

# Lecture 8: Fisher information

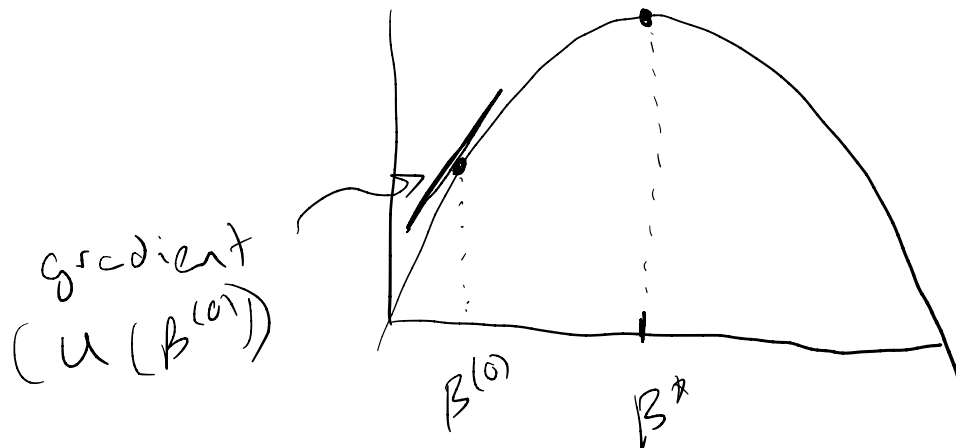
# Recap: Newton's method

To find  $\beta^*$  such that  $U(\beta^*) = 0$ , when there is no closed-form solution we use Newton's method:

- Begin with an initial guess  $\beta^{(0)}$
- Iteratively update:  $\beta^{(r+1)} = \beta^{(r)} - \mathbf{H}^{-1}(\beta^{(r)})U(\beta^{(r)})$
- Stop when the algorithm converges

Intuition:

log likelihood



Hessian tells us about curvature of the log-likelihood

# Some intuition about Hessians

**Example:** Suppose that  $\beta = (\beta_0, \beta_1)^T \in \mathbb{R}^2$ , and

$$\ell(\beta) = -\beta_0^2 - 100\beta_1^2$$

Calculate

$$\begin{aligned} \mathbf{U}(\beta) &= \begin{bmatrix} \frac{\partial \ell}{\partial \beta_0} \\ \frac{\partial \ell}{\partial \beta_1} \end{bmatrix} & \mathbf{H}(\beta) &= \begin{bmatrix} \frac{\partial^2 \ell}{\partial \beta_0^2} & \frac{\partial^2 \ell}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 \ell}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 \ell}{\partial \beta_1^2} \end{bmatrix} \\ &= \begin{bmatrix} -2\beta_0 \\ -200\beta_1 \end{bmatrix} & &= \begin{bmatrix} -2 & 0 \\ 0 & -200 \end{bmatrix} \end{aligned}$$

$$L(\beta) = -\beta_0^2 - 100\beta_1^2$$

$$H(\beta) = - \begin{bmatrix} 2 & 0 \\ 0 & 200 \end{bmatrix}$$

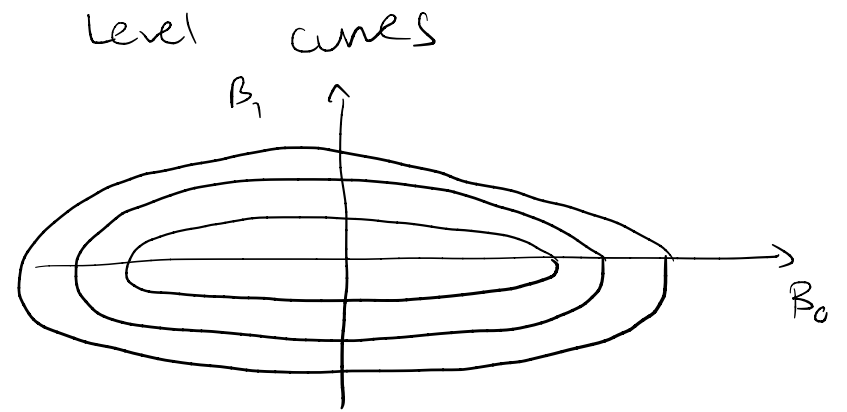
$$H^{-1}(\beta) = - \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{200} \end{bmatrix}$$

Suppose start with

$$\beta^{(0)} = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$$

$$\beta^{(1)} = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} + \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{200} \end{bmatrix} \begin{bmatrix} -2(0.1) \\ -200(0.1) \end{bmatrix}$$

$$= \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} - \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



