

Fisher scoring

Recap: Fisher information

Def: Let $\ell(\theta|Y)$ be a log likelihood,
and $u(\theta) = \frac{\partial \ell}{\partial \theta}$. The Fisher information
is
$$\hat{I}(\theta) = \text{Var}(u(\theta)|\theta)$$

$$\mathbb{E}[g(Y)] = \int g(y) f(y) dy$$

Properties

$$g(Y) = U(\theta)$$

Theorem: Under appropriate regularity conditions,

$$\mathbb{E}[U(\theta) | \theta] = 0$$

$$\text{Proof: } \mathbb{E}[U(\theta) | \theta] = \mathbb{E}\left[\frac{\partial}{\partial \theta} \ell(\theta | Y) | \theta\right] = \mathbb{E}\left[\frac{\partial}{\partial \theta} \log f(Y | \theta) | \theta\right]$$

$$= \mathbb{E}\left[\frac{1}{f(Y | \theta)} \cdot \frac{\partial}{\partial \theta} f(Y | \theta) | \theta\right] = \int \frac{\partial}{\partial \theta} f(y | \theta) \frac{1}{f(y | \theta)} dy$$

$$= \int \frac{\partial}{\partial \theta} f(y | \theta) dy \quad \begin{matrix} \downarrow \\ \text{required regularity:} \\ \text{swap derivative and integral} \end{matrix}$$

$$= \frac{\partial}{\partial \theta} \int f(y | \theta) dy \quad (\text{see Casella \& Berger 2.4})$$

$$= \frac{\partial}{\partial \theta} 1 = 0$$

$$\frac{\partial}{\partial \theta} \log f(\gamma|\theta) = \frac{\frac{\partial}{\partial \theta} f(\gamma|\theta)}{f(\gamma|\theta)}$$

Properties

Theorem: Under appropriate regularity conditions,

$$I(\theta) = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta | \mathbf{Y}) \mid \theta \right]$$

Proof:

$$\begin{aligned}
 & -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta | \mathbf{Y}) \mid \theta \right] = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \log f(\gamma|\theta) \mid \theta \right] \\
 &= -\mathbb{E} \left[\frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} \log f(\gamma|\theta) \right) \mid \theta \right] = -\mathbb{E} \left[\frac{\partial}{\partial \theta} \left(\frac{\frac{\partial}{\partial \theta} f(\gamma|\theta)}{f(\gamma|\theta)} \right) \mid \theta \right] \\
 &= -\mathbb{E} \left[\frac{\frac{\partial^2}{\partial \theta^2} f(\gamma|\theta)}{f(\gamma|\theta)} \right] - \left(\frac{\frac{\partial}{\partial \theta} f(\gamma|\theta)}{f(\gamma|\theta)} \right)^2 \mid \theta \\
 &= -\mathbb{E} \left[\frac{\frac{\partial^2}{\partial \theta^2} f(\gamma|\theta)}{f(\gamma|\theta)} \mid \theta \right] + \mathbb{E} \left[\left(\frac{\partial}{\partial \theta} \log f(\gamma|\theta) \right)^2 \mid \theta \right] \\
 &\quad \underbrace{\frac{\partial^2}{\partial \theta^2} \int f(y|\theta) dy = 0}_{\text{Property}}
 \end{aligned}$$

= $\mathbb{E}[(u(\theta))^2 \mid \theta]$ ($\mathbb{E}[u(\theta) \mid \theta] = 0$)

= $\text{Var}(u(\theta) \mid \theta)$ ($= 0$)

Fisher information vs. Hessian

For logistic regression:

$$\mathcal{L}(\beta) = \mathbf{x}^T \mathbf{w} \mathbf{x} = -H(\beta)$$

under regularity conditions, $\mathcal{L}(\beta) = -\mathbb{E}[H(\beta)]$

\mathcal{L} and $-H$ are not always the same

Example: $y_1, \dots, y_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

$$\frac{\partial^2}{\partial p^2} \ell(p|y) = -\frac{\sum_i y_i}{p^2} - \frac{(n - \sum_i y_i)}{(1-p)^2}$$

depends
on the
observed data

$$\mathcal{L}(p) = -\mathbb{E}\left[\frac{\partial^2}{\partial p^2} \ell(p|Y)\right] = \frac{n}{p(1-p)}$$

does not
depend on the
observed data

In general, $H(\theta)$ can be more work to calculate
 \Rightarrow typically use $\mathcal{L}(\theta)$ in place of $H(\theta)$

Newton's method: $\beta^{(r)} - \underbrace{H^{-1}(\beta^{(r)})}_{\text{replace w/ } \lambda^{-1}(\beta^{(r)})} u(\beta^{(r)})$

Fisher scoring

1) Start with initial guess $\beta^{(0)}$

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

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

IRLS for logistic regression

$$\begin{aligned}\beta^{(r+1)} &= \beta^{(r)} + (X^T w^{(r)} X)^{-1} X^T (y - p^{(r)}) \\&= \underbrace{(X^T w^{(r)} X)^{-1}}_{\text{identity matrix}} X^T w^{(r)} \beta^{(r)} + (X^T w^{(r)} X)^{-1} X^T (y - p^{(r)}) \\&= (X^T w^{(r)} X)^{-1} X^T w^{(r)} \underbrace{\left(X\beta^{(r)} + (w^{(r)})^{-1} (y - p^{(r)}) \right)}_{Z^{(r)}} \\&\Rightarrow \beta^{(r+1)} = (X^T w^{(r)} X)^{-1} X^T w^{(r)} Z^{(r)} \quad \nwarrow \text{working responses at iteration } r\end{aligned}$$

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

actually doing weighted least squares with weights $w^{(r)}$, responses $Z^{(r)}$, explanatory variables X

Suppose $\gamma = X\beta + \varepsilon$

$\varepsilon \sim N(0, W^{-1})$ (allow non-constant variance)

$$\underbrace{w^{\frac{1}{2}}\gamma}_{\gamma_w} = \underbrace{w^{\frac{1}{2}}X\beta}_{X_w} + \underbrace{w^{\frac{1}{2}}\varepsilon}_{\varepsilon_w}$$

$$\gamma_w = X_w\beta + \varepsilon_w$$

$$W = \text{diag}(w_1, \dots, w_n) \quad W^{\frac{1}{2}} = \text{diag}(\sqrt{w_i})$$

$$\varepsilon_w \sim N(0, I)$$

$$\text{var}(W^{\frac{1}{2}}\varepsilon) = W^{\frac{1}{2}}W^{-1}W^{\frac{1}{2}} = I$$

$$\begin{aligned}\hat{\beta} &= (X_w^T X_w)^{-1} X_w^T \gamma_w \\ &= (X^T W X)^{-1} X^T W \gamma\end{aligned}$$

A preview of Fisher information properties