

# COMPS NOTES

## Table of contents

<b>Probability Foundation</b>	<b>4</b>
MGF-Moment Generate Function . . . . .	4
Independent . . . . .	8
Conditional . . . . .	9
Definition . . . . .	9
Conditional Expectation & Variance ... . . . .	10
Total Variance . . . . .	11
Hierarchical Model . . . . .	11
Multivariate Distribution . . . . .	12
Bivariate Normal Distribution . . . . .	12
Multinomial . . . . .	15
Transformation/Location Scale Family . . . . .	16
Univariate Transformation . . . . .	16
Location Scale Family . . . . .	16
Multivariate Transformation . . . . .	20
Order Statistics . . . . .	21
Discrete Ordered Statistics . . . . .	22
Continuous Ordered Statistics . . . . .	22
<b>Asymptotic Probability</b>	<b>24</b>
Necessary limits . . . . .	24
Exponential . . . . .	24
Necessary Inequality . . . . .	24
Markov's Inequality . . . . .	24
Chebychev's Inequality . . . . .	25
Holder's Inequality . . . . .	25
Cauchy's Inequality . . . . .	25
Jenson's Inequality . . . . .	26
Boferroni's Inequality . . . . .	26
Convergence . . . . .	27
Converge in Probability . . . . .	27
Converge in Distribution . . . . .	27
Converge Almost Surely (related to Strong Law of Large Number) . . . . .	29

Law of Large Number . . . . .	29
Weak Law of Large Number . . . . .	29
Strong Law of Large Number . . . . .	30
CLT . . . . .	31
Delta Methods . . . . .	33
Taylor Series . . . . .	33
Delta Method for Univariate . . . . .	34
Delta Method for Multivariate . . . . .	36
<b>Statistics &amp; Estimator</b>	<b>36</b>
Sufficient Statistics . . . . .	36
Definition . . . . .	36
Factorization Theorem: How to find Sufficient Statistics . . . . .	36
Minimum Sufficient Statistics and Complete Statistics . . . . .	39
MSS . . . . .	39
Ancillary Statistics . . . . .	44
Complete Statistics . . . . .	46
Estimator . . . . .	51
MLE . . . . .	51
Methods of Mommment . . . . .	51
Linear Regression . . . . .	52
UMVUE . . . . .	52
CRLB-Definition . . . . .	52
Rao-Blackwell-Finding UMVUE( <b>Must be Unbiased</b> ) . . . . .	60
Lehmann-scheffe-Finding unique UMVUE . . . . .	65
Consistency, Asymptotic Variance, efficiency of MLE . . . . .	68
<b>Statistical Inference</b>	<b>74</b>
Hypothesis Testing . . . . .	74
Power/Size of the Test . . . . .	74
Most Powerful Test . . . . .	77
Likelihood Ratio Test . . . . .	84
Large Sample Test Method . . . . .	91
Confidence Interval . . . . .	96
Confidence Interval: Definition, Coverage Probability, Coverage Coefficient and Length . . . . .	96
Unbiased CI . . . . .	100
Find CI by Pivoting CDF . . . . .	102
Find CI by MLE asymptotic . . . . .	105
Find CI by Inverting Test . . . . .	109
P-value . . . . .	113
One Sample Testing . . . . .	116
Two Sample Testing . . . . .	116
Assumption . . . . .	116
Two sample t-test . . . . .	116

<b>Bayesian</b>	<b>118</b>
Posterior . . . . .	118
Credible Interval . . . . .	118
<b>Regression</b>	<b>118</b>
Linear Regression . . . . .	118
Assumption . . . . .	118
Algreba Form . . . . .	119
Projection . . . . .	119
Interpretation . . . . .	120
Multivariate Regression . . . . .	120
Confounding/Percision Variable . . . . .	120
Spline . . . . .	121
Weighted Regression . . . . .	123
Saturated Model . . . . .	123
Multivariate Regression . . . . .	124
Logistic Regression . . . . .	124
Ordinal Regression . . . . .	124
Poisson Regression . . . . .	124
Longitudinal Regressoin . . . . .	124
Survival Regression . . . . .	124
<b>Integral</b>	<b>124</b>
Substitution . . . . .	124
<b>Distribution</b>	<b>124</b>
Discrete . . . . .	124
Poisson . . . . .	124
Geometric . . . . .	124
<b>Simulation</b>	<b>124</b>
<b>COMPS Practice</b>	<b>125</b>

## Probability Foundation

### MGF-Moment Generate Function

- Use to find moment, if question asked to find second moment and you find yourself using pdf, usually wrong and much easier using MGF
- Use MGF to prove converge in distribution, see comps 2022

#### Theorem

**Definition 2.3.6** Let  $X$  be a random variable with cdf  $F_X$ . The moment generating function (mgf) of  $X$  (or  $F_X$ ), denoted by  $M_X(t)$ , is

$$M_X(t) = Ee^{tX}$$

provided that the expectation exists for  $t$  in some neighborhood of 0. That is, there is an  $h > 0$  such that, for all  $t$  in  $-h < t < h$ ,  $Ee^{tX}$  exists. If the expectation does not exist in a neighborhood of 0, we say that the moment generating function does not exist. More explicitly, we can write the mgf of  $X$  as

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx \quad \text{if } X \text{ is continuous,}$$

or

$$M_X(t) = \sum_x e^{tx} P(X = x) \quad \text{if } X \text{ is discrete.}$$

It is very easy to see how the mgf generates moments. We summarize the result in the following theorem.

**Theorem 2.3.7** If  $X$  has mgf  $M_X(t)$ , then

$$EX^n = M_X^{(n)}(0)$$

where we define

$$M_X^{(n)}(0) = \left. \frac{d^n}{dt^n} M_X(t) \right|_{t=0}$$

That is, the  $n$ th moment is equal to the  $n$ th derivative of  $M_X(t)$  evaluated at  $t = 0$ .

**Theorem 2.3.12 (Convergence of mgfs)** Suppose  $\{X_i, i = 1, 2, \dots\}$  is a sequence of random variables, each with mgf  $M_{X_i}(t)$ . Furthermore, suppose that

$$\lim_{i \rightarrow \infty} M_{X_i}(t) = M_X(t), \quad \text{for all } t \text{ in a neighborhood of } 0,$$

and  $M_X(t)$  is an mgf. Then there is a unique cdf  $F_X$  whose moments are determined by  $M_X(t)$  and, for all  $x$  where  $F_X(x)$  is continuous, we have

$$\lim_{i \rightarrow \infty} F_{X_i}(x) = F_X(x)$$

That is, convergence, for  $|t| < h$ , of mgfs to an mgf implies convergence of cdfs.

**Theorem 2.3.15** For any constants  $a$  and  $b$ , the mgf of the random variable  $aX + b$  is given by

$$M_{aX+b}(t) = e^{bt} M_X(at)$$

### Example: Use MGF to show poison Convergence

$$P(X = x) \approx P(Y = x)$$

for large  $n$  and small  $np$ . We now show that the mgfs converge, lending credence to this approximation. Recall that

$$M_X(t) = [pe^t + (1 - p)]^n$$

For the Poisson( $\lambda$ ) distribution, we can calculate (see Exercise 2.33)

$$M_Y(t) = e^{\lambda(e^t - 1)}$$

and if we define  $p = \lambda/n$ , then  $M_X(t) \rightarrow M_Y(t)$  as  $n \rightarrow \infty$ . The validity of the approximation in (2.3.9) will then follow from Theorem 2.3.12.

We first must digress a bit and mention an important limit result, one that has wide applicability in statistics. The proof of this lemma may be found in many standard calculus texts.

**Lemma 2.3.14** Let  $a_1, a_2, \dots$  be a sequence of numbers converging to  $a$ , that is,  $\lim_{n \rightarrow \infty} a_n = a$ . Then

$$\lim_{n \rightarrow \infty} \left(1 + \frac{a_n}{n}\right)^n = e^a$$

Returning to the example, we have

$$M_X(t) = [pe^t + (1 - p)]^n = \left[1 + \frac{1}{n}(e^t - 1)(np)\right]^n = \left[1 + \frac{1}{n}(e^t - 1)\lambda\right]^n,$$

because  $\lambda = np$ . Now set  $a_n = a = (e^t - 1)\lambda$ , and apply Lemmas 2.3.14 to get

$$\lim_{n \rightarrow \infty} M_X(t) = e^{\lambda(e^t - 1)} = M_Y(t)$$

the moment generating function of the Poisson.

