

# STAT8027 Final Review

## Outline

C.C

---

### Catalogue

- Introduction
- Properties of Estimators
- Methods of Estimation
- Hypothesis Testing
- Interval Estimation
- Bayesian Inference
- Decision Theory
- Non-parametric Methods
- Computationally Intensive Methods

---

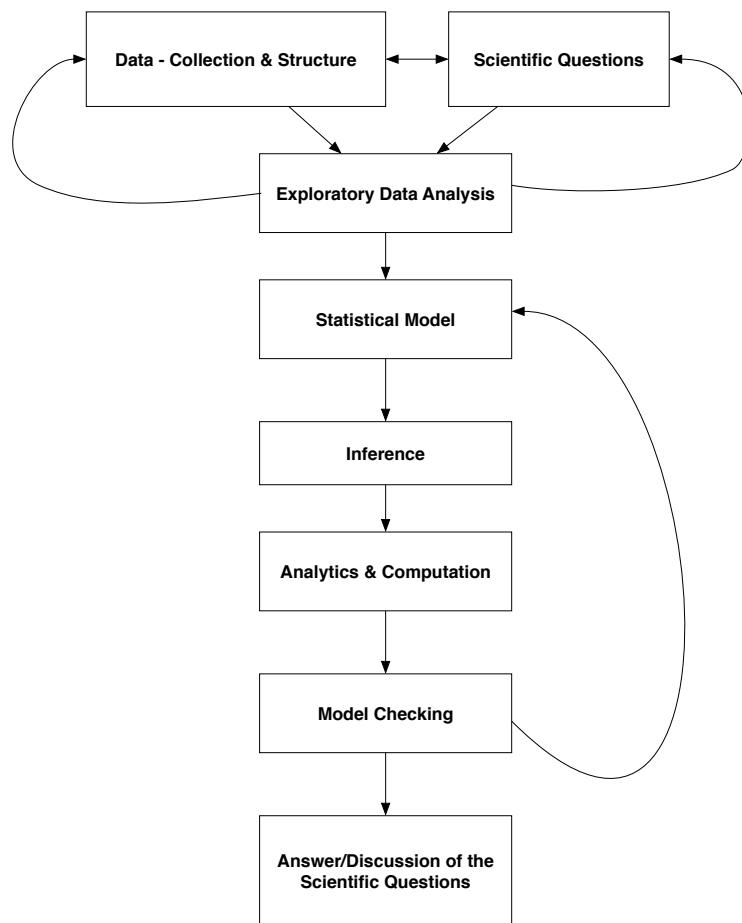
### Reference

1. STAT8027 Lecture slides  
Lecture: Dr. Anton Westveld  
ANU - RSFAS
2. *Statistical Inference* (second edition)  
Paul Garthwaite, Ian Jolliffe, and Byron Jones  
Oxford Science Publication

---

# Introduction

## Workflow



Areas:

1. Point estimation
2. Interval estimation
3. Hypothesis testing

Estimation Methods:

1. Least-squares Estimation
2. Method of moments
3. Maximum Likelihood Estimation (MLE)
4. Bayesian
5. Non-parametric
6. ...

## 1. Generation Random Variables

### 1.1 Monte Carlo integration

- Many quantities of statistical analyses can be expressed as the expectation of a function of a random variable  $E[h(X)]$ .
- Let  $f(X|\theta)$  denote the density of  $X$
- Let  $\mu$  denote the expectation of  $h(X)$ .
- Then when an iid sample  $X_1, \dots, X_n$  is obtained from  $f(X; \theta)$ , we can approximate  $\mu$  by a sample average:

$$\hat{\mu}_{MC} = \frac{1}{n} \sum_{i=1}^n h(X_i) \rightarrow \int h(x)f(x)dx = \mu$$

14.

$$\hat{\sigma}_{MC}^2 = \frac{1}{n-1} \sum_{i=1}^n (h(X_i) - \hat{\mu}_{MC})^2 \rightarrow \sigma^2$$

### 1.2 Generating Random Samples

The probability inverse transform (Tutorial 0~1)

L12a

## 2. Theorems between distributions (Proof matters!)

**Theorem R1:** If  $Z$  is a standard normal random variable, the  $U = Z^2$  is a  $\chi^2$  distribution with 1 degree of freedom.

**Theorem R2:** If  $U_1, \dots, U_n$  are independent and  $X_1 \sim \chi^2_1$  then

$$U_1 + U_2 \sim \chi^2_{n+m} \quad \sum_{i=1}^n U_i \sim \chi^2_n$$

**Theorem Extend**

**3. If**

- $Z \sim \text{normal}(0, 1)$
- $U \sim \chi^2_n$
- $Z$  and  $U$  are independent, then:

$T = Z/\sqrt{U/n}$  is a t distribution with  $n$  degrees of freedom

$$f(t) = \frac{\Gamma[(n+1)/2]}{\sqrt{n\pi}\Gamma(n/2)} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2}$$

**4. If**

- $U \sim \chi^2_m$
- $V \sim \chi^2_n$
- $U$  and  $V$  are independent then:  $W = \frac{U/m}{V/n} \sim F(m, n)$

**5. If**

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{normal}(\mu, \sigma^2)$ , then

1.  $\bar{X} \sim \text{normal}(\mu, \sigma^2/n)$
2.  $\bar{X}$  and  $S^2$  are independent
3.  $(n-1)S^2/\sigma^2 \sim \chi^2_{n-1}$

**6.**

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

---

# Properties of Estimators

## Properties and MVUE

### 1. Unbiased

**Definition 2.1:**  $\hat{\theta} = T(X_1, \dots, X_n)$  is an unbiased estimator for  $\theta$  if  $E[T(\mathbf{X})] = \theta$ . The bias of an estimator is defined as:

$$\text{bias}(\hat{\theta}) = E[T(\mathbf{X})] - \theta$$

$$\begin{aligned} 1.1 \quad \text{MSE}(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\ E[(\hat{\theta} - \theta)^2] &= E[(\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta)^2] \\ &= E[(\hat{\theta} - E(\hat{\theta}))^2 + 2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta) + (E(\hat{\theta}) - \theta)^2] \\ &= E[(\hat{\theta} - E(\hat{\theta}))^2] + 2(E(\hat{\theta}) - \theta)E[(\hat{\theta} - E(\hat{\theta}))] + E[(E(\hat{\theta}) - \theta)^2] \\ &= E[(\hat{\theta} - E(\hat{\theta}))^2] + 0 + E[(E(\hat{\theta}) - \theta)^2] \\ &= E[(\hat{\theta} - E(\hat{\theta}))^2] + (E(\hat{\theta}) - \theta)^2 \\ &= V(\hat{\theta}) + \text{Bias}(\hat{\theta})^2 \end{aligned}$$

### 2. Weakly consistent

**Definition 2.1:** An estimator  $\hat{\theta}$  is **weakly consistent** if

$$P(|\hat{\theta} - \theta| > \epsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

for any  $\epsilon > 0$ .

### 3. Efficiency

An unbiased estimator is said to be **efficient** if it has the minimum possible variance; the **efficiency** of an unbiased estimator is the ratio of the minimum possible variance to the variance of the estimator.

#### 3.1 Cramer-Rao Inequality [lower bound]

**Section 2.4 (Cramer-Rao Inequality [lower bound]):** Let  $X_1, \dots, X_n$  be a random sample from a distribution family with density function  $f_X(x; \theta)$  where  $\theta$  is a scalar parameter. Also, let  $T = t(X_1, \dots, X_n)$  be an unbiased estimator for  $\tau(\theta)$ . Then under certain regularity (smoothness) conditions:

$$\text{Var}(T) \geq \frac{\{\tau'(\theta)\}^2}{ni(\theta)} = \{\tau'(\theta)\}^2 I(\theta)^{-1}$$

- $\tau'(\theta) = \frac{d}{d\theta}\tau(\theta)$
- Where  $I(\theta) = ni(\theta)$  is called the **Expected Fisher Information**.
- $I(\theta) = E\left[\left(\frac{\partial I(\theta)}{\partial \theta}\right)^2\right] = -E\left[\frac{\partial^2 I(\theta)}{\partial \theta^2}\right]$

**C-R Inequality Extended:** Let  $X_1, \dots, X_n$  be a sample [note we don't have to have iid] with pdf  $f(\mathbf{x}|\theta)$  and let  $T(\mathbf{X})$  be an estimator [doesn't have to be unbiased] then based on regularity conditions we have:

$$V[T(\mathbf{X})] \geq \frac{\left[ \frac{\partial}{\partial \theta} E[T(\mathbf{X})] \right]^2}{E \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{x}; \theta) \right)^2 \right]} = \frac{\left[ \frac{\partial}{\partial \theta} E[T(\mathbf{X})] \right]^2}{I(\theta)}$$

- If  $E[T(\mathbf{X})] = \tau(\theta)$ , so  $T(\mathbf{X})$  is an unbiased estimator for  $\tau(\theta)$ ,
- If we have iid samples:

$$V[T(\mathbf{X})] \geq \frac{[\tau'(\theta)]^2}{ni(\theta)}$$

### Regularity conditions

- $\frac{\partial}{\partial \theta} \ln \{f(x|\theta)\}$  exists for all  $x$  and  $\theta$ ;
- interchange of integration and differentiation is permissible;
- The expectation  $i(\theta) = E \left[ \left( \frac{\partial}{\partial \theta} \ln \{f(x; \theta)\} \right)^2 \right]$ , where  $X$  is a generic random variable having distribution with density  $f(x|\theta)$ , is finite for all  $\theta \in \Theta$ .

