# Maximum likelihood estimation for regression models

Ciaran Evans

# Maximum likelihood estimation and linear regression

Let $(\mathbf{x}_1, Y_1), ..., (\mathbf{x}_n, Y_n)$ be iid samples from the model

$$Y_i | \mathbf{x}_i \sim N(\mu_i, \underline{\sigma^2})$$

$$\mu_i = \mathbf{x}_i^T \widehat{\beta} = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$$

where the distribution of $\mathbf{x}_i$ does not depend on $\beta$ or $\sigma^2$.

$$L(\beta, \sigma^2 | \mathbf{y}, \mathbf{X}) = \prod_{i=1}^{n} f(x_i, Y_i \mid \beta, \sigma^2) = \prod_{i=1}^{n} f(x_i) f(Y_i | x_i, \beta, \sigma^2)$$

does not involve $\beta$ or $\sigma^2$ $\longrightarrow$

does involve $\beta$ and $\sigma^2$ $\longleftarrow$

$$= \left( \prod_{i=1}^{n} f(x_i) \right) \left( \prod_{i=1}^{n} f(Y_i | x_i, \beta, \sigma^2) \right)$$

$$\propto \prod_{i=1}^{n} f(Y_i | x_i, \beta, \sigma^2)$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ -\frac{1}{2\sigma^2} (Y_i - x_i^T \beta)^2 \right\} \quad \swarrow SSE$$

$$= \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp\left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^{n} (Y_i - x_i^T \beta)^2 \right\}$$

$\Rightarrow$ choosing $\beta$ to maximize $L$ is equivalent to minimizing $SSE$!

# Maximum likelihood estimation and logistic regression

Let $(\mathbf{x}_1, Y_1), ..., (\mathbf{x}_n, Y_n)$ be iid samples from the model

$$Y_i | \mathbf{x}_i \sim Bernoulli(p_i)$$

$$\log\left(\frac{p_i}{1 - p_i}\right) = \mathbf{x}_i^T \beta$$

$$f(Y_i | x_i, \beta) = p_i^{Y_i}(1 - p_i)^{1 - Y_i}$$

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

where the distribution of $\mathbf{x}_i$ does not depend on $\beta$.

$$L(\beta | \mathbf{y}, \mathbf{X}) \propto \prod_{i=1}^{n} f(Y_i | \mathbf{x}_i, \beta) = \prod_{i=1}^{n} p_i^{Y_i}(1 - p_i)^{1 - Y_i}$$

(up to a constant)

$$= \prod_{i=1}^{n} \left(\frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}\right)^{Y_i} \left(\frac{1}{1 + e^{x_i^T \beta}}\right)^{1 - Y_i}$$

$$\ell(\beta | \mathbf{y}, \mathbf{X}) = \sum_{i=1}^{n} \left\{ Y_i \log\left(\frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}\right) + (1 - Y_i)\log\left(\frac{1}{1 + e^{x_i^T \beta}}\right) \right\}$$

$$= \sum_{i=1}^{n} \left\{ Y_i \, x_i^T \beta - \log\left(1 + e^{x_i^T \beta}\right) \right\}$$

# Maximizing

$$\ell(\beta|\mathbf{y}, \mathbf{X}) \stackrel{\text{(up to a constant)}}{=} \sum_{i=1}^{n} \left\{ Y_i \mathbf{x}_i^T \beta - \log(1 + e^{\mathbf{x}_i^T \beta}) \right\}$$

$$\left[ \begin{array}{c} \frac{\partial \ell}{\partial \beta_0} \\ \frac{\partial \ell}{\partial \beta_1} \\ \vdots \\ \frac{\partial \ell}{\partial \beta_u} \end{array} \right.$$

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^{n} \left\{ \underbrace{\frac{\partial}{\partial \beta} Y_i x_i^T \beta}_{Y_i x_i} - \underbrace{\frac{\partial}{\partial \beta} \log\left(1 + e^{x_i^T \beta}\right)}_{} \right\}$$

$$\left( \frac{\partial}{\partial u} a^T u = a \right)$$

$$- \frac{1}{1 + e^{x_i^T \beta}} \cdot e^{x_i^T \beta} \cdot x_i \qquad \text{(chain rule)}$$

$$= \sum_{i=1}^{n} \left\{ Y_i x_i - \boxed{\frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}} x_i \right\} \qquad P_i$$

$$= \sum_{i=1}^{n} (Y_i - p_i) x_i \qquad = X^T (y - p)$$

$$X = \text{design matrix} \qquad \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \qquad \begin{pmatrix} P_1 \\ \vdots \\ P_n \end{pmatrix}$$

# Score

**Definition (score):** Let $\mathbf{y} = (Y_1, ..., Y_n)$ be a sample of $n$ observations from some distribution with parameter vector $\boldsymbol{\theta}$. Let $L(\boldsymbol{\theta}|\mathbf{y})$ be the likelihood function, and $\ell(\boldsymbol{\theta}|\mathbf{y}) = \log L(\boldsymbol{\theta}|\mathbf{y})$ the log-likelihood.