The Poisson approximation can be quite good even for moderate  $p$  and  $n$ . In Figure 2.3.3 we show a binomial mass function along with its Poisson approximation, with  $\lambda = np$ . The approximation appears to be satisfactory.

**Example: Rao's Blackwell 6.8**

- If the Blackwellization of  $W$  yields an estimator different from  $W$ , then the new estimator has a smaller MSE.

**Example 6.8** Suppose  $X_1, X_2, \dots, X_n \sim \text{iid Poisson}(\theta)$ . Let  $\tau(\theta) = P\{X = 0\} = e^{-\theta}$ .

- Consider the statistic  $W = \mathbf{1}_{\{0\}}(X_1) \sim \text{Bernoulli}(\tau(\theta))$
- $E_\theta[W] = \tau(\theta)$ , that is,  $W$  is unbiased for  $\tau(\theta)$ .
- $\text{Var}_\theta[W] = \tau(\theta)(1 - \tau(\theta))$
- From previous work we know that  $T = \sum X_i$  is sufficient for  $\theta$ . Note that  $T \sim \text{Poisson}(n\theta)$ .
- Find a better unbiased estimator, in terms of the variance.
- Confirm this estimator is unbiased.
- Is this estimator UMVUE?

$$\begin{aligned}\Phi(s) &= E[W|T = s] = P(X_1 = 0|T = s) \\ &= \frac{P(X_1 = 0, \sum_1 X_i = s)}{P(T = s)} \\ &= \frac{P(X_1 = 0) * P(\sum_2 X_i = s)}{P(T = s)} \\ &= \frac{\underbrace{Poi((n-1)\theta)}}{P(T = s)} \\ &= \left(\frac{n-1}{n}\right)^s \\ \Phi(T) &= \left(\frac{n-1}{n}\right)^{\sum X_i}\end{aligned}$$

This is a better estimator by Rao-blackwell.

$$\begin{aligned}E[\Phi(T)] &= E\left[\left(\frac{n-1}{n}\right)^{\sum X_i}\right] = E\left[\exp\left(\log\left(\frac{n-1}{n}\right) \sum X_i\right)\right] \\ &\text{by moment generate function : } M_T(t) = \exp(n\theta(e^t - 1)) \\ &= M\left(\log\left(\frac{n-1}{n}\right)\right) = \exp\left[n\theta\left(\frac{n-1}{n} - 1\right)\right] \\ &= \exp(-\theta)\end{aligned}$$

This shows that  $\Phi$  is unbiased.

## Independent

- $f(x, y) = f(x)f(y)$
- $f(x|y) = f(x)$
- $E[xy] = E[x]E[y]$

**Example: Show  $\bar{X} \perp S^2$**

Without loss of generality, we assume i.i.d sample  $X_1, \dots, X_n \sim N(0, 1)$

Let's look at the first case  $X_1$ :

Recall that  $\sum_1 (X_i - \bar{X}) = n \frac{\sum_1 X_i}{n} - n\bar{X} = 0$ , we have:

$$\begin{aligned} \sum_2 (X_i - \bar{X}) &= -(X_1 - \bar{X}) \\ \Rightarrow (X_1 - \bar{X})^2 &= [\sum_2 (X_i - \bar{X})]^2 \end{aligned}$$

Replace it into sample variance we have:

$$\begin{aligned} S^2 &= \frac{1}{n-1} \sum_1 (X_i - \bar{X})^2 \\ &= \frac{1}{n-1} [(X_1 - \bar{X})^2 + \sum_2 (X_i - \bar{X})^2] \\ &= \frac{1}{n-1} \{ [\sum_2 (X_i - \bar{X})]^2 + \sum_2 (X_i - \bar{X})^2 \} \\ &\text{is a function of } X_i - \bar{X}, i \in \{2, \dots, n\} \end{aligned}$$

Next, we need to show jointly  $\bar{X}$  and  $X_i - \bar{X}, i \in \{2, \dots, n\}$  are independent.

We do that by finding the joint pdf and

Now, let  $Y_1 = \bar{X}, Y_2 = X_2 - \bar{X}, \dots, Y_n = X_n - \bar{X}$ , we have

$$\begin{aligned} \sum_2 Y_i &= \sum_2 X_i - \sum_2 \bar{X} \\ &= \sum_1 \bar{X} - X_1 - \sum_2 \bar{X} \\ &= \bar{X} - X_1 \\ \Rightarrow \begin{cases} X_1 = Y_1 - \sum_2 Y_i \\ X_i = Y_i + Y_1 \end{cases} \end{aligned}$$

So the jacobian is



$$\begin{bmatrix} \frac{\partial X_1}{\partial Y_1} & \dots & \frac{\partial X_n}{\partial Y_1} \\ \dots & \dots & \dots \\ \frac{\partial X_1}{\partial Y_n} & \dots & \frac{\partial X_n}{\partial Y_n} \end{bmatrix} = \begin{bmatrix} 1 & -1 & -1 & \dots & -1 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ 1 & \dots & \dots & \dots & 0 \\ 1 & 0 & 0 & \dots & 1 \end{bmatrix}$$

We can show that  $|J| = n$ ,

By i.i.d, we have

$$f_{X_1, \dots, X_n} = (2\pi)^{-n/2} \exp\left\{-\frac{1}{2} \sum X_i^2\right\}$$

Replacing  $Y_i$  we have

$$\begin{aligned} f_{Y_1, \dots, Y_n} &= n(2\pi)^{-n/2} \exp\left\{-\frac{1}{2}[(Y_1 - \sum_2 Y_i)^2 + \sum_2 (Y_i + Y_1)^2]\right\} \\ &= c * \exp\left\{-0.5[(Y_1^2 - 2Y_1 \sum_2 Y_i + (\sum_2 Y_i)^2) + \sum_2 (Y_1^2 + 2Y_1 Y_i + Y_i^2)]\right\} \\ &= c * \exp(-0.5nY_1^2) * \exp(-0.5[(\sum_2 Y_i)^2 + \sum_2 Y_i^2]) \end{aligned}$$

\*We left out the indicator function since they don't have interaction indicator function\*

By factorization,  $Y_1 = \bar{X} \perp \sum_2 Y_i$ , since a function of independent R.V is also independent, we have  $\bar{X} \perp S^2$  ■

## Conditional

### Definition

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

- if  $A \perp B$ ,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A)$$

- if  $A \subset B$ ,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)}{P(B)}$$

## Conditional Expectation & Variance ...

$$E[g(X)|Y] = \int g(x) * f_{X|Y}(x)dx \Rightarrow h(Y)$$

- Conditional something of  $Y$  should be a function of  $Y$
- $E[X|X] = X$  itself:

Proof:

$$E[X|X] = \int x f(x|X)dx$$

where  $f(x|X)$  is the conditional probability density function of  $X$  given  $X$ .

However, since we already know the value of  $X$ , the conditional probability density function of  $X$  given  $X$  is just a Dirac delta function centered at  $X$ , which is defined as:

$$f(x|X) = \delta(x - X)$$

where  $\delta(x - X)$  is the Dirac delta function, which is zero for all values of  $x$  except  $x = X$ , where it is infinite.

