$$= \sum_{i=1}^n (y_i x_{ij} - y_i x_{ij} p_i - x_{ij} p_i + y_i x_{ij} p_i) \quad (11)$$

$$= \sum_{i=1}^n (x_{ij} (y_i - p_i)) \quad (12)$$

where $j = 0, \dots, d$ and d is the number of parameters. Each of the partial derivatives is then set to zero. In matrix form, equation (12) is written as

$$\mathbf{X}^T (\mathbf{y} - \mathbf{p}) = 0 \quad (13)$$

Now, the second derivatives with respect to β are given by

$$\frac{\partial^2}{\partial \beta_j \partial \beta_k} \ln L(\beta) = \sum_{i=1}^n \left(\frac{-x_{ij} x_{ik} e^{x_i \beta}}{(1 + e^{x_i \beta})^2} \right) \quad (14)$$

$$= \sum_{i=1}^n (-x_{ij} x_{ik} (p_i (1 - p_i))) \quad (15)$$

If v_i is defined as $p_i (1 - p_i)$ and $\mathbf{V} = \text{diag}(v_1, \dots, v_n)$ then the Hessian matrix is written as

$$\mathbf{H}(\beta) = \frac{\partial^2}{\partial^2 \beta} \ln L(\beta) = -\mathbf{X}^T \mathbf{V} \mathbf{X} \quad (16)$$

Over-fitting the training data may arise in LR [17], especially when the data are very high dimensional and/or sparse. One of the approaches to reduce over-fitting is through quadratic *regularization*, known also as *ridge regression*, which introduces a penalty for large values of β and to obtain better generalization [19]. The regularized log likelihood is defined as

$$\ln L(\beta) = \sum_{i=1}^n \left(y_i \ln \left(\frac{e^{x_i \beta}}{1 + e^{x_i \beta}} \right) + (1 - y_i) \ln \left(\frac{1}{1 + e^{x_i \beta}} \right) \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (17)$$

$$= \sum_{i=1}^n \left(y_i \left(\ln \left(\frac{e^{x_i \beta}}{1 + e^{x_i \beta}} \right) - \ln \left(\frac{1}{1 + e^{x_i \beta}} \right) \right) + \ln \left(\frac{1}{1 + e^{x_i \beta}} \right) \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (18)$$

$$= \sum_{i=1}^n \left(y_i \left(\ln(e^{x_i \beta}) - \ln(1 + e^{x_i \beta}) - \ln 1 + \ln(1 + e^{x_i \beta}) \right) + \ln 1 - \ln(1 + e^{x_i \beta}) \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (19)$$

$$= \sum_{i=1}^n \left(y_i (\ln(e^{x_i \beta}) - 0) + 0 - \ln(1 + e^{x_i \beta}) \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (20)$$

$$= \sum_{i=1}^n \left(\ln(e^{x_i \beta})^{y_i} - \ln(1 + e^{x_i \beta}) \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (21)$$

$$= \sum_{i=1}^n \ln \left(\frac{e^{y_i x_i \beta}}{1 + e^{x_i \beta}} \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (22)$$

$$= - \sum_{i=1}^n \ln \left(e^{-y_i x_i \beta} (1 + e^{x_i \beta}) \right) - \frac{\lambda}{2} \|\beta\|^2 \quad (23)$$

where the regularization (penalty) term $\frac{\lambda}{2} \|\beta\|^2$ was added to formula (5). The objective is then to find the Maximum Likelihood Estimate (MLE), $\hat{\beta}$, which maximizes the log likelihood. For binary outputs, the loss function or the deviance DEV is the negative log likelihood and is given by the formula

$$DEV(\hat{\beta}) = -2 \ln L(\hat{\beta}) \quad (24)$$

4. Iteratively Re-weighted Least Squares (IRLS)

An alternative to numerically maximizing the LR maximum likelihood equations is the *iteratively re-weighted least squares* (IRLS) technique. This technique uses the Newton-Raphson algorithm to solve the LR score equations. Each iteration finds the *weighted least squares* (WLS) estimates for a given set of weights, which are used to construct a new set of weights [18]. The gradient and the Hessian are obtained by differentiating the regularized likelihood in (22) with respect to β .

