

Asymptotics & optimality (Module 11)



Statistics (MAST20005) &
Elements of Statistics
(MAST90058)

School of Mathematics and Statistics
University of Melbourne

Semester 2, 2020

Aims of this module

- Explain some of the theory that we skipped in previous modules
- Show why the MLE is usually a good (or best) estimator
- Explain some related important theoretical concepts

Outline

Likelihood theory

Asymptotic distribution of the MLE

Cramér–Rao lower bound

Sufficient statistics

Factorisation theorem

Optimal tests

Previous claims (from modules 2 & 4)

The MLE is asymptotically:

- unbiased
- efficient (has the optimal variance)
- normally distributed



Can use the 2nd derivative of the log-likelihood (the 'observed information function') to get a standard error for the MLE.

Motivating example (non-zero binomial)

- Consider a factory producing items in batches. Let θ denote the proportion of defective items. From each batch 3 items are sampled at random and the number of defectives is determined. However, records are only kept if there is at least one defective.
- Let Y be the number of defectives in a batch.
- Then $Y \sim \text{Bi}(3, \theta)$,

$$\Pr(Y = y) = \binom{3}{y} \theta^y (1 - \theta)^{3-y}, \quad y = 0, 1, 2, 3$$

- But we only take an observation if $Y > 0$, so the pmf is

$$\Pr(Y = y \mid Y > 0) = \frac{\binom{3}{y} \theta^y (1 - \theta)^{3-y}}{1 - (1 - \theta)^3}, \quad y = 1, 2, 3$$


 $\Pr(Y=0)$

$$\ln L(\theta) = (x_1 + x_2 + 3x_3) \ln \theta \\ + \ln(1-\theta) (2x_1 + x_2) \\ - n \ln [1 - (1-\theta)^3]$$

$$\frac{\partial \ln L(\theta)}{\partial \theta} = \frac{(x_1 + 2x_2 + 3x_3)}{\theta} \\ - \frac{2x_1 + x_2}{1-\theta}$$

$$-(x_1 + x_2 + x_3) \frac{3(1-\theta)^2}{1-(1-\theta)^3} > 0$$

$$(x_1 + 2x_2 + 3x_3)(1-\theta)[1 - (1-\theta)^3]$$

$$-(2x_1 + x_2)\theta[1 - (1-\theta)^3]$$

$$-3(x_1 + x_2 + x_3)(1-\theta)^2\theta = 0$$

$$(x_1 + 2x_2 + 3x_3)[(1-\theta) - (1-\theta)^4] - (x_1 + x_2)\theta - (x_1 + 2x_2 + 3x_3)\theta(1-\theta)^3 = 0$$

$$t(\theta^3 - 3\theta^2 + 2\theta) - (2x_1 + x_2)\theta = 0$$

- Let X_i be the number of times we observe i defectives and let $n = X_1 + X_2 + X_3$ be the total number of observations.

- The likelihood is, $\text{no } \theta \text{ inside} \rightarrow \text{coefficient means order not matter}$

$$L(\theta) = \frac{n!}{x_1! x_2! x_3!} \left(\frac{3\theta(1-\theta)^2}{1-(1-\theta)^3} \right)^{x_1} \left(\frac{3\theta^2(1-\theta)}{1-(1-\theta)^3} \right)^{x_2} \left(\frac{\theta^3}{1-(1-\theta)^3} \right)^{x_3}$$

- This simplifies to,

$$L(\theta) \propto \frac{\theta^{x_1+2x_2+3x_3} (1-\theta)^{2x_1+x_2}}{(1-(1-\theta)^3)^n}$$

- After taking logarithms and derivatives, the MLE is found to be the smaller root of

$$t\theta^2 - 3t\theta + 3(t-n) = 0$$

where $t = x_1 + 2x_2 + 3x_3$.

- This gives:

$$\hat{\theta} = \frac{3t - \sqrt{-3t^2 + 12tn}}{2t}$$

$$\hat{\theta} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad \text{since } \hat{\theta} < 1$$

$$\hat{\theta} = \frac{3t \pm \sqrt{9t^2 - 12t(t-n)}}{2t} = \frac{3t \pm \sqrt{3t^2 + 12tn}}{2t}$$

$$\therefore \hat{\theta} = \frac{3t - \sqrt{3t^2 + 12tn}}{2t}$$

$$\tau(\theta^2 - 3\theta + 2) - (2x_1 + x_2) = 0$$

$$\tau(\theta^2 - 3\theta + 3(t-n)) = 0$$

- We now have the MLE...
- ...but finding its sampling distribution is not straightforward!
- In general, finding the exact distribution of a statistic is often difficult.
- We've used the Central Limit Theorem to approximate the distribution of the sample mean.
- Gave us approximate CIs for a population mean μ of the form,

$$\bar{x} \pm \Phi^{-1} \left(1 - \frac{\alpha}{2}\right) \times \frac{s}{\sqrt{n}}$$

- Similar results hold more generally for MLEs (and other estimators)

Definitions

- Start with the log-likelihood:

$$\ell(\theta) = \ln L(\theta)$$

- Taking the first derivative gives the **score function** (also known simply as the **score**). Let's call it U ,

$$U(\theta) = \frac{\partial \ell}{\partial \theta}$$

- Note: we solve $U(\hat{\theta}) = 0$ to get the MLE

- Taking the second derivative, and then it's negative, gives the observed information function (also known simply as the observed information). Let's call it V .