(a small change in  $\beta_1$  makes a bigger difference to  $L(\beta)$  than a small change in  $\beta_0$ )

Contrast: gradient ascent / gradient descent

$\beta^{(n)}$  current guess

$$X^T W X$$

new guess:  $\beta^{(n+1)} = \beta^{(n)} + \alpha \nabla U(\beta^{(n)})$

step size

previous example:

$$\begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} + \alpha \begin{bmatrix} -2(0.1) \\ -200(0.1) \end{bmatrix}$$

(probably converge more slowly than Newton's method,  
b/c we don't use information about second derivative)

Trade-off:

gradient ascent/descent:

- more steps
- doesn't require matrix inversion of second derivatives

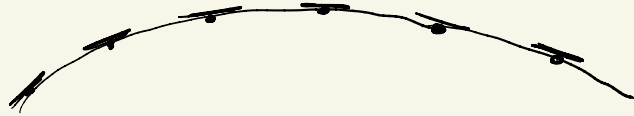
Newton's method

- fewer steps
- steps are more computationally expensive

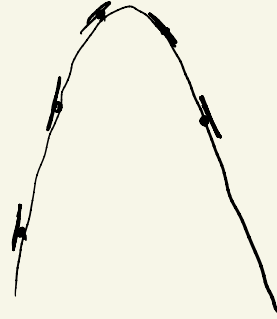
One more perspective

↙ harder to maximize

↘ easier to maximize



slopes  $u(\beta)$  are close to 0



slopes  $u(\beta)$  are far from 0

⇒ suggest look @ variability in  $u(\beta)$  (as a function of observed data)

Logistic regression:  $u(\beta) = X^T(Y-p)$

$X$  design matrix

$Y$  vector of responses

$p$  = vector of probabilities

$$\text{Var}(u(\beta) | X) = \text{Var}(X^T(Y-p) | X)$$

$$= X^T \text{Var}(Y-p | X) X$$

$$= X^T \text{Var}(Y | X) X$$

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

# Fisher information

Def: Let  $\ell(\theta|Y)$  be a log likelihood, and  $u(\theta) = \frac{\partial \ell}{\partial \theta}$

The Fisher information is

$$\mathcal{I}(\theta) = \text{Var}(u(\theta) | \theta)$$

i.e. The variance of  $u(\theta)$ , given  $\theta$  is the true parameter

# Example: Bernoulli sample

$$\text{var}(Y+c) = \text{var}(Y)$$

$$\text{var}(Y_i) = p(1-p)$$

Suppose that  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ .

$$L(p|Y) = p^{\sum_i Y_i} (1-p)^{n - \sum_i Y_i}$$

$$\ell(p|Y) = \left(\sum_i Y_i\right) \log p + \left(n - \sum_i Y_i\right) \log(1-p)$$

$$u(p) = \frac{\partial}{\partial p} \ell(p|Y) = \frac{\sum_i Y_i}{p} - \frac{(n - \sum_i Y_i)}{1-p}$$

$$\text{var}(u(p)|p) = \text{var}\left(\frac{\sum_i Y_i}{p} - \frac{(n - \sum_i Y_i)}{(1-p)} \mid p\right)$$

$$= \text{var}\left(\frac{(\sum_i Y_i)(1-p) - (n - \sum_i Y_i)p}{p(1-p)} \mid p\right)$$

$$= \text{var}\left(\frac{\sum_i Y_i}{p(1-p)} \mid p\right) = \frac{1}{p^2(1-p)^2} \sum_i \text{var}(Y_i) = \frac{n p(1-p)}{p^2(1-p)^2} = \frac{n}{p(1-p)}$$



$$I(p) = \frac{n}{p(1-p)}$$

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

$$\hat{p} = \frac{1}{n} \sum_i Y_i$$

$$\text{Var}(\hat{p}) = \frac{1}{n^2} \cdot n p(1-p) = \frac{p(1-p)}{n} = I^{-1}(p)$$

Linear:  
regression

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} = I^{-1}(\beta)$$

Logistic  
regression

$$\text{Var}(\hat{\beta}) = (X^T W X)^{-1} = I^{-1}(\beta)$$

$$W = \begin{bmatrix} p_1(1-p_1) & & \\ & p_2(1-p_2) & \\ & & \ddots & p_n(1-p_n) \end{bmatrix}$$

# Properties

Under certain regularity conditions, we have:

$$① \quad \mathbb{E}[u(\theta) | \theta] = 0$$

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

## Example: Bernoulli sample

Suppose that  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p_i)$ .