$$\frac{\partial}{\partial \beta} \ln L(\beta) = \mathbf{X}^T (\mathbf{y} - \mathbf{p}) - \lambda \beta = 0 \quad (25)$$

$$\frac{\partial^2}{\partial^2 \beta} \ln L(\beta) = -\mathbf{X}^T \mathbf{V} \mathbf{X} - \lambda \mathbf{I} = 0 \quad (26)$$

where \mathbf{I} is a $d \times d$ identity matrix. Now that the first and second derivatives are obtained, the Newton-Raphson update formula on the $(c+1)$ -th iteration is given by

$$\hat{\beta}^{(c+1)} = \hat{\beta}^{(c)} + (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T (\mathbf{y} - \mathbf{p}) - \lambda \hat{\beta}^{(c)}) \quad (27)$$

Since $\hat{\beta}^{(c)} = (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I}) \hat{\beta}^{(c)}$, then (27) can be rewritten as

$$\hat{\beta}^{(c+1)} = (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{V} \mathbf{X} \hat{\beta}^{(c)} + (\mathbf{y} - \mathbf{p})) \quad (28)$$

$$= (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{V} \mathbf{z}^{(c)} \quad (29)$$

where $\mathbf{z}^{(c)} = \mathbf{X} \hat{\beta}^{(c)} + \mathbf{V}^{-1} (\mathbf{y} - \mathbf{p})$ and is referred to as the adjusted response [4].

If the matrix $(\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I})$ were dense, the iterative computation could become unacceptably slow [20]. This necessitates the need for a “trade off” between convergence speed and accurate Newton direction [21]. The method which provides such a trade-off is known as the truncated Newton’s method.

4.1. Truncated Regularized Iteratively Re-weighted Least Squares (TR-IRLS)

The weighted least squares (WLS) sub problem is then

$$(\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I}) \hat{\boldsymbol{\beta}}^{(c+1)} = \mathbf{X}^T \mathbf{V} \mathbf{z}^{(c)} \quad (30)$$

is a linear system of d equations and variables, and solving it is equivalent to minimizing the quadratic function

$$\frac{1}{2} \hat{\boldsymbol{\beta}}^{(c+1)} (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I}) \hat{\boldsymbol{\beta}}^{(c+1)} - \hat{\boldsymbol{\beta}}^{(c+1)} \mathbf{X}^T \mathbf{V} \mathbf{z}^{(c)} \quad (31)$$

Komarek and Moore [22] used *truncated-regularized iteratively-reweighted least squares* (TR-IRLS) technique that implemented a modified linear Conjugate Gradient (CG) to approximate the Newton direction in solving the IRLS for LR.

Conjugate gradient (CG) is an iterative minimization algorithm. The main advantage of the CG method is that it guarantees convergence in at most d steps [21]. The TR-IRLS algorithm consists of two loops. Algorithm 1 represents the outer loop finds the solution to the WLS problem and is terminated when the relative difference of deviance between two consecutive iterations is no larger than a specified threshold ε_1 . Algorithm 2 represents the inner loop, which solves the WLS sub problems in Algorithm 1 through the linear CG method, which is the Newton direction. Algorithm 2 is terminated when the residual

$$\mathbf{r}^{(c+1)} = (\mathbf{X}^T \mathbf{V} \mathbf{X} + \lambda \mathbf{I}) \hat{\boldsymbol{\beta}}^{(c+1)} - \mathbf{X}^T \mathbf{V} \mathbf{z}^{(c)} \quad (32)$$

is no greater than a specified threshold ε_2 [6].