### Fisher Information

- The Fisher information, or expected Fisher information, or the information number is:

$$I(\theta) = E \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{x}; \theta) \right)^2 \right] = -E \left[ \left( \frac{\partial^2}{\partial \theta^2} \log f(\mathbf{x}; \theta) \right) \right]$$

- For one data point we have:

$$i(\theta) = E \left[ \left( \frac{\partial}{\partial \theta} \log f(x; \theta) \right)^2 \right]$$

- For iid data:

$$ni(\theta) = I(\theta)$$

(Principles of Data Reduction, L03a / p2)

## 4. Sufficiency

**Sufficiency Principle:** If  $T(X_1, \dots, X_n)$  is a sufficient statistic for  $\theta$ , then any inference about  $\theta$  should depend on the sample  $\mathbf{X}$  only through  $T(X_1, \dots, X_n)$ .

**Definition 2.5:** A statistic  $T(X_1, \dots, X_n)$  is **sufficient** for  $\theta$  if the conditional distribution of the sample  $X_1, \dots, X_n$  given  $T(X_1, \dots, X_n)$  does not depend on  $\theta$ .

**Theorem 2.1 (the factorization theorem/criterion):** Suppose  $X_1, \dots, X_n$ , form a random sample from  $f(x; \theta)$ . Then  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  if and only if there exists two non-negative functions  $K_1$  and  $K_2$ , such that the likelihood  $L(\theta, \mathbf{x})$  can be written:

$$L(\theta; \mathbf{x}) = K_1 [t(\mathbf{x}); \theta] \ K_2 [\mathbf{x}]$$

### 4.1 Minimal Sufficient

**Definition 2.6:** A sufficient statistic  $T(\mathbf{X})$  is called a **minimal sufficient statistic** if, for any other sufficient statistic  $T'(\mathbf{X})$ ,  $T(\mathbf{X})$  is a function of  $T'(\mathbf{X})$ .

- Not easy to use the definition to find a minimal sufficient statistic!

**Lemma 2.3:** Let  $f(\mathbf{x}; \theta)$  be the pdf or pmf of a sample  $\mathbf{X}$ . Suppose there exists a function  $T(\mathbf{x})$  such that, for every two sample points  $\mathbf{x}$  and  $\mathbf{y}$  the ratio

$$L(\theta; \mathbf{x})/L(\theta; \mathbf{y})$$

is constant as function of  $\theta$  [note: this can be a vector] if and only if

$$T(\mathbf{x}) = T(\mathbf{y}).$$

Then  $T(\mathbf{X})$  is a minimal sufficient statistic.

**Example:**

- Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} n(\mu, \sigma^2)$ , with both  $\mu, \sigma^2$  unknown.
- Let  $\mathbf{x}$  and  $\mathbf{y}$  be two sample points.
- Let  $(\bar{x}, s_x^2)$  and  $(\bar{y}, s_y^2)$  be the sample means and sample variances for the samples  $\mathbf{x}$  and  $\mathbf{y}$ .

$$\begin{aligned} \frac{f(\mathbf{x}|\mu, \sigma^2)}{f(\mathbf{y}|\mu, \sigma^2)} &= \frac{(2\pi\sigma^2)^{-n/2} \exp(-[n(\bar{x}-\mu)^2 - (n-1)s_x^2]/(2\sigma^2))}{(2\pi\sigma^2)^{-n/2} \exp(-[n(\bar{y}-\mu)^2 - (n-1)s_y^2]/(2\sigma^2))} \\ &= \exp([-n(\bar{x}^2 - \bar{y}^2) + 2n\mu(\bar{x} - \bar{y}) - (n-1)(s_x^2 - s_y^2)]/(2\sigma^2)) \end{aligned}$$

- This ratio will not depend on  $\mu$  and  $\sigma^2$  if and only if  $\bar{x} = \bar{y}$  and  $s_x^2 = s_y^2$ .
- $(\bar{X}, S^2)$  are minimally sufficient for  $\mu, \sigma^2$ .

## 4.2 Rao-Blackwell Theorem (How sufficiency helps finding MVUEs)

- Let  $W$  be any unbiased estimator of  $\tau(\theta)$ .
- Let  $T$  be a sufficient statistic for  $\theta$ .
- Define  $\phi(T) = E[W|T]$ .
- Then

$$E[\phi(T)] = \tau(\theta)$$

$$V[\phi(T)] \leq V[W]$$

- So if we have unbiased estimator and condition it on a sufficient statistic, our new statistic  $\phi(T)$  has the same or smaller variance!!

### **Corollary to Theorem**

If an MVUE for  $\theta$  exists, then there must be a function of the minimum sufficient statistic for  $\theta$  which is an MVUE.

## 5. Complete

**Definition 2.9:** Let  $f_T(t; \theta)$  be a family of pdfs or pmfs for a statistic  $T(x)$ . The family of probability distributions is called **complete** if

$$E[h(T)] = \int h(t)f_T(t)dt = 0$$

for all  $\theta$  implies that  $P(h(T) = 0) = 1$

for all  $\theta$ .

### **Scheffe Theorem**

**Lemma 2.6:** Let  $X_1, \dots, X_n$  be a random sample from a distribution with density function  $f(x; \theta)$ . If  $T = T(\mathbf{X})$  is a complete and sufficient statistic, and  $\phi(T)$  is an unbiased estimator of  $\tau(\theta)$ , then  $\phi(T)$  is the unique MVUE of  $\tau(\theta)$ .

## How to find MVUEs?

- How to find MVUEs? It seems we have an approach:

- Find or construct a sufficient and complete statistic  $T$ .
- Find an unbiased estimator  $W$  for  $\tau(\theta)$ .
- Compute  $\phi(T) = E[W|T]$ , then  $\phi(T)$  is the UMVUE.

- Or:

- Find or construct a sufficient and complete statistic  $T$ .
- Find a function  $h(T)$ , where  $E[h(T)] = \tau(\theta)$  (i.e. it is unbiased).
- Then  $h(T)$  is the MVUE.

## Exponential Families

**Definition 2.7:** We say that a random variable belongs to the  $k$ -parameter exponential family of distributions if its pdf can be written in the following form:

$$f(x; \theta) = \exp \left( \sum_{j=1}^k A_j(\theta) B_j(x) + C(x) + D(\theta) \right)$$

$$f(x; \theta) = C^*(x) D^*(\theta) \exp \left( \sum_{j=1}^k A_j(\theta) B_j(x) \right)$$

- If we define:

$$\phi = (\phi_1, \dots, \phi_k) = A(\theta) = \{A_1(\theta), \dots, A_k(\theta)\}$$

then  $\phi$  is referred to as the canonical parameter for the exponential family and the density function can be written in the form:

$$f(x; \theta) = \exp \left\{ \sum_{j=1}^k \phi_j B_j(x) + C(x) + D(\phi) \right\}$$

- Note:

$$\theta = A^{-1}(\phi)$$

$$D(\phi) = D\{A^{-1}(\phi)\}$$

### Exponential families and MVUE

**Lemma 2.4:** If the usual regularity condition hold, then a vector of  $k$  sufficient statistics  $T$  exists for a vector of parameters  $\theta$  if and only if the distribution of  $X$  belong to the  $k$ -parameter exponential family.

**Proof:**

$$f(x; \theta) = \exp \left\{ \sum_{j=1}^k A_j(\theta) \left( \sum_{i=1}^n B_j(x_i) \right) + nD(\theta) + \sum_{i=1}^n C(x_i) \right\}$$

- Let  $t = (\sum_{i=1}^n B_1(x_i), \dots, \sum_{i=1}^n B_k(x_i))$ .

- $K_1 = \exp \left\{ \sum_{j=1}^k A_j(\theta) t_j + nD(\theta) \right\}$
- $K_2 = \exp \left\{ \sum_{i=1}^n C(x_i) \right\}$

**Lemma 2.5:** Under the same conditions as Lemma 2.4,  $T$  is also minimal sufficient.

**Lemma 2.8:** Under the same conditions as Lemma 2.4,  $T$  is also complete.

- For exponential families, we can easily find complete and sufficient statistics.
- All that is needed to find  $h(T)$  such that  $E[h(T)] = \tau(\theta)$ , then we have the unique MVUE!