$$V(\theta) = -\frac{\partial U}{\partial \theta} = -\frac{\partial^2 \ell}{\partial \theta^2}$$

- This represents the curvature of the log-likelihood. Greater curvature \Rightarrow narrower likelihood around a certain value \Rightarrow the likelihood is more informative.

more curvature
↓

more precise estimate

Fisher information

randomness of data
↓

function of data is random
↓ same for U, V

has probability distribution
↓ sampling distribution

- All of the above are functions of the data (and parameters). Therefore they are random variables and have sampling distributions.
- For example, we can show that $\mathbb{E}(U(\theta)) = 0$.
- An important quantity is $I(\theta) = \mathbb{E}(V(\theta))$, which is the Fisher information function (or just the Fisher information). It is also known as the expected information function (or simply as the expected information).
- Many results are based on the Fisher information.
- For example, we can show that $\text{var}(U(\theta)) = I(\theta)$.
- More importantly, it arises in theory about the distribution of the MLE.

Asymptotic distribution

- The following is a key result:

$$\hat{\theta} \approx N\left(\theta, \frac{1}{I(\theta)}\right) \quad \text{as } n \rightarrow \infty$$

- unbiased asymptotically*
- It requires some conditions for it to hold. The main one being that the parameter should not be defining a boundary of the sample space (e.g. like in the boundary problem examples we've looked at).
 - Let's see a proof...

Asymptotic distribution (derivation)

- Assumptions:

- X_1, \dots, X_n is a random sample from $f(x, \theta)$
- Continuous pdf, $f(x, \theta)$
- θ is not a boundary parameter

- Suppose the MLE satisfies:

$$U(\hat{\theta}) = \frac{\partial \ln L(\hat{\theta})}{\partial \theta} = 0$$

Note: this requires that θ is not a boundary parameter.

- Taylor series approximation for $U(\hat{\theta})$ about θ :

taylor expansion

$$U(t) = U(\theta) + \frac{(t-\theta)U'(\theta)}{1}$$

first order taylor expansion

set $t = \hat{\theta}$

$$\begin{aligned} 0 = U(\hat{\theta}) &= \frac{\partial \ln L(\hat{\theta})}{\partial \theta} \approx \frac{\partial \ln L(\theta)}{\partial \theta} + (\hat{\theta} - \theta) \frac{\partial^2 \ln L(\theta)}{\partial \theta^2} \\ &= U(\theta) - (\hat{\theta} - \theta)V(\theta) \end{aligned}$$

$$\theta = U(\hat{\theta}) = U(\theta) + (\hat{\theta} - \theta)$$

$$= \frac{\partial \ln L(\theta)}{\partial \theta} + (\hat{\theta} - \theta) \frac{\partial \ln L(\theta)}{\partial \theta}$$

since $V(\theta) = -U'(\theta)$

$$\Rightarrow U(\theta) - (\hat{\theta} - \theta) V(\theta) = 0$$

- We can write this as:

$$V(\theta)(\hat{\theta} - \theta) \approx U(\theta)$$

- Remember that we have a random sample (iid rvs), so we have,

$$U(\theta) = \frac{\partial \ln L(\theta)}{\partial \theta} = \sum_{i=1}^n \frac{\partial \ln f(X_i, \theta)}{\partial \theta}$$

- Since the X_i are iid so are:

$$U_i = \frac{\partial \ln f(X_i, \theta)}{\partial \theta}, \quad i = 1, \dots, n.$$

u_i iid.

- And the same for:

$$V_i = -\frac{\partial^2 \ln f(X_i, \theta)}{\partial \theta^2}, \quad i = 1, \dots, n.$$

v_i iid

- Determine $\mathbb{E}(U_i)$ by integration by substitution and exchanging the order of integration and differentiation,

$$\begin{aligned}\mathbb{E}(U_i) &= \int_{-\infty}^{\infty} \frac{\partial \ln f(x, \theta)}{\partial \theta} f(x, \theta) dx = \int_{-\infty}^{\infty} \frac{\partial f(x, \theta)}{\partial \theta} \frac{f(x, \theta)}{f(x, \theta)} dx \\ &= \int_{-\infty}^{\infty} \frac{\partial f(x, \theta)}{\partial \theta} dx = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f(x, \theta) dx = \frac{\partial}{\partial \theta} 1 = 0\end{aligned}$$

- To get the variance of U_i , we start with one of the above results,

$$\int_{-\infty}^{\infty} \frac{\partial \ln f(x, \theta)}{\partial \theta} f(x, \theta) dx = 0$$

- Taking another derivative of both sides gives,

$$\int_{-\infty}^{\infty} \left\{ \frac{\partial^2 \ln f(x, \theta)}{\partial \theta^2} f(x, \theta) + \frac{\partial \ln f(x, \theta)}{\partial \theta} \frac{\partial f(x, \theta)}{\partial \theta} \right\} dx = 0$$

- But,

$$\frac{\partial f(x, \theta)}{\partial \theta} = \frac{\partial \ln f(x, \theta)}{\partial \theta} f(x, \theta)$$

- Combining the previous two equations gives,

$$\int_{-\infty}^{\infty} \left\{ \frac{\partial \ln f(x, \theta)}{\partial \theta} \right\}^2 f(x, \theta) dx = - \int_{-\infty}^{\infty} \frac{\partial^2 \ln f(x, \theta)}{\partial \theta^2} f(x, \theta) dx$$