5. Logistic Regression in Rare Events Data

King and Zeng [23] recommend two methods of estimation for choice-based sampling, *prior correction* and *weighting*. This study uses weighting method.

Under pure endogenous sampling, the conditioning is on \mathbf{X} rather than \mathbf{y} [24, 25], and the joint distribution of \mathbf{y} and \mathbf{X} in the sample is

$$f_s(\mathbf{y}, \mathbf{X} | \boldsymbol{\beta}) = P_s(\mathbf{X} | \mathbf{y}, \boldsymbol{\beta}) P_s(\mathbf{y}) \quad (33)$$

where $\boldsymbol{\beta}$ is the unknown parameter to be estimated. Yet, since \mathbf{X} is a matrix of exogenous variables, then the conditional probability of \mathbf{X} in the sample is equal to that in the population, or $P_s(\mathbf{X} | \mathbf{y}, \boldsymbol{\beta}) = P(\mathbf{X} | \mathbf{y}, \boldsymbol{\beta})$. However, the conditional probability in the population is

$$P(\mathbf{X} | \mathbf{y}, \boldsymbol{\beta}) = \frac{f(\mathbf{y}, \mathbf{X} | \boldsymbol{\beta})}{P(\mathbf{y})} \quad (34)$$

but

$$f(\mathbf{y}, \mathbf{X} | \boldsymbol{\beta}) = P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) P(\mathbf{X}) \quad (35)$$

and hence, substituting and rearranging yields

$$f_s(\mathbf{y}, \mathbf{X} | \boldsymbol{\beta}) = \frac{P_s(\mathbf{y})}{P(\mathbf{y})} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) P(\mathbf{X}) \quad (36)$$

$$= \frac{H}{Q} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) P(\mathbf{X}) \quad (37)$$

where $\frac{H}{Q} = \frac{P_s(\mathbf{y})}{P(\mathbf{y})}$, with H representing the proportions in the sample and Q the proportions in the population. The likelihood is then

$$L_{Endogenous} = \prod_{i=1}^n \frac{H_i}{Q_i} P(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\beta}) P(\mathbf{x}_i) \tag{38}$$

where $\frac{H_i}{Q_i} = \left(\frac{\bar{y}}{\tau}\right) y_i + \left(\frac{1-\bar{y}}{1-\tau}\right) (1 - y_i)$, with \bar{y} is the proportion of events in the sample and τ is the proportion of events in the population. When dealing with REs and imbalanced data, the likelihood needs to be maximized.

The log-likelihood for LR can then be rewritten as

$$\ln L(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \sum_{i=1}^n \frac{Q_i}{H_i} \ln P(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\beta}) \tag{39}$$

$$= \sum_{i=1}^n \frac{Q_i}{H_i} \ln \left(\frac{e^{y_i \mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{y_i \mathbf{x}_i \boldsymbol{\beta}}} \right) \tag{40}$$

$$= \sum_{i=1}^n w_i \ln \left(\frac{e^{y_i \mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{y_i \mathbf{x}_i \boldsymbol{\beta}}} \right) \tag{41}$$

where $w_i = \frac{Q_i}{H_i}$. If the proportion of events in the sample is more than that in the population, then the ratio $\frac{Q_i}{H_i}$ is less than one, and hence the events are given less weight, while the non-events would be given more weight if their proportion in the sample is less than that in the population. The above estimator, however, is not fully efficient, because the information matrix equality does not hold. This is demonstrated as

$$-E \left[\frac{Q}{H} \nabla_{\boldsymbol{\beta}}^2 P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) \right] \neq E \left[\left(\frac{Q}{H} \nabla_{\boldsymbol{\beta}} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) \right) \left(\frac{Q}{H} \nabla_{\boldsymbol{\beta}} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) \right)^T \right] \tag{42}$$

and for the LR model it is

$$- \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{Q_i}{H_i} \right) p_i (1 - p_i) \mathbf{x}_i \mathbf{x}_j \right] \neq \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{Q_i}{H_i} \right)^2 p_i (1 - p_i) \mathbf{x}_i \mathbf{x}_j \right] \tag{43}$$