---

## Methods of Estimation

### Method of Moments (MOM)

#### 1. MOM

- Typically the population moments are implicitly defined by parameters  $\theta = (\theta_1, \dots, \theta_K)$ .
- Equate the sample and population moments:

$$\begin{aligned}\mu_1(\theta_1, \dots, \theta_K) &= \hat{\mu}_1(x_1, \dots, x_n) \\ &\vdots \quad \vdots \\ \mu_K(\theta_1, \dots, \theta_K) &= \hat{\mu}_K(x_1, \dots, x_n)\end{aligned}$$

- The estimator  $T(\mathbf{X}) = \tilde{\theta}$  is the value for  $\theta$  which solves the system of  $K$  equations.

#### 2. Central Moments

$$\mu'_k = E_\theta(\{X - E_\theta(X)\}^k)$$

And sample moments:

$$\hat{\mu}'_k = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^k$$

for  $k = 2, \dots, K$ . Set  $\mu_1 = \hat{\mu}_1 = \bar{x}$ .

#### 3. Generalized MOM

$$\begin{aligned}E_\theta(g_1(X)) &= \frac{1}{n} \sum_{i=1}^n g_1(x_i) \\ &\vdots \\ E_\theta(g_k(X)) &= \frac{1}{n} \sum_{i=1}^n g_k(x_i)\end{aligned}$$

Note: If we set  $g_i(x) = x^i$  then we recover the standard method of moments.

# Maximum Likelihood Estimation

*"This is the **best known, most widely used, most intuitive, most important**, . . . of estimation procedures."*

## 1. Definition:

**Definition 3.1:** The MLE is the value  $\hat{\theta}$  which maximizes  $L(\theta; \mathbf{x})$ .

- If the likelihood is differentiable in  $(\theta_i)$ , possible candidates for the MLE are the values  $(\theta_1, \dots, \theta_k)$  that solve

$$\frac{\partial}{\partial \theta_i} L(\theta; \mathbf{x}) = 0, \quad i = 1, \dots, k$$

- Possible: local vs. global maximum, extrema may occur on the boundary and thus the first derivative may not be 0, . . .
- All the good points and bad points of optimizing a function are here!

The log likelihood function is

$$\ell(\theta) = \log[L(\theta)]$$

Score Function is the gradient of  $\ell(\theta)$

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

## 2. General approach for checking 2-parameter MOM

1. First-order partial derivatives (score equations) at  $\hat{\theta}_1, \hat{\theta}_2$  are zero:

$$\begin{aligned} \left. \frac{\partial}{\partial \theta_1} H(\theta_1, \theta_2) \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2} &= 0 \\ \left. \frac{\partial}{\partial \theta_2} H(\theta_1, \theta_2) \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2} &= 0 \end{aligned}$$

2. At least one second-order partial derivative is negative:

$$\begin{aligned} \left. \frac{\partial^2}{\partial \theta_1^2} H(\theta_1, \theta_2) \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2} &< 0 \\ \left. \frac{\partial^2}{\partial \theta_2^2} H(\theta_1, \theta_2) \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2} &< 0 \end{aligned}$$

3. The determinant of the matrix of second order partial derivatives (the Hessian matrix) is positive.

$$\begin{vmatrix} \frac{\partial^2}{\partial \theta_1^2} & \frac{\partial^2}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2}{\partial \theta_2^2} \end{vmatrix} \Big|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2} > 0$$

### 3. Newton-Raphson (N-R) Method

- N-R is an extremely fast root finding approach, but is sensitive to starting values.
- Again let's consider a log likelihood function  $\ell(\theta|x)$ .
- Let  $U(\theta) = \ell'(\theta|x)$  denote first derivative of  $\ell(\theta)$ .
- Let  $H(\theta) = \ell''(\theta|x)$  denote the derivative of  $\ell(\theta)$ .
- Let  $\theta_0$  be an initial estimate of  $\theta$ .
- Let  $\hat{\theta}$  be the MLE.
- Let's do a Taylor series expansion of  $U(\theta)$  around  $\theta_0$

$$U(\theta) = U(\theta_0) + (\theta - \theta_0)H(\theta_0) + \dots$$

- At  $\theta = \hat{\theta}$  we know  $U(\hat{\theta}) = 0$ , so we have

$$0 = U(\theta_0) + (\hat{\theta} - \theta_0)H(\theta_0) + \dots$$

- Let's say the one-step approximation is reasonable enough.

$$\hat{\theta} = \theta_0 - H^{-1}(\theta_0)U(\theta_0)$$

- This suggests that  $\hat{\theta}$  is well approximated by  $\theta_1$ :

$$\theta_1 = \theta_0 - H^{-1}(\theta_0)U(\theta_0)$$

- We can then get an improved estimate:

$$\theta_2 = \theta_1 - U(\theta_1)H^{-1}(\theta_1)$$

- We can continue with  $\theta_3, \theta_4, \dots$  until convergence is achieved.

$$|\theta_k - \theta_{k-1}| < \epsilon = 1e-07$$

#### 3.1 Multivariate N-R Method

- We can naturally extend the N-R approach to a multivariate setting.
- Let  $U(\theta)$  denote the vector of first partial derivatives of  $\ell(\theta)$ .
- Let  $H(\theta)$  denote the matrix of second partial derivatives of  $\ell(\theta)$ .

#### 3.2 Fisher Scoring

- The method of “scoring” is a simple modification of the Newton-Raphson method.
- The Hessian  $H(\theta)$  is replaced by its expectation.

$$E[H(\theta)] = -I(\theta)$$

where  $I(\theta)$  is Fisher's information matrix.

$$I(\theta) = E \left[ \left( \frac{\partial \ell(\theta|x)}{\partial \theta_i} \right) \left( \frac{\partial \ell(\theta|x)}{\partial \theta_j} \right) \right] = -E \left[ \frac{\partial^2 \ell(\theta|x)}{\partial \theta_i \partial \theta_j} \right]$$

$$\theta_{t+1} = \theta_t + I^{-1}(\theta_t)U(\theta_t)$$

- A great advantage is that  $E[H(\theta)]$  is guaranteed to be positive definite, thus eliminating some possible convergence issue with the Newton-Raphson approach.

## 4. Expectation - Maximization (EM) Algorithm

- Presentation adapted from CB & *Computational Statistics*.
- The EM algorithm is a general algorithm to find MLEs when some of the data are missing (or the problem can be set in a manner that there are missing data).
- Suppose we observe all of the data  $\mathbf{y} = \{y_1, \dots, y_n\}$ , then all we do to find the MLE is maximize:  $\ell(\theta; \mathbf{y})$
- Suppose we don't observe all the  $\mathbf{y}$ s then based on the notation by Donald Rubin we have  $\mathbf{y} = (\mathbf{y}_{obs}, \mathbf{y}_{miss})$ .
 
$$\begin{aligned} f(\mathbf{y}; \theta) &= f(\mathbf{y}_{obs}, \mathbf{y}_{miss}; \theta) \\ &= k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta) g(\mathbf{y}_{obs}; \theta) \end{aligned}$$
- This leads to:  $g(\mathbf{y}_{obs}; \theta) = \frac{f(\mathbf{y}; \theta)}{k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)}$ 

$$\begin{aligned} \log [g(\mathbf{y}_{obs}; \theta)] &= \log [f(\mathbf{y}_{obs}, \mathbf{y}_{miss}; \theta)] - \log [k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)] \\ \ell_{obs}(\theta; \mathbf{y}_{obs}) &= \ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss}) - \log [k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)] \end{aligned}$$
- As  $\mathbf{y}_{miss}$  is missing, we replace the right side of the equation with its expectation:

$$\begin{aligned} \ell_{obs}(\theta; \mathbf{y}_{obs}) &= E \left\{ \ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss}) \middle| \theta', \mathbf{y}_{obs} \right\} \\ &\quad - E \left\{ \log [k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)] \middle| \theta', \mathbf{y}_{obs} \right\} \end{aligned}$$

### 4.1 E-M Steps

- The EM algorithm seeks to maximize  $\ell(\theta; \mathbf{y}_{obs})$  with respect to  $\theta$  through the following process:
 

*don't want missing data so integrate it out*
- 1. **E step:** Calculate the expectation of the complete likelihood conditional on the observed data and the current value of  $\theta$ :

$$\begin{aligned} Q(\theta | \theta^{(r)}) &= E \left\{ \ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss}) \middle| \theta^{(r)}, \mathbf{y}_{obs} \right\} \quad \text{conditional} \\ &= \int [\ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss})] k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta) d\mathbf{y}_{miss} \end{aligned}$$

- 2. **M step:** Maximize  $Q(\theta | \theta^{(r)})$  with respect to  $\theta$ . Set  $\theta^{(r+1)}$  equal to the maximizer of  $Q$ .
- 3. Return to the E step unless a stopping criterion has been reached.