- In other words,

$$\mathbb{E}(U_i^2) = \mathbb{E}(V_i)$$

previous slides

$$\begin{aligned} \frac{\partial \ln f(x, \theta)}{\partial \theta} f(x, \theta) &= \frac{1}{f(x, \theta)} \cdot \frac{\partial f(x, \theta)}{\partial \theta} \cdot f(x, \theta) \\ &= \frac{\partial f(x, \theta)}{\partial \theta} \end{aligned}$$

- Since $\mathbb{E}(U_i) = 0$ we also have $\mathbb{E}(U_i^2) = \text{var}(U_i)$, so we can conclude,

$$\text{var}(U_i) = \mathbb{E}(V_i)$$

- Thus $U = \sum_i U_i$ is the sum of iid rvs with mean 0 and this variance.
- Thus,

$$\text{var}(U) = n \mathbb{E}(V_i)$$

- Also, since $V = \sum_i V_i$, we can conclude that,

$$\mathbb{E}(V) = n \mathbb{E}(V_i)$$

- Note that this is just the Fisher information, i.e.

$$\mathbb{E}(V) = \text{var}(U) = I(\theta)$$

- Looking back at,

$$V(\theta)(\hat{\theta} - \theta) \approx U(\theta)$$

We want to know what happens to U and V as the sample size gets large.

- U has mean 0 and variance $I(\theta)$.
- Central Limit Theorem $\Rightarrow U \approx N(0, I(\theta))$.
- V has mean $I(\theta)$.
- Law of Large Numbers $\Rightarrow V \rightarrow I(\theta)$
- Putting these together gives, as $n \rightarrow \infty$,

$$I(\theta)(\hat{\theta} - \theta) \sim N(0, I(\theta))$$

$$V(\theta)(\hat{\theta} - \theta) \approx N(0, I(\theta))$$

$$I(\theta)(\hat{\theta} - \theta) \approx N(0, I(\theta))$$

$$\hat{\theta} \sim N(\theta, \frac{1}{I(\theta)})$$

- Equivalently,

$$\hat{\theta} \sim N\left(\theta, \frac{1}{I(\theta)}\right)$$

- This is a very powerful result. For large (or even modest) samples we do not need to find the exact distribution of the MLE but can use this approximation.
- In other words, as a standard error of the MLE we can use:

$$se(\hat{\theta}) = \frac{1}{\sqrt{I(\hat{\theta})}}$$

$E(V(\hat{\theta})) = I(\theta)$ expected version

if we know $I(\theta)$, or otherwise replace it with its realised
(observed) version.

$$se(\hat{\theta}) = \frac{1}{\sqrt{V(\hat{\theta})}} \rightarrow \text{observed version}$$

- Furthermore, we use the normal distribution to construct approximate confidence intervals.

Example (exponential distribution)

- X_1, \dots, X_n random sample from

$$f(x | \theta) = \frac{1}{\theta} e^{-x/\theta}, \quad 0 < x < \infty, \quad 0 < \theta < \infty$$

- MLE is \bar{X} . *one observation contribution*
- $\ln f(x | \theta) = -\ln \theta - x/\theta$, so

$$U_i(\theta) = \frac{\partial}{\partial \theta} \ln f(x | \theta) = -\frac{1}{\theta} + \frac{x}{\theta^2}$$

$$V_i(\theta) = -\frac{\partial^2}{\partial \theta^2} \ln f(x | \theta) = -\frac{1}{\theta^2} + \frac{2x}{\theta^3}$$

- Since $\mathbb{E}(X) = \theta$,

$$I_i(\theta) = \mathbb{E}(V_i(\theta)) = \mathbb{E}\left(-\frac{1}{\theta^2} + \frac{2X}{\theta^3}\right) = -\frac{1}{\theta^2} + \frac{2\theta}{\theta^3} = \frac{1}{\theta^2}$$

$$I(\theta) = \sum_{i=1}^n I_i(\theta)$$

- Then $I(\theta) = n/\theta^2$ and $\hat{\theta} \approx N(\theta, \theta^2/n)$
- Suppose we observe $n = 20$ and $\bar{x} = 3.7$. An approximate 95% CI is,

$$3.7 \pm 1.96 \sqrt{\frac{3.7^2}{20}} = (2.1, 5.3)$$

approximation

$$\Pr\left(-1.96 < \frac{\hat{\theta} - \theta}{\frac{\theta}{\sqrt{n}}} < 1.96\right)$$

$$\Rightarrow \frac{\hat{\theta}}{\sqrt{20}} \Rightarrow \frac{\hat{\theta}}{\sqrt{20}} \text{ to plug in}$$

$$\Rightarrow \Pr\left(-1.96 < \frac{3.7 - \theta}{\frac{\theta}{\sqrt{20}}} < 1.96\right)$$

Example (Poisson distribution)

$$f(x | \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$
$$\ln f(x | \lambda) = x \ln \lambda - \lambda - \ln x!$$
$$v_i(\lambda) = \frac{\lambda^x}{x!}$$
$$v_i(\lambda) = \frac{x}{\lambda}$$
$$z_i(u) = e(v_i(u)) = \frac{x}{\lambda}$$
$$z(u) = \sum z_i(u) = \frac{x}{\lambda}$$
$$z(u) = \frac{x}{\lambda}$$

- Same arguments hold for discrete distributions, e.g. $P_n(\lambda)$.

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

We have seen $\hat{\lambda} = \bar{X}$.