Let $\mathbf{A} = \frac{1}{n} \sum_{i=1}^n \left(\frac{Q_i}{H_i} \right) p_i (1 - p_i) \mathbf{x}_i \mathbf{x}_j$, and $\mathbf{B} = \frac{1}{n} \sum_{i=1}^n \left(\frac{Q_i}{H_i} \right)^2 p_i (1 - p_i) \mathbf{x}_i \mathbf{x}_j$, then the asymptotic variance matrix of the estimator $\boldsymbol{\beta}$ is given by the *sandwich estimate*, such that $\mathbf{V}(\boldsymbol{\beta}) = \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}$ [26, 27, 28].

King and Zeng [23] extended the small-sample bias corrections, as described by McCullagh and Nelder [29] to include the weighted likelihood (41). According to McCullagh and Nelder [29], and later Cordeiro and McCullagh [30], the bias vector is given by

$$\text{bias}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V} \boldsymbol{\xi} \tag{44}$$

where \mathbf{Q}_{ii} are the diagonal elements of $\mathbf{Q} = \mathbf{X}(\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T$ which is the approximate covariance matrix of the logistic link function. The second-order bias-corrected estimator is then

$$\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - \mathbf{B}(\hat{\boldsymbol{\beta}}) \tag{45}$$

As for the variance matrix $\mathbf{V}(\tilde{\boldsymbol{\beta}})$ of $\tilde{\boldsymbol{\beta}}$, it is estimated using

$$\mathbf{v}(\tilde{\boldsymbol{\beta}}) = \left(\frac{n}{n+d} \right)^2 \mathbf{v}(\hat{\boldsymbol{\beta}}) \tag{46}$$

Since $\left(\frac{n}{n+d} \right)^2 < 1$, then $\mathbf{V}(\tilde{\boldsymbol{\beta}}) < \mathbf{V}(\hat{\boldsymbol{\beta}})$, and hence both the variance and the bias are now reduced.

6. Rare Event Weighted Logistic Regression (RE-WLR)

Now the formulation of the weighted LR suggested by King and Zeng [23] is applied to the WLR model in (41), the weighted likelihood can be rewritten as

$$L_w(\beta) = \prod_{i=1}^n (p_i)^{w_1 y_i} (1 - p_i)^{w_0 (1 - y_i)} \tag{47}$$

where $w_1 = \frac{\tau}{\bar{y}}$, $w_0 = \frac{1-\tau}{1-\bar{y}}$, $p_i = \frac{1}{1+e^{-\eta_i}}$, and $(1 - p_i) = \frac{e^{-\eta_i}}{1+e^{-\eta_i}}$.

Now,

$$p_i = \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} = p_i^{w_1} \tag{48}$$

and hence,

$$p_i' = \frac{\partial p_i}{\partial \eta_i} = w_1 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1 - 1} (-1) (1 + e^{-\eta_i})^{-2} (-e^{-\eta_i}) \tag{49}$$

$$= w_1 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} (1 + e^{-\eta_i}) (1 + e^{-\eta_i})^{-2} (e^{-\eta_i}) \tag{50}$$

$$= w_1 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\frac{e^{-\eta_i}}{1 + e^{-\eta_i}} \right) \tag{51}$$

$$= w_1 p_i^{w_1} (1 - p_i) \tag{52}$$

$$p_i'' = \frac{\partial^2 p_i}{\partial \eta_i^2} = w_1^2 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\frac{e^{-\eta_i}}{1 + e^{-\eta_i}} \right)^2 + w_1 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\frac{-e^{-\eta_i}}{1 + e^{-\eta_i}} + \frac{e^{-\eta_i}(-e^{-\eta_i})}{(1 + e^{-\eta_i})^2} \right) \tag{53}$$