## 5. Invariance Property of MLEs

**Lemma 3.2:** Suppose that  $\theta$  and  $\eta$  represent two alternative parameterizations for some probability distribution and that  $\eta$  is a (1-1) function of  $\theta$ , so that we can write  $\eta = \mathbf{g}(\theta)$ ,  $\theta = \mathbf{h}(\eta)$  for appropriate functions  $\mathbf{g}(\cdot)$ ,  $\mathbf{h}(\cdot)$ .

- If  $\hat{\theta}$  is the MLE of  $\theta$  then  $\hat{\eta} = \mathbf{g}(\hat{\theta})$  is the MLE for  $\eta$
- If the mapping is (1-1) we simply note:  

$$\eta = \tau(\theta) \rightarrow \tau^{-1}(\eta) = \theta$$
- Define our likelihood based on the reparameterization ( $\theta = \tau^{-1}(\eta)$ ):

$$L^*(\eta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \tau^{-1}(\eta)) = L(\tau^{-1}(\eta); \mathbf{x}) = L(\theta; \mathbf{x})$$

- We find the supremum of likelihood

$$\sup_{\eta} L^*(\eta; \mathbf{x}) = \sup_{\eta} L(\tau^{-1}(\eta); \mathbf{x}) = \sup_{\theta} L(\theta; \mathbf{x})$$

to see that the maximum of  $L^*(\eta; \mathbf{x})$  is when  $\eta = \tau(\theta) = \tau(\hat{\theta})$ .

### 5.1 not one-to-one case

- However, many functions of interest are not one-to-one:  $\theta \rightarrow \theta^2$ .
- We proceed by defining the **induced likelihood function** of  $L^*$  for  $\tau(\theta)$

$$L^*(\eta; \mathbf{x}) = \sup_{\theta: \tau(\theta)=\eta} L(\theta; \mathbf{x})$$

- The value  $\hat{\eta}$  that maximizes  $L^*(\eta; \mathbf{x})$  will be called the MLE of  $\eta$ .

## 6. Asymptotic of MLE

### Score Function

**Lemma:** Let  $X_1, \dots, X_n \stackrel{\text{indep}}{\sim} f(x; \theta)$  and let  $\hat{\theta}$  be the MLE of  $\theta$ . Under regularity conditions of  $f(x; \theta)$  and thus  $L(\theta; \mathbf{x})$  (under appropriate smoothness conditions), we can state:

$$W = \frac{1}{\sqrt{n}} \ell'(\theta; \mathbf{x}) \xrightarrow{D} \text{normal}(0, i(\theta))$$

### MLE

**Lemma 3.3:** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x; \theta)$ . Let  $\hat{\theta}$  be the MLE of  $\theta$ . Under regularity conditions of  $f(x; \theta)$  and thus  $L(\theta; \mathbf{x})$  (under appropriate smoothness conditions), we have:

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{D} \text{normal}(0, i(\theta)^{-1})$$

## 6.1 Delta Method

**Theorem:** Let  $Y_n$  be a sequence of random variables such that:

$$\sqrt{n}(Y_n - \theta) \xrightarrow{D} \text{normal}(0, \sigma^2)$$

For a given function  $g$  and a specific value  $\theta$ , suppose that  $g'(\theta)$  exists and is not 0, then

$$\sqrt{n}(g(Y_n) - g(\theta)) \xrightarrow{D} \text{normal}(0, \sigma^2[g'(\theta)]^2)$$

### Extension

**Lemma:** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x; \theta)$ . Let  $\hat{\theta}$  be the MLE of  $\theta$  and let  $\tau(\theta)$  be a continuous function of  $\theta$ . Under regularity conditions (i.e. under appropriate smoothness conditions) of  $f(x; \theta)$  and thus  $L(\theta; \mathbf{x})$ , we have:

$$\sqrt{n}(\tau(\hat{\theta}) - \tau(\theta)) \xrightarrow{D} \text{normal}(0, \nu(\theta))$$

Where  $\nu(\theta) = \frac{[\tau'(\theta)]^2}{I(\theta)}$  is the Cramer-Rao lower bound for a single data point.

$$\tau(\hat{\theta}) \stackrel{\text{d}}{\sim} \text{normal}\left(\tau(\theta), \frac{[\tau'(\theta)]^2}{I(\theta)}\right)$$

## 6.2 Asymptotic of MLE

So asymptotically, MLEs are:

1. unbiased;
2. achieve the Cramer-Rao lower bound (efficient);
3. asymptotically normally distributed.

- We can also note that MLEs are consistent estimators.
- Because these estimators achieve (1-3) they are \*asymptotically efficient! **best asymptotically normal (BAN) estimators**

---

# Hypothesis Testing

## 1. Definition

**Definitions:** Suppose that  $X_1, \dots, X_n$  represent a simple random sample from a parametric family with density function  $f(x; \theta)$  for some parameter  $\theta \in \Theta$ .

- A statistical hypothesis is simply a subset of the parameter space,  $\Theta$ .
- Any statistical hypothesis of interest, often termed the **null hypothesis**, is associated with a competing **alternative hypothesis**.
- A null hypothesis and its alternative form a partition of the parameter space  $\Theta$  consisting of the sets  $\Theta_0$  and

$$\Theta_1 = \Theta_0^c \cap \Theta$$

**Definition:** A **hypothesis testing procedure** or **hypothesis test** is a rule that specifies:

1. For which sample values the decision is made to accept  $H_0$ .
2. For which sample values  $H_0$  is rejected and  $H_1$  is accepted as true.

## 2. Rejection Region ( $C$ )

The subset of the space for which  $H_0$  will be rejected is called the **rejection region** or **critical region** ( $C$ ). The complement of the rejection region is called the **acceptance region**.

we can define a statistical test in terms of a **rejection region** ( $C$ ) which is just a set for some statistic  $T(X_1, \dots, X_n)$ :

$$C = \{X \in \mathcal{X} : T(X) < k\}$$

## 3. Type I and Type II errors

- Type I Error: Reject  $H_0$  given that it is true. Thus the observations fall in the **rejection region**  $C$  when in fact that null hypothesis,  $H_0$ , is true.
- Type II Error: Do not Reject  $H_0$  when it is false. Thus the observed data values fall outside the **rejection region** when in fact the null hypothesis is false.

		Decision
Truth	Accept $H_0$	Accept $H_1$
$H_0$	Correct Decision	Type I Error
$H_1$	Type II Error	Correct Decision

There is a strong relationship between Type I and Type II errors. Note that for a given value of  $\theta$ , only one type of error can occur (since for any given  $\theta$ ,  $H_0$  either is or is not true).

### 3.1 Significance Level & Power

**Definition 4.2:** The probability of a Type I error,  $\alpha$  in a test of hypotheses is called the **size** or **significance level** of the test. The complement of the probability of a Type II error

$$\eta(\theta) = 1 - \beta, \quad \text{is the } \mathbf{power} \text{ of the test.}$$

- **Power =  $1 - P(\text{Type II Error})$**

$$= 1 - P(\mathbf{X} \in C^c | H_1 \text{ is true}) = P(\mathbf{X} \in C | H_1 \text{ is true})$$

- Given that  $H_1$  is true, what is the probability I reject  $H_0$

- It is standard to focus on tests which have sizes 0.05 or 0.01.
- If we focus on tests with rejection regions of the form  $C = \{X < k_\alpha\}$ , we can choose  $k_\alpha$  such that:

$$\max_{\theta \in \Theta_0} \eta(\theta) = \max_{\theta \geq 1000} 1 - \exp(-k_\alpha/\theta) = 1 - \exp(-k_\alpha/1000) = \alpha$$

*Maximize the Power w.r.t  $\theta \in \Theta_0$*

### 3.2 UMP Test

**Definition 4.3** A test which minimizes  $\beta$  for fixed  $\alpha$  is called a **most powerful or best test of size  $\alpha$** .

**Definition 4.5** Suppose we have a test of size  $\alpha$ ; then it is **uniformly most powerful (UMP)** of size  $\alpha$  if its power function  $\eta(\theta)$  is such that  $\eta(\theta) \geq \eta^*(\theta)$  for all  $\theta \in \Omega - \omega$ , where  $\eta^*(\theta)$  is the power function of any other size- $\alpha$  test.

