

# ST2132 Variance of ML Estimators

Semester 1 2022/2023

# Large-sample variance of ML estimator

Let  $\hat{\theta}_n$  be the ML estimator of  $\theta \in \Theta \subset \mathbb{R}$ , based on IID RV's  $X_1, \dots, X_n$  with density  $f(x|\theta)$ .

- ▶ In virtually all applications, as  $n \rightarrow \infty$ ,

$$\text{var}(\hat{\theta}_n) \approx \frac{\mathcal{I}(\theta)^{-1}}{n}$$

$\mathcal{I}(\theta): p \times p$

where the Fisher information  $\mathcal{I}(\theta)$  is determined by  $f(x|\theta)$ .

- ▶ For large  $n$ , Monte Carlo may not be required to estimate the SE of an ML estimate. Bootstrap is still needed.
- ▶ The MOM estimate is given by a formula more often than the ML estimate, but the opposite is true for SE, for large  $n$ .

# Fisher information: Poisson

$X \sim \text{Poisson}(\lambda)$ :

$$f(x) = \frac{\lambda^x e^{-\lambda}}{x!}, x = 0, 1, \dots$$

→  $\log f(X) = X \log \lambda - \lambda - \log X!$

$$\frac{d \log f(X)}{d\lambda} = \frac{X}{\lambda} - 1$$

$$\frac{d^2 \log f(X)}{d\lambda^2} = -\frac{X}{\lambda^2}$$

The Fisher information is

$$\mathcal{I}(\lambda) = -E \left( \frac{d^2 \log f(X)}{d\lambda^2} \right) = \frac{1}{\lambda}$$

# Fisher information: Bernoulli

$X \sim \text{Bernoulli}(p)$ :

$$f(x) = p^x(1-p)^{(1-x)}, \quad x = 0, 1$$

$$\log f(X) = X \log p + (1-X) \log(1-p)$$

$$\frac{d \log f(X)}{dp} = \frac{X}{p} - \frac{1-X}{1-p} \quad \frac{d^2 \log f(X)}{dp^2} = -\frac{X}{p^2} - \frac{1-X}{(1-p)^2}$$

The Fisher information is

$$\mathcal{I}(p) = -E \left( \frac{d^2 \log f(X)}{dp^2} \right) = \frac{1}{p} + \frac{1}{1-p} = \frac{1}{p(1-p)}$$

# Fisher information: normal

$$X \sim N(\mu, \sigma^2), \theta = (\mu, \sigma).$$

$$E(X - \mu)^2 = \sigma^2 \quad \frac{\partial}{\partial \sigma} \left( \frac{X - \mu}{\sigma^2} \right)$$

$$\left[ \begin{array}{c} \frac{\partial}{\partial \mu} \\ \frac{\partial}{\partial \sigma} \end{array} \right] f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}, \quad -\infty < x < \infty$$

$$\log f(X) = -\frac{\log 2\pi}{2} - \log \sigma - \frac{(X - \mu)^2}{2\sigma^2}$$

$$\frac{d \log f(X)}{d\theta} =$$

$$\left[ \begin{array}{cc} \frac{\partial^2}{\partial \mu^2} & \frac{\partial^2}{\partial \mu \partial \sigma} \\ \frac{\partial^2}{\partial \sigma \partial \mu} & \frac{\partial^2}{\partial \sigma^2} \end{array} \right]$$

$$\cdot \left[ \begin{array}{c} \frac{X - \mu}{\sigma^2} \\ -\frac{1}{\sigma} + \frac{(X - \mu)^2}{\sigma^3} \end{array} \right]$$

$$\mathcal{I}(\theta) = \left[ \begin{array}{cc} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{array} \right]$$

$$\frac{d^2 \log f(X)}{d\theta^2} =$$

$$\left[ \begin{array}{cc} -\frac{1}{\sigma^2} & -\frac{2(X - \mu)}{\sigma^3} \\ -\frac{2(X - \mu)}{\sigma^3} & \frac{1}{\sigma^2} - \frac{3(X - \mu)^2}{\sigma^4} \end{array} \right]$$

# Definition: Fisher information

- ▶ Let  $X$  have density  $f(x|\theta)$ ,  $\theta \in \Theta \subset \mathbb{R}^p$ . The Fisher information is the  $p \times p$  matrix

$$\mathcal{I}(\theta) = -\mathbb{E} \left[ \frac{d^2 \log f(X)}{d\theta^2} \right]$$

$$\frac{d \log f(X)}{d\theta}$$

$p \times 1$

- ▶  $\mathcal{I}(\theta)$  is symmetric, with  $(i, j)$ -entry

$$-\mathbb{E} \left[ \frac{\partial^2 \log f(X)}{\partial \theta_i \partial \theta_j} \right]$$

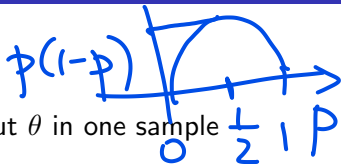
which is

$$-\int_{-\infty}^{\infty} \frac{\partial^2 \log f(x)}{\partial \theta_i \partial \theta_j} f(x) dx, \quad \sum_x \frac{\partial^2 \log f(x)}{\partial \theta_i \partial \theta_j} f(x)$$

according to whether  $X$  is continuous or discrete.

