

Problem set 2

Notation

$$y_i \quad i = 1, \dots, N$$

(scalar)

$$x_i = \begin{pmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ki} \end{pmatrix} \quad \begin{matrix} \text{p.d.f.} \\ f(y_i | x_i, \theta) \end{matrix}$$

$$\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_k \end{pmatrix}$$

$$L(\theta) = \prod_{i=1}^N f(y_i | x_i, \theta)$$

likelihood

$$L(\theta) = \sum_{i=1}^n \ln f(y_i | x_i, \theta)$$

log likelihood

$$f_i(\theta) \equiv f(y_i | x_i, \theta)$$

$$\ell_i(\theta) \equiv \ln f_i(\theta)$$

$$\text{Then } \mathcal{L}(\theta) = \sum_{i=1}^n \ell_i(\theta)$$

$$\mathcal{L}(\theta) \rightarrow \max_{\theta \in \Theta}$$

FOC

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_1} \stackrel{!}{=} 0$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_2} \stackrel{!}{=} 0$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_k} \stackrel{!}{=} 0$$

$$\begin{aligned} \underset{\theta}{\nabla} \mathcal{L}(\theta) &= \left[\frac{\partial \mathcal{L}(\theta)}{\partial \theta_1} \dots \frac{\partial \mathcal{L}(\theta)}{\partial \theta_k} \right] \\ &= \frac{\partial \mathcal{L}(\theta)}{\partial \theta} \end{aligned}$$

$$\nabla_{\theta} \left(\sum \ell_i \right) = \sum \underbrace{\nabla_{\theta} \ell_i(\theta)}_{\text{score for } i\text{th observation}}$$

$$s(y_i | x_i, \theta) = s_i(\theta) = \nabla_{\theta} \ell_i(\theta) \quad \text{① transpose}$$

$K \times 1$

$$= \nabla_{\theta} \ln f_i(\theta)' = \frac{1}{f_i(\theta)} \nabla f_i(\theta)' \stackrel{!}{=} 0$$

$$= \frac{1}{f_i(\theta)} \begin{bmatrix} \partial f_i(\theta) / \partial \theta_1 \\ \vdots \\ \partial f_i(\theta) / \partial \theta_K \end{bmatrix}$$

Why ∇f is a row?

$$f(x) \approx f(x_0) + \nabla f(x_0) (\underline{x - x_0})$$

Taylor expansion

$$E[s_i(\theta) | x_i] = E\left(\frac{1}{f_i(\theta)} \nabla_{\theta} f_i(\theta)\right)$$

$$= \int_{-\infty}^{+\infty} \frac{1}{f_i(\theta)} \nabla_{\theta} f_i(\theta) * f_i(\theta_0) dy_i$$

Note : $f_i(\theta_0)$ is the true density function. We take expectation wrt to true density

So for $s_i(\theta_0)$ we have

$$\begin{aligned} E(s_i(\theta_0) | x_i) &= \int_{-\infty}^{\infty} \frac{1}{f_i(\theta_0)} \nabla_{\theta} f_i(\theta_0) \cancel{f_i(\theta_0)} dy \\ &= \int_{-\infty}^{\infty} \nabla_{\theta} f_i(\theta_0) dy; \end{aligned}$$

For discrete y_i

Example : $y_i = \begin{cases} 1 & \text{w.p. } P(x_i, \theta) \\ 0 & \text{w.p. } 1 - P(x_i, \theta) \end{cases}$

Then

$$\begin{aligned} E(s_i(\theta_0) | x_i) &= \frac{1}{P(x_i, \theta_0)} \nabla_{\theta} P(x_i, \theta_0) \cdot P(x_i, \theta_0) \\ &+ \frac{1}{1 - P(x_i, \theta_0)} (-\nabla_{\theta} P(x_i, \theta_0)) \\ &\times (1 - P(x_i, \theta_0)) \\ E(s_i(\theta_0) | x_i) &= \int_{-\infty}^{+\infty} \nabla_{\theta} f_i(\theta) dy \end{aligned}$$

j indexes possible values of y_i

OR

$$E(s_i(\theta_0) | x_i) = \sum_j \nabla_{\theta} P_j(\theta_0) = \sum_j \nabla_{\theta} P_j(y_i | x_i, \theta_0)$$