Substituting this into the expression for the conditional expectation, we get:

$$E[X|X] = \int x \delta(x - X)dx$$

Since the Dirac delta function is zero everywhere except at  $x = X$ , this integral evaluates to:

$$E[X|X] = X * \int \delta(x - X)dx$$

where the integral evaluates to 1 since the Dirac delta function integrates to 1 over its support. Therefore, we have:

$$E[X|X] = X$$

- $Var(X|X) = 0$ , I know mind blowing

Proof:

$$Var(X|X) = E[(X - E[X|X])^2|X]$$

Since  $E[X|X] = X$ , we have:

$$Var(X|X) = E[(X - X)^2|X] = E[0|X] = 0$$

## Total Variance

- You failed to show it in homework

$$\text{Var}(X) = \text{Var}(E[X|Y]) + E[\text{Var}(X|Y)]$$

Proof:

$$\begin{aligned} \text{Var}(X) &= E[X - E[X]]^2 \\ &= E\{X - E[X|Y] + E[X|Y] - E[X]\}^2 \\ &= E[X - E[E[X|Y]]]^2 + E[E[X|Y] - E[X]]^2 \\ &\quad + 2E[(X - E[X|Y])(E[X|Y] - E[X])] \end{aligned}$$

The last term in this expression is equal to 0, however, which can easily be seen by iterating the expectation:

$$E([X - E(X | Y)][E(X | Y) - EX]) = E(E\{[X - E(X | Y)][E(X | Y) - EX] | Y\})$$

In the conditional distribution  $X | Y$ ,  $X$  is the random variable. So in the expression

$$E\{[X - E(X | Y)][E(X | Y) - EX] | Y\}$$

$E(X | Y)$  and  $EX$  are constants. Thus,

$$\begin{aligned} E\{[X - E(X | Y)][E(X | Y) - EX] | Y\} &= (E(X | Y) - EX)(E\{[X - E(X | Y)] | Y\}) \\ &= (E(X | Y) - EX)(E(X | Y) - E(X | Y)) \\ &= (E(X | Y) - EX)(0) \\ &= 0. \end{aligned}$$

Thus, from (4.4.6), we have that  $E((X - E(X | Y))(E(X | Y) - EX)) = E(0) = 0$ . Referring back to equation (4.4.5), we see that

$$\begin{aligned} E([X - E(X | Y)]^2) &= E(E\{[X - E(X | Y)]^2 | Y\}) \\ &= E(\text{Var}(X | Y)) \end{aligned}$$

and

$$E([E(X | Y) - EX]^2) = \text{Var}(E(X | Y))$$

establishing (4.4.4).

## Hierarchical Model

- [Amber's Homework 6\(Pretty Bad\)](#)

$$E[X] = E[E[X|Y]] = E[E[E[X|Y, Z]]] = \dots$$

### Example

$$\begin{aligned} X | Y &\sim \text{binomial}(Y, p) \\ Y | \Lambda &\sim \text{Poisson}(\Lambda) \\ \Lambda &\sim \text{exponential}(\beta) \end{aligned}$$

where the last stage of the hierarchy accounts for the variability across different mothers. The mean of  $X$  can easily be calculated as

$$\begin{aligned} EX &= E(E(X | Y)) \\ &= E(pY) \\ &= E(E(pY | \Lambda)) \\ &= E(p\Lambda) && \text{(as before)} \\ &= p\beta, && \text{(exponential expectation)} \end{aligned}$$

## Multivariate Distribution

### Bivariate Normal Distribution

Let  $\mathbf{x}_1$  be the first partition and  $\mathbf{x}_2$  the second. Now define  $\mathbf{z} = \mathbf{x}_1 + \mathbf{A}\mathbf{x}_2$  where  $\mathbf{A} = -\Sigma_{12}\Sigma_{22}^{-1}$ . Now we can write

$$\begin{aligned} \text{cov}(\mathbf{z}, \mathbf{x}_2) &= \text{cov}(\mathbf{x}_1, \mathbf{x}_2) + \text{cov}(\mathbf{A}\mathbf{x}_2, \mathbf{x}_2) \\ &= \Sigma_{12} + \mathbf{A}\text{var}(\mathbf{x}_2) \\ &= \Sigma_{12} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{22} \\ &= 0 \end{aligned}$$

Therefore  $\mathbf{z}$  and  $\mathbf{x}_2$  are uncorrelated and, since they are jointly normal, they are independent. Now, clearly  $E(\mathbf{z}) = \mu_1 + \mathbf{A}\mu_2$ , therefore it follows that

$$\begin{aligned} E(\mathbf{x}_1 | \mathbf{x}_2) &= E(\mathbf{z} - \mathbf{A}\mathbf{x}_2 | \mathbf{x}_2) \\ &= E(\mathbf{z} | \mathbf{x}_2) - E(\mathbf{A}\mathbf{x}_2 | \mathbf{x}_2) \\ &= E(\mathbf{z}) - \mathbf{A}\mathbf{x}_2 \\ &= \mu_1 + \mathbf{A}(\mu_2 - \mathbf{x}_2) \\ &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \mu_2) \end{aligned}$$

which proves the first part. For the covariance matrix, note that

$$\begin{aligned} \text{var}(\mathbf{x}_1 | \mathbf{x}_2) &= \text{var}(\mathbf{z} - \mathbf{A}\mathbf{x}_2 | \mathbf{x}_2) \\ &= \text{var}(\mathbf{z} | \mathbf{x}_2) + \text{var}(\mathbf{A}\mathbf{x}_2 | \mathbf{x}_2) - \mathbf{A}\text{cov}(\mathbf{z}, -\mathbf{x}_2) - \text{cov}(\mathbf{z}, -\mathbf{x}_2)\mathbf{A}' \\ &= \text{var}(\mathbf{z} | \mathbf{x}_2) \\ &= \text{var}(\mathbf{z}) \end{aligned}$$

Now we're almost done:

$$\begin{aligned}
 \text{var}(\mathbf{x}_1 \mid \mathbf{x}_2) &= \text{var}(\mathbf{z}) = \text{var}(\mathbf{x}_1 + \mathbf{A}\mathbf{x}_2) \\
 &= \text{var}(\mathbf{x}_1) + \mathbf{A} \text{var}(\mathbf{x}_2) \mathbf{A}' + \mathbf{A} \text{cov}(\mathbf{x}_1, \mathbf{x}_2) + \text{cov}(\mathbf{x}_2, \mathbf{x}_1) \mathbf{A}' \\
 &= \Sigma_{11} + \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{22} \Sigma_{22}^{-1} \Sigma_{21} - 2 \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \\
 &= \Sigma_{11} + \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} - 2 \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \\
 &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}
 \end{aligned}$$

So

$$f_{X,Y} = N\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}\right)$$

And

$$f_{X|Y} = N\left\{\mu_x + \rho\sqrt{\frac{\sigma_x}{\sigma_y}}(Y - \mu_y), \sigma_x^2(1 - \rho^2)\right\}$$

#### Example exercise 4.45

Show that if  $(X, Y) \sim \text{bivariate normal}(\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2, \rho)$ , then the following are true.

(a) The marginal distribution of  $X$  is  $n(\mu_X, \sigma_X^2)$  and the marginal distribution of  $Y$  is  $n(\mu_Y, \sigma_Y^2)$ . (b) The conditional distribution of  $Y$  given  $X = x$  is

$$n(\mu_Y + \rho(\sigma_Y/\sigma_X)(x - \mu_X), \sigma_Y^2(1 - \rho^2))$$

(c) For any constants  $a$  and  $b$ , the distribution of  $aX + bY$  is

$$n(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\rho\sigma_X\sigma_Y)$$

The mean is easy to check,

$$E(aX + bY) = aEX + bEY = a\mu_X + b\mu_Y$$

as is the variance,

$$\text{Var}(aX + bY) = a^2 \text{Var} X + b^2 \text{Var} Y + 2ab \text{Cov}(X, Y) = a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\rho\sigma_X\sigma_Y$$

To show that  $aX + bY$  is normal we have to do a bivariate transform. One possibility is  $U = aX + bY, V = Y$ , then get  $f_{U,V}(u, v)$  and show that  $f_U(u)$  is normal. We will do this in the standard case. Make the indicated transformation and write  $x = \frac{1}{a}(u - bv), y = v$  and obtain

$$|J| = \begin{vmatrix} 1/a & -b/a \\ 0 & 1 \end{vmatrix} = \frac{1}{a}$$

Then

$$f_{UV}(u, v) = \frac{1}{2\pi a \sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[ \left( \frac{1}{a}(u-bv) \right)^2 - 2\frac{b}{a}(u-bv)+v^2 \right]}.$$

Now factor the exponent to get a square in  $u$ . The result is

$$-\frac{1}{2(1-\rho^2)} \left[ \frac{b^2 + 2\rho ab + a^2}{a^2} \right] \left[ \frac{u^2}{b^2 + 2\rho ab + a^2} - 2 \left( \frac{b + a\rho}{b^2 + 2\rho ab + a^2} \right) uv + v^2 \right]$$

