# Logistic and Probit Regression Computational Algorithms

Authored by: Kyle Kolsti, PhD

Kaity Jones

Alex McBride

6 August 2021



HSCoBP@afit.edu 937-255-3636 x 4768

The HS CoBP core mission involves leading a consortium of government, industry, and academic experts to assess future homeland threats to inform strategic plans across nine DHS T&E capability areas.

# **Table of Contents**

Executive Summary3
Introduction
Logistic Regression Development3
Data Structure3
Model4
Likelihood4
Firth Bias Correction5
Determining Parameters5
Newton Iterations5
Matrix Formulas7
Uncertainty
Example9
Conclusion10
Works Cited10

## **Executive Summary**

This report provides the computational steps to employ logistic regression and probit regression under the construct of a generalized linear model (GLM). Matrix equations and pseudo-code facilitate the development of scripts to determine the coefficients of maximum likelihood through Newton-Raphson iteration method.

Keywords: logistic regression, maximum likelihood, Newton-Raphson

## Introduction

Many systems encountered in the Department of Homeland Security (DHS) are black box systems which provide only binary response data to the test team during Test and Evaluation (T&E). Logistic regression is a widely-used method to analyze binary response data that provides the probability of observing one of the two response values given certain values of a continuous factor (Natoli, Burke, & Oimoen, 2020). A common example is a sensor where the probability of detection depends on the range to the item of interest. Probit regression is a related method applicable to the same scenario that may be encountered is some fields such as medical research. Both logistic and probit regression can be employed under the construct of a generalized linear model (GLM) which is available in most statistical analysis software packages. On occasion, a test team may need to script these methods rather than relying on commercially available software products. The following best practice provides step by step computations for practitioners to improve understanding of these important analytical methods and to implement them in custom code if necessary.

# **Logistic Regression Development**

#### **Data Structure**

Data in a system with a binary outcome can be expressed as a vector of n binary outcomes Y ( $Y_i = 0$  or 1). For this report we will define  $Y_i = 1$  as a successful outcome such as a detection. Each outcome  $Y_i$  was observed at a value of the single predictor variable, the continuous factor x. (The matrix methods developed here easily extend to multiple factors if needed.) The  $n \times 2$  design matrix X contains 1's in the first column and the values of x used in the trials in the second column. Thus the first trial was performed with x set to  $x_1$  and the observed outcome was  $Y_1$ , and so forth.

$$\boldsymbol{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \boldsymbol{Y} = \begin{cases} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{cases}$$
(1)

#### Model

To model binary response data, we begin by defining  $\pi_i = \pi(x_i)$  as the prediction of the probability of observing Y = 1 at  $x = x_i$ . We also define two model parameters  $\beta_i$  such that

$$y_i = y(x_i) = \beta_0 + \beta_1 x_i \tag{2}$$

In matrix form,  $y = X\beta$  where the parameter vector  $\beta = {\beta_0, \beta_1}^T$ . The generalized linear model (GLM) relates  $\pi$  to y through a function called the link function. For binary responses, two common link functions are the logit link function, based on the odds, and the probit link function, which is based on the normal distribution (Agresti, 2002). The link function g and inverse link function  $g^{-1}$  for logit and probit regression are shown below in Table 1. The symbol  $\Phi$  is the cumulative distribution function for the standard normal distribution, which has a mean of zero and a standard deviation of one.

Link type	Link Function $y = g(\pi)$	Inverse Link Function $\pi = g^{-1}(y)$
Logit	$y_i = \log\left(\frac{\pi_i}{1 - \pi_i}\right)$	$\pi_i = \frac{1}{1 + e^{-y_i}}$
Probit	$y_i = \Phi^{-1}(\pi_i)$	$\pi_i = \Phi(y_i)$

