

Lecture 8: Fisher Information

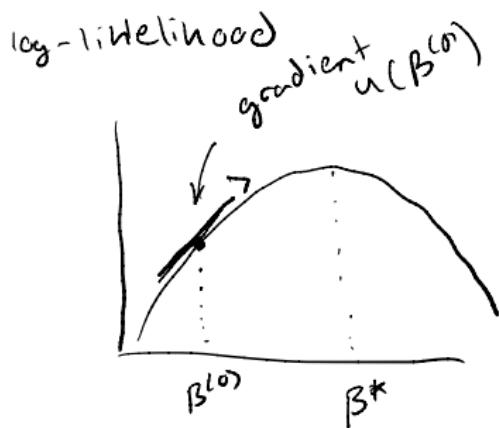
Ciaran Evans

Recap: Newton's method

To find β^* 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:



use gradient to
walk uphill

Hessian tells us
how far along
gradient to move

Some intuition about Hessians

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

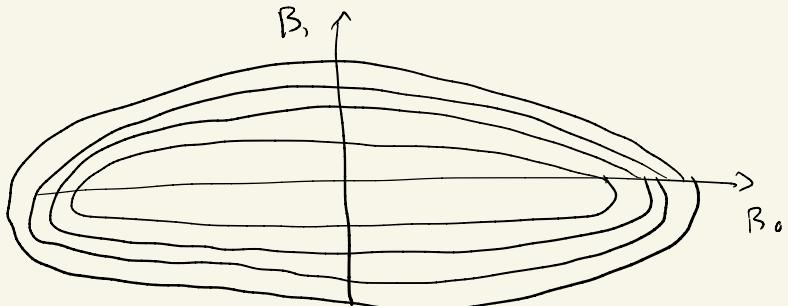
$$\ell(\beta) = -\beta_0^2 - 100\beta_1^2$$

Calculate

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

$$= \begin{bmatrix} -2\beta_0 \\ -200\beta_1 \end{bmatrix} \quad \begin{bmatrix} -2 & 0 \\ 0 & -200 \end{bmatrix}$$

Level curves $\ell(\beta) = -\beta_0^2 - 100\beta_1^2$ (maximum at $(0,0)$)



a small change in β_1
makes a bigger difference to
 $\ell(\beta)$ than a small change
in β_0

Two different approaches to moving along gradient

gradient ascent: step size

$$\beta^{(r+1)} = \beta^{(r)} + \alpha \nabla \ell(\beta^{(r)})$$

Example: $\beta^{(r)} = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$

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

$$\begin{aligned} \alpha &= \frac{1}{200} \\ \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} &+ \frac{1}{200} \begin{bmatrix} -2(0.1) \\ -200(0.1) \end{bmatrix} \\ &= \begin{bmatrix} 0.1 - 0.001 \\ 0 \end{bmatrix} \end{aligned}$$

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

$\alpha = \frac{1}{200}$
 $\Rightarrow \text{update} = \begin{bmatrix} 0.099 \\ 0 \end{bmatrix}$

$$\alpha = \frac{1}{2}$$

$$\begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} -2(0.1) \\ -200(0.1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0.1 - 10 \end{bmatrix} = \begin{bmatrix} 0 \\ -9.9 \end{bmatrix}$$

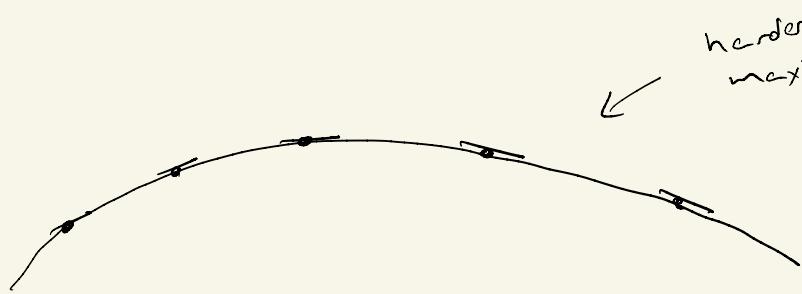
Using Hessian:

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

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

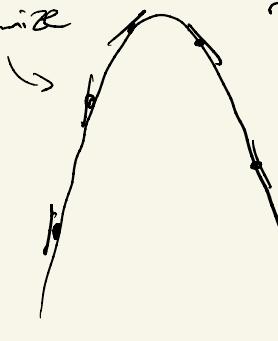
Newton's method is faster (fewer iterations) by using info about second derivative

one more perspective



harder to maximize

easier to maximize



matrix A ,
random vector γ
 $\text{var}(A\gamma)$
 $= A \text{var}(\gamma) A^T$

Slopes $u(B)$ close to 0

less variability in slopes

(most values close to 0)

Slope $u(B)$ are far from 0

more variability in slopes
(some values near 0, some values far from 0)

⇒ look at Variance of $u(B)$ (as a function of the data)

Logistic regression : $u(B) = X^T (\gamma - p)$

$$\begin{aligned} \text{var}(u(B)) &= \text{var}(X^T (\gamma - p)) = X^T \text{var}(\gamma - p) X \\ &= X^T \text{var}(\gamma) X - X^T \begin{bmatrix} p_1(1-p_1) & 0 & \cdots & 0 \\ 0 & p_2(1-p_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & p_n(1-p_n) \end{bmatrix} X = X^T W X \\ \text{var}(\gamma_i) &= p_i(1-p_i) \end{aligned}$$

$$= -H(B)$$

Fisher information

Def : Let $\ell(\theta|\gamma)$ be a log-likelihood given data γ , and θ some parameter of interest.

$$\text{Let } u(\theta) = \frac{\partial \ell}{\partial \theta}$$

The Fisher information is

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

Variance of score $u(\theta)$, if θ is the true parameter

Thm : Under certain regularity conditions,

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

Example: Bernoulli sample

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

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

$$L(p|y) = p^{\sum_i y_i} (1-p)^{n - \sum_i y_i}$$

$$\ell(p|y) = (\sum_i y_i) \log p + (n - \sum_i y_i) \log (1-p)$$

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

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

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

$$= \text{var} \left(\frac{\sum_i y_i}{p(1-p)} \right) = \frac{1}{p^2(1-p)^2} \sum_i \text{var}(y_i) = \frac{np(1-p)}{p^2(1-p)^2}$$

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

$$\begin{aligned}
 \text{MLE: } \hat{p} &= \frac{1}{n} \sum_i y_i \\
 \text{var}(\hat{p}) &= \text{var}\left(\frac{1}{n} \sum_i y_i\right) \\
 &= \frac{1}{n^2} \sum_i \text{var}(y_i) \\
 &= \frac{1}{n^2} \cdot n p(1-p) \\
 &= \frac{p(1-p)}{n} \\
 \text{var}(\hat{p}) &= \Sigma^{-1}(p)
 \end{aligned}$$

Linear regression: $\text{var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = \Sigma^{-1}(\beta)$

Preview: variance of MLEs is closely connected
to Σ^{-1}

Properties