Note that this is joint bivariate normal form since  $\mu_U = \mu_V = 0, \sigma_u^2 = 1, \sigma_v^2 = a^2 + b^2 + 2ab\rho$  and

$$\rho^* = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = \frac{\text{E}(aXY + bY^2)}{\sigma_U \sigma_V} = \frac{a\rho + b}{\sqrt{a^2 + b^2 + 2ab\rho}}$$

thus

$$(1 - \rho^{*2}) = 1 - \frac{a^2\rho^2 + ab\rho + b^2}{a^2 + b^2 + 2ab\rho} = \frac{(1 - \rho^2) a^2}{a^2 + b^2 + 2ab\rho} = \frac{(1 - \rho^2) a^2}{\sigma_u^2}$$

where  $a\sqrt{1-\rho^2} = \sigma_U\sqrt{1-\rho^{*2}}$ . We can then write

$$f_{UV}(u, v) = \frac{1}{2\pi\sigma_U\sigma_V\sqrt{1-\rho^{*2}}} \exp \left[ -\frac{1}{2\sqrt{1-\rho^{*2}}} \left( \frac{u^2}{\sigma_U^2} - 2\rho \frac{uv}{\sigma_U\sigma_V} + \frac{v^2}{\sigma_V^2} \right) \right]$$

which is in the exact form of a bivariate normal distribution. Thus, by part a),  $U$  is normal.

#### Example exercise 4.46

4.46 ( A derivation of the bivariate normal distribution) Let  $Z_1$  and  $Z_2$  be independent  $n(0, 1)$  random variables, and define new random variables  $X$  and  $Y$  by

$$X = a_X Z_1 + b_X Z_2 + c_X \text{ and } Y = a_Y Z_1 + b_Y Z_2 + c_Y$$

where  $a_X, b_X, c_X, a_Y, b_Y$ , and  $c_Y$  are constants. (a) Show that

$$\begin{aligned} EX &= c_X, & \text{Var } X &= a_X^2 + b_X^2, \\ EY &= c_Y, & \text{Var } Y &= a_Y^2 + b_Y^2, \\ \text{Cov}(X, Y) &= a_X a_Y + b_X b_Y. \end{aligned}$$

(b) If we define the constants  $a_X, b_X, c_X, a_Y, b_Y$ , and  $c_Y$  by

$$\begin{aligned} a_X &= \sqrt{\frac{1+\rho}{2}}\sigma_X, & b_X &= \sqrt{\frac{1-\rho}{2}}\sigma_X, & c_X &= \mu_X, \\ a_Y &= \sqrt{\frac{1+\rho}{2}}\sigma_Y, & b_Y &= -\sqrt{\frac{1-\rho}{2}}\sigma_Y, & c_Y &= \mu_Y, \end{aligned}$$

where  $\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2$ , and  $\rho$  are constants,  $-1 \leq \rho \leq 1$ , then show that

$$\begin{aligned} EX &= \mu_X, & \text{Var } X &= \sigma_X^2, \\ EY &= \mu_Y, & \text{Var } Y &= \sigma_Y^2, \\ \rho_{XY} &= \rho. \end{aligned}$$

(c) Show that  $(X, Y)$  has the bivariate normal pdf with parameters  $\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2$ , and  $\rho$ .

Let  $D = a_X b_Y - a_Y b_X = -\sqrt{1-\rho^2}\sigma_X\sigma_Y$  and solve for  $Z_1$  and  $Z_2$ ,

$$\begin{aligned} Z_1 &= \frac{b_Y(X - c_X) - b_X(Y - c_Y)}{D} = \frac{\sigma_Y(X - \mu_X) + \sigma_X(Y - \mu_Y)}{\sqrt{2(1+\rho)}\sigma_X\sigma_Y} \\ Z_2 &= \frac{\sigma_Y(X - \mu_X) + \sigma_X(Y - \mu_Y)}{\sqrt{2(1-\rho)}\sigma_X\sigma_Y}. \end{aligned}$$

Then the Jacobian is

$$J = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_1}{\partial y} \\ \frac{\partial z_2}{\partial x} & \frac{\partial z_2}{\partial y} \end{pmatrix} = \begin{pmatrix} \frac{b_Y}{D} & \frac{-b_X}{D} \\ \frac{-a_Y}{D} & \frac{a_X}{D} \end{pmatrix} = \frac{a_X b_Y}{D^2} - \frac{a_Y b_X}{D^2} = \frac{1}{D} = \frac{1}{-\sqrt{1-\rho^2}\sigma_X\sigma_Y}$$

and we have that

$$\begin{aligned} f_{X,Y}(x, y) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(\sigma_Y(x-\mu_X) + \sigma_X(y-\mu_Y))^2}{2(1+\rho)\sigma_X^2\sigma_Y^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(\sigma_Y(x-\mu_X) + \sigma_X(y-\mu_Y))^2}{2(1-\rho)\sigma_X^2\sigma_Y^2}} \frac{1}{\sqrt{1-\rho^2}\sigma_X\sigma_Y} \\ &= (2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2})^{-1} \exp\left(-\frac{1}{2(1-\rho^2)}\left(\frac{x-\mu_X}{\sigma_X}\right)^2\right) \\ &\quad - 2\rho \frac{x-\mu_X}{\sigma_X} \left(\frac{y-\mu_Y}{\sigma_Y}\right) + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2, \quad -\infty < x < \infty, -\infty < y < \infty, \end{aligned}$$

a bivariate normal pdf.

## Multinomial

- [Amber's Homework 5](#)

## **Transformation/Location Scale Family**

- [Amber's Homework 8](#)
- [Amber's Homework 7](#)

## **Univariate Transformation**

### **Discrete Random Variable**

No jacobian

### **Continuous Random Variable**

## **Location Scale Family**



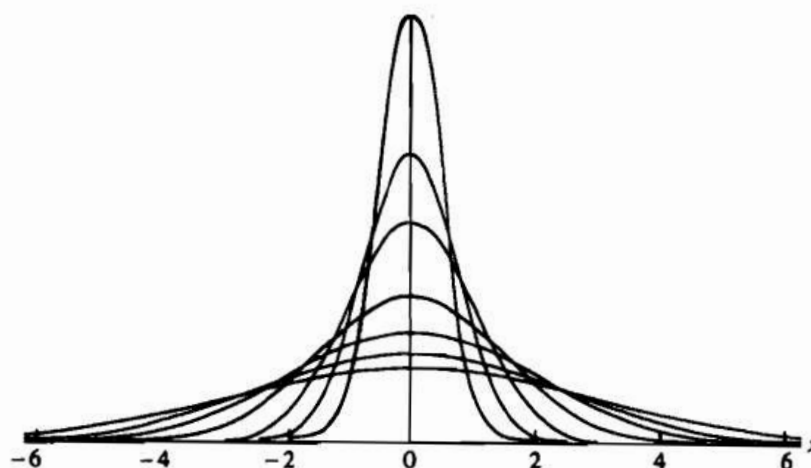


Figure 3.5.3. Members of the same scale family

The other two types of families to be discussed in this section are scale families and location-scale families.

**Definition 3.5.4** Let  $f(x)$  be any pdf. Then for any  $\sigma > 0$ , the family of pdfs  $(1/\sigma)f(x/\sigma)$ , indexed by the parameter  $\sigma$ , is called the *scale family with standard pdf  $f(x)$*  and  $\sigma$  is called the *scale parameter* of the family.

The effect of introducing the scale parameter  $\sigma$  is either to stretch ( $\sigma > 1$ ) or to contract ( $\sigma < 1$ ) the graph of  $f(x)$  while still maintaining the same basic shape of the graph. This is illustrated in Figure 3.5.3. Most often when scale parameters are used,  $f(x)$  is either symmetric about 0 or positive only for  $x > 0$ . In these cases the stretching is either symmetric about 0 or only in the positive direction. But, in the definition, any pdf may be used as the standard.