Notice that

$$\int_{-\infty}^{+\infty} f_i(\theta) dy = 1$$

$$\sum_j P_j(\theta) = 1$$

for any θ

$$\nabla_{\theta} \int f_i(\theta) dy = \nabla_{\theta}(1) = 0$$

If $f_i(\cdot)$ is "good", can bring
 ✤ under \int

$$\int_{-\infty}^{\infty} \nabla_{\theta} f_i(\theta) dy = 0 = E(s_i(\theta_0) | x_i)$$

$E(s_i(\theta_0) | x_i)$ does not depend on x_i

$$\Rightarrow \boxed{E(s_i(\theta_0)) = 0} \quad (1)$$

Second derivative

$$\nabla_{\theta}^2 L(\theta) = \begin{pmatrix} \frac{\partial^2 L(\theta)}{\partial \theta_1^2} & \frac{\partial^2 L(\theta)}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 L(\theta)}{\partial \theta_1 \partial \theta_k} \\ \vdots & \ddots & & \vdots \\ \frac{\partial^2 L(\theta)}{\partial \theta_k \partial \theta_1} & \frac{\partial^2 L(\theta)}{\partial \theta_k \partial \theta_2} & \cdots & \frac{\partial^2 L(\theta)}{\partial \theta_k^2} \end{pmatrix}$$

$$= \nabla_{\theta} (\nabla_{\theta} L(\theta)') = \nabla_{\theta} (\sum s_i(\theta))$$

$$= \sum_{i=1}^N \underbrace{\nabla s_i(\theta)}_{H_i(\theta)} \quad \begin{matrix} \text{Hessian matrix for} \\ \text{i-th obs.} \end{matrix}$$

$$H_i(\theta) = \nabla_{\theta} \left[\frac{1}{f_i(\theta)} \nabla_{\theta} f_i(\theta) \right]$$

scalar

$K \times 1$

$$= -\frac{1}{f_i(\theta)^2} \nabla_{\theta} f_i(\theta)' \nabla_{\theta} f_i(\theta) = S_i(\theta) S_i(\theta)$$

since

$S_i(\theta) = \frac{1}{f_i(\theta)} \nabla_{\theta} f_i(\theta)'$

$\nabla_{\theta} \left(\frac{1}{f_i(\theta)} \right)$

$$+ \frac{1}{f_i(\theta)} \nabla_{\theta}^2 f_i(\theta)$$

$\nabla_{\theta} (\nabla_{\theta} f_i(\theta))$

The results here are correct but the derivation is a bit sloppy.
See the added sheet

$$H_i(\theta) = -S_i(\theta) S_i(\theta)' + \frac{1}{f_i(\theta)} \nabla_{\theta}^2 f_i(\theta)$$

$$E(S_i(\theta) S_i(\theta)' | X_i) = V(S_i(\theta) | X_i)$$

variance-covariance matrix of s_i

$$= -E(H_i(\theta) | X_i) + E\left[\frac{1}{f_i(\theta)} \nabla_{\theta}^2 f_i(\theta) | X_i\right]$$

= 0 → same trick as for $E\left[\frac{1}{f_i(\theta)} \nabla_{\theta} f_i(\theta) | X_i\right] = 0$

Additional sheet

This is the rigorous and more detailed derivation of the stuff in the purple box above.

$$\begin{aligned}
 H_i(\theta) &= \nabla_{\theta} \left[\frac{1}{f_i(\theta)} \nabla_{\theta} f_i(\theta) \right]' = \nabla_{\theta} \left[\nabla_{\theta} f_i(\theta) \frac{1}{f_i(\theta)} \right] \\
 &\quad \left[(AB)' = B' A' \right] \xrightarrow{\text{scalar}} \\
 &= \nabla_{\theta}^2 f_i(\theta) \frac{1}{f_i(\theta)} + \nabla_{\theta} f_i(\theta)' \nabla_{\theta} \left[\frac{1}{f_i(\theta)} \right] \\
 &= \nabla_{\theta}^2 f_i(\theta) \frac{1}{f_i(\theta)} + \nabla_{\theta} f_i(\theta)' \left(-\frac{1}{f_i(\theta)^2} \nabla_{\theta} f_i(\theta) \right) \\
 &\quad \text{chain rule} \\
 &= \frac{1}{f_i(\theta)} \nabla_{\theta}^2 f_i(\theta) - \frac{1}{f_i(\theta)^2} \nabla_{\theta} f_i(\theta)' \nabla_{\theta} f_i(\theta) \\
 &= -\frac{1}{f_i(\theta)} \nabla_{\theta}^2 f_i(\theta) - s_i(\theta) s_i(\theta)'
 \end{aligned}$$