## 4. Essential Nature of a Hypothesis Test (Experimental Design, Hoff 2009)

Given  $H_0, H_1$  and data  $\mathbf{x} = \{x_1, \dots, x_n\}$ :

1. From the data, compute a relevant test statistic  $T(\mathbf{x})$ : The test statistic  $T(\mathbf{x})$  should be chosen so that it can differentiate between  $H_0$  and  $H_1$  in ways that are scientifically relevant. Typically,  $T(\mathbf{x})$  is chosen so that

$$T(\mathbf{x}) \text{ is probably } \begin{cases} \text{small under } H_0 \\ \text{large under } H_1 \end{cases}$$

2. Obtain a null distribution: A probability distribution over the possible outcomes of  $T(\mathbf{X})$  under  $H_0$ . Here,  $X = \{X_1, \dots, X_n\}$  are potential experimental results that could have happened under  $H_0$ .
3. Compute the p-value: The probability under  $H_0$  of observing a test statistic  $T(\mathbf{X})$  as or more extreme than the observed statistic  $t(\mathbf{x})$ .

$$\text{p-value} = P(T(\mathbf{X}) \geq t(\mathbf{x}) | H_0)$$

If the p-value is small  $\Rightarrow$  evidence against  $H_0$

If the p-value is large  $\Rightarrow$  not evidence against  $H_0$

## 5. Neyman-Pearson Set-up (simple case)

- Suppose that  $X_1, \dots, X_n$  are a sample from a population characterized by a probability model with density function  $f(x; \theta)$  for  $\theta \in \Theta$  where  $\Theta = \{\theta_0, \theta_1\}$ .

$$H_0 : \theta = \theta_0$$

$$H_1 : \theta = \theta_1$$

- Consider the likelihood-ratio:

$$\lambda(\mathbf{x}) = \frac{L(\theta_0; \mathbf{x})}{L(\theta_1; \mathbf{x})}$$

- The test we shall define has a critical region of the form

$$C = \{\lambda(\mathbf{x}) \leq k\}$$

- The ratio of the likelihood for any given sample at each of the two possible parameter values is precisely a relative measure of how plausible the two hypotheses are.
- In other words, when  $\lambda(\mathbf{x})$  is very small, this is strong evidence that the observations arose from the alternative hypothesis rather than the null hypothesis.
- It should seem intuitively reasonable that the likelihood ratio is a good method of distinguishing between samples which support the null hypothesis versus samples which support the alternative hypothesis.
- From what we have done, we know for a given  $\alpha$  we could compare the power  $\beta(\theta)$ .
- We would like to find a uniformly most powerful test .  $\eta(\theta) \geq \eta(\theta^*)$
- It turns out that N-P tests lead to UMP tests.
- Suppose that  $H_0$  and  $H_1$  are simple hypotheses and that the test that rejects  $H_0$  Note:  $0 \leq \lambda \leq 1$ , and  $\lambda$  will be close to 1 if  $H_0$  is true. significance level Where  $0 \leq k \leq 1$ .
- Lemma 4.2:** Then any other test for which the significance level is less than or equal to  $\alpha$  has power less than or equal to that of the likelihood ratio test.

Examples compare  $H_1 = H_1 >$  and  $H_1 <$ , they have the same result as the result of  $H_1 =$  (applying NP), which is UMP.

## 5. Maximum Likelihood Ratio Test

$$H_0 : \theta \in \omega \text{ versus } H_1 : \theta \in \Omega - \omega$$

$$\lambda(\mathbf{x}) = \frac{\max_{\theta \in \omega} L(\theta; \mathbf{x})}{\max_{\theta \in \Omega} L(\theta; \mathbf{x})}$$

- $\max_{\theta \in \omega} L(\theta; \mathbf{x})$  is a restricted maximization.
- $\max_{\theta \in \Omega} L(\theta; \mathbf{x})$  is a unrestricted maximization.

We construct a test of the form:

$$C = \{\mathbf{x} : \lambda(\mathbf{x}) \leq k\}$$

Note:  $0 \leq \lambda \leq 1$ , and  $\lambda$  will be close to 1 if  $H_0$  is true.  
Where  $0 \leq k \leq 1$ .

## 5.1 Asymptotic of MLRT

**Theorem** For testing  $H_0 : \theta \in \omega$  versus  $H_1 : \theta \in \Omega - \omega$   
 suppose  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x; \theta)$  and  $\hat{\theta}$  is the MLE of  $\theta$  and  $f(x; \theta)$   
 satisfies the regularity conditions (smoothness).

Then under  $H_0$ , as  $n \rightarrow \infty$ ,

$$-2\log[\lambda(\mathbf{x})] \xrightarrow{D} \chi_1^2$$

**Theorem Extend**  $-2\log(\lambda) \xrightarrow{D} \chi_\nu^2$

where  $\nu = \#$  number of constraints set in  $H_0$ .

- Another way to think about it is: Let  $p$  be the number of parameters estimated (are free) under  $H_1$ . And let  $p_0$  be the the number of parameters estimated (are free) under  $H_0$ .
- Then  $\nu = p - p_0$ .

## 5.2 Properties of MLRTs

The MLRT

1. is asymptotically most powerful unbiased;
2. is asymptotically similar;
3. is asymptotically efficient.

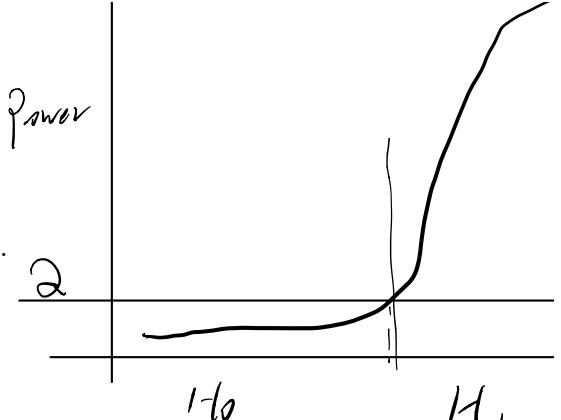
### 6. Unbiased Test

Suppose that we wish to test:

$$H_0 : \theta \in \Theta_0 \quad \text{vs.} \quad H_1 : \theta \in \Theta_1.$$

A test of size  $\alpha$  is said to be **unbiased** if

$$\eta(\theta) \geq \alpha \text{ for all } \theta \in \Theta_1.$$



### 7. Similar Test

Suppose that we wish to test:

$$H_0 : \theta \in \Theta_0 \quad \text{vs.} \quad H_1 : \theta \in \Theta_1.$$

A test of size  $\alpha$  is said to be **similar** if

$$\eta(\theta) = \alpha \text{ for all } \theta \in \Theta_0.$$

### 8. Efficiency

Suppose that we have ~~to~~ possible tests of  $H_0$  vs.  $H_1$ , where both tests are simple.

If  $n_1$  and  $n_2$  are the minimum possible sample sizes for tests 1 and 2  
 for which we can achieve a size  $\alpha$  and power  $\geq \eta$ , then the **relative efficiency** of test 1 compared to test 2 is:

$$n_2/n_1.$$

## 9. Two other tests

### 9.1 **The Score Test**

$$\mathbf{u}(\boldsymbol{\theta}) = \left( \frac{\partial \ell}{\partial \theta_1}, \frac{\partial \ell}{\partial \theta_2}, \dots, \frac{\partial \ell}{\partial \theta_k} \right)^T$$

- Suppose that we wish to test:

$$H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0 \quad \text{vs.} \quad H_1 : \Omega - \{\boldsymbol{\theta}_0\}.$$

- Test statistic:

$$\mathbf{u}(\boldsymbol{\theta})^T \mathbf{I}_{\boldsymbol{\theta}_0}^{-1} \mathbf{u}(\boldsymbol{\theta}) \stackrel{\sim}{=} \chi_{df=k}^2$$

- Note: We don't have to determine the MLEs!

## 9.2 The Wald Test

- Suppose that we wish to test:

$$H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0 \quad \text{vs.} \quad H_1 : \Omega - \{\boldsymbol{\theta}_0\}.$$

- Test statistic:

$$(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \mathbf{I}_{\hat{\boldsymbol{\theta}}} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \stackrel{\sim}{=} \chi_{df=k}^2$$

- The MLRT, the Score Test, and the Wald Test are asymptotically equivalent!
-

# Interval Estimation / Confidence set

## 1. Definition

Suppose that  $\{f(x; \theta); \theta \in \Omega\}$  define a family of distributions

If  $S_X$  is a subset of  $\Omega$ , depending on  $X$ , such that

$$P(X : S_X \ni \theta) = 1 - \alpha$$

then  $S_X$  is a **confidence set** for  $\theta$  with **confidence coefficient**  $1 - \alpha$ .

### Construction Methods

- parametric "exact" intervals
- parametric asymptotic intervals

### Some general approaches

- Inverting a test statistic
- Pivotal Quantities
- Pivoting the CDF

## 2. Hypothesis Testing & Confidence Sets/Intervals

- There is a strong relationship between hypothesis testing and interval estimation. In general, every confidence set corresponds to a test and vice versa.