$$= w_1^2 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\frac{e^{-\eta_i}}{1 + e^{-\eta_i}} \right)^2 + w_1 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\frac{e^{-\eta_i}(-1 - e^{-\eta_i} + e^{-\eta_i})}{(1 + e^{-\eta_i})^2} \right) \tag{54}$$

$$= w_1^2 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\frac{e^{-\eta_i}}{1 + e^{-\eta_i}} \right)^2 + w_1 \left(\frac{1}{1 + e^{-\eta_i}} \right)^{w_1} \left(\left(\frac{e^{-\eta_i}}{1 + e^{-\eta_i}} \right) \left(-\frac{1}{1 + e^{-\eta_i}} \right) \right) \tag{55}$$

$$= w_1 p_i^{w_1} (1 - p_i) w_1 (1 - p_i) + w_1 p_i^{w_1} (1 - p_i) (-p_i) \tag{56}$$

$$= w_1 p_i^{w_1} (1 - p_i) (w_1 - w_1 p_i) + w_1 p_i^{w_1} (1 - p_i) (-p_i) \tag{57}$$

$$= w_1 p_i^{w_1} (1 - p_i) (w_1 - (1 + w_1) p_i) \tag{58}$$

Finally, the bias vector for WLR can now be rewritten as

$$\mathbf{B}(\hat{\beta}) = (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \boldsymbol{\xi} \tag{59}$$

where the i-th element of the vector $\boldsymbol{\xi}$ is now

$$\xi_i = 0,5 Q_{ii} ((1 + w_1) p_i - w_1) \tag{60}$$

with Q_{ii} as the diagonal elements of \mathbf{Q} , which is now $\mathbf{Q} = \mathbf{X}(\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T$, and $\mathbf{D} = \text{diag}(v_i w_i)$ for $i = 1 \dots n$. The bias-corrected estimator becomes

$$\tilde{\beta} = \hat{\beta} - \mathbf{B}(\hat{\beta}) \tag{61}$$

Iteratively re-weighted least squares (IRLS) method is used to find the MLE of β , which uses Newton-Raphson algorithm to solve LR score equations. Each iteration finds the *weighted least squares* (WLS)

estimates for a given set of weights, which are used to construct a new set of weights. For WLR, the gradient and Hessian are obtained by differentiating the regularized weighted log-likelihood,

$$\ln L_w(\boldsymbol{\beta}) = \sum_{i=1}^n w_i \ln \frac{e^{y_i \mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} - \frac{\lambda}{2} \|\boldsymbol{\beta}\|^2 \quad (62)$$

with respect to $\boldsymbol{\beta}$. In matrix form, the gradient is

$$\frac{\partial}{\partial \beta_j} \ln L_w(\boldsymbol{\beta}) = \sum_{i=1}^n w_i \left(\left(\frac{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}}{e^{y_i \mathbf{x}_i \boldsymbol{\beta}}} \right) \left(\frac{y_i \mathbf{x}_{ij} e^{y_i \mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} + \frac{e^{y_i \mathbf{x}_i \boldsymbol{\beta}} (-\mathbf{x}_{ij} e^{\mathbf{x}_i \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i \boldsymbol{\beta}})^2} \right) \right) - \lambda \beta_j \quad (63)$$

$$= \sum_{i=1}^n w_i \left(\left(\frac{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}}{e^{y_i \mathbf{x}_i \boldsymbol{\beta}}} \right) \left(\frac{y_i \mathbf{x}_{ij} e^{y_i \mathbf{x}_i \boldsymbol{\beta}} (1 + e^{\mathbf{x}_i \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i \boldsymbol{\beta}})^2} + \frac{e^{y_i \mathbf{x}_i \boldsymbol{\beta}} (-\mathbf{x}_{ij} e^{\mathbf{x}_i \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i \boldsymbol{\beta}})^2} \right) \right) - \lambda \beta_j \quad (64)$$