Note that the order of multiplication matters!

$$V(S_i(\theta_0) | x_i) = E(S_i(\theta_0) S_i(\theta_0)' | x_i) \\ = -E(H_i(\theta_0) | x_i)$$

$$E(S_i(\theta_0) S_i(\theta_0)') = -E(H_i(\theta_0)) \quad (2)$$

1(B) Why the log?

→ computational convenience

→ $\log \Rightarrow$ sums \Rightarrow can apply LLN
CLT

→ It can be shown that

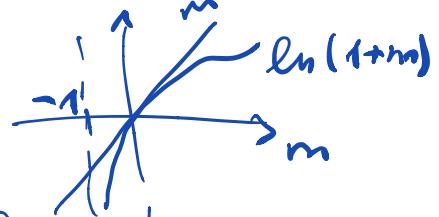
θ_0 maximizes $E(\mathcal{L}(\theta))$

Slide 55

$$\underbrace{\forall \theta \in \Theta, E[\ln f(y_i | x_i, \theta_0)]}_{\text{for all}} \geq E[\ln f(y_i | x_i, \theta)]$$

To show that, we are going to use

- $\ln(1+m) \leq m$



- $f(z), g(z)$ density functions

$$S_f = \{z \mid f(z) > 0\}$$

$$\int_{S_f} f(z) dz \geq \int_{S_f} g(z) dz$$

$$J(f, g) = \int_{S_f} \ln \frac{f(z)}{g(z)} f(z) dz \geq 0$$

$$= - \int_{S_f} \ln \frac{g(z)}{f(z)} f(z) dz$$

$$\geq - \int_{S_f} \left(\ln \frac{g(z)}{f(z)} f(z) + \underbrace{f(z) - g(z)} \right) dz$$

$$= - \int \left[\ln \left(\frac{g}{f} \right) - \underbrace{\left(\frac{g}{f} - 1 \right)}_m \right] dz$$

$$= - \int [\ln(1+m) - m] dz \geq 0$$

$\gamma(f, g)$ is called the Kullback-Leibler distance / information criterion

$$\gamma(f, g) = \int_{S_f} \ln \frac{f(z)}{g(z)} f(z) dz$$

$$= \int_{S_f} \ln f(z) \cdot f(z) dz - \int_{S_f} \ln g(z) \cdot f(z) dz$$

$$= E_f \ln f(z) - E_f (\ln g(z)) \geq 0$$

Substitute $f_i(\theta)$ for f

$f_i(\theta)$ for g

$$\Rightarrow E[\ln f_i(y_i | x_i, \underline{\theta_0})] \geq E[\ln f_i(y_i | x_i, \underline{\theta})]$$

$$\hat{\theta}_0 = \arg \max_{\theta \in \Theta} \mathbb{E}[\ln f_i(\theta)]$$

1c Asymptotic distribution of

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \mathcal{L}(\theta)$$

$$\mathcal{L}(\theta) = \sum_{i=1}^N l_i(\theta) = \sum_{i=1}^N \ln f_i(\theta)$$

If we have a sample of size N , maximization does not depend on N .

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N l_i(\theta)$$

Our observations are i.i.d.

$$\forall \theta \quad \frac{1}{N} \sum_{i=1}^N l_i(\theta) \xrightarrow{\text{P}} E[l_i(\theta)] \\ = E(\ln f_i(\theta))$$

pointwise convergence

It would seem logical

$$\operatorname{argmax}_{\theta} \frac{1}{N} \sum_i l_i(\theta) \xrightarrow{\text{P}} \operatorname{argmax}_{\theta} E(l_i(\theta))$$

Unfortunately, pointwise convergence is not enough.