# Interpretation

- $\mathcal{I}(\theta)$  measures the information about  $\theta$  in one sample  $X \sim f(x|\theta)$ .



$$\mathcal{I}(\lambda) = \frac{1}{\lambda}$$

- Let  $X \sim \text{Poisson}(\lambda)$ . The larger  $\lambda$ , the less information in  $X$ .

- For which value of  $p$  does a Bernoulli( $p$ ) sample have the least information?

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

- Guess: How much information about  $p$  is in  $n$  IID Bernoulli( $p$ ) samples?

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

- ▶ IID  $X_1, \dots, X_n$  with density  $f(x|\theta)$  can be regarded as a sample from  $\mathbf{X} = (X_1, \dots, X_n)$  with random joint density

$$g(\mathbf{X}|\theta) = f(X_1|\theta) \cdots f(X_n|\theta)$$

- ▶ The information in  $\mathbf{X}$ ,

$$-E \left[ \frac{d^2 \log g(\mathbf{X})}{d\theta^2} \right]$$

is  $n\mathcal{I}(\theta)$ , where  $\mathcal{I}(\theta)$  is the information in any one of the  $X$ 's.



# Binomial

$X \sim \text{Binomial}(n, p)$ :  $f(x) = \binom{n}{x} p^x (1-p)^{n-x}$ ,  $x = 0, 1, \dots, n$ .

$$\log f(X) = \log \binom{n}{X} + X \log p + (n - X) \log(1 - p)$$

$$\frac{d \log f}{dp} = \frac{X}{p} - \frac{n-X}{1-p}$$

$$\frac{d^2 \log f}{dp^2} = -\frac{X}{p^2} - \frac{n-X}{(1-p)^2}$$

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

# Bernoulli vs binomial

- ▶ One binomial( $n, p$ ) sample has the same information about  $p$  as  $n$  IID Bernoulli( $p$ ) samples.

*binomial(1, p)*

- ▶ Surprising? The binomial sample only tells us the number of successes,  $\sum_{i=1}^n X_i$ , while the  $n$  Bernoulli samples tell us the full sequence of successes and failures:  $X_1, \dots, X_n$ .

- ▶ Unsurprising? The ML estimators of  $p$  are identical.

- ▶ Similarly, the multinomial( $n, (p_1, \dots, p_r)$ ) has the same information about  $(p_1, \dots, p_r)$  as  $n$  IID mult(1,  $(p_1, \dots, p_r)$ ) samples.

*poss. values:  $(1, 0, \dots, 0) \dots (0, \dots, 0, 1)$*

# HWE trinomial distribution

$$\log f(\mathbf{X}) = c + (2X_1 + X_2) \log(1-\theta) + (X_2 + 2X_3) \log \theta$$

$$\frac{d}{d\theta} = -\frac{2X_1 + X_2}{1-\theta} - \frac{X_2 + 2X_3}{\theta} \quad \frac{d^2}{d\theta^2} = -\frac{2X_1 + X_2}{(1-\theta)^2} - \frac{X_2 + 2X_3}{\theta^2}$$

$\mathbf{X} = (X_1, X_2, X_3) \sim \text{multinomial}(n, \mathbf{p})$ , where

$$p_1 = (1-\theta)^2, p_2 = 2\theta(1-\theta), p_3 = \theta^2, \quad 0 < \theta < 1$$

$$I(\theta) = \frac{2n}{1-\theta} + \frac{2n}{\theta} \quad \mathcal{I}(\theta) = \frac{2n}{\theta(1-\theta)}$$

$$E(X_2 + 2X_3) = 2n\theta$$

$$E(2X_1 + X_2) = 2n - 2n\theta \quad \leftarrow$$

Name two other sets of samples (number of IID samples, and the distribution) that have the same information.