Table 1: Link and inverse link functions

#### Likelihood

There is no closed-form solution to determine the parameters  $\beta$  in logistic or probit regression as there is for standard linear regression, so the model best fit must be determined using a nonlinear solver with the Maximum Likelihood Estimate (MLE) method. The likelihood of a pair of parameter values  $\beta$  is the product of the probabilities of obtaining each of the observed outcomes in Y. For example, for a single trial at some value of x, if  $\pi(\beta) = 0.7$  then the probability of obtaining Y = 1 is 0.7 and the probability of obtaining Y = 0 is 0.3; thus if actual observation were Y = 1, the likelihood of  $\beta$  for that trial is 0.7. The formula for the likelihood of a pair of parameter values  $\beta$  as a product of the probabilities for each trial is

$$\prod_{i=1}^{n} \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i} \tag{3}$$

For better computational behavior we will deal with the natural logarithm or "log likelihood." Recall that since  $Y_i$  is either 0 or 1, one of the two terms inside the summation will vanish for each data point.

$$L(\beta_0, \beta_1) = \sum_{i=1}^n \{ Y_i \log(\pi_i) + (1 - Y_i) \log(1 - \pi_i) \}$$
(4)

Now we can fit the statistical model to the data using MLE; i.e., find  $\beta_0$  and  $\beta_1$  such that the log likelihood L is at its maximum value. There is a wealth of numerical methods that can be used to minimize – L, which provides the same solution as maximizing L. These methods, which include the conjugate gradient method, steepest descent method, BFGS (Broyden–Fletcher–Goldfarb–Shanno algorithm, and differential evolution, are available for most scripting languages. As an alternative, the following section will provide a root-finding method which may be computationally faster than these minimization procedures.

## **Firth Bias Correction**

If there is any value of x below which all the Y values are 0 and above which all the Y values are 1 (or vice versa), the likelihood function does not have a maximum value and the iterative procedure therefore fails. This behavior is what is known as complete separation (Zorn, 2005). An example of a data set which has complete separation is shown in Figure 1.



Figure 1: Example data set exhibiting complete separation.

The Firth correction is a commonly used and effective way to address this situation. It is performed by adding an additional term to the log likelihood equation. The term equals one half of the natural logarithm of the determinant of the expected Fisher Information Matrix, *I* (Zorn, 2005), which will be defined in subsequent sections of this paper.

$$L(\boldsymbol{\beta}) = \sum_{i=1}^{n} \{ Y_i \log(\pi_i) + (1 - Y_i) \log(1 - \pi_i) \} + \frac{1}{2} \log|\boldsymbol{I}(\boldsymbol{\beta})|$$
(5)

# **Determining Parameters**

### **Newton Iterations**

This section will provide the step by step instructions for fitting the GLM using the Newton-Raphson method (Agresti, 2002). Like standing at the top of a hill, the maximum likelihood occurs where the slope in every direction is zero. Mathematically, the solution for  $\beta$  is where the partial derivatives of the

log likelihood function with respect to the parameters  $\beta_0$  and  $\beta_1$  are all equal to zero. The vector of the partial derivatives is called the gradient vector. Since there are two parameters in  $\beta$ , the gradient vector is of shape 2 × 1.

$$\nabla L(\beta_0, \beta_1) = \begin{cases} \frac{\partial L}{\partial \beta_0} \\ \frac{\partial L}{\partial \beta_1} \end{cases}$$
(6)

If a guess for  $\beta$  does not satisfy  $\nabla L(\beta) = 0$ , an update has to be made to  $\beta$  to nudge it closer to the solution. This update requires computation of the symmetric 2 × 2 Hessian matrix, H. The Hessian describes the local curvature of the surface, as indicated by the second partial derivatives that comprise it.

$$\boldsymbol{H}(\boldsymbol{\beta}_{0},\boldsymbol{\beta}_{1}) = \begin{bmatrix} \frac{\partial^{2}L}{\partial\boldsymbol{\beta}_{0}^{2}} & \frac{\partial^{2}L}{\partial\boldsymbol{\beta}_{0}\partial\boldsymbol{\beta}_{1}} \\ \frac{\partial^{2}L}{\partial\boldsymbol{\beta}_{0}\partial\boldsymbol{\beta}_{1}} & \frac{\partial^{2}L}{\partial\boldsymbol{\beta}_{1}^{2}} \end{bmatrix}$$
(7)

Table 2 contains the steps for the Newton-Raphson iterations. The gradient vector and Hessian matrix are fully defined in the next section. After obtaining the solution, the model fit may be assessed using the methods detailed in (Allison, 2014).