Recall

$$x_n \xrightarrow{\text{P}} a \text{ if}$$

$$\forall \varepsilon > 0, P(|x_n - a| > \varepsilon) \xrightarrow{\cdot} 0$$

Example

$$f_n(x) = n x^n (1-x), x \in [0,1]$$

$$\lim_{n \rightarrow \infty} \underset{x \rightarrow 0}{\overset{\infty}{\underset{\rightharpoonup}{\lim}}} n x^n (1-x) = [\infty, 0]$$

$$= \lim_{n \rightarrow \infty} \frac{n(1-x)}{\left(\frac{1}{x}\right)^n} = \left[\frac{\infty}{\infty} \right]$$

L'Hopital's rule

$$\lim_{x \rightarrow \infty} \frac{h(x) \rightarrow \infty}{m(x) \rightarrow \infty}$$

$$= \lim_{x \rightarrow \infty} \frac{h'(x)}{m'(x)}$$

Analogously, for sequences

$$\lim_{n \rightarrow \infty} \frac{h(n)}{m(n)} = \lim_{n \rightarrow \infty} \frac{h(n+1) - h(n)}{m(n+1) - m(n)}$$

Numerator : $n(1-x)$

$$[(n+1) - n](1-x) = 1-x$$

Denominator : $(\frac{1}{x})^n$

$$\left(\frac{1}{x}\right)^{n+1} - \left(\frac{1}{x}\right)^n = \left(\frac{1}{x}\right)^n \left(\frac{1-x}{x}\right)$$

$$\lim_{n \rightarrow \infty} \frac{n(1-x)}{\left(\frac{1}{x}\right)^n} = \lim_{n \rightarrow \infty} \frac{\frac{1-x}{x}}{\left(\frac{1}{x}\right)^n \frac{1-x}{x}} =$$

$$= \lim_{n \rightarrow \infty} \frac{x}{(1/x)^n} = \underline{0}$$

$$f_n(x) = n x^n (1-x)$$

$$\frac{d f_n(x)}{dx} = n^2 x^{n-1} (1-x) - n x^n \stackrel{!}{=} 0$$

$$x^* = \frac{n}{n+1} \rightarrow \boxed{1}$$

What we need is uniform convergence:

$$\max_{\theta \in \Theta} \left(N^{-1} \sum_i \ell_i(\theta) - \mathbb{E} \ell_i(\theta) \right) \xrightarrow{P} 0$$

$\rightarrow \Theta$ is closed and bounded

$\rightarrow \ell_i(\theta)$ continuous for all values of y and x

$\rightarrow \underline{|\ell_i(\theta)| < \delta(y, x_i)}, \quad \delta(\cdot) > 0$

$$\underline{\mathbb{E}(\delta(\cdot)) < \infty}$$

Uniform convergence +
identification assumption

$$\hat{\theta}_0 = \arg \max_{\theta \in \Theta} E(\ln f_i(\theta))$$

$\hat{\theta}_0$ is unique

$$\Rightarrow \hat{\theta} \xrightarrow{P} \theta_0$$

$\hat{\theta}$ is consistent

Asymptotic normality

Multivariate normal distr.

P.d.f.

$$g(x) = \frac{\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))}{(\sqrt{2\pi})^n \sqrt{\det \Sigma}}$$

μ is the mean

Σ is the V-Cov matrix

$$\mathcal{N}(\mu, \Sigma)$$

$$E(S_i(\theta_0)) = 0$$

$$E(S_i(\theta_0) S_i(\theta_0)') = -E(H_i(\theta_0))$$

Taylor expansion of $S(\cdot)$ around θ_0

$$\sum_{i=1}^n S_i(\hat{\theta}) = \underbrace{\sum_i S_i(\theta_0)}_{=0 \text{ (FOC)}} + \underbrace{\sum_i H_i(\theta_0)(\hat{\theta} - \theta_0)}_{+ R(\hat{\theta})}$$

$R(\hat{\theta})$ converges to 0 faster

than $\hat{\theta} \rightarrow \theta_0$