

# **The Classification Problem, Machine learning demonstrated**

Written by Brian Lesko, this is a hand written demonstration of Statistical Machine Learning theories largely originating from the book, *An Introduction to Statistical Learning*, by Gareth James.

**In part for the course offering at Ohio State University, Statistics and Machine Learning 6500,  
Problem set 2**

2/14/23

For The Logistic regression model  
 (with one feature,  $x$  ( $P=1$ ) Predictor)  
 $\rightarrow \log \frac{p(x)}{1-p(x)} = \beta_0 + \beta_1 x$

the maximum likelihood yields the coefficients  $\beta_0$  and  $\beta_1$  by

where the likelihood function measures the likelihood of observing a particular set of data

where  $l(\beta_0, \beta_1)$

$$= \sum_{i=1}^n [y_i(\beta_0 + \beta_1 x_i) - \log(1 + e^{\beta_0 + \beta_1 x_i})]$$

note: with one class assigned a probability of success, this can be called a bernoulli trial, which results in a binomial distribution (also called bernoulli)

maximizing the likelihood function requires the gradient and hessian of the likelihood function,  $l(\beta_0, \beta_1)$

the hessian is a matrix of form

$$\begin{bmatrix} \frac{\partial^2 l}{\partial \beta_0^2} & \frac{\partial^2 l}{(\partial \beta_0)(\partial \beta_1)} \\ \frac{\partial^2 l}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 l}{\partial \beta_1^2} \end{bmatrix}$$

the gradient has form

$$\begin{bmatrix} \frac{\partial l}{\partial \beta_0} \\ \frac{\partial l}{\partial \beta_1} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n y_i - \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \\ \sum_{i=1}^n x_i \left( y_i - \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \end{bmatrix}$$

→ the gradient can be used to find the parameters numerically with a gradient descent algorithm

however, finding the hessian reveals the curvature of the surface revealing if the estimate is a maximum or minimum as well as the stability

(+) not stable  
(-) stable

the hessian can also be used for estimating the standard errors of the coefficient estimates and the efficiency of the solution

# fisher info matrix and Cramer-Rao lower bound

a) The hessian

$$\begin{bmatrix} \frac{\partial^2 l}{\partial \beta_0^2} & \frac{\partial^2 l}{(\partial \beta_0)(\partial \beta_1)} \\ \frac{\partial^2 l}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 l}{\partial \beta_1^2} \end{bmatrix} = \begin{matrix} \text{from below} \\ -\sum_{i=1}^n \frac{1}{\beta_0^2} & -\sum_{i=1}^n \frac{1}{\beta_0^2} \\ -\sum_{i=1}^n \frac{x_i}{\beta_1^2} & -\sum_{i=1}^n \frac{x_i}{\beta_1^2} \end{matrix}$$

$$\begin{aligned} \frac{\partial^2 l}{\partial \beta_0^2} &= -\frac{\partial}{\partial \beta_0} \left( \sum_{i=1}^n \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \quad \text{quotient rule} \\ &= -\left[ (1 + e^{\beta_0 + \beta_1 x_i}) - (e^{\beta_0 + \beta_1 x_i}) \right] / \beta_0^2 \quad \leftarrow \text{sum left out} \\ &= -\sum_{i=1}^n 1 / \beta_0^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 l}{\partial \beta_1 \partial \beta_0} &= \frac{\partial}{\partial \beta_1} \left( \frac{\partial l}{\partial \beta_0} \right) = -\frac{\partial}{\partial \beta_1} \left( \sum_{i=1}^n \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \\ &= -\left[ x_i (1 + e^{\beta_0 + \beta_1 x_i}) - x_i e^{\beta_0 + \beta_1 x_i} \right] / \beta_1^2 \\ &= \sum_{i=1}^n -x_i / \beta_1^2 = -\sum_{i=1}^n \frac{x_i}{\beta_1^2} \end{aligned}$$

$$\beta_0 + \beta_1 x_i$$

$$\begin{aligned}
\frac{\partial^2 l}{(\partial \beta_0)(\partial \beta_1)} &= \frac{\partial}{\partial \beta_0} \left( \frac{\partial l}{\partial \beta_1} \right) = \frac{\partial}{\partial \beta_0} \sum_{i=1}^n x_i \left( y_i - \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \\
&= - \frac{\partial}{\partial \beta_0} \sum_{i=1}^n \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \\
&= \sum_{i=1}^n - \left[ \left( 1 + e^{\beta_0 + \beta_1 x_i} \right) - \left( e^{\beta_0 + \beta_1 x_i} \right) \right] / \beta_0^2 \\
&= - \sum_{i=1}^n 1 / \beta_0^2 = - \sum_{i=1}^n \frac{1}{\beta_0^2}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial^2 l}{\partial \beta_1^2} &= \frac{\partial}{\partial \beta_1} \left( \frac{\partial l}{\partial \beta_1} \right) = \frac{\partial}{\partial \beta_1} \sum_{i=1}^n x_i \left( y_i - \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \\
&= - \frac{\partial}{\partial \beta_1} \sum_{i=1}^n \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \quad \xrightarrow{\text{P}(x)} \\
&= \sum_{i=1}^n \left[ x_i (1 + e^{\beta_0 + \beta_1 x_i}) - x_i (e^{\beta_0 + \beta_1 x_i}) \right] / \beta_1^2 \\
&= \sum_{i=1}^n - x_i / \beta_1^2 = - \sum_{i=1}^n \frac{x_i}{\beta_1^2}
\end{aligned}$$

b) the hessian can be simplified  
(now negative)

$$\begin{bmatrix} \sum_{i=1}^n \frac{1}{\beta_0^2} & \sum_{i=1}^n \frac{1}{\beta_0^2} \\ \sum_{i=1}^n \frac{x_i}{\beta_1^2} & \sum_{i=1}^n \frac{x_i}{\beta_1^2} \end{bmatrix} = X^T W X$$

$$n \times 2$$

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \end{bmatrix}$$

$$n \times n$$

$$W = \begin{bmatrix} P(x_1)(1-P(x_1)) & 0 & 0 & \dots \\ 0 & P(x_2)(1-P(x_2)) & 0 & \dots \\ 0 & 0 & P(x_3)(1-P(x_3)) & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

$$P(x_i)(1-P(x_i)) =$$

$$\text{where } P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

...

c) the hessian of LSR

Gradient

$$\frac{\partial l}{\partial \beta_0} = \sum -2(y_i - (\beta_0 + \beta_1 x_i))$$

$$\frac{\partial l}{\partial \beta_1} = \sum -2(x_i (y_i - (\beta_0 + \beta_1 x_i)))$$

$$\partial \beta_1 = \frac{1}{n} \sum_{i=1}^n x_i$$

hessian

$$\frac{\partial^2 l}{\partial \beta_0^2} = \sum 2$$

$$\frac{\partial^2 l}{\partial \beta_1^2} = \sum 2x_i^2$$

$$\frac{\partial^2 l}{\partial \beta_0 \partial \beta_1} = \sum 2x_i$$

how is this similar in terms of  $x$ ?

...