Several of the families introduced in Section 3.3 either are scale families or have scale families as subfamilies. These are the gamma family if  $\alpha$  is a fixed value and  $\beta$  is the scale parameter, the normal family if  $\mu = 0$  and  $\sigma$  is the scale parameter, the exponential family, and the double exponential family if  $\mu = 0$  and  $\sigma$  is the scale parameter. In each case the standard pdf is the pdf obtained by setting the scale parameter equal to 1. Then all other members of the family can be shown to be of the form in Definition 3.5.4.

**Definition 3.5.5** Let  $f(x)$  be any pdf. Then for any  $\mu$ ,  $-\infty < \mu < \infty$ , and any  $\sigma > 0$ , the family of pdfs  $(1/\sigma)f((x - \mu)/\sigma)$ , indexed by the parameter  $(\mu, \sigma)$ , is called the *location-scale family with standard pdf  $f(x)$* ;  $\mu$  is called the *location parameter* and  $\sigma$  is called the *scale parameter*.

The effect of introducing both the location and scale parameters is to stretch ( $\sigma > 1$ ) or contract ( $\sigma < 1$ ) the graph with the scale parameter and then shift the graph so that the point that was above 0 is now above  $\mu$ . Figure 3.5.4 illustrates this transformation of  $f(x)$ . The normal and double exponential families are examples of location-scale families. Exercise 3.39 presents the Cauchy as a location-scale family.

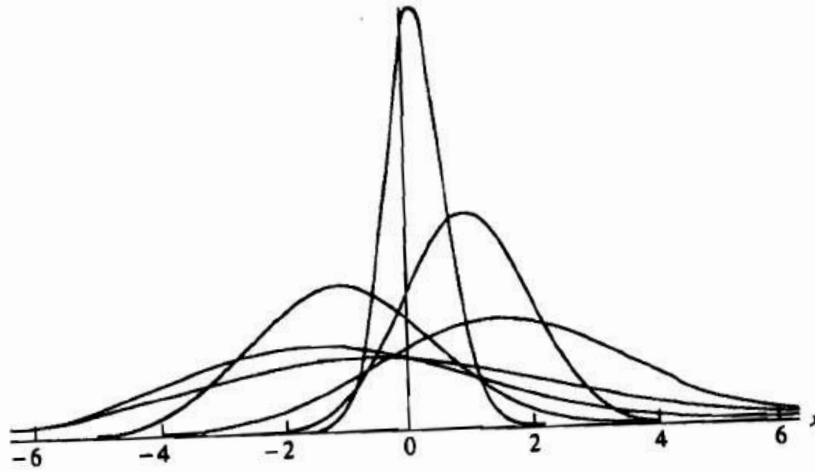


Figure 3.5.4. Members of the same location-scale family

The following theorem relates the transformation of the pdf  $f(x)$  that defines a location-scale family to the transformation of a random variable  $Z$  with pdf  $f(z)$ . As mentioned earlier in the discussion of location families, the representation in terms of  $Z$  is a useful mathematical tool and can help us understand when a location-scale family might be appropriate in a modeling context. Setting  $\sigma = 1$  in Theorem 3.5.6 yields a result for location (only) families, and setting  $\mu = 0$  yields a result for scale (only) families.

**Theorem 3.5.6** *Let  $f(\cdot)$  be any pdf. Let  $\mu$  be any real number, and let  $\sigma$  be any positive real number. Then  $X$  is a random variable with pdf  $(1/\sigma)f((x - \mu)/\sigma)$  if and only if there exists a random variable  $Z$  with pdf  $f(z)$  and  $X = \sigma Z + \mu$ .*

**Proof:** To prove the “if” part, define  $g(z) = \sigma z + \mu$ . Then  $X = g(Z)$ ,  $g$  is a monotone function,  $g^{-1}(x) = (x - \mu)/\sigma$ , and  $|(d/dx)g^{-1}(x)| = 1/\sigma$ . Thus by Theorem 2.1.5, the pdf of  $X$  is

$$f_X(x) = f_Z(g^{-1}(x)) \left| \frac{d}{dx} g^{-1}(x) \right| = f\left(\frac{x - \mu}{\sigma}\right) \frac{1}{\sigma}.$$

To prove the “only if” part, define  $g(x) = (x - \mu)/\sigma$  and let  $Z = g(X)$ . Theorem 2.1.5 again applies:  $g^{-1}(z) = \sigma z + \mu$ ,  $|(d/dz)g^{-1}(z)| = \sigma$ , and the pdf of  $Z$  is

$$f_Z(z) = f_X(g^{-1}(z)) \left| \frac{d}{dz} g^{-1}(z) \right| = \frac{1}{\sigma} f\left(\frac{(\sigma z + \mu) - \mu}{\sigma}\right) \sigma = f(z).$$

Also,

$$\sigma Z + \mu = \sigma g(X) + \mu = \sigma \left( \frac{X - \mu}{\sigma} \right) + \mu = X. \quad \square$$



An important fact to extract from Theorem 3.5.6 is that the random variable  $Z = (X - \mu)/\sigma$  has pdf

$$f_Z(z) = \frac{1}{1} f\left(\frac{z - 0}{1}\right) = f(z).$$

That is, the distribution of  $Z$  is that member of the location-scale family corresponding to  $\mu = 0, \sigma = 1$ . This was already proved for the special case of the normal family in Section 3.3.

Often, calculations can be carried out for the “standard” random variable  $Z$  with pdf  $f(z)$  and then the corresponding result for the random variable  $X$  with pdf  $(1/\sigma)f((x - \mu)/\sigma)$  can be easily derived. An example is given in the following, which is a generalization of a computation done in Section 3.3 for the normal family.

**Theorem 3.5.7** *Let  $Z$  be a random variable with pdf  $f(z)$ . Suppose  $EZ$  and  $\text{Var } Z$  exist. If  $X$  is a random variable with pdf  $(1/\sigma)f((x - \mu)/\sigma)$ , then*

$$EX = \sigma EZ + \mu \quad \text{and} \quad \text{Var } X = \sigma^2 \text{Var } Z.$$

*In particular, if  $EZ = 0$  and  $\text{Var } Z = 1$ , then  $EX = \mu$  and  $\text{Var } X = \sigma^2$ .*

**Proof:** By Theorem 3.5.6, there is a random variable  $Z^*$  with pdf  $f(z)$  and  $X = \sigma Z^* + \mu$ . So  $EX = \sigma EZ^* + \mu = \sigma EZ + \mu$  and  $\text{Var } X = \sigma^2 \text{Var } Z^* = \sigma^2 \text{Var } Z$ .  $\square$

For any location-scale family with a finite mean and variance, the standard pdf  $f(z)$  can be chosen in such a way that  $EZ = 0$  and  $\text{Var } Z = 1$ . (The proof that this choice can be made is left as Exercise 3.40.) This results in the convenient interpretation of  $\mu$  and  $\sigma^2$  as the mean and variance of  $X$ , respectively. This is the case for the usual definition of the normal family as given in Section 3.3. However, this is not the choice for the usual definition of the double exponential family as given in Section 3.3. There,  $\text{Var } Z = 2$ .

Probabilities for any member of a location-scale family may be computed in terms of the standard variable  $Z$  because

$$P(X \leq x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = P\left(Z \leq \frac{x - \mu}{\sigma}\right).$$

Thus, if  $P(Z \leq z)$  is tabulated or easily calculable for the standard variable  $Z$ , then probabilities for  $X$  may be obtained. Calculations of normal probabilities using the standard normal table are examples of this.