# Multinomial data

- $\mathbf{X} \sim \text{multinomial}(n, (p_1, \dots, p_r))$ ,  $\theta = (p_1, \dots, p_{r-1})$ .

$$p_r = 1 - p_1 - \dots - p_{r-1}$$

$$\log f(\mathbf{X}) = \sum_{i=1}^r X_i \log p_i$$

$$+ C \frac{n!}{x_1! \dots x_r!} p_1^{x_1} \dots p_r^{x_r}$$

$$\rightarrow \frac{\partial \log f(\mathbf{X})}{\partial p_i} = \frac{X_i}{p_i} - \frac{X_r}{p_r} \quad 1 \leq i \leq r-1$$

$$\left\{ \begin{aligned} \frac{\partial^2 \log f(\mathbf{X})}{\partial p_i^2} &= -\frac{X_i}{p_i^2} - \frac{X_r}{p_r^2} \\ \frac{\partial^2 \log f(\mathbf{X})}{\partial p_i \partial p_j} &= -\frac{X_r}{p_r^2} \end{aligned} \right. \quad 1 \leq i \leq r-1$$

$$\frac{\partial^2 \log f(\mathbf{X})}{\partial p_i \partial p_j} = -\frac{X_r}{p_r^2} \quad 1 \leq i \neq j \leq r-1$$

- $(i, j)$ -entry of  $\mathcal{I}(\theta)$ :

$$r=3 \quad p_3 = 1 - p_1 - p_2$$

$$\frac{\partial}{\partial p_1} = \frac{X_1}{p_1} - \frac{X_3}{p_3}$$

$$\frac{n}{p_i} + \frac{n}{p_r}, \quad i=j$$

$$\frac{n}{p_r}, \quad i \neq j$$

$$\frac{\partial}{\partial p_2} = \frac{X_2}{p_2} - \frac{X_3}{p_3}$$

$$X_1 \log p_1 + X_2 \log p_2 + X_3 \log p_3$$

# Multinomial parameterisation

- ▶ The general trinomial distribution has two parameters:  $p_1, p_2$ , or any other set of two probabilities. The HWE trinomial is defined by a single parameter.
- ▶ They are opposite ends of a collection of multinomial distributions with the same  $n$  and  $r$ , but of varying number of independent parameters.
- ▶ Any such distribution can be described as

$$\text{Multinomial}(n, \mathbf{p}(\theta))$$

$$\mathbf{p}(\theta) = (p_1(\theta), \dots, p_r(\theta))$$

$$\theta \in \Theta \subset \mathbb{R}^k, 1 \leq k \leq r-1.$$

## Summary: Fisher information

- ▶ Let  $X$  have mass/density  $f(x|\theta)$ ,  $\theta \in \Theta \subset \mathbb{R}^p$ . The Fisher information at  $\theta$  in  $X$  is the  $p \times p$  matrix

$$-E \left[ \frac{d^2 \log f(X)}{d\theta^2} \right]$$

$X$  can be a random vector.

- ▶ Information in  $n$  IID samples is  $n$  times that in one sample.
- ▶ A binomial( $n, p$ ) sample has the same information as  $n$  IID Bernoulli( $p$ ) samples. Similarly, a multinomial( $n, \mathbf{p}$ ) has the same information as  $n$  IID multinomial( $1, \mathbf{p}$ ) samples.

# Variance of ML estimator

- ▶  $X_1, \dots, X_n$  IID Poisson( $\lambda$ ). What is  $\mathcal{I}(\lambda)$ , the information in any  $X_i$ ? What is the variance of the ML estimator  $\hat{\lambda} = \bar{X}$ ?

$$\mathcal{I}(\lambda) = \frac{1}{\lambda} \quad \frac{\mathcal{I}(\lambda)^{-1}}{n} \quad \text{var}(\bar{X}) = \frac{\lambda}{n}$$

- ▶  $X_1, \dots, X_n$  IID Bernoulli( $p$ ). What is  $\mathcal{I}(p)$ , the information in any  $X_i$ ? What is the variance of the ML estimator  $\hat{p} = \bar{X}$ ?

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

- ▶ What is your conjecture?

$$\text{var}(\hat{\theta}_n) = \frac{\mathcal{I}(\theta)^{-1}}{n}$$

# Variance of ML estimator: normal

- $X_1, \dots, X_n$  are IID normal( $\mu, \nu$ ). Let  $\theta = (\mu, \nu)$ . What is  $\mathcal{I}(\theta)$ , the information in any  $X_i$ ?

$$\mathcal{I}(\mu, \nu) = \begin{bmatrix} \frac{1}{\nu} & 0 \\ 0 & \frac{1}{2\nu^2} \end{bmatrix} \quad \frac{\mathcal{I}(\mu, \nu)^{-1}}{n} = \begin{bmatrix} \frac{\nu}{n} & 0 \\ 0 & \frac{2\nu^2}{n} \end{bmatrix}$$

- What is the variance of the ML estimator  $\hat{\theta}_n = (\bar{X}, \hat{\sigma}^2)$ ?

