

Maximum likelihood estimation for logistic regression

Recap: Newton's method

Want to find β^* st $u(\beta^*) = 0$
= Hessian = $\frac{\partial u(\beta^{(r)})}{\partial \beta^{(r)}}$

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

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

3) stop when $\beta^{(r)} \approx \beta^{(r+1)}$

$\gamma_c \neq \gamma_0$

Example

Solve(...)

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).

$$\beta^{(c+1)} = \begin{bmatrix} -3.1 \\ 0.9 \end{bmatrix} + \begin{bmatrix} 17.834 & 53.218 \\ 53.218 & 180.718 \end{bmatrix}^{-1} \begin{bmatrix} 9.16 \\ 31.91 \end{bmatrix}$$

$$= \begin{bmatrix} -3.21 \\ 1.11 \end{bmatrix}$$

Actual MLE : $\begin{bmatrix} -3.36 \\ 1.17 \end{bmatrix}$ So we got closer!

Newton's method for logistic regression

$$\beta^{(r+1)} = \beta^{(r)} - (H(\beta^{(r)}))^{-1} u(\beta^{(r)})$$

$$u(\beta) = X^T(Y - P)$$

$$X = \begin{bmatrix} 1 & x_{i1} & \dots \\ 1 & x_{i2} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$$\Rightarrow u(\beta^{(r)}) = X^T(Y - P^{(r)})$$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad P = \begin{bmatrix} p_1 \\ \vdots \\ p_n \end{bmatrix}$$

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

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

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

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

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

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

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

$$p_i = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} = g(f(\beta))$$

$$\Rightarrow \frac{\partial p_i}{\partial \beta} = \frac{e^{\beta^T x_i}}{(1 + e^{\beta^T x_i})^2} \cdot X_i \left\{ \begin{array}{l} g(u) = \frac{e^u}{1 + e^u} \\ f(\beta) = \beta^T x_i \\ \frac{\partial f}{\partial \beta} = x_i \\ g'(u) = \frac{e^u}{(1 + e^u)^2} \end{array} \right.$$

$$H(\beta) = - \begin{bmatrix} p_1(1-p_1)x_1 & p_2(1-p_2)x_2 & \dots & p_n(1-p_n)x_n \end{bmatrix} X$$

$$= -[x_1 \dots x_n] \begin{bmatrix} p_1(1-p_1) & p_2(1-p_2) & \dots & p_n(1-p_n) \end{bmatrix} X = \boxed{-X^T W X}$$

W = diag($(p_i(1-p_i))$)

Rules for matrix derivatives

if x and $v(x)$ are vectors, and A is a matrix which does not depend on x , then

$$\frac{\partial A v(x)}{\partial x} = \frac{\partial v(x)}{\partial x} A^T$$

If x is a vector, $f(x) \in \mathbb{R}$, and $g(f(x)) \in \mathbb{R}$, then

$$\frac{\partial g(f(x))}{\partial x} = g'(f(x)) \frac{\partial f(x)}{\partial x}$$

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

\uparrow
Var(ε)

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

$$W = \text{diag}(\lambda_i)$$

Newton's method for logistic regression

$$U(\beta) = X^T(Y - P) \quad H(\beta) = -X^T W X$$

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

$$2) \quad \beta^{(r+1)} = \beta^{(r)} + (X^T W^{(r)} X)^{-1} X^T (Y - P^{(r)})$$

3) Stop when $\beta^{(r+1)} \approx \beta^{(r)}$

Checking the solution is a unique maximum

Newton's method finds β^* such that $U(\beta^*) = 0$. How do we know that β^* maximizes the likelihood?