- $\ln f(x | \lambda) = x \ln \lambda - \lambda - \ln(x!)$, so

$$\frac{\partial \ln f(x | \lambda)}{\partial \lambda} = \frac{x}{\lambda} - 1, \quad \text{and} \quad \frac{\partial^2 \ln f(x | \lambda)}{\partial \lambda^2} = -\frac{x}{\lambda^2}$$

- Thus

$$-\mathbb{E}\left(-\frac{X}{\lambda^2}\right) = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$$

$$v_i(\lambda) = \frac{x}{\lambda^2}$$

$$E(v_i(\lambda)) = E\left(\frac{x}{\lambda^2}\right) = \frac{1}{\lambda^2} E(x)$$

$$I(\lambda) = E(v(\lambda)) = \frac{n}{\lambda}$$

$$= \pi^2 \lambda = \frac{1}{\lambda}$$

- Then $\hat{\lambda} \approx N(\lambda, \lambda/n)$

✓

- Suppose we observe $n = 40$ and $\bar{x} = 2.225$. An approximate 90% CI is,

$$2.225 \pm 1.645 \sqrt{\frac{2.225}{40}} = (1.837, 2.612)$$

✓

$$\text{var}(\hat{\theta}) = \frac{1}{I(\lambda)} = \frac{\lambda}{n}$$

Cramér–Rao lower bound

- How good can our estimator get?
- Suppose we know that it is unbiased.
- What is the minimum variance we can achieve?
- Under similar assumptions to before (esp. the parameter must not define a boundary), we can find a lower bound on the variance
- This is known as the Cramér–Rao lower bound
- It is equal to the asymptotic variance of the MLE.
- In other words, if we take any unbiased estimator T , then

$$\text{var}(T) \geq \frac{1}{I(\theta)}$$

↗ unbiased estimator ↗

⇒ MLE is asymptotically
the best estimator

Cramér–Rao lower bound (proof)

- Let T be an unbiased estimator of θ
- Consider its covariance with the score function

$$E(U) = 0$$

$$\text{cov}(T, U) = \mathbb{E}(TU) - \mathbb{E}(T)\mathbb{E}(U) = \mathbb{E}(TU)$$

$$= \int T \frac{\partial \ln L}{\partial \theta} L d\mathbf{x} \xrightarrow{\text{chain rule}} \int T \frac{\partial L}{\partial \theta} d\mathbf{x}$$

$$= \frac{\partial}{\partial \theta} \int TL d\mathbf{x} = \frac{\partial}{\partial \theta} \mathbb{E}(T) \xrightarrow{\text{unbiased estimator}} \frac{\partial}{\partial \theta} \theta = 1 \quad E(TU)$$

$$\frac{\partial \ln L}{\partial \theta} = 0 \quad \uparrow \quad = \int T(\mathbf{x}) \frac{\partial}{\partial \theta} \ln L(\mathbf{x}; \theta) f(\mathbf{x}) d\mathbf{x}$$

$$\text{cor}(T, U) = \frac{\text{cov}(T, U)}{\sqrt{\text{var}(T)} \sqrt{\text{var}(U)}} \leq 1$$

$$\Rightarrow 1 = \text{cov}(T, U)^2 \leq \text{var}(T) \text{var}(U)$$

$$\text{var}(T) \geq \frac{1}{\text{var}(U)} = \frac{1}{I(\theta)}$$

since we run
for all observations

Implications of the Cramér–Rao lower bound

you can get
biased estimator
With smaller variance

- If an unbiased estimator attains this bound, then it is best in the sense that it has minimum variance compared with other unbiased estimators.
- Therefore, MLEs are approximately (or exactly) optimal for large sample size because:
 - They are asymptotically unbiased
 - Their variance meets the Cramér–Rao lower bound asymptotically

Efficiency

- We can compare any unbiased estimator against the lower bound
- We define the **efficiency** of the unbiased estimator T as its variance relative to the lower bound,

with lower bound if 3?

→ 其實解是怎樣
scale it to [0,1]

$$\text{eff}(T) = \frac{1/I(\theta)}{\text{var}(T)} = \frac{1}{I(\theta) \text{var}(T)}$$

- Note that $0 \leq \text{eff}(T) \leq 1$
- If $\text{eff}(T) \approx 1$ we say that T is an **efficient** estimator

var(T) 的分子

∴ $\text{var}(T) = 0 \Rightarrow \text{eff}(T) = 0$

Example (exponential distribution)

- Sampling from an exponential distribution
- We saw that $I(\theta) = n/\theta^2$
- Therefore, the Cramér–Rao lower bound is θ^2/n . $\frac{1}{I(\theta)} = \frac{\theta^2}{n}$
- Any unbiased estimator must have variance at least as large as this.
- The MLE in this case is the sample mean, $\hat{\theta} = \bar{X}$
- Therefore, $\text{var}(\hat{\theta}) = \text{var}(X)/n = \theta^2/n$
- So the MLE is efficient (for all sample sizes!)

$\text{eff}(\hat{\theta}) = 1$
always asymptotically efficient
 \Rightarrow in this case is efficient

Outline

Likelihood theory

Asymptotic distribution of the MLE

Cramér–Rao lower bound

Sufficient statistics

Factorisation theorem

Optimal tests

Sufficiency: a starting example

- We toss a coin 10 times
- Want to estimate the probability of heads, θ
- $X_i \sim \text{Be}(\theta)$
- Suppose we use $\hat{\theta} = \frac{1}{2}(X_1 + X_2)$
- Only uses the first 2 coin tosses
- Clearly, we have not used all of the available information!