$$\text{var}(\hat{\theta}_n) = \begin{bmatrix} \frac{\nu}{n} & 0 \\ 0 & \frac{2\nu^2}{n} \end{bmatrix}$$

- Revise conjecture?

$$\text{var}(\hat{\theta}_n) \approx \frac{\mathcal{I}(\theta)^{-1}}{n} \text{ for large } n$$



# Gamma distribution

$X \sim \text{Gamma}(\alpha, \lambda)$ ,  $\theta = (\alpha, \lambda)$ .

$$\log f(X) = \alpha \log \lambda + (\alpha - 1) \log X - \lambda X - \log \Gamma(\alpha)$$

$$\frac{\partial \log f(X)}{\partial \alpha} = \log \lambda + \log X - \psi(\alpha), \quad \frac{\partial \log f(X)}{\partial \lambda} = \frac{\alpha}{\lambda} - X$$

where  $\psi(\alpha) = \frac{d}{d\alpha} \log \Gamma(\alpha) = \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}$  is the digamma function.

$$\frac{\partial^2 \log f(X)}{\partial \alpha^2} = -\psi'(\alpha), \quad \frac{\partial^2 \log f(X)}{\partial \lambda^2} = -\frac{\alpha}{\lambda^2}$$

$$\frac{\partial^2 \log f(X)}{\partial \alpha \partial \lambda} = \frac{\partial^2 \log f(X)}{\partial \lambda \partial \alpha} = \frac{1}{\lambda}$$

$$\mathcal{I}(\theta) = -\mathbb{E} \left[ \frac{d^2 \log f(X)}{d\theta^2} \right] = \begin{bmatrix} \psi'(\alpha) & -\frac{1}{\lambda} \\ -\frac{1}{\lambda} & \frac{\alpha}{\lambda^2} \end{bmatrix}$$

## Revisiting rainfall again

- ▶ Based on 227 realisations, the ML estimate of  $\theta = (\alpha, \lambda)$  is  $(0.44, 1.96)$ . Let  $\hat{\alpha}$  and  $\hat{\lambda}$  be the ML estimators. Assuming 227 is large enough,

$$\text{var}(\hat{\alpha}, \hat{\lambda}) \approx \frac{\mathcal{I}(\alpha, \lambda)^{-1}}{227}$$

- ▶ Bootstrap approximation:

$$\text{var}(\hat{\alpha}, \hat{\lambda}) \approx \frac{\mathcal{I}(0.44, 1.96)^{-1}}{227} \approx \begin{bmatrix} 0.0011 & 0.0051 \\ 0.0051 & 0.0610 \end{bmatrix}$$

- ▶ SE in 0.44 is  $\text{SD}(\hat{\alpha}) \approx \sqrt{0.0011} \approx 0.03$ .  
SE in 1.96 is  $\text{SD}(\hat{\lambda}) \approx \sqrt{0.0610} \approx 0.25$ .  
They are very close to the Monte Carlo approximations 0.03 and 0.26 (Parameter Estimation II slide 20).

- ▶  $X_1, \dots, X_n$  IIR RV's with density  $f(x|\theta)$ . ML estimator is  $\hat{\theta}$ .  
As  $n \rightarrow \infty$ ,

$$\text{var}(\hat{\theta}) \approx \frac{\mathcal{I}(\theta)^{-1}}{n}$$

where  $\mathcal{I}(\theta)$  is the Fisher information in any  $X_i$  (one sample).

- ▶ For large  $n$ , SE in ML estimate can be approximated without Monte Carlo.
- ▶ Some technical conditions are required, but they are true in almost all applications. For technical details, see Rice's book.

## Special case: multinomial

The previous result applies to multinomial data, provided the Fisher information is suitably defined.

$\mathbf{X} \sim \text{Multinomial}(n, \mathbf{p}(\theta))$ ,  $\theta \in \Theta \subset \mathbb{R}^k$  with  $1 \leq k \leq r - 1$ . For large  $n$ ,

$$\text{var}(\hat{\theta}) \approx \frac{\mathcal{I}(\theta)^{-1}}{n}$$

where  $\mathcal{I}(\theta)$  is the information in  $\text{Multinomial}(1, \mathbf{p}(\theta))$ .

- ▶  $X \sim \text{Binomial}(n, p)$ ,  $\text{var}(\hat{p}) = \frac{p(1-p)}{n}$ ,  $\frac{1}{p(1-p)}$  the information in  $\text{Binomial}(1, p) = \text{Bernoulli}(p)$ .
- ▶ HWE trinomial:  $\text{var}(\hat{\theta}) = \frac{\theta(1-\theta)}{2n}$ ,  $\frac{1}{\theta(1-\theta)}$  the information in  $\text{Trinomial}(1, \mathbf{p}(\theta))$ .