

Lecture 10 (or whatever): Fisher information and the MLE

- We'll go over the univariate case of Fisher Information today.
- we start with a 1-parameter family of densities

$$\mathcal{F} = \{ f_{\theta}(x), \underbrace{\theta \in \Omega}_{\text{Recall: Parameter Space}}, \underbrace{x \in \mathcal{X}}_{\text{Sample Space}} \}$$

→ We'll consider just the continuous case for simplicity.

- Recall that we defined the log-likelihood function as:

$$\ell_x(\theta) = \log f_{\theta}(x)$$

- we'll call its derivative wrt θ the score function

$$\dot{\ell}_x(\theta) = \frac{\partial}{\partial \theta} \log f_{\theta}(x) = \frac{\dot{f}_{\theta}(x)}{f_{\theta}(x)}$$

Recall (chain rule)

How higher/lower the likelihood value of our sample gets as the true θ varies.

$$\frac{\partial}{\partial x} \log(f(x)) = \frac{f'(x)}{f(x)}$$

- First, let's compute the expectation of the score function.

$$\text{Recall } E(\theta) = \int_{\mathcal{X}} \theta \cdot \underbrace{f(x)}_{\text{density}} dx$$

$$E[\dot{\ell}_x(\theta)] = \int_{\mathcal{X}} \dot{\ell}_x(\theta) \cdot \underbrace{f_{\theta}(x)}_{\text{density}} dx = \int_{\mathcal{X}} \dot{f}_{\theta}(x) dx = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} f_{\theta}(x) dx$$

- Assuming $f_{\theta}(x)$ continuous & continuously differentiable.

$$\int \frac{\partial}{\partial \theta} f_0(x) dx = \frac{\partial}{\partial \theta} \int f_0(x) dx = \frac{\partial}{\partial \theta} \cdot 1 = 0$$

First result today

$$E[\dot{\ell}_x(\theta)] = 0$$

• OK, got the expectation. Now, let's talk about the variance of $\dot{\ell}_x(\theta)$.

Recall $V(X) = \int [x - E(X)]^2 \cdot f_0(x) dx$

Since $E[\dot{\ell}_x(\theta)] = 0$

$$I_0 \triangleq V(\dot{\ell}_x(\theta)) = \int [\dot{\ell}_x(\theta)]^2 \cdot f_0(x) \cdot dx$$

Define the Fisher Information as the variance of the score function.

Now, we'll show that:

Then, $\hat{\theta}_{MLE} \sim N(\theta, \frac{1}{I_0})$

(Sketch of Proof): Let $\ddot{\ell}_x(\theta) = \frac{\partial^2}{\partial \theta^2} \log f_0(x)$

Recall: Quotient differentiation rule $\frac{\partial}{\partial x} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - h'(x)g(x)}{[h(x)]^2}$

$$= \frac{\ddot{f}_0(x) \cdot f_0(x) + \dot{f}_0(x) \cdot \dot{f}_0(x)}{[f_0(x)]^2} = \frac{\ddot{f}_0(x)}{f_0(x)} - \left(\frac{\dot{f}_0(x)}{f_0(x)} \right)^2$$

Now, $\ddot{\ell}_x(\theta)$ has expectation

$$E[\ddot{\ell}_x(\theta)] = \int \frac{\ddot{f}_0(x)}{f_0(x)} \cdot f_0(x) dx - \int \left(\frac{\dot{f}_0(x)}{f_0(x)} \right)^2 \cdot f_0(x) dx$$

$$= \int \ddot{f}_0(x) dx - \int [\dot{\ell}_x(\theta)]^2 f_0(x) dx$$

$$\begin{aligned}
 &= \frac{\partial^2}{\partial \theta^2} \int f_0(x) dx - I_0 \\
 &= \frac{\partial^2}{\partial \theta^2} g(\theta) - I_0
 \end{aligned}$$

Definition of Fisher Information

- The expectation of the second derivative of the log-likelihood function is equal to the negative of the Fisher information.