Motivation

- Point estimation reduces the whole sample to a few statistics.
- Different methods of estimation can yield different statistics.
- Is there a preferred reduction?
- Toss a coin with probability of heads θ 10 times.
Observe T H T H T H H T T T.
- Intuitively, knowing we have 4 heads in 10 tosses is all we need.
- But are we missing something? Does the length of the longest run give extra information?

Definition

- Intuition: want to find a statistic so that any **other** statistic provides no additional information about the value of the parameter
- Definition: the statistic $T = g(X_1, \dots, X_n)$ is **sufficient** for an underlying parameter θ if the conditional probability distribution of the data (X_1, \dots, X_n) , given the statistic $g(X_1, \dots, X_n)$, does not depend on the parameter θ .
- Sometimes need more than one statistic, e.g. T_1 and T_2 , in which case we say they are **jointly sufficient** for θ

data | statistic
doesn't depend on θ

⇒ if does depend on θ , know the value of statistic is not sufficient

⇒ extra thing we can know from data

Example (binomial)

- The pmf is, $f(x | p) = p^x(1 - p)^{1-x}$, $x = 0, 1$
- The likelihood is,

$$\prod_{i=1}^n f(x_i | p) = p^{\sum x_i} (1 - p)^{n - \sum x_i}$$

of heads

- Let $Y = \sum X_i$, we have that $Y \sim \text{Bi}(n, p)$ and then,

$$\begin{aligned} & \Pr(X_1 = x_1, \dots, X_n = x_n \mid Y = y) \\ &= \frac{\Pr(X_1 = x_1, \dots, X_n = x_n)}{\Pr(Y = y)} \quad \begin{array}{l} \rightarrow \text{particular sequence} \\ \rightarrow \text{order doesn't matter} \end{array} \\ &= \frac{p^{x_1}(1 - p)^{1-x_1} \dots p^{x_n}(1 - p)^{1-x_n}}{\binom{n}{y} p^y (1 - p)^{n-y}} = \frac{1}{\binom{n}{y}} \quad \begin{array}{l} \text{doesn't depend on } p \\ \Rightarrow \text{the special ordering is uniform} \end{array} \end{aligned}$$

- Given $Y = y$, the conditional distribution of X_1, \dots, X_n does not depend on p .
- Therefore, Y is sufficient for p .

Factorisation theorem

- Let X_1, \dots, X_n have joint pdf or pmf $f(x_1, \dots, x_n | \theta)$
- $Y = g(x_1, \dots, x_n)$ is sufficient for θ if and only if

*factorise
the likelihood*

$$f(x_1, \dots, x_n | \theta) = \phi\{g(x_1, \dots, x_n) | \theta\} h(x_1, \dots, x_n)$$

*sufficient statistics & parametric
data*

- ϕ depends on x_1, \dots, x_n only through $g(x_1, \dots, x_n)$ and h doesn't depend on θ .

Example (binomial)

- The pdf is, $f(x | p) = p^x(1 - p)^{1-x}$, $x = 0, 1$
- The likelihood is,

$$\prod_{i=1}^n f(x_i | p) = p^{\sum x_i} (1 - p)^{n - \sum x_i}$$

- So $y = \sum x_i$ is sufficient for p , since we can factorise the likelihood into:

$$\phi(y, p) = p^y(1 - p)^{n-y} \quad \text{and} \quad h(x_1, \dots, x_n) = 1$$

- So in the coin tossing example, the total number of heads is sufficient for θ .

Example (Poisson)

- X_1, \dots, X_n random sample from a Poisson distribution with mean λ .
- The likelihood is,

$$\prod_{i=1}^n f(x_i | \lambda) = \frac{\lambda^{\sum x_i} e^{-n\lambda}}{x_1! \dots x_n!} = (\lambda^{n\bar{x}} e^{-n\lambda}) \left(\frac{1}{x_1! \dots x_n!} \right)$$

- We see that \bar{X} is sufficient for λ .

Exponential family of distributions

- We often use distributions which have pdfs of the form:

$$f(x \mid \theta) = \exp\{K(x)p(\theta) + S(x) + q(\theta)\}$$

- This is called the **exponential family**.
- Let X_1, \dots, X_n be iid from an exponential family. Then $\sum_{i=1}^n K(X_i)$ is sufficient for θ .
- To prove this note that the joint pdf is

$$\begin{aligned} & \exp \left\{ p(\theta) \sum K(x_i) + \sum S(x_i) + nq(\theta) \right\} \\ &= \left[\exp \left\{ p(\theta) \sum K(x_i) + nq(\theta) \right\} \right] \exp \left\{ \sum S(x_i) \right\} \end{aligned}$$

- The factorisation theorem then shows sufficiency.

Example (exponential)

- The pdf is,

$$f(x | \theta) = \frac{1}{\theta} e^{-x/\theta} = \exp \left[x \left(-\frac{1}{\theta} \right) - \ln \theta \right], \quad 0 < x < \infty$$

- This is of the form

$$f(x | \theta) = \exp\{K(x)p(\theta) + S(x) + q(\theta)\}$$

- So $K(x) = x$ and $\sum X_i$ is sufficient for θ (and so is $\bar{X} = \sum X_i/n$).

$$\sum k(x_i) = \sum x_i$$

Sufficiency and MLEs

- If there exist sufficient statistics, the MLE will be a function of them.
- Factorise the likelihood:

$$L(\theta) = f(x_1, \dots, x_n | \theta) = \phi\{g(x_1, \dots, x_n) | \theta\} h(x_1, \dots, x_n)$$