**Lemma 5.1:** Suppose that  $\bar{C}(\theta_0)$  is the acceptance region for a test of size  $\alpha$ :

$$H_0 : \theta = \theta_0 \quad \text{vs} \quad H_1 : \theta \in \Omega - \bar{C}(\theta_0).$$

Then a **confidence set** for  $\theta$  with **confidence coefficient**  $(1 - \alpha)$ , is given by

$$S_X = \{\theta_0 : X \in \bar{C}(\theta_0)\}$$

## 3. Pivotal quantity

**Definition 5.2:** A random variable  $g(X, \theta)$  is a **pivotal quantity (or a pivot)** if the distribution of  $g(X, \theta)$  is independent of all parameters.

$$\bar{X} \sim N(\theta, \frac{1}{n}) \Rightarrow \frac{\bar{X} - \theta}{\sqrt{\frac{1}{n}}} \sim N(0, 1) \quad \text{with } \bar{X} - \theta \sim N(0, 1) \text{ also pivot.}$$

we can construct a pivot  $g(x, \theta)$  so that the distribution of  $g$  is not dependent on  $\theta$ , thus we can use d.s of  $g$  to get the interval of  $\theta$

## 4. MLEs & Asymptotic

$$\text{MLE} \Rightarrow \hat{\theta} \sim \text{normal}(\theta, I(\theta)^{-1})$$

$$\frac{\hat{\theta} - \theta}{1/\sqrt{I(\theta)}} \sim \text{normal}(0, 1)$$

We have a pivotal quantity. Based on the same approach as before we can construct an asymptotic  $100(1 - \alpha)\%$  confidence interval as:

$$\left[ \hat{\theta} - z_{\alpha/2} \frac{1}{\sqrt{I(\hat{\theta})}}, \hat{\theta} + z_{\alpha/2} \frac{1}{\sqrt{I(\hat{\theta})}} \right]$$

If we are interested in a function of  $\theta$ , say  $\tau(\theta)$ , then we have:

$$\tau(\hat{\theta}) \sim \text{normal} \left( \tau(\theta), \frac{[\tau'(\theta)]^2}{I(\theta)} \right)$$

$$\frac{\tau(\hat{\theta}) - \tau(\theta)}{\sqrt{\frac{[\tau'(\theta)]^2}{I(\theta)}}} \sim \text{normal}(0, 1)$$

We can construct an asymptotic  $100(1 - \alpha)\%$  confidence interval as:

$$\text{CRLB}(\tau(\theta)) = \frac{(\tau'(\theta))^2}{I(\theta)}$$

$$\left[ \tau(\hat{\theta}) - z_{\alpha/2} \frac{\tau'(\hat{\theta})}{\sqrt{I(\hat{\theta})}}, \tau(\hat{\theta}) + z_{\alpha/2} \frac{\tau'(\hat{\theta})}{\sqrt{I(\hat{\theta})}} \right] \in \sqrt{\text{CRLB}(\tau(\theta))}$$

19 / 43

### 4.1 Invertible Interval

interval construction, as we have done it, is **not** functionally equivalent!!

- Can we come up with an approach which does possess the equivariance property?
  - Yes, as long as the functional transformation in question is invertible.
  - Let's consider an **asymptotic likelihood-based confidence interval procedure** which is **parameterization equivariant**.
  - Specifically, this means that if we find a confidence region,  $C$ , for  $\theta$  based on **this new procedure** and transform all of its values [which we sometimes denote as  $\tau(C) = \{\tau(\theta) : \theta \in C\}$ ] then we will arrive at the same confidence region as if we had applied our new procedure to the parameter  $\tau$  directly.

## 5. Asymptotic MLRT Interval Estimation (L09b/p26)

Did we have to use the asymptotic result of the LRT for our interval.  
No, but it is more straightforward.

## 6. CDF Method

- Let  $T$  be a statistic with a continuous cdf  $F_T(t; \theta)$ .
- Let  $\alpha_1 + \alpha_2 = \alpha$  with  $0 < \alpha < 1$ .
- Suppose that for each  $t \in T$ , the functions  $\theta_L(t)$  and  $\theta_U(t)$  can be defined as:

- If  $F_T(t; \theta)$  is a decreasing function of  $\theta$  for each  $t$ , define  $\theta_L(t)$  and  $\theta_U(t)$  by:

$$F_T(t; \theta_U(t)) = \alpha_1 \quad F_T(t; \theta_L(t)) = 1 - \alpha_2$$

- If  $F_T(t; \theta)$  is an increasing function of  $\theta$  for each  $t$ , define  $\theta_L(t)$  and  $\theta_U(t)$  by:

$$F_T(t; \theta_L(t)) = \alpha_1 \quad F_T(t; \theta_U(t)) = 1 - \alpha_2$$

Then the interval  $[\theta_L(t), \theta_U(t)]$  is a  $1 - \alpha$  confidence interval for  $\theta$ .

- We can prove that  $F_T(t; \theta)$  is monotone in  $\theta$ . See C&B.

**Example:** Consider  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Unif}(0, \theta)$ .

- So we have the following CDF for  $X$ :

$$F_X(x; \theta) = \frac{x}{\theta} \mathbb{I}_{(0 \leq x \leq \theta)}$$

- We know the MLE for  $\theta$  is  $T = \max(X_1, \dots, X_n)$

$$\begin{aligned} F_T(t; \theta) &= \Pr(T \leq t) = \Pr\{\max(X_1, \dots, X_n) \leq t\} \\ &= \Pr\{X_1 \leq t, \dots, X_n \leq t\} \\ &= \Pr\{X_1 \leq t\} \times \dots \times \Pr\{X_n \leq t\} \\ &= \{F_X(t; \theta)\}^n \\ &= \frac{t^n}{\theta^n} \mathbb{I}_{(0 \leq t \leq \theta)} \end{aligned}$$

- Note:  $F_T(t; \theta)$  is a decreasing function for  $\theta$ . Let  $\alpha_1 = \alpha_2 = \alpha/2$ . We have:

$$\begin{aligned} F_T(t; \theta_U(t)) &= \alpha/2 \\ \left(\frac{t}{\theta_U}\right)^n &= \alpha/2 \\ \theta_U &= t(\alpha/2)^{-(1/n)} \end{aligned}$$

$$\begin{aligned} F_T(t; \theta_L(t)) &= 1 - \alpha/2 \\ \left(\frac{t}{\theta_L}\right)^n &= 1 - \alpha/2 \\ \theta_L &= t(1 - \alpha/2)^{-(1/n)} \end{aligned}$$

# Bayesian Inference

## 1. Introduction

- Bayesian inference is the application of Bayes' rule to the problem of 'guessing' an unknown quantity(ies).
- The key distinction between Bayesian and sampling theory statistics is the issue of what is to be regarded as random and what is to be regarded as fixed.
  - To a Bayesian, parameters are random and data, once observed, are fixed.
  - To a sampling theorist (classical statistician), data are random even after being observed (how inference is done), but parameters are fixed.

## 2. Advantages

- Bayesian methods provide the user with the ability to formally incorporate prior information (Carlin and Louis, 2009).
- If  $p(y|\theta)$  and  $p(\theta)$  represent a rational person's beliefs then Bayes' rule is an optimal method for updating those beliefs.
- If  $p(y|\theta)$  and  $p(\theta)$  approximately represents beliefs then  $p(\theta|y)$  is also approximate and may be useful (Hoff 2009).
- Inferences are conditional on the actual data (Carlin and Louis, 2009).
- Bayesian answers are more easily interpretable by non-statisticians (Carlin and Louis, 2009).
- All Bayesian analyses follow directly from the posterior; no separate theories of estimation, testing, multiple comparisons, etc. are needed (Carlin and Louis, 2009).
- In many complicated statistical problems there are no obvious non-Bayesian methods of estimation or inference (Hoff 2009).

## 3. Bayesian Inference

$$\begin{aligned}\pi(\theta|x) &= \frac{f(x|\theta)\pi(\theta)}{\int_{\Theta} f(x|\theta)\pi(\theta)d\theta} \\ &= \frac{f(x|\theta)\pi(\theta)}{m(x)} \text{ marginal dist for } X \\ &\propto f(x|\theta)\pi(\theta)\end{aligned}$$

$L(\theta|x)$

- $\pi(\theta|x)$  is the posterior distribution for  $\theta$ .
- $f(x|\theta)$  is the joint sampling distribution for  $X$  or the likelihood for  $\theta$ .
- $\pi(\theta)$  is the prior distribution for  $\theta$ .
- $m(x)$  is the marginal distribution for  $X$ .

It is important to remember though that in a Bayesian framework  $\theta$  is random. In the frequentist framework it is fixed!

- Bayesian approach:

- We start with a prior belief about a situation (represented through a probability distribution parameterized by  $\theta$ ).
- We observe data  $x$ .
- Given the data  $x$ , we update our beliefs about  $\theta$  via Bayes' rule and obtain the posterior distribution.