### 3.6 Inequalities and Identities

Statistical theory is literally brimming with inequalities and identities—so many that entire books are devoted to the topic. The major work by Marshall and Olkin (1979) contains many inequalities using the concept of majorization. The older work by Hardy, Littlewood, and Polya (1952) is a compendium of classic inequalities. In this section and in Section 4.7 we will mix some old and some new, giving some idea of the

## Multivariate Transformation

Year	Question	How it goes
2021 T	2	Bad, very important

---

### Exercise CB 5.17

### Homework 4.16

(a)  $X \sim \text{Geometric}(p), Y \sim \text{Geometric}(p), Y \perp X$ , Show that  $U = \min(X, Y), V = X - Y$  is independent:

$$f_{X,Y} = p(1-p)^{x-1}p(1-p)^{y-1}I_{1,2,\dots}(x)I_{1,2,\dots}(y)$$

- when  $x < y$

$$\begin{cases} u = x \\ v = x - y \end{cases} \Rightarrow \begin{cases} x = u \\ y = u - v, v < 0 \text{ b/c } x < y \end{cases}$$

$$\begin{cases} u = x \in \{1, 2, \dots\} \\ v = x - y \in \{-1, -2, \dots\} \end{cases}$$

### IT'S DISCRETE, NO JACOBIAN

$$\begin{aligned} P(U = u, V = v) &= p(1-p)^{u-1}p(1-p)^{u-v-1}I_{1,2,\dots}(u)I_{-1,-2,\dots}(v) \\ &= p^2(1-p)^{2u-1}I_{1,2,\dots}(u)(1-p)^{-v-1}I_{-1,-2,\dots}(v) \end{aligned}$$

- when  $x > y$

$$\begin{aligned} P(U = u, V = v) &= p(1-p)^{u-1}p(1-p)^{u-v-1}I_{1,2,\dots}(u)I_{-1,-2,\dots}(v) \\ &= p^2(1-p)^{2u-1}I_{1,2,\dots}(u)(1-p)^{-v-1}I_{1,2,\dots}(v) \end{aligned}$$

- when  $x = y$

$$u = x = y, v = 0$$

$$\begin{aligned} P(U = u, V = v) &= p(1-p)^{u-1}p(1-p)^{u-v-1}I_{1,2,\dots}(u)I_{-1,-2,\dots}(v) \\ &= p^2(1-p)^{2u-2}I_{1,2,\dots}(u)I_{-1,-2,\dots}(v) \end{aligned}$$

see that indicator function is necessary to show independent, don't be lazy

So  $U \perp V$  by factorization.

(b) Find the distribution of  $Z = X/(X + Y)$ , where we define  $Z = 0$  if  $X + Y = 0$ .

Set:

$$\begin{aligned} \begin{cases} V = X \\ Z = X/(X + Y) \end{cases} &\Rightarrow \begin{cases} X = V \\ ZX + ZY = X \end{cases} \\ \Rightarrow \begin{cases} V = X \\ Y = \frac{X-ZX}{Z} \end{cases} &\Rightarrow \begin{cases} V = X \\ Y = V(1 - Z)/Z \end{cases} \end{aligned}$$

$$V \in \{1, 2, \dots\}, Z \in \{0, 1\}$$

$$\begin{aligned} P_{V,Z} &= p(1-p)^{v-1}p(1-p)^{v(1-z)/z-1}I_{1,\dots}(v)I_{(0,1)}(z) \\ &\Rightarrow P_Z = \sum_v^{\infty} P_{V=v,Z} \\ &= \frac{p^2}{1-p} \sum_v (1-p)^{v/Z} \\ &= \frac{p^2}{1-p} \sum_v [(1-p)^{1/Z}]^v \end{aligned}$$

Now, let  $Y = V - 1$ , such that  $Y \in \{0, 1, \dots\}$

By  $\sum_{v=0}^{\infty} r^k = \frac{1}{1-r}$ :

$$\begin{aligned} P_Z &= \frac{p^2}{1-p} (1-p)^{1/Z} \sum_v [(1-p)^{1/Z}]^{v-1} \\ &= \frac{p^2}{1-p} (1-p)^{1/Z} \sum_v [(1-p)^{1/Z}]^y \\ &= \frac{p^2}{1-p} (1-p)^{1/Z} \frac{1}{1 - (1-p)^{1/Z}} I_{(0,1)}(z) \end{aligned}$$

## Homework7 4.19

### Order Statistics

If the  $k$ th order statistics is less than  $X$ , then it means that at least  $k-1$  observations are less than  $X$ .

Based on this we can have the cdf for the  $k$ th order statistics:

### Discrete Ordered Statistics

Given a random sample  $X_1, X_2, \dots, X_n$  from a discrete distribution with pmf  $f(x)$ ,

The probability that there are exact  $k$  observations less than  $x$  is:

$$P(X_{(k)} = x) = \binom{n}{k} P(X_1 < x)^k (1 - P(X_1 < x))^{n-k}$$

The probability that at least  $k$  observations are less than  $x$  is:

$$P(X_{(k)} \leq x) = \sum_{i=k}^n \binom{n}{i} F_X(x)^i (1 - F_X(x))^{n-i}$$

### Continuous Ordered Statistics

Given a random sample  $X_1, X_2, \dots, X_n$  from a continuous distribution with cdf  $F(x)$ ,

the pdf of the  $k$ th order statistics is:

$$f_{X_{(k)}}(x) = \frac{n!}{(k-1)!(n-k)!} F_X(x)^{k-1} (1 - F_X(x))^{n-k} f_X(x)$$

**Just draw a line and you can write down the joint pdf**

### Joint Ordered Statistics

Theorem 5.4.6 Let  $X_{(1)}, \dots, X_{(n)}$  denote the order statistics of a random sample,  $X_1, \dots, X_n$ , from a continuous population with cdf  $F_X(x)$  and pdf  $f_X(x)$ . Then the joint pdf of  $X_{(i)}$  and  $X_{(j)}$ ,  $1 \leq i < j \leq n$ , is

$$(5.4.7) \quad f_{X_{(i)}, X_{(j)}}(u, v) = \frac{n!}{(i-1)!(j-1-i)!(n-j)!} f_X(u) f_X(v) [F_X(u)]^{i-1} \\ \times [F_X(v) - F_X(u)]^{j-1-i} [1 - F_X(v)]^{n-j}$$

for  $-\infty < u < v < \infty$

[Some Simple Examples](#)

For general  $n$ , we have

$$\boxed{f_{X_{(1)}, X_{(2)}, \dots, X_{(n)}}(x_1, x_2, \dots, x_n)} \quad \begin{aligned} & \text{“=”} \quad P(X_{(1)} = x_1, X_{(2)} = x_2, \dots, X_{(n)} = x_n) \\ & = \quad \boxed{n! f(x_1) f(x_2) \cdots f(x_n)} \end{aligned}$$

which holds for  $\boxed{x_1 < x_2 < \cdots < x_n}$  with all  $x_i$  in the support for the original distribution. The joint pdf is zero otherwise.

### The Formalities:

The joint cdf,

$$P(X_{(1)} \leq x_1, X_{(2)} \leq x_2, \dots, X_{(n)} \leq x_n),$$

is a little hard to work with. Instead, we consider something similar:

$$P(y_1 < X_{(1)} \leq x_1, y_2 < X_{(2)} \leq x_2, \dots, y_n < X_{(n)} < x_n)$$

for values  $y_1 < x_1 \leq y_2 < x_2 \leq y_3 < x_3 \leq \cdots \leq y_n < x_n$ .

This can happen if

$$y_1 < X_1 \leq x_1, \quad y_2 < X_2 \leq x_2, \quad \dots, \quad y_n < X_n < x_n,$$

or if

$$y_1 < X_5 \leq x_1, \quad y_2 < X_3 \leq x_2, \quad \dots, \quad y_n < X_{n-2} < x_n,$$

or...

Because of the constraints on the  $x_i$  and  $y_i$ , these are disjoint events. So, we can add these  $n!$  probabilities, which will all be the same, together to get

$$P(y_1 < X_{(1)} \leq x_1, \dots, y_n < X_{(n)} < x_n) = n! P(y_1 < X_1 \leq x_1, \dots, y_n < X_n < x_n).$$

Note that

$$P(y_1 < X_1 \leq x_1, \dots, y_n < X_n < x_n) \stackrel{\text{indep}}{=} \prod_{i=1}^n P(y_i < X_i \leq x_i) = \prod_{i=1}^n [F(x_i) - F(y_i)].$$

So,

$$P(y_1 < X_{(1)} \leq x_1, \dots, y_n < X_{(n)} < x_n) = n! \prod_{i=1}^n [F(x_i) - F(y_i)] \quad (2)$$

The left-hand side is

$$\int_{y_n}^{x_n} \int_{y_{n-1}}^{x_{n-1}} \cdots \int_{y_1}^{x_1} f_{X_{(1)}, X_{(2)}, \dots, X_{(n)}}(u_1, u_2, \dots, u_n) du_1 du_2 \cdots du_n.$$

Taking derivatives  $\frac{d}{dx_1} \frac{d}{dx_2} \cdots \frac{d}{dx_n}$  gives

$$f_{X_{(1)}, X_{(2)}, \dots, X_{(n)}}(x_1, x_2, \dots, x_n)$$

Differentiating both sides of (2) with respect to  $x_1, x_2, \dots, x_n$  gives us

$$f_{X_{(1)}, X_{(2)}, \dots, X_{(n)}}(x_1, x_2, \dots, x_n) = n! f(x_1) f(x_2) \cdots f(x_n)$$

which holds for  $x_1 < x_2 < \cdots, x_n$  and all  $x_i$  in the support of the original distribution. The pdf is zero otherwise.