$$= \sum_{i=1}^n w_i \left(\left(\frac{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}}{e^{y_i \mathbf{x}_i \boldsymbol{\beta}}} \right) \left(\frac{\mathbf{x}_{ij} e^{y_i \mathbf{x}_i \boldsymbol{\beta}} (y_i (1 + e^{\mathbf{x}_i \boldsymbol{\beta}}) - e^{\mathbf{x}_i \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i \boldsymbol{\beta}})^2} \right) \right) - \lambda \beta_j \quad (65)$$

$$= \sum_{i=1}^n w_i \left(\mathbf{x}_{ij} \left(\frac{y_i (1 + e^{\mathbf{x}_i \boldsymbol{\beta}}) - e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} \right) \right) - \lambda \beta_j \quad (66)$$

$$= \sum_{i=1}^n w_i \mathbf{x}_{ij} (y_i - p_i) - \lambda \beta_j \quad (67)$$

$$\nabla_{\boldsymbol{\beta}} \ln L_w(\boldsymbol{\beta}) = \mathbf{X}^T \mathbf{W} (\mathbf{y} - \mathbf{p}) - \lambda \boldsymbol{\beta} \quad (68)$$

where $\mathbf{W} = \text{diag}(w_i)$ and \mathbf{p} is the probability vector. The Hessian with respect to $\boldsymbol{\beta}$ is then

$$\frac{\partial^2}{\partial \beta_j \partial \beta_k} \ln L_w(\boldsymbol{\beta}) = \sum_{i=1}^n -w_i \mathbf{x}_{ij} \mathbf{x}_{ik} \left(\frac{e^{\mathbf{x}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i \boldsymbol{\beta}}} + \frac{e^{\mathbf{x}_i \boldsymbol{\beta}} (-e^{\mathbf{x}_i \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i \boldsymbol{\beta}})^2} \right) - \lambda \mathbf{I} \quad (69)$$

$$= \sum_{i=1}^n -w_i \mathbf{x}_{ij} \mathbf{x}_{ik} \left(\frac{e^{\mathbf{x}_i \boldsymbol{\beta}} (1 + e^{\mathbf{x}_i \boldsymbol{\beta}} - e^{\mathbf{x}_i \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i \boldsymbol{\beta}})^2} \right) - \lambda \mathbf{I} \quad (70)$$

$$= \sum_{i=1}^n (-\mathbf{x}_{ij} \mathbf{x}_{ik} w_i (p_i (1 - p_i))) - \lambda \mathbf{I} \quad (71)$$

$$\nabla_{\boldsymbol{\beta}}^2 \ln L_w(\boldsymbol{\beta}) = -\mathbf{X}^T \mathbf{D} \mathbf{X} - \lambda \mathbf{I} \quad (72)$$

The first and second derivatives are obtained, the Newton-Raphson methods iterates via

$$\hat{\boldsymbol{\beta}}^{(c+1)} = \hat{\boldsymbol{\beta}}^{(c)} - (\nabla_{\boldsymbol{\beta}}^2 \ln L_w(\boldsymbol{\beta}))^{-1} (\nabla_{\boldsymbol{\beta}} \ln L_w(\boldsymbol{\beta})) \quad (73)$$

$$\hat{\boldsymbol{\beta}}^{(c+1)} = \hat{\boldsymbol{\beta}}^{(c)} + (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} (\mathbf{y} - \mathbf{p}) - \lambda \hat{\boldsymbol{\beta}}^{(c)}) \quad (74)$$