As  $\pi(\theta|x)$  is a probability distribution function, we can consider a number of summaries for  $\theta$  after we observe the data:

- Posterior mode: the value of  $\theta$  which maximizes the posterior distribution
- Posterior median
- Posterior mean  $\hat{\theta}_B = E\{[\theta|x]\} = \int_{\Theta} \theta \pi(\theta|x) d\theta$
- We can also consider functions of  $\theta$ ,  $\tau(\theta)$ :

$$\widehat{\tau(\theta)}_B = E\{[\tau(\theta)|x]\} = \int_{\Theta} \tau(\theta) \pi(\theta|x) d\theta$$

### 3.1 Conjugate **Family**

**Definition (Section 6.4.2):** Let  $\mathcal{F}$  denote the class of pdfs or pmfs  $f(x|\theta)$ . A class  $\mathcal{P}$  of prior distributions is a **conjugate family** for  $\mathcal{F}$  if the posterior distribution is in the class  $\mathcal{P}$  for all  $f \in \mathcal{F}$ , all priors in  $\mathcal{P}$ , and all  $x \in \mathcal{X}$ .

- In other words, if the prior distribution is in the same family of distributions as the posterior, then it is a conjugate prior distribution.
- For  $f(x|\theta)$  examine the kernel with regard to  $\theta$  and see if you recognize the distribution. This will be the conjugate distribution.

**Example:** Let  $X_1, \dots, X_n \sim \text{Poisson}(\lambda)$ .

- Determine a conjugate prior distribution and then determine the posterior distribution for  $\lambda$ .

$$f(x|\lambda) = \frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod x_i!} \\ \Rightarrow \lambda^{\sum x_i} e^{-n\lambda}$$

- This is a kernel for a gamma distribution. Thus we have the following conjugate prior:

$$\lambda \sim \text{gamma}(a, b)$$

$$p(\lambda|x) \propto p(x|\lambda)p(\lambda) \\ \propto \left[ \lambda^{\sum x_i} e^{-n\lambda} \right] \left[ \lambda^{a-1} e^{-\lambda/b} \right] \\ = \lambda^{\sum x_i + a - 1} e^{-\lambda(n+1/b)} \\ = \lambda^{a^*-1} e^{-\lambda/b^*}$$

$$p(\lambda|x) \propto p(x|\lambda)p(\lambda) \\ \propto \left[ \lambda^{\sum x_i} e^{-n\lambda} \right] \left[ \lambda^{a-1} e^{-\lambda/b} \right] \\ = \lambda^{\sum x_i + a - 1} e^{-\lambda(n+1/b)} \\ = \lambda^{a^*-1} e^{-\lambda/b^*}$$

## 4. Bayesian Testing

consider the following ratio of the posterior model probabilities

$$\text{Ratio} = \frac{\pi(M_2|\mathbf{x})}{\pi(M_1|\mathbf{x})} = \underbrace{\frac{f(\mathbf{x}|M_2)}{f(\mathbf{x}|M_1)} \times \frac{\pi(M_2)}{\pi(M_1)}}_{BF(M_2; M_1) \times \frac{\pi(M_2)}{\pi(M_1)}}$$

Where  $BF(M_2; M_1)$  is called the **Bayes factor**.

- Typically  $\pi(M_2) = \pi(M_1)$ , so the ratio of the posterior probabilities is the Bayes factor.
- The Bayes factor looks like a likelihood ratio. However, the difference is that  $\theta$  has been integrated out in both the numerator and denominator, so we have the marginal distribution of the data given the model.
- If  $f(\mathbf{x}|M_2) > f(\mathbf{x}|M_1)$  or  $\frac{f(\mathbf{x}|M_2)}{f(\mathbf{x}|M_1)} > 1$  then we have support for  $M_2$  against  $M_1$ .
- Jeffreys, H. (1961 appendix B) suggested the following:

$BF(M_2; M_1) = B_{21}$	Evidence against model 1 ( $H_0$ )
1 to 3.2	Not worth more than a bare mention
3.2 to 10	Substantial
10 to 100	Strong
> 100	Decisive

## 5. Bayesian Interval Estimation

In Bayesian estimating, we have the whole distribution for  $\pi(\theta|\mathbf{x}) = \frac{L(\theta|\mathbf{x})\pi(\theta)}{m(\mathbf{x})}$

- To obtain an interval we simply consider:  $P\pi(\theta|\mathbf{x})(C) = \int_C \pi(\theta|\mathbf{x})d\theta = 1 - \alpha$

### 5.1 Choices of Reject Region (C)

- Equal tailed:  $\int_{-\infty}^{\theta_L} \pi(\theta|\mathbf{x})d\theta = \alpha/2, \quad \int_{\theta_U}^{\infty} \pi(\theta|\mathbf{x})d\theta = \alpha/2$
- Smallest length: We can choose  $C$  to minimize  $\theta_U - \theta_L$ .
- Highest posterior density region (HPD): We define  $C$  to be that set with posterior probability  $1 - \alpha$  which satisfies the criterion:

$$\theta_1 \in C \quad \text{and} \quad \pi(\theta_2|\mathbf{x}) > \pi(\theta_1|\mathbf{x}) \Rightarrow \theta_2 \in C$$

$C$  contains the values of  $\theta$  which have the highest posterior density values, so that we can determine HPD regions as the set:

$$C = \{\theta \in \Theta : \pi(\theta|\mathbf{x}) > c_\alpha\}$$

If the posterior is unimodal then this will be the smallest length interval!

## 6. **Properties** of Bayesian Inference

### 6.1 Sufficient

**Definition 7.1:** A statistic  $T(\mathbf{X})$  is **sufficient** for  $\theta$  if and only if the posterior distribution of  $\theta$  given  $\mathbf{X}$  is the same as the posterior distribution of  $\theta$  given  $T(\mathbf{X})$ .

**Proof:** Note that Definitions 2.5 and 7.1 are the same!

Suppose that  $T(\mathbf{X})$  satisfies Definition 2.5. Then:

$$f(\mathbf{x}; \theta) = g(\mathbf{x}|t, \theta)h(t|\theta) = g(\mathbf{x}|t)h(t|\theta)$$

- The posterior is

$$\begin{aligned} p(\theta|\mathbf{x}) &\propto f(\mathbf{x}; \theta)p(\theta) \\ &\propto h(t|\theta)p(\theta) \\ &\propto p(\theta|t) \end{aligned}$$

- Now assume that  $T(\mathbf{X})$  satisfies Definition 7.1.

! key idea

$$\begin{aligned} f(\mathbf{x}|\theta) &= \frac{p(\theta|\mathbf{x})h(\mathbf{x})}{p(\theta)} \\ &= \frac{p(\theta|t)h(\mathbf{x})}{p(\theta)} \\ &= K_1[t|\theta] K_2[\mathbf{x}] \end{aligned}$$

- From the **factorization theorem**, it follows that  $T(\mathbf{X})$  is a sufficient statistic.

### 6.2 Asymptotic (to Normal)

$$\underline{p(\theta|\mathbf{y})} \approx \text{normal}\left(\hat{\theta}, [I(\hat{\theta})]^{-1}\right)$$

( rough idea -> L10b/p33)

---

# Decision Theory

## 1. Basic Elements

1. A number of 'actions' are possible; we must decide which to take.
2. A number of 'states of nature' are possible; we don't know in general which will occur.
3. The relative desirability of the various actions for each state of nature can be quantified.
4. Prior information may be available regarding the relative probabilities of the various states of nature.
5. Data may be available which will add to our knowledge of the relative probabilities of the states of nature.

## 2. General Process

- Let  $\theta$  denote the true state of nature.
- Suppose we are able to observe some data . . . a draw from the random variable  $\mathbf{X}$ , whose distribution depends on  $\theta$ . Sometimes no data are available.
- A decision procedure  $\delta$ , specifies which action to take for each value of  $\mathbf{X}$ .
- If we observe  $\mathbf{X} = \mathbf{x}$ , then we adopt the procedure  $\delta(\mathbf{x})$ .
- Whether  $\delta(\mathbf{x})$  was a good choice depends on the loss function, which measures the loss from  $\delta(\mathbf{x})$  when  $\theta$  holds.

$$L_S(\theta, \delta(\mathbf{x}))$$

Note: The negative of a loss function is a utility function.

- Frequentists want to minimize expected loss.

$$\min_{\delta} \int L_S(\theta, \delta(\mathbf{x})) f(\mathbf{x}; \theta) d\mathbf{x}$$

- Bayesians want to minimize posterior expected loss.

$$\min_{\delta} \int L_S(\theta, \delta(\mathbf{x})) p(\theta | \mathbf{x}) d\theta$$

## 3. Risk Function

$$R(\theta, \delta(\mathbf{x})) = \int L_S(\theta, \delta(\mathbf{x})) L(\theta; \mathbf{x}) d\mathbf{x}.$$

This is the expected loss with respect to the joint distribution (i.e. likelihood).