#### Table 2: Pseudo-code for Newton-Raphson iterations

- 1. Set k = 0.
  - 2. Start with initial guess for the vector of parameters,  $\beta_k$ . It is usually sufficient to make an initial guess of  $\beta_0 = \{0,0\}^T$ .
  - 3. Calculate the gradient vector,  $\nabla L(\boldsymbol{\beta}_k)$
  - 4. Calculate the Hessian,  $H(\boldsymbol{\beta}_k)$ .
  - 5. Calculate the update vector for the parameters:  $\Delta \beta = -H^{-1} \nabla L$ . It is better computationally to avoid inverting the Hessian matrix and instead solve the linear system:

$$H\Delta\beta = -\nabla L$$

- 6. Update the parameter vector to the latest guess.  $\beta_{k+1} = \beta_k + \Delta \beta$
- 7. Check convergence. One way to declare convergence is when the largest change magnitude in  $\beta$  (hence biggest absolute value in  $\Delta\beta$ ) was smaller than a tolerance. Another way is to check that the  $L_2$  norm of the residual vector  $\nabla L$  is below a preset tolerance; i.e., both values of the residual vector are nearly zero. a. If converged: STOP.
  - b. If not converged: k = k + 1; Go to Step 3 and perform another Newton iteration.

## **Matrix Formulas**

This section provides the recipes for the matrices used in the Newton iteration algorithm. For reference the matrix sizes are  $X_{n\times 2}$ ,  $Y_{n\times 1}$ ,  $\pi_{n\times 1}$ ,  $Y_{n\times 1}$ ,  $H_{2\times 2}$ ,  $\nabla L_{2\times 1}$ ,  $W_{n\times n}$ , and  $Z_{n\times n}$  where n is the number of trials that were conducted. The constituent matrices are calculated using differently for logistic and probit regression, as shown in Table 4. Note that the function f is the standard normal probability distribution function ( $\mu = 0$ ,  $\sigma = 1$ ). The formula for Z may be found in (Agresti, 2002), Problem 6.32.

Link Function	Matrix Definitions	
Logistic	$\boldsymbol{W} = \text{diag}\{\pi_i(1-\pi_i)\}$	
	$\boldsymbol{Z} = identity(n  imes n)$	
Probit	$\boldsymbol{W} = \text{diag}\left\{\frac{[f(y_i)]^2}{\pi_i(1-\pi_i)}\right\}$	
	$\boldsymbol{Z} = \operatorname{diag} \left\{ \frac{f(y_i)}{\pi_i (1 - \pi_i)} \right\}$	

Table 3: Definitions of the matrices *W* and *Z*.

The Hessian matrix is calculated using Equation 8. In truth the matrix is different depending upon whether the Firth correction is used; however, convergence to the solution is not significantly affected either way, so for simplicity Equation 8 will be used in either case.

$$H = -X^T W X \tag{8}$$

The gradient vector is calculated using Equation 9,

$$\operatorname{grad} = \nabla L = X^T Z (Y - \pi) + \nabla L_{Firth}$$
(9)

If not using Firth bias correction, the Firth gradient correction vector  $\nabla L_{Firth}$  vanishes. If using the Firth bias correction,  $\nabla L_{Firth}$  is calculated using Equations 10 through 12. The Fisher Information Matrix I is defined in Equation 10. To take the square root of the diagonal matrix W, simply take the square root of each element along the diagonal. The bolded **0.5** indicates a column vector of size  $n \times 1$  with all elements equal to 0.5.

$$I = -H = X^T W X \tag{10}$$

$$\widehat{H} = W^{1/2} X I^{-1} X^T W^{1/2}$$
(11)

$$\nabla L_{Firth} = X^T \widehat{H}(0.5 - \pi)$$
(12)

### **Uncertainty**

Uncertainty is captured by the  $2 \times 2$  covariance matrix  $\Sigma$ , which is the inverse of the observed Fisher Information Matrix, I.

$$\Sigma = I^{-1} = (-H)^{-1}$$
(13)

Because of a special property of the logit link function (it is called "canonical"), the observed Hessian is the same as the expected Hessian we have used up to this point (Cao, 2013). Therefore, the covariance matrix is easy to compute using Equation 10. Unfortunately, the probit link is not canonical and for the purposes of the covariance matrix, Equation 10 cannot be used. For more a more accurate probit covariance matrix, one technique is to employ numerical approximations of the second partial derivatives of the log likelihood function to approximate the Hessian. This procedure is computationally expensive; fortunately, it only must be done once and the matrix is only  $2 \times 2$  so the expense is reasonable in this case.