Done with the definitions. Now suppose $\underline{X} = (X_1, \dots, X_n)$ is sample from $f_0(X)$. Then

$$\dot{\ell}_{\underline{X}}(\theta) = \sum_{i=1}^n \dot{\ell}_{X_i}(\theta)$$

Because $\dot{\ell}_{\underline{X}}(\theta)$ is in log space

Further,
$$-\ddot{\ell}_{\underline{X}}(\theta) = -\sum_{i=1}^n \ddot{\ell}_{X_i}(\theta)$$

- Since $\hat{\theta}^{MLE}$ for the full sample \underline{X} satisfies maximizing condition $\dot{\ell}_{\underline{X}}(\hat{\theta}^{MLE})$. Then, we can get Taylor Series approximation...

↑
Recall this is just a local linearization

$$0 = \dot{\ell}_{\underline{X}}(\hat{\theta}^{MLE}) \approx \dot{\ell}_{\underline{X}}(\theta) + \ddot{\ell}_{\underline{X}}(\theta) \cdot (\hat{\theta}^{MLE} - \theta)$$

$$\Rightarrow \theta \ddot{\ell}_{\underline{X}}(\theta) \approx \dot{\ell}_{\underline{X}}(\theta) + \ddot{\ell}_{\underline{X}}(\theta) \cdot \hat{\theta}^{MLE}$$

$$\Rightarrow \hat{\theta}^{MLE} \approx \frac{\theta \cdot \ddot{\ell}_{\underline{X}}(\theta) - \dot{\ell}_{\underline{X}}(\theta)}{\ddot{\ell}_{\underline{X}}(\theta)}$$

$$\Rightarrow \hat{\theta}^{MLE} \approx \theta - \frac{\dot{\ell}_{\underline{X}}(\theta)}{\ddot{\ell}_{\underline{X}}(\theta)}$$

$$\ell_x(\theta)$$

Finally, recall $\ell_x(\theta) = \sum \ell_{xi}(\theta)$, so it's a sum of a function of random variables.

while we like sums of random variables, we love means of random variables (why?)

Central Limit Theorem ^{the version}: If you take the mean of a large # of random variables, it will approximately follow a Normal distr.

By CLT... $\frac{\dot{\ell}_x(\theta)}{n} \sim N(0, \underbrace{I_0/n}_{\text{variance}})$

Then, it is convenient to write

$$\hat{\theta}_{MLE} \approx \underline{\theta} - \frac{\dot{\ell}_x(\theta)/n}{\ddot{\ell}_x(\theta)/n}$$

As $n \rightarrow \infty$, we know the numerator $\sim N(0, \frac{I_0}{n})$

Further $\frac{\ddot{\ell}_x(\theta)}{n} \rightarrow \frac{E_0[\ddot{\ell}_x(\theta)]}{n} = \underline{-I_0}$

Using the denominator as a constant and recalling $V(c \cdot X) = c^2 \cdot V(X)$

we get that $\frac{\dot{\ell}_x(\theta)/n}{\ddot{\ell}_x(\theta)/n} \sim N(0, \frac{I_0/n}{I_0})$
 $\sim N(0, \frac{1}{n \cdot I_0})$

Thus, $\hat{\theta}_{MLE} \xrightarrow{\text{Approx}} N(\theta, \frac{1}{n \cdot I_0})$

This is the Fisher Information
for n iid observations.

$$E\left(\theta - \frac{\dot{\ell}_x(\theta)}{\dot{\ell}_x'(\theta)}\right) = E(\theta) - E\left(\frac{\dot{\ell}_x(\theta)}{\dot{\ell}_x'(\theta)}\right) \\ = \theta - 0 = \theta$$

Equivalent to writing

$$\hat{\theta}^{MLE} \sim N\left(\theta, \frac{1}{I_0(n)}\right)$$