## 4. Inadmissible

A procedure  $\delta_1$  is **inadmissible** if there exists another procedure  $\delta_2$  such that

$$R(\theta, \delta_1) \geq R(\theta, \delta_2) \quad \forall \theta,$$

with strict inequality for some  $\theta$ .

**" $\delta_1$  is never better than  $\delta_2$  and sometimes worse"**

5. **Minimax** Procedure       $\max_{\theta} R(\theta, \delta)$     is minimized.

6. **Bayes** Procedure       $\int R(\theta, \delta) p(\theta) d\theta$     is minimized.

Expected risk with regard to the prior distribution minimized.

### Extended

$$\begin{aligned}\int R(\theta, \delta) p(\theta) d\theta &= \int_{\Theta} \left[ \int_{\mathcal{X}} L_S(\theta, \delta) L(\theta; \mathbf{x}) d\mathbf{x} \right] p(\theta) d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} L_S(\theta, \delta) \frac{L(\theta; \mathbf{x}) p(\theta)}{p(\mathbf{x})} p(\mathbf{x}) d\mathbf{x} d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} L_S(\theta, \delta) p(\theta|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} d\theta \\ &= \int_{\mathcal{X}} p(\mathbf{x}) \left[ \int_{\Theta} L_S(\theta, \delta) p(\theta|\mathbf{x}) d\theta \right] d\mathbf{x}\end{aligned}$$

For any value of  $\mathbf{x}$  we should minimize  $\int_{\Theta} L_S(\theta, \delta) p(\theta|\mathbf{x}) d\theta$ , which is **posterior expected loss**.

- Under certain conditions:

For example, if  $\theta$  is discrete and only takes a number of finite values then (2) holds. Additionally if  $P(\theta) > 0$  for all  $\theta$  then (1) holds.

1. A Bayes procedure is necessarily admissible.
2. Every admissible procedure is a Bayes procedure for some prior distribution.

- The link between Bayes and minimax procedures:

1. A Bayes procedure with constant risk for  $\theta$  is minimax.
2. A minimax procedure is generally a Bayes procedure for some prior distribution. In particular, the so called **least favourable prior distribution**.

## 7. Point Estimation

1. The possible actions are possible estimators.
2. Possible states of nature correspond to the true value of  $\theta$ .
3. Loss is some function of the estimator  $\hat{\theta}$  and  $\theta$ :  $L_S(\theta, \delta = \hat{\theta})$ .
4. Prior information is quantified by means of a prior distribution  $p(\theta)$ .
5. The available data, are our draws from  $f(\mathbf{x}|\theta)$ .

- Some simple loss functions:

1. Zero-one loss:  $L_S(\theta, \delta = \hat{\theta}) = \begin{cases} 0, & |\hat{\theta} - \theta| < b, \\ a, & |\hat{\theta} - \theta| \geq b \end{cases}$  where  $a, b > 0$ .

2. Absolute error loss:  $L_S(\theta, \delta = \hat{\theta}) = a|\hat{\theta} - \theta|$

3. Squared Error Loss (quadratic loss):  $L_S(\theta, \delta = \hat{\theta}) = a(\hat{\theta} - \theta)^2$

To minimize, we set  $\hat{\theta} = \bar{\theta}$ .

- Absolute error loss:  $E[|\theta - \hat{\theta}|] = \int |\theta - \hat{\theta}| p(\theta|\mathbf{x}) d\theta$ .

- This is minimized when  $\hat{\theta} = \text{median}$ .

(James-Stein Estimator L11a/p32)

# Non-Parametric Methods

## 1. Empirical Distribution Function

- The empirical distribution function  $\hat{F}$  is the CDF of a new discrete random variable, say  $X^*$ .
- It can be shown that  $\hat{F}$  is a sufficient statistic for  $F$  (based on a random sample), so
- $\hat{F}$  and  $X^*$  mimics the relationship between  $F$  and the  $X$ .
- This leads to studying  $(\hat{F}, X^*)$  to learn about  $(F, X)$ .

$$\begin{aligned}\hat{F}(x) &= \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{(x_i \leq x)} \\ E\{\hat{F}(x)\} &= \frac{1}{n} \sum_{i=1}^n Pr(X_i \leq x) = \frac{1}{n} \sum_{i=1}^n F(x) = F(x) \\ Var\{\hat{F}\} &= \frac{1}{n} F(x)(1 - F(x))\end{aligned}$$

we are interested in  $\theta(F) = E_F(X)$ :

$$\hat{\theta} = \theta(\hat{F}) = E_{\hat{F}}(X) = \sum_{x \in \mathcal{X}} x p_{\hat{F}}(x) = \sum_{i=1}^n x_i \frac{1}{n} = \bar{x}$$

## 2. Bootstrap

- We create  $B$  “bootstrap” (re-sampled) data sets:  $\{X_{1,1}^*, \dots, X_{n,1}^*\}, \dots, \{X_{1,B}^*, \dots, X_{n,B}^*\}$
- With each sample we estimate  $\hat{\theta}^b = \theta(\hat{F}_b^*)$ .
- Since we know the “true” distribution  $\hat{F}$ , we can determine exactly the bias and variance of  $\theta(\hat{F}^*)$ .

$$\begin{aligned}\hat{B}_B &= E_{\hat{F}}\{\theta(\hat{F}^*)\} - \theta(\hat{F}) \approx \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^* - \hat{\theta} \\ \hat{Var}_B\{\theta(\hat{F}^*)\} &\approx \frac{1}{B-1} \sum_{b=1}^B \left( \hat{\theta}_b^* - \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^* \right)^2\end{aligned}$$

### 2.1

#### Disadvantages

- As implemented, the bootstrap method yields a different answer every time (of course, the differences will be very small if  $B$  is large).
- Another drawback is that if  $\theta$  is complicated to calculate (perhaps because it is implicitly defined as the solution to an equation, just as the MLE was) then computing its value for each of  $B$  re-sampled data sets is computationally quite expensive and time consuming.
- If we truly believe the parametric structure we have set up, then the parametric estimators have nice optimal properties.

#### Advantages

- The bootstrap is a very flexible and widely applicable approach which deserves more attention than it currently gets among statistical practitioners
- The bootstrap can even be extended to circumstances beyond the *iid* setting on which we have focused here.
- But a word of warning on complicated (non iid settings):

- We cannot always guarantee that using the bootstrap paradigm (replacing  $F$  by  $\hat{F}$  and  $\hat{F}$  by  $\hat{F}^*$ ) to estimate bias and variance will yield valid estimates.

## 2.2 Bootstrap Interval Estimation

### 2.21 Asymptotically

The estimated standard deviation was:

$$\hat{\sigma}_B(\hat{\theta}) = \sqrt{\underbrace{\frac{1}{B-1}}_{\text{B-1}} (\hat{\theta}_b^* - \bar{\hat{\theta}}^*)^2}$$

We saw MLEs are asymptotically normal, and in fact many estimators are, we could just use that idea:

$$[\hat{\theta} - z_{\alpha/2} \hat{\sigma}_B(\hat{\theta}), \hat{\theta} + z_{\alpha/2} \hat{\sigma}_B(\hat{\theta})]$$

### 2.2.2 Empirical

We can use the empirical quantiles of the “bootstrap distribution”. So we have:

$$P_{\hat{F}}(\hat{\theta}^* \leq \hat{\theta}_L^*) = \alpha/2, \quad P_{\hat{F}}(\hat{\theta}^* \leq \hat{\theta}_U^*) = 1 - \alpha/2$$

Using B re-samples, calculate  $\hat{Q}_b$  for each re-sample and approximate  $q_L$  and  $q_U$  with (where  $\alpha_1 + \alpha_2 = \alpha$ ) using the empirical distribution of  $\hat{Q}$ :

$$\hat{q}_L = \hat{Q}_{\alpha_1} \quad \hat{q}_U = \hat{Q}_{1-\alpha_2} \quad \text{from } \text{Bootstrap}$$
$$[\theta(\hat{F}) - \hat{q}_U \hat{\sigma}(F), \theta(\hat{F}) + \hat{q}_L \hat{\sigma}(F)]$$

### 2.2.3 Properties of Intervals

- Shortest intervals for a given confidence or credibility (eg. 95%).
- Range respecting
- Parameterization equivariance (We would like our interval construction procedures to transform appropriately if we change our focus from  $\tau = \tau(\theta)$  to  $\gamma = \gamma(\tau) = \gamma\{\tau(\theta)\}$ )

## Non-Parametric Testing

- We will consider two types of tests:
  - permutation/randomization tests
  - bootstrap tests

(L11b)