The standard deviations of the MLE estimates of  $\beta_0$  and  $\beta_1$  are the square roots of the diagonal of the covariance matrix.

The confidence intervals on the predicted probability of success  $\pi$  at a set of desired prediction points  $x_p$  can be calculated using the covariance matrix. Assume an  $m \times 2$  prediction design matrix  $X_p$  where m is the number of points at which to calculate the confidence limits (it may be a dense grid in x for plotting purposes). As before, the first column of  $X_p$  is all ones; the second column is the x coordinates of interest,  $x_p$ . The  $m \times 1$  vector of standard error, **SE**, is

$$SE = \sqrt{\operatorname{diag}(X_p \Sigma X_p^T)}$$
(14)

Calculate  $\hat{y}_j = \hat{\beta}_0 + \hat{\beta}_1 x_{pj}$  at the *m* prediction points with the best fit parameters from the Newton-Raphson MLE solution  $\hat{\beta}$ . (In matrix form the equation is  $\hat{y} = X_p \hat{\beta}$ ). Calculate the *z* value based on the desired confidence using the standard normal distribution – for example, *z* = 1.28 for 95% confidence and *z* = 1.96 for 95% confidence (2-sided interval). The lower and upper confidence intervals are computed using the link function defined in Table 1.

$$LCL_{j} = g^{-1}(\hat{y}_{j} - zSE_{j})$$

$$UCL_{j} = g^{-1}(\hat{y}_{j} + zSE_{j})$$
(15)

The resulting confidence limits should be in the range 0 to 1 except for round-off errors which should be controlled using min/max statements.

## Example

This example is provided for code verification. The fabricated data consists of 12 observations (n = 12). Figure 2 shows a plot of the predicted probability  $\pi(x)$  and its 80% confidence intervals for logistic regression, logistic regression with the Firth bias correction, and probit regression.

X values: 0.8, 0.9, 1.2, 1.7, 1.8, 1.9, 2.0, 2.1, 2.7, 2.9, 3.3, 3.3

Y values: 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1





Figure 1: Graph of the MLE solutions (solid lines) and corresponding confidence intervals (dashed lines) for logistic regression with and without Firth correction and probit regression.

# Conclusion

Logistic and probit regression are effective tools for predicting the probability of success for a system when the probability changes with a continuous variable, x; for example, a probability of detection that depends on the range to the item of interest. This paper provided the steps for constructing the matrices and implementing Newton-Raphson iterations to obtain a solution based on the test data. Formulas were also provided for estimating the uncertainty in the solution's parameters and in the probability of success at any value of x. These recipes can be applied in scripts to automate analysis of a large number of data sets.

# **Works Cited**

Agresti, A. (2002). Categorical Data Analysis. John Wiley & Sons, Inc.

- Allison, P. (2014). Measures of Fit for Logistic Regression. *Proceedings of the SAS® Global Forum 2014 Conference. Paper 1485-2014.* Cary, NC: SAS Institute Inc.
- Cao, X. (2013). *Relative Performance of Expected and Observed Fisher Information in Covariance Estimation for Maximum Likelihood Estimates.* Dissertation, Johns Hopkins University.
- Natoli, C., Burke, S., & Oimoen, S. (2020). *Categorical Data in a Designed Experiment Part 3: Logistic Regression.* Scientific Test and Analysis Techniques (STAT COE) Report 10-2020.
- Zorn, C. (2005). A Solution to Separation in Binary Response Models. *Political Analysis, 13*(2), 157-170. doi:10.1093/pan/mpi009

HS CoBP-Report-03-2021