### Exercise CB 5.24

Show that  $X_{(1)}/X_{(n)}$  and  $X_{(n)}$  are independent, where  $f_X(x) = 1/\theta$

Proof: Use  $f_X(x) = 1/\theta, F_X(x) = x/\theta, 0 < x < \theta$ . Let  $Y = X_{(n)}, Z = X_{(1)}$ . Then, from Theorem 5.4.6,

$$f_{Z,Y}(z,y) = \frac{n!}{0!(n-2)!0!} \frac{1}{\theta} \frac{1}{\theta} \left(\frac{z}{\theta}\right)^0 \left(\frac{y-z}{\theta}\right)^{n-2} \left(1 - \frac{y}{\theta}\right)^0 = \frac{n(n-1)}{\theta^n} (y-z)^{n-2}, 0 < z < y < \theta$$

Solutions Manual for Statistical Inference Now let  $W = Z/Y, Q = Y$ . Then  $Y = Q, Z = WQ$ , and  $|J| = q$ . Therefore

$$f_{W,Q}(w,q) = \frac{n(n-1)}{\theta^n} (q - wq)^{n-2} q = \frac{n(n-1)}{\theta^n} (1-w)^{n-2} q^{n-1}, 0 < w < 1, 0 < q < \theta.$$

The joint pdf factors into functions of  $w$  and  $q$ , and, hence,  $W$  and  $Q$  are independent.

### Exercise CB 5.27

### Exercise CB 5.42

### Exercise CB 7.49

Example: 2021-T 3a-b

## Asymptotic Probability

### Amber's Homework9

### Necessary limits

#### Exponential

$$\lim_{n \rightarrow \infty} \left(1 + \frac{a}{n}\right)^n = e^a$$

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

### Necessary Inequality

#### Markov's Inequality

$$P(X \geq s) \leq E(X)/s$$



### Chebychev's Inequality

Let  $X$  be a random variable and let  $g(z)$  be a nonnegative function. Then, for any  $r > 0$ ,

$$P(g(X) \geq r) \leq \frac{Eg(X)}{r}.$$

Proof:

$$\begin{aligned} Eg(X) &= \int_{-\infty}^{\infty} g(x) f_X(x) dx \\ &\geq \int_{\{x: g(x) \geq r\}} g(x) f_X(x) dx \quad (g \text{ is nonnegative}) \\ &\geq r \int_{\{x: g(x) \geq r\}} f_X(x) dx \\ &= rP(g(X) \geq r). \quad (\text{definition}) \end{aligned}$$

Example

### Holder's Inequality

Let  $X, Y$  be R.V and  $q > 1, p > 1, \frac{1}{q} + \frac{1}{p} = 1$

$$\begin{aligned} |E[XY]| &\leq E(|XY|) \leq [E(|X|^p)]^{\frac{1}{q}} [E(|Y|^q)]^{\frac{1}{p}} \\ \Rightarrow \sum |X_i Y_i| &\leq [(|X|^p)]^{\frac{1}{q}} [E(|Y|^q)]^{\frac{1}{p}} \end{aligned}$$

### Cauchy's Inequality

Extend from Holder,

$$|E(XY)| \leq E(|XY|) \leq E[(X^2)]^{\frac{1}{2}} E[(Y^2)]^{\frac{1}{2}}$$

Use this to show that  $|cov(X, Y)| \leq \sigma_x \sigma_y$

$$\begin{aligned} \Rightarrow |cov(X, Y)| &= |E[(X - \mu_x)(Y - \mu_y)]| \\ &\leq E[|(X - \mu_x)(Y - \mu_y)|] \\ &\leq E[(X - \mu_x)^2]^{\frac{1}{2}} * E[(Y - \mu_y)^2]^{\frac{1}{2}} \\ &= \sigma_x \sigma_y \end{aligned}$$

$$\Rightarrow |\rho(x, y)| \leq 1$$

### Jenson's Inequality

Theorem 4.7.7 (Jensen's Inequality)

For any random variable  $X$ , if  $g(x)$  is a convex function, then

$$Eg(X) \geq g(EX)$$

Equality holds if and only if, for every line  $a + bx$  that is tangent to  $g(x)$  at  $x = EX$ ,  $P(g(X) = a + bX) = 1$

Proof:

To establish the inequality, let  $l(x)$  be a tangent line to  $g(x)$  at the point  $g(EX)$ . (Recall that  $EX$  is a constant.) Write  $l(x) = a + bx$  for some  $a$  and  $b$ . The situation is illustrated in Figure 4.7.2.

Now, by the convexity of  $g$  we have  $g(x) \geq a + bx$ . Since expectations preserve inequalities,

$$\begin{aligned} Eg(X) &\geq E(a + bX) \\ &= a + bEX \quad \left( \begin{array}{l} \text{linearity of expectation,} \\ \text{Theorem 2.2.5} \end{array} \right) \\ &= l(EX) \quad \text{(definition of } l(x)) \\ &= g(EX), \quad \text{(} l \text{ is tangent at } EX) \end{aligned}$$

as was to be shown. If  $g(x)$  is linear, equality follows from properties of expectations (Theorem 2.2.5). For the "only if" part see Exercise 4.62 .

One immediate application of Jensen's Inequality shows that  $EX^2 \geq (EX)^2$ , since  $g(x) = x^2$  is convex. Also, if  $x$  is positive, then  $1/x$  is convex; hence  $E(1/X) \geq 1/EX$ , another useful application.

To check convexity of a twice differentiable function is quite easy.

The function  $g(x)$  is convex if  $g''(x) \geq 0$ , for all  $x$ , and  $g(x)$  is concave if  $g''(x) \leq 0$ , for all  $x$ .

Jensen's Inequality applies to concave functions as well. If  $g$  is concave, then  $Eg(X) \leq g(EX)$ .

### Boferroni's Inequality

$$P(A \cap B) \geq P(A) + P(B) - 1$$

## Convergence

### Converge in Probability

#### Definition:

A sequence of random variables,  $X_1, X_2, \dots$ , converges in probability to a random variable  $X$  if, for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0 \quad \text{or, equivalently,} \quad \lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1.$$

The  $X_1, X_2, \dots$  in Definition 5.5.1 (and the other definitions in this section) are typically not independent and identically distributed random variables, as in a random sample. The distribution of  $X_n$  changes as the subscript changes, and the convergence concepts discussed in this section describe different ways in which the distribution of  $X_n$  converges to some limiting distribution as the subscript becomes large.

#### Theorem 5.5.4

Suppose that  $X_1, X_2, \dots$  converges in probability to a random variable  $X$  and that  $h$  is a continuous function. Then  $h(X_1), h(X_2), \dots$  converges in probability to  $h(X)$ .

#### Exercise CB 5.32

#### Exercise CB 5.33

### Converge in Distribution

We have already encountered the idea of convergence in distribution in Chapter 2. Remember the properties of moment generating functions (mgfs) and how their convergence implies convergence in distribution (Theorem 2.3.12).

#### Definition 5.5.10

A sequence of random variables,  $X_1, X_2, \dots$ , converges in distribution to a random variable  $X$  if

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

at all points  $x$  where  $F_X(x)$  is continuous.

