

Lecture 7: Maximum likelihood estimation for logistic regression

Ciaran Evans

Recall: Newton's method

Score equation: $u(\beta) = \frac{\partial l}{\partial \beta} \stackrel{\text{set}}{=} 0$

1) Initial guess $\beta^{(0)}$

$$\text{Hessian} = \frac{\partial^2 l}{\partial \beta^2}$$

2) Update : $\beta^{(r+1)} = \beta^{(r)} - \underbrace{(\text{H}(\beta^{(r)}))^{-1}}_{\leftarrow} u(\beta^{(r)})$

3) Stop when algorithm converges

e.g. $\|\beta^{(r+1)} - \beta^{(r)}\| < \varepsilon$

or $\|l(\beta^{(r+1)}) - l(\beta^{(r)})\| < \delta$

Newton's method for logistic regression

$$u(\beta) = \frac{\partial \ell}{\partial \beta} = X^T(Y - P)$$

$$H(\beta) = \frac{\partial^2}{\partial \beta^2} X^T(Y - P)$$

$$= - \frac{\partial^2}{\partial \beta^2} X^T P$$

$$= \left(- \frac{\partial^2}{\partial \beta^2} P \right) X$$

$$\frac{\partial P}{\partial \beta} = \begin{bmatrix} \frac{\partial p_1}{\partial \beta} & \frac{\partial p_2}{\partial \beta} & \dots & \frac{\partial p_n}{\partial \beta} \end{bmatrix} \in \mathbb{R}$$

So: need $\frac{\partial p_i}{\partial \beta} \in \mathbb{R}^{n+1}$

X = design matrix

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$$

$$P = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix}$$

$$p_i = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$$

$$\underbrace{\# \beta_s}_{m \times n^{(x+1)}} \times n^{(x)} \leftarrow \underbrace{\# \text{obs}}_{n^{(y)}}$$

$$p_i = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$$

$\beta^T x_i \in \mathbb{R}$ (scaled) $\beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_n \end{pmatrix} \quad x_i = \begin{pmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{in} \end{pmatrix}$

$$p_i = \frac{e^u}{1 + e^u} \quad u = \beta^T x_i$$

$$= g(h(\beta))$$

$$g(u) = \frac{e^u}{1 + e^u}$$

$$h(\beta) = \beta^T x_i$$

$$g'(u) = \frac{e^u}{(1 + e^u)^2} \quad \frac{\partial}{\partial \beta} \beta^T x_i = x_i$$

$$\frac{\partial p_i}{\partial \beta} = \frac{e^{\beta^T x_i}}{(1 + e^{\beta^T x_i})^2} \quad x_i = p_i(1 - p_i)x_i$$

$$H(\beta) = \left(-\frac{\partial p_i}{\partial \beta} \right) X = - \underbrace{\left[p_1(1-p_1)x_1, p_2(1-p_2)x_2, \dots, p_n(1-p_n)x_n \right]}_{W = \text{diag}(p_i(1-p_i))} X$$

chain rule for single variable calculus:

$$\frac{\partial}{\partial x} g(h(x)) = \underbrace{g'(h(x))}_{\text{Scalar piece}} \underbrace{h'(x)}_{\text{vector piece}}$$

chain rule for matrix calculus:

If x is a vector and $h(x) \in \mathbb{R}$ and $g(h(x)) \in \mathbb{R}$. Then

$$\frac{\partial g(h(x))}{\partial x} = \underbrace{g'(h(x))}_{\text{Scalar piece}} \underbrace{\frac{\partial h}{\partial x}}_{\text{vector piece}}$$

$$\underbrace{[p_1(1-p_1)x_1, p_2(1-p_2)x_2, \dots, p_n(1-p_n)x_n]}_{\text{X}} \times$$

$$X^T = [x_1 \ x_2 \ \dots \ x_n]$$

i^{th} column = x_i

= i^{th} observation in dataset

j^{th} row = j^{th} variable

$$[x_1 \ x_2 \ \dots \ x_n] \underbrace{?}_{\text{?}} = [p_1(1-p_1)x_1 \ \dots \ p_n(1-p_n)x_n]$$

