

Lecture 6: Maximum likelihood estimation for logistic regression

Newton's method for logistic regression

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

1) Initial guess $\beta^{(0)}$ Hessian = $\frac{\partial u}{\partial \beta} = H(\beta)$

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

3) Stop when algorithm converged
(usually when $l(\beta)$ stops changing, or
changes by a very small amount)

Example

multiplication: $\gamma_0 * \delta_0$
inverse: solve(...)

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

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

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

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).

$$\beta^{(r+1)} = ?$$

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

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

(Actual MLE (from `glm` function) is

$$\begin{bmatrix} -3.36 \\ 1.17 \end{bmatrix},$$

So we got closer!)

Newton's method for logistic regression

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

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

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

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

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

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

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

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

$\downarrow z_1$

$$w_{\text{cont}} \quad \frac{\partial p}{\partial \beta} = \left[\begin{array}{cccc} \frac{\partial p_1}{\partial \beta} & \frac{\partial p_2}{\partial \beta} & \cdots & \frac{\partial p_n}{\partial \beta} \end{array} \right] \in \mathbb{R}^{(k+1) \times n}$$

↑ ↑ ↗ ↑
obs # BS

↑
observations

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

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

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

$$= p_i(1 - p_i) x_i$$

$$H(\beta) = \left(-\frac{\partial^2 p}{\partial \beta^2} \right) X = -\begin{bmatrix} p_1(1-p_1)x_1 \\ -X^T W X \end{bmatrix}$$

Columns of $\frac{\partial p}{\partial \beta}$

$$g(u) = \frac{e^u}{1 + e^u} \quad h(\beta) = \beta^T x_i$$

Chain rule for matrix derivatives:
 If x is vector, and $h(x) \in \mathbb{R}$
 and $g(h(x)) \in \mathbb{R}$, then
 $\frac{\partial g(h(x))}{\partial x} = g'(h(x)) \frac{\partial h(x)}{\partial x}$

$$p_2(1-p_2)x_2 \cdots p_n(1-p_n)x_n] X$$

$$W = \text{diag}(p_1(1-p_1), p_2(1-p_2), \dots, p_n(1-p_n))$$

comparison:

Linear regression: $H(\beta) = -\frac{1}{\sigma^2} X^\top X$
 $\nwarrow \text{Var}(\varepsilon)$

Poisson: $H(\beta) = -X^\top W X$
 $W = \text{diag}(\lambda_i)$

Preview to GLMs:

in general, $H(\beta) = -X^\top W X$
 $W = \text{diag}\left(\frac{\text{Var}(\gamma_i)}{\phi^2}\right)$ $\phi = \text{dispersion parameter}$

Bernoulli: $\phi = 1$

Poisson: $\phi = 1$

Normal: $\phi = \sigma^2$

Checking the solution is a unique maximum

Newton's method finds β^* st $u(\beta^*) = 0$

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 of β (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!



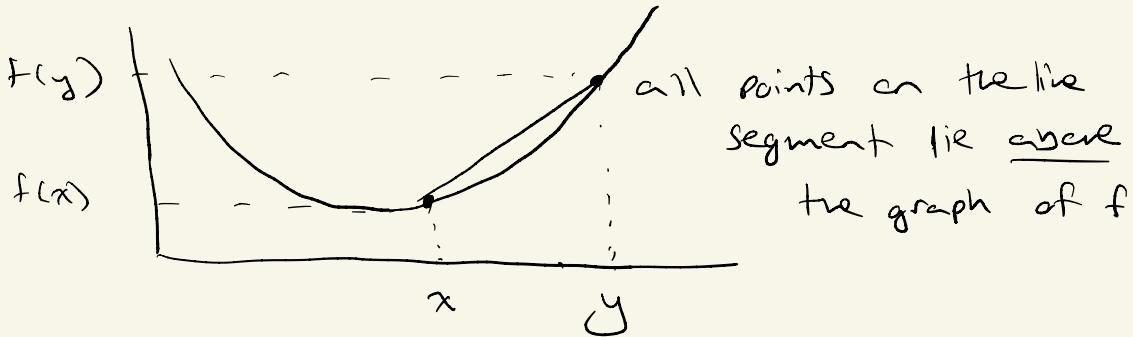
Convex functions

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

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

Convex combination

e.g.



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 $\forall x \in \mathbb{R}^n$

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

Claim: for logistic regression, $-H(\beta)$ is positive semi-definite

$$\text{Pf: } H(\beta) = -X^T W X \Rightarrow -H(\beta) = X^T W X$$

$$\text{Let } v \in \mathbb{R}^{n+1} \quad (\beta \in \mathbb{R}^{n+1})$$

$$v^T X^T W X v = (Xv)^T W (Xv)$$

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

$$s^T W s = \sum_{i=1}^n \underbrace{p_i(1-p_i)}_{\in (0,1)} s_i^2 \geq 0$$

$p_i \in (0,1)$

$$\geq 0$$

$\Rightarrow -H(\beta)$ is PSD, so $-L(\beta | X, Y)$ is convex

$\rightarrow L(\beta | X, Y)$ is concave \Rightarrow Solving $u(\beta)$ gives MLE //