Since $\hat{\boldsymbol{\beta}}^{(c)} = (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} (\mathbf{y} - \mathbf{p}) - \lambda \hat{\boldsymbol{\beta}}^{(c)})$, then (74) can be rewritten as

$$\hat{\boldsymbol{\beta}}^{(c+1)} = (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I}) \hat{\boldsymbol{\beta}}^{(c)} + (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} (\mathbf{y} - \mathbf{p}) - \lambda \hat{\boldsymbol{\beta}}^{(c)}) \quad (75)$$

$$\hat{\boldsymbol{\beta}}^{(c+1)} = (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{D} (\mathbf{X} \hat{\boldsymbol{\beta}}^{(c)} + (\mathbf{X}^T)^{-1} \mathbf{D}^{-1} \lambda \hat{\boldsymbol{\beta}}^{(c)} + \mathbf{D}^{-1} \mathbf{W} (\mathbf{y} - \mathbf{p}) - (\mathbf{X}^T)^{-1} \mathbf{D}^{-1} \lambda \hat{\boldsymbol{\beta}}^{(c)}) \quad (76)$$

$$\hat{\boldsymbol{\beta}}^{(c+1)} = (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{D} (\mathbf{X} \hat{\boldsymbol{\beta}}^{(c)} + \mathbf{D}^{-1} \mathbf{W} (\mathbf{y} - \mathbf{p})) \quad (77)$$

The Newton–Raphson update with respect to β on the $(c+1)$ th iteration is

$$\hat{\beta}^{(c+1)} = (\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{z}^{(c)} \quad (78)$$

where $\mathbf{z}^{(c)} = \mathbf{X} \hat{\beta}^{(c)} + \mathbf{D}^{-1} \mathbf{W}(\mathbf{y} - \mathbf{p})$ is the adjusted dependent variable or the adjusted response.

The weighted least squares (WLS) sub problem is then

$$(\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I}) \hat{\beta}^{(c+1)} = \mathbf{X}^T \mathbf{D} \mathbf{z}^{(c)} \quad (79)$$

which is a system of linear equations with a matrix \mathbf{X} , a vector of adjusted responses \mathbf{z} , and a weight matrix \mathbf{D} . Both the weights and the adjusted response vector are dependent on $\hat{\beta}^{(c)}$, which is the current estimate of the parameter vector. Specifying an initial estimate $\hat{\beta}^{(0)}$ for $\hat{\beta}$ can be solved iteratively, giving a sequence of estimates that converges to the MLE of $\hat{\beta}$. This iterative process can be done using the CG method. Like the TR-IRLS algorithm, in order to avoid the long computations that the CG may suffer from, a limit can be placed on the number of CG iterations, thus creating an approximate or truncated Newton direction [6].

7. Conclusion

Rare Event Weighted Logistic Regression (RE-WLR) is a combination of LR rare events correction [23] and TR-IRLS [22]. Rare event correction in LR is applied by setting weighting to LR with the result that Weighted Logistic Regression (WLR) will be formed. Regularization was added to reduce over-fitting. The estimation of β was performed by using the method of *maximum likelihood* (ML), but there was no closed form solution to maximize $\ln L(\beta)$ with respect to β . Therefore WLR *maximum likelihood estimates* (MLE) was obtained by using IRLS method of Newton-Raphson algorithm to solve WLR score equations. Each of iterations found the *weighted least squares* (WLS) estimating to give set of weights, which were used to construct a new set of weights. The gradient and the Hessian were obtained by differentiating the regularized weighted likelihood with respect to β . In order to solve large optimization problems, Truncated-Newton method was applied.