- We find the MLE by maximizing $\phi\{g(x_1, \dots, x_n) | \theta\}$ which is a function of the sufficient statistics and θ
- So the MLE must be a function of the sufficient statistics

bayesian
not of $n(\theta)$...
we usually throw $n(x_1, \dots, x_n)$
the left part is about
 $\phi\{g(x_1, \dots, x_n)\}$

Importance of sufficiency

- Why are sufficient statistics important?
- Once the sufficient statistics are known there is no additional information on the parameter in the sample
- Samples that have the same values of the sufficient statistic yield the same estimates
- The optimal estimators/tests are based on sufficient statistics (such as the MLE)
- A lot of statistical theory is based on them
- Easy to find the sufficient statistics in some special cases (e.g. exponential family)

Disclaimer

- But... the concept of sufficiency relies on knowing the population distribution
- So, it is mostly important for theoretical work. → assume population distribution
- In practice, we want to also look at all aspects of our data
- That is, we should go beyond any putative sufficient statistics, as a sanity check of our assumptions (e.g. QQ plots).

Outline

Likelihood theory

Asymptotic distribution of the MLE

Cramér–Rao lower bound

Sufficient statistics

Factorisation theorem

Optimal tests

Previous claims (from module 8)

- The likelihood ratio test (LRT) gives the optimal test
- The likelihood ratio has a known distribution

Neyman–Pearson lemma

- Comparing **simple** hypotheses:

$$H_0: \theta = \theta_0 \quad \text{versus} \quad H_1: \theta = \theta_1$$

- The **Neyman–Pearson lemma** states that the **most powerful test**,
for a given significance level, is the LRT
- (Proof of lemma not shown)

↓
likelihood ratio test

Uniformly most powerful tests

- Now consider a **composite alternative** hypothesis,

$$H_1: \theta \in A_1$$

- If the same test (from the LRT) is most powerful for all $\theta_1 \in A_1$, then we say it is **uniformly most powerful** for $\theta_1 \in A_1$.
- If the form of the LRT **differs** for different values of θ_1 , then any given one will only be the best for particular values of θ_1 .
- If so, then we do not have a uniformly best test.
- But any given test might still be a reasonably good test for other values of θ_1

eg. two sided test $(\theta \neq \theta_0)$ → may be split into 2 parts
one side test $(\theta > \theta_0)$ have uniformly most powerful test for either $\theta > \theta_0$
uniformly most powerful

Asymptotic distribution of the likelihood ratio*

- Consider the test,

$$H_0: \theta = \theta_0 \quad \text{versus} \quad H_1: \theta \neq \theta_0$$

diff

$$\begin{aligned} H_0: \theta = \theta_0 & \# \text{ free par} = 1 \\ H_1: \theta \neq \theta_0 & \# \text{ freq par} = n > 1 \end{aligned}$$

- The likelihood ratio is,

$$\lambda = \frac{L_0}{L_1} = \frac{L(\theta_0)}{L(\hat{\theta})}$$

reject $H_0: \lambda \leq k$

$$-2 \ln \lambda \geq -2 \ln k$$

↓

(1- α) quantile
of χ^2

- The function $2 \ln(\lambda)$ asymptotically follows a χ^2_1 distribution
- This can be used to set up approximate hypothesis tests
- Is often used to formally compare different models

eg different regression models.

2015 exam (MAST20005), question 7

best test

LRT ?

Suppose Y has the binomial distribution $\text{Bi}(n, \theta)$. Show that a best region for testing the null hypothesis $H_0: \theta = \theta_0$ against $H_1: \theta > \theta_0$ is $\{y: y/n > c\}$. Then find c so that the test has significance $\alpha \approx 0.05$ when n is large.

$$Y \sim \text{Bi}(n, \theta)$$

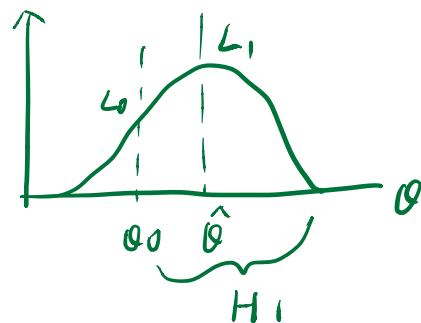
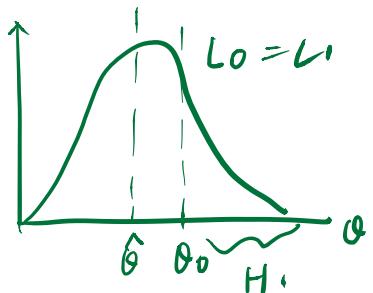
$$\text{under } H_0 \quad Y \sim \text{Bi}(n, \theta_0)$$