Recall: $[x_1 \ x_2 \ \dots \ x_n]$ $\begin{matrix} \uparrow \\ \uparrow \\ \uparrow \end{matrix}$ $\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = a_1x_1 + a_2x_2 + \dots + a_nx_n$
 vectors
 = linear combination
 of columns of
 matrix

$$[x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} p_1(1-p_1) \\ 0 \\ \vdots \\ 0 \end{bmatrix} = p_1(1-p_1)x_1$$

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots \end{pmatrix}$$

i^{th} row = i^{th} individual

(i^{th} observation)

j^{th} column = j^{th} variable

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix}$$

$$x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{in} \end{pmatrix}$$

$$[x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} p_1(1-p_1) \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = p_1(1-p_1)x_1$$

$$[x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} 0 \\ p_2(1-p_2) \\ 0 \\ \vdots \\ 0 \end{bmatrix} = p_2(1-p_2)x_2$$

w

$$\underbrace{[x_1 \ x_2 \ \dots \ x_n]}_{X^T} \begin{bmatrix} p_1(1-p_1) & 0 & 0 & \cdots & 0 \\ 0 & p_2(1-p_2) & 0 & & 0 \\ \vdots & 0 & p_3(1-p_3) & & \vdots \\ 0 & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & & p_n(1-p_n) \end{bmatrix}$$

$$\Rightarrow X^T w X = [p_1(1-p_1)x_1 \ p_2(1-p_2)x_2 \ \cdots \ p_n(1-p_n)x_n]$$

Comparison:

Linear regression: $H(\beta) = -\frac{1}{\sigma^2} X^T X$

$\uparrow \text{Var}(\varepsilon_i)$

Newton's method
(if we used it)
converges in a
single step

$$-(H(\beta))^{-1} = \sigma^2 (X^T X)^{-1}$$

$$u(\beta) = \frac{1}{\sigma^2} X^T (Y - \mu) = \frac{1}{\sigma^2} X^T (Y - X\beta)$$

$$\Rightarrow -(H(\beta))^{-1} u(\beta) = \sigma^2 (X^T X)^{-1} \frac{1}{\sigma^2} X^T (Y - X\beta)$$

minimize quadratic function of β

$$\beta^{(0)} + (X^T X)^{-1} X^T (Y - X\beta^{(0)})$$

$$u(\beta) = \frac{1}{\sigma^2} X^T (Y - X\beta)$$

linear function of β

$$= \beta^{(0)} + (X^T X)^{-1} X^T Y - \underbrace{(X^T X)^{-1} X^T X \beta^{(0)}}_I$$

\Rightarrow closed form solution

\Rightarrow first-order Taylor expansion is exact
(no approximation)

$$= \beta^{(0)} + (X^T X)^{-1} X^T Y - \beta^{(0)}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Linear regression: $H(\beta) = -\frac{1}{\sigma^2} X^T X$

Logistic regression: $H(\beta) = -X^T w X$

$$w = \text{diag} \left(\underbrace{p_i(1-p_i)}_{\text{var}(Y_i|X_i)} \right)$$

Poisson regression: $H(\beta) = -X^T w X$

$$w = \text{diag} \left(\lambda_i \right) \uparrow \text{var}(Y_i|X_i)$$

Preview to GLMS: for many models,

$$H(\beta) = -X^T w X$$

\uparrow related to variance of Y_i

Example

Suppose that $\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_i$, and we have

$$\beta^{(r)} = \begin{bmatrix} -3.1 \\ 0.9 \end{bmatrix}, \quad U(\beta^{(r)}) = \begin{bmatrix} 9.16 \\ 31.91 \end{bmatrix},$$

$$\mathbf{H}(\beta^{(r)}) = - \begin{bmatrix} 17.834 & 53.218 \\ 53.218 & 180.718 \end{bmatrix}$$

Use Newton's method to calculate $\beta^{(r+1)}$ (you may use R or a calculator, you do not need to do the matrix arithmetic by hand).

Checking the solution is a unique maximum