References

- [1] J. Han and M. Kamber. *Data Mining Concepts and Techniques*. Academic Press, 2001.
- [2] J. Pacheco, S. Casado, and L. Núñez. A variable selection method based on Tabu search for logistic regression models. *European Journal of Operational Research*, 199:506–511, 2009.
- [3] G.-X. Yuan, C.-H. Ho, C.-J. Lin. Recent advances of large-scale linear classification. *Proc. IEEE 100* (9), page 2584–2603, 2012.
- [4] T. Hastie, R. Tibshirani and J. Friedman. *The Elements of Statistical Learning*. Springer Verlag, 2001.
- [5] P. Karsmakers, K. Pelckmans and J.A.K. Suykens. Multi-class kernel logistic regression: a fixed-size implementation. *Int. Joint Conf. Neural Network*, page 1756–1761, 2007.
- [6] M. Maalouf and M. Siddiqi. Weighted logistic regression for large-scale imbalanced and rare events data. *Journal of Knowledge-Based Systems*, 59:141–148, 2014.
- [7] D. Lewis and J. Carlett. Heterogeneous Uncertainly Sampling for Supervised Learning. *Proceedings of ICML-94, 11th International Conference on Machine Learning*, Eds: Cohen, W. dan Hirsh, H., Morgan Kaufmann, San Fransisco, page 148–156, 1994.
- [8] M. Kubat, S. Matwin, and R. Holte. Machine Learning for the Detection of Oil Spills In Satellite Radar Images. *Machine Learning*, 30:195-215, 1998.
- [9] I. Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. *Artificial Intelligence in Medicine*, 23:89–109, 2001.
- [10] Chawla, et.al. SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence and Research*, 16:321-357, 2002.

- [11] G. Wu and E. Chang. Class-Boundary Alignment for Imbalanced Dataset Learning. *ICML 2003 Workshop on Learning from Imbalanced Data Sets II*, Washington, 2003.
- [12] C.S. Hilar. Designing Expert System for Fraud Detection in Private Telecommunication Networks. *Expert Systems with Applications*, 36 (9):11559-11569, 2009.
- [13] R. Longadge, S. Dongre, and L. Malik. Class Imbalance Problem in Data Mining: Review. *International Journal of Computer Science and Network*, Vol. 2, 2013.
- [14] Maalouf and Trafalis. Robust weighted kernel logistic regression in imbalanced and rare events data. *Computational Statistics and Data Analysis 2011*, 55:168–183, 2010.
- [15] S.P. Rahayu. *Logistic regression methods for classification of imbalanced data*. Tesis Ph.D, University Malaysia Pahang (UMP), Pahang, 2012.
- [16] Wang, Xu, and Zhou. Large Unbalanced Credit Scoring Using Lasso-Logistic Regression Ensemble. *PLoS ONE 10(2): e0117844*, 2015.
- [17] D.W. Hosmer and S. Lemeshow. *Applied Logistic Regression, second ed.*. Wiley, 2000.
- [18] P. Garthwaite, I. Jolliffe, and J. Byron, J. *Statistical Inference*. Oxford University Press, 2002.
- [19] C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [20] P. Komarek. *Logistic regression for data mining and high-dimensional classification*. Ph.D. thesis, Carnegie Mellon University, 2004.
- [21] J.M. Lewis, S. Lakshminarayanan, S. Dhall. *Dynamic Data Assimilation: A Least Squares Approach*. Cambridge University Press, 2006.
- [22] P. Komarek and A. Moore, A. *Making logistic regression a core data mining tool: A practical investigation of accuracy, speed, and simplicity*. Tech. rep., Carnegie Mellon University, 2005.
- [23] G. King and L. Zeng. Logistic regression in rare events data. *Political Analysis*, 9:137–163, 2001.
- [24] A.C. Cameron, and P.K. Trivedi. *Microeconometrics: Methods and Applications*. Cambridge University Press, 2005.
- [25] M. Milgate, J. Eatwell, and P.K. Newma. *Econometrics*. W. W. Norton & Company (1990)
- [26] T. Amemiya. *Advanced Econometrics*. Harvard University Press, 1985.
- [27] Y. Xie and C.F. Manski. The logit model and response-based samples. *Sociological Methods & Research*, 17:283–302, 1989.
- [28] C.F. Manski and S.R. Lerman. The estimation of choice probabilities from choice based samples. *Econometrica*, 45:1977–1988, 1977,
- [29] P. McCullagh and J. Nelder. *Generalized Linear Model*. Chapman and Hall/CRC, 1989.
- [30] G.M. Cordeiro and P. McCullagh. Bias correction in generalized linear models, *Journal of Royal Statistical Society*, 53(3):629–643, 1991.