$$\begin{cases} H_0: \theta = \theta_0 \\ H_1: \theta > \theta_0 \end{cases}$$

$$L(\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

$$L_0 = L(\theta_0) = \binom{n}{y} \theta_0^y (1-\theta_0)^{n-y}$$

$$L_1: \hat{\theta} = \frac{y}{n}$$



$$L_1 = L(\theta_0) = L_0$$

$$L_1 = L(\hat{\theta}) > L_0$$

likelihood ratio = 1

likelihood ratio < 1

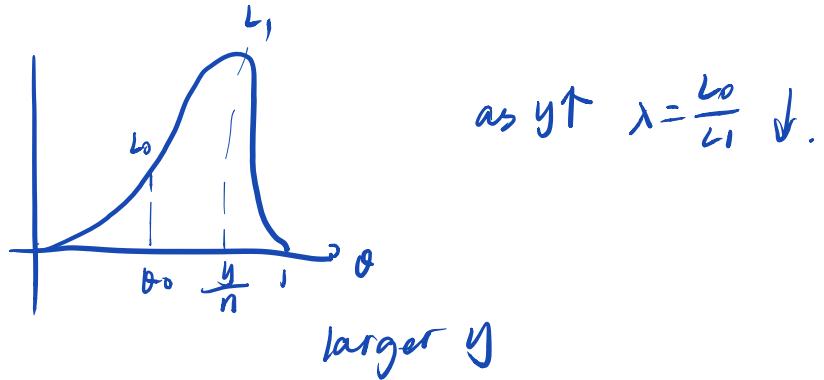
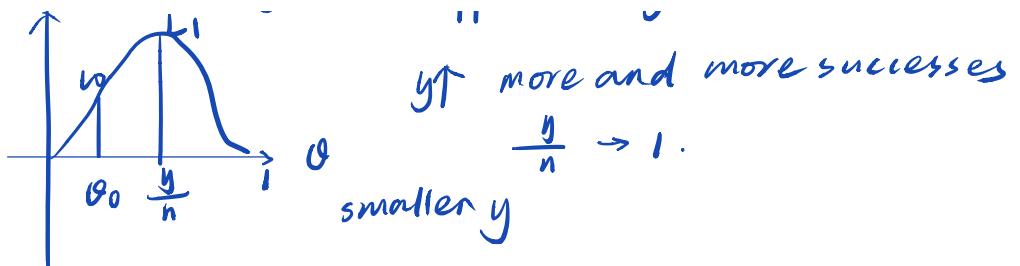
will never reject H_0 ($\lambda=1$)

can ignore.

n, θ_0 , is fixed. function of y .

$$\text{when } \hat{\theta} > \theta_0 \quad \lambda = \frac{L_0}{L_1} = \frac{\binom{n}{y} \theta_0^y (1-\theta_0)^{n-y}}{\binom{n}{\hat{y}} \hat{\theta}^y (1-\hat{\theta})^{n-y}}$$

draw λ (based on y), what happens as y varies.



$$(\lambda \leq k \Leftrightarrow \frac{y}{n} \geq c.)$$

Show λ decrease as y increase.

$$g(y) = \log(\lambda(y)) = y \log \alpha_0 + (n-y) \log(1-\alpha_0) - y \log\left(\frac{y}{n}\right) - (n-y) \log\left(\frac{n-y}{n}\right).$$

$$= y \left[\log \alpha_0 - \log y + \log n \right] + (n-y) \left[\log(1-\alpha_0) - \log(n-y) + \log n \right]$$

$$\begin{aligned} g'(y) &= \log \alpha_0 - (\log y + 1) + \log n \\ &\quad + \left[-\log(1-\alpha_0) + \cancel{c} \cancel{\log(n-y)} + (n-y) \cdot \frac{1}{ny} - \log n \right] \\ &= \cancel{\log \alpha_0} - \cancel{\log y} + \cancel{\log n} - \log(1-\alpha_0) + \log(n-y) + \cancel{-\log n} \\ &= \log \alpha_0 - \log y + \log(n-y) - \log(1-\alpha_0) \end{aligned}$$

$$= \log \left(\frac{\theta_0}{1-\theta_0} \cdot \frac{n-y}{y} \right).$$

$$= \log \left(\frac{\theta_0}{1-\theta_0} \cdot \frac{1-\hat{\theta}}{\hat{\theta}} \right).$$

$$= \log \left(\frac{\theta_0 / 1-\theta_0}{\hat{\theta} / 1-\hat{\theta}} \right).$$

log odds ratio

$$\text{since } \hat{\theta} > \theta_0, \quad \frac{\hat{\theta}}{1-\hat{\theta}} > \frac{\theta_0}{1-\theta_0}$$



odds will also bigger

$$\log \left(\frac{\theta_0 / 1-\theta_0}{\hat{\theta} / 1-\hat{\theta}} \right) < 1$$

$$g'(y) < 0$$

$\Rightarrow g(y)$ is decreasing

$\Rightarrow \lambda$ is decreasing

$$\Rightarrow \lambda \leq k \Leftrightarrow \frac{y}{n} \geq 0 \quad \text{binomial}$$

$$n \text{ large} \Rightarrow \frac{k}{n} \approx N(\theta, \frac{\theta(1-\theta)}{n}).$$

$$z = \frac{\frac{y}{n} - \theta_0}{\sqrt{\frac{\theta_0(1-\theta_0)}{n}}} \approx N(0,1) \text{ under } H_0.$$

test with $\alpha = 0.05 \quad z \geq \Phi^{-1}(0.95) \text{. reject } H_0.$

$$\frac{y}{n} - \theta_0 \geq \Phi^{-1}(0.95) \cdot \sqrt{\frac{\theta_0(1-\theta_0)}{n}} \Rightarrow \frac{y}{n} \geq \theta_0 + \Phi^{-1}(0.95) \cdot \sqrt{\frac{\theta_0(1-\theta_0)}{n}}.$$

2014 exam (MAST20005), question 1

Let X_1, \dots, X_n be a random sample from the probability density function:

$$f(x | \theta) = \frac{x}{\theta^2} e^{-x/\theta}, \quad 0 < x < \infty, \quad 0 < \theta < \infty.$$

- (a) Determine a sufficient statistic for θ .
- (b) Write the log-likelihood function and the score function.
- (c) Determine the maximum likelihood estimator of θ .
- (d) Give the Cramér-Rao lower bound of unbiased estimators of θ .
Hint: If X follows a Gamma(α, β) distribution with pdf $(x^{\alpha-1} e^{-x/\beta}) / (\beta^\alpha \Gamma(\alpha))$, with $x, \alpha, \beta > 0$, then $E(X) = \alpha\beta$.
- (e) A random sample of size $n = 35$ gave $\bar{x} = 10.5$. Determine the maximum likelihood estimate of θ and an approximate 95% confidence interval for θ .