Important property: If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and differentiable, then x^* is a global minimizer of f if and only if

$$\frac{\partial f}{\partial x} \Big|_{x=x^*} = 0$$

\Rightarrow if $-\ell(\beta | X, Y)$ is a convex function (so $\ell(\beta | X, Y)$ is concave), then β^* maximizes $\ell(\beta | X, Y)$ if and only if $U(\beta^*) = 0$

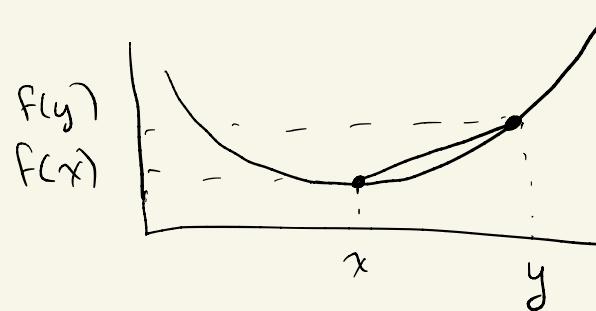
\Rightarrow Newton's method gives us the MLE if $-\ell(\beta | X, Y)$ is convex!

Convex functions

Def: $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if $\forall x, y \in \mathbb{R}^n$ and $0 \leq \lambda \leq 1$

$$f(\underbrace{\lambda x + (1-\lambda)y}_{\text{convex combination}}) \leq \lambda f(x) + (1-\lambda)f(y)$$

e.g.



all points on the line segment lie above the graph of f

Theorem: A twice-differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if the Hessian $H(x)$ is positive-semidefinite $\forall x \in \mathbb{R}^n$

i.e. $y^T H(x) y \geq 0 \quad \forall y \in \mathbb{R}^n$

Claim : for logistic regression, $-H(\beta)$ is positive semi-definite (so $-\ell(\beta|X, Y)$ is convex)

Pf : $H(\beta) = -X^T W X$ $W = \begin{bmatrix} p_1(1-p_1) \\ p_2(1-p_2) \\ \vdots \\ p_n(1-p_n) \end{bmatrix}$

 $\Rightarrow -H(\beta) = X^T W X \in \mathbb{R}^{n+1 \times n+1}$

Let $v \in \mathbb{R}^{n+1}$

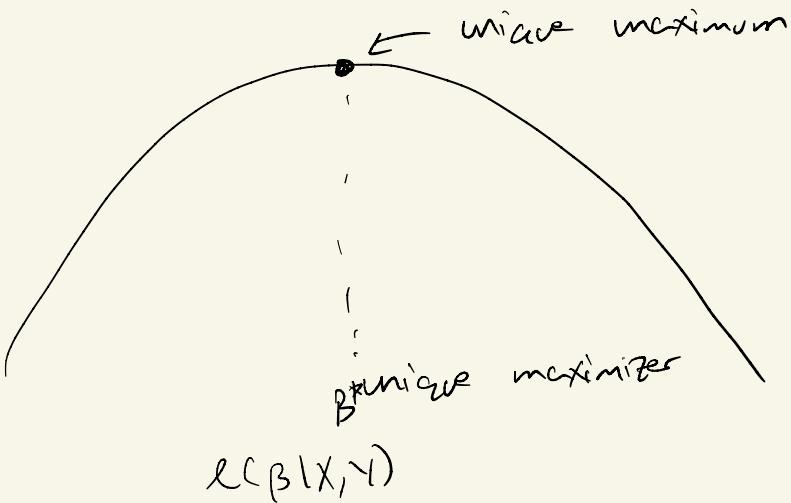
$$-v^T H(\beta) v = v^T X^T W X v = (Xv)^T W (Xv)$$

Let $s = Xv \in \mathbb{R}^n$

$$\begin{aligned} \Rightarrow -v^T H(\beta) v &= s^T W s \\ &= \sum_{i=1}^n p_i(1-p_i) s_i^2 \end{aligned}$$

$$\begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix} \begin{bmatrix} p_1(1-p_1) & & & \\ & \ddots & & \\ & & p_n(1-p_n) & \end{bmatrix} \begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix} \geq 0$$

$\Rightarrow -\ell(\beta|X, Y)$ is convex, so solving $U(\beta) = 0$ gives MLE of β



Concave function



Convex function

Some intuition about Hessians

Fisher information

Properties

Example