The **score**, which we will denote $U(\boldsymbol{\theta})$, is the gradient of the log-likelihood with respect to $\boldsymbol{\theta}$:

$$U(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \ell(\boldsymbol{\theta}|\mathbf{y}).$$

**Example:** For logistic regression: $U(\beta) = \mathbf{X}^T(\mathbf{y} - \mathbf{p})$

**Question:** How would I solve $\mathbf{X}^T(\mathbf{y} - \mathbf{p}) = 0$?

Challenge: p is a nonlinear function of β

In fact, there is no closed form solution for β!

# Newton's method

We want to find $\beta^*$ such that $U(\beta^*) = \mathbf{0}$. Issue: no closed form solution!

**Idea:** Approximate $U(\beta^*)$ with a first-order Taylor expansion:

Taylor expansion of $g(x)$ around $a$:
$$g(x) \approx g(a) + g'(a)(x-a)$$

$$U(\beta^*) \approx u(\beta) + \left( \frac{\partial u(\beta)}{\partial \beta} \right) (\beta^* - \beta) \qquad \text{(if } \beta \text{ is close enough to } \beta^*\text{)}$$

Let $\beta^{(0)}$ be an initial guess for $\beta^*$

Want to iteratively update and improve this initial guess

$$u(\beta^*) \approx u(\beta^{(0)}) + \left( \frac{\partial u}{\partial \beta} \Big|_{\beta = \beta^{(0)}} \right) (\beta^* - \beta^{(0)})$$

$$\overset{\shortparallel}{0}$$

$$\Rightarrow \quad \beta^* \approx \beta^{(0)} - \left( \frac{\partial u}{\partial \beta} \Big|_{\beta = \beta^{(0)}} \right)^{-1} u(\beta^{(0)})$$

we can evaluate this!

$$u(\beta) = \frac{\partial}{\partial \beta} \ell(\beta | y, x)$$

$$H(\beta) = \frac{\partial u}{\partial \beta} = \frac{\partial^2 \ell(\beta | y, x)}{\partial \beta^2} \quad \text{(Hessian matrix)}$$

second derivative!

# Newton's method

- Want $\beta^*$ such that $U(\beta^*) = \mathbf{0}$
- Begin with initial estimate ~~$\beta^{(0)}$~~ $\beta^{(0)}$
- Iterative updates:

$$\beta^{(r+1)} = \beta^{(r)} - \left(\mathbf{H}(\beta^{(r)})\right)^{-1} U(\beta^{(r)})$$

$$U(\beta) = \frac{\partial \ell}{\partial \beta} \qquad \text{(gradient)} \qquad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_u \end{pmatrix} \in \mathbb{R}^{(u+1)}$$

$$\in \mathbb{R}^{u+1}$$

$$H(\beta) = \frac{\partial^2 \ell}{\partial \beta^2} \in \mathbb{R}^{(u+1) \times (u+1)} \qquad \text{(Hessian)}$$

$$= \begin{bmatrix} \frac{\partial^2 \ell}{\partial \beta_0^2} & \frac{\partial^2 \ell}{\partial \beta_0 \partial \beta_1} & \cdots & \frac{\partial^2 \ell}{\partial \beta_0 \partial \beta_u} \\ \vdots & & & \\ \frac{\partial^2 \ell}{\partial \beta_u \partial \beta_0} & & \cdots & \frac{\partial^2 \ell}{\partial \beta_u^2} \end{bmatrix}$$

# The Hessian

$$U(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} \ell(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) = \mathbf{X}^T(\mathbf{y} - \mathbf{p})$$

$$\mathbf{H}(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} U(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{X}^T(\mathbf{y} - \mathbf{p})$$

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

(chain rule)

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

$$p = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix} \in \mathbb{R}^n$$

$$\frac{\partial p}{\partial \beta} = \left[ \frac{\partial p_1}{\partial \beta} \quad \frac{\partial p_2}{\partial \beta} \quad \cdots \quad \frac{\partial p_n}{\partial \beta} \right] \in \mathbb{R}^{(k+1) \times n}$$

$$\beta \in \mathbb{R}^{k+1}$$

So : need to find $\frac{\partial p_i}{\partial \beta}$

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

# Putting everything together

Want to maximize the log likelihood $\ell(\beta|\mathbf{y}, \mathbf{X})$.