The following R output may help.

```
> z <- c(0.95, 0.975, 0.99, 0.995)
> qnorm(z)
[1] 1.644854 1.959964 2.326348 2.575829
```

factorisation theorem

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n f(x_i | \theta) = \prod_{i=1}^n \frac{x_i}{\theta^2} e^{-x_i/\theta} \\ &= \frac{\prod_{i=1}^n x_i}{\theta^{2n}} e^{-\frac{1}{\theta} \sum_{i=1}^n x_i} \\ &= \prod_{i=1}^n x_i \left(\frac{1}{\theta^n} e^{-\frac{1}{\theta} \sum_{i=1}^n x_i} \right) \end{aligned}$$

$$(a) f(x | \theta) = \exp \{ k(x) p(\theta) + S(x) + q(\theta) \}$$

$$\begin{aligned} \text{since } f(x | \theta) &= \frac{x}{\theta^2} e^{-x/\theta} = \exp \left\{ -\frac{x}{\theta} + \ln \left(\frac{x}{\theta^2} \right) \right\} \\ &= \exp \left\{ -\frac{1}{\theta} x + \ln x - 2 \ln \theta \right\} \end{aligned}$$

$$\text{so } K(x) = x, \quad p(\theta) = -\frac{1}{\theta}, \quad S(x) = \ln x, \quad q(\theta) = -2 \ln \theta$$

$\therefore \sum_{i=1}^n k(x_i) = \sum_{i=1}^n x_i$ is sufficient statistic for θ .

$$(b) f(x | \theta) = \frac{x}{\theta^2} e^{-x/\theta}$$

$$L(\theta) = \prod_{i=1}^n f(x_i | \theta) = \prod_{i=1}^n \frac{x_i}{\theta^2} e^{-x_i/\theta} = \left(\frac{1}{\theta^2} \right)^n \prod_{i=1}^n x_i e^{-\frac{\sum_{i=1}^n x_i}{\theta}}$$

$$\ln L(\theta) = -2n \ln \theta + \ln \left(\prod_{i=1}^n x_i \right) - \frac{1}{\theta} \sum_{i=1}^n x_i$$

$$\frac{\partial \ln L(\theta)}{\partial \theta} = -\frac{2n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n x_i = 0$$

$$(c) \quad 0 = -\frac{2n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n x_i \quad \Rightarrow \quad \hat{\theta} = \frac{\sum_{i=1}^n x_i}{2n} = \frac{\bar{x}}{2}$$

$\hat{\theta} = \frac{\bar{x}}{2}$ estimator

$$(d) \quad V = -\frac{\partial U}{\partial \theta} = -\frac{2n}{\theta^2} + 2 \frac{1}{\theta^3} \sum_{i=1}^n x_i$$

$$\begin{aligned} I(\theta) &= E(V(\theta)) = E\left(2 \frac{\sum_{i=1}^n x_i - n\theta}{\theta^3}\right) \\ &= \frac{2}{\theta^3} E\left(\sum_{i=1}^n x_i - n\theta\right) \end{aligned}$$

Hint: If X follows a $\text{Gamma}(\alpha, \beta)$ distribution with pdf $(x^{\alpha-1} e^{-x/\beta})/(\beta^\alpha \Gamma(\alpha))$, with $x, \alpha, \beta > 0$, then $E(X) = \alpha\beta$.

$$f(x|\theta) = \frac{x}{\theta^2} e^{-\frac{x}{\theta}}$$

$$\alpha = 2 \quad \beta = \theta.$$

$\sim \text{Gamma}(2, \theta)$.

$$E(x) = 2\theta.$$

$$\therefore I(\theta) = \frac{2}{\theta^3} (2n\theta - n\theta) = \frac{2}{\theta^3} \cdot n\theta = \frac{2n}{\theta^2}$$

$$\text{Lower bound } \overline{I(\theta)} = \frac{\theta^2}{2n}$$

- (e) A random sample of size $n = 35$ gave $\bar{x} = 10.5$. Determine the maximum likelihood estimate of θ and an approximate 95% confidence interval for θ .

The following R output may help.

```
> z <- c(0.95, 0.975, 0.99, 0.995)
> qnorm(z)
[1] 1.644854 1.959964 2.326348 2.575829
```

$$\begin{aligned} \text{MLE} \quad \hat{\theta} &= \frac{1}{2} \bar{x} = 5.25 \\ \hat{\theta} &\sim N(\theta, \frac{1}{I(\theta)}) \quad I(\hat{\theta}) = 2.54. \end{aligned}$$

$$\text{se}(\hat{\theta}) = \sqrt{\frac{1}{I(\hat{\theta})}} = 0.6275 \quad \hat{\theta} \pm 1.96 \text{se}(\hat{\theta})$$

$$5.25 \pm 1.96 \times 0.6275 \quad (4.0205, 6.487)$$

1

2