### Slutsky's Theorem

Theorem 5.5.17 (Slutsky's Theorem) If  $X_n \rightarrow X$  in distribution and  $Y \rightarrow a$ , where  $a$  is constant, in probability, then

$$Y_n X_n \rightarrow^d a + X \text{ in distribution}$$

$$Y_n + X_n \rightarrow^d a + X \text{ in distribution}$$

### Continuous Mapping Theorem

#### Example 5.5.11 (Maximum of uniforms)

If  $X_1, X_2, \dots$  are iid uniform(0,1) and  $X_{(n)} = \max_{1 \leq i \leq n} X_i$ , let us examine if (and to where)  $X_{(n)}$  converges in distribution. As  $n \rightarrow \infty$ , we expect  $X_{(n)}$  to get close to 1 and, as  $X_{(n)}$  must necessarily be less than 1, we have for any  $\varepsilon > 0$ ,

$$\begin{aligned} P(|X_{(n)} - 1| \geq \varepsilon) &= P(X_{(n)} \geq 1 + \varepsilon) + P(X_{(n)} \leq 1 - \varepsilon) \\ &= 0 + P(X_{(n)} \leq 1 - \varepsilon). \end{aligned}$$

Next using the fact that we have an iid sample, we can write

$$P(X_{(n)} \leq 1 - \varepsilon) = P(X_i \leq 1 - \varepsilon, i = 1, \dots, n) = (1 - \varepsilon)^n$$

which goes to 0. So we have proved that  $X_{(n)}$  converges to 1 in probability. However, if we take  $\varepsilon = t/n$ , we then have

$$P(X_{(n)} \leq 1 - t/n) = (1 - t/n)^n \rightarrow e^{-t}$$

which, upon rearranging yields

$$P(n(1 - X_{(n)}) \leq t) \rightarrow 1 - e^{-t}$$

that is, the random variable  $n(1 - X_{(n)})$  converges in distribution to an exponential(1) random variable.

#### Exercise CB 5.18

#### Exercise CB 5.23

**Converge Almost Surely (related to Strong Law of Large Number)****Definition 5.5.6**

A sequence of random variables,  $X_1, X_2, \dots$ , converges almost surely to a random variable  $X$  if, for every  $\epsilon > 0$ ,

$$P\left(\lim_{n \rightarrow \infty} |X_n - X| < \epsilon\right) = 1$$

Notice the similarity in the statements of Definitions 5.5.1 and 5.5.6. Although they look similar, they are very different statements, with Definition 5.5.6 much stronger. To understand almost sure convergence, we must recall the basic definition of a random variable as given in Definition 1.4.1. A random variable is a real-valued function defined on a sample space  $S$ . If a sample space  $S$  has elements denoted by  $s$ , then  $X_n(s)$  and  $X(s)$  are all functions defined on  $S$ . Definition 5.5.6 states that  $X_n$  converges to  $X$  almost surely if the functions  $X_n(s)$  converge to  $X(s)$  for all  $s \in S$  except perhaps for  $s \in N$ , where  $N \subset S$  and  $P(N) = 0$ .

**Example 5.5.7 (Almost sure convergence)**

Let the sample space  $S$  be the closed interval  $[0, 1]$  with the uniform probability distribution. Define random variables  $X_n(s) = s + s^n$  and  $X(s) = s$ . For every  $s \in [0, 1)$ ,  $s^n \rightarrow 0$  as  $n \rightarrow \infty$  and  $X_n(s) \rightarrow s = X(s)$ . However,  $X_n(1) = 2$  for every  $n$  so  $X_n(1)$  does not converge to  $1 = X(1)$ . But since the convergence occurs on the set  $[0, 1)$  and  $P([0, 1)) = 1$ ,  $X_n$  converges to  $X$  almost surely.

**Law of Large Number****Weak Law of Large Number****Definition**

Let  $X_1, X_2, \dots$  be iid random variables with  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 < \infty$ . Define  $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$ . Then, for every  $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \epsilon) = 1$$

that is,  $\bar{X}_n$  converges in probability to  $\mu$ .

Proof: The proof is quite simple, being a straightforward application of Chebychev's Inequality. We have, for every  $\epsilon > 0$ ,

$$P(|\bar{X}_n - \mu| \geq \epsilon) = P((\bar{X}_n - \mu)^2 \geq \epsilon^2) \leq \frac{E(\bar{X}_n - \mu)^2}{\epsilon^2} = \frac{\text{Var } \bar{X}_n}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}$$

Hence,  $P(|\bar{X}_n - \mu| < \epsilon) = 1 - P(|\bar{X}_n - \mu| \geq \epsilon) \geq 1 - \sigma^2/(n\epsilon^2) \rightarrow 1$ , as  $n \rightarrow \infty$ .

The Weak Law of Large Numbers (WLLN) quite elegantly states that, under general conditions, the sample mean approaches the population mean as  $n \rightarrow \infty$ . In fact, there are more general versions of the WLLN, where we need assume only that the mean is finite. However, the version stated in Theorem 5.5.2 is applicable in most practical situations.

The property summarized by the WLLN, that a sequence of the “same” sample quantity approaches a constant as  $n \rightarrow \infty$ , is known as consistency.

### Example 5.5.3 (Consistency of $S^2$ )

Suppose we have a sequence  $X_1, X_2, \dots$  of iid random variables with  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 < \infty$ . If we define

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

can we prove a WLLN for  $S_n^2$ ? Using Chebychev’s Inequality, we have

$$P(|S_n^2 - \sigma^2| \geq \epsilon) \leq \frac{E(S_n^2 - \sigma^2)^2}{\epsilon^2} = \frac{\text{Var } S_n^2}{\epsilon^2}$$

and thus, a sufficient condition that  $S_n^2$  converges in probability to  $\sigma^2$  is that  $\text{Var } S_n^2 \rightarrow 0$  as  $n \rightarrow \infty$ .

If  $S_n^2$  is a consistent estimator of  $\sigma^2$ , then by Theorem 5.5.4, the sample standard deviation  $S_n = \sqrt{S_n^2} = h(S_n^2)$  is a consistent estimator of  $\sigma$ . Note that  $S_n$  is, in fact, a biased estimator of  $\sigma$  (see Exercise 5.11), but the bias disappears asymptotically.

### Strong Law of Large Number

**Theorem 5.5.9 (Strong Law of Large Numbers)** Let  $X_1, X_2, \dots$  be iid random variables with  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 < \infty$ , and define  $X_n = (1/n) \sum_{i=1}^n X_i$ . Then, for every  $c > 0$ ,

$$P\left(\lim_{n \rightarrow \infty} |\hat{X}_n - \mu| < c\right) = 1$$

that is,  $X_n$  converges almost surely to  $\mu$ . For both the Weak and Strong Laws of Large Numbers we had the assumption of a finite variance. Although such an assumption is true (and desirable) in most applications, it is, in fact, a stronger assumption than is needed. Both the weak and strong laws hold without this assumption. The only moment condition needed is that  $E|X_i| < \infty$  (see Resnick 1999, Chapter 7, or Billingsley 1995, Section 22).

**Exercise 5.38**

The following extensions of the inequalities established in Exercise 3.45 are useful in establishing a SLLN (see Miscellanea 5.8.4). Let  $X_1, X_2, \dots, X_n$  be iid with mgf  $M_X(t)$ ,  $-h < t < h$ , and let  $S_n = \sum_{i=1}^n X_i$  and  $\bar{X}_n = S_n/n$  (a) Show that  $P(S_n > a) \leq e^{-at} [M_X(t)]^n$ , for  $0 < t < h$ , and  $P(S_n \leq a) \leq e^{-at} [M_X(t)]^n$ , for  $-h < t < 0$  (b) Use the facts that  $M_X(0) = 1$  and  $M'_X(0) = E(X)$  to show that, if  $E(X) < 0$ , then there is a  $0 < c < 1$  with  $P(S_n > a) \leq c^n$ . Establish a similar bound for  $P(S_n \leq a)$  (c) Define  $Y_i = X_i - \mu - \varepsilon$  and use the above argument, with  $a = 0$ , to establish that  $P(\bar{X}_n - \mu > \varepsilon) \leq c^n$ . (d) Now define  $Y_i = -X_i + \mu - \varepsilon$ , establish an equality similar to part (c), and combine the two to get

$$P(|\bar{X}_n - \mu| > \varepsilon) \leq 2c^n \text{ for some } 0 < c < 1$$

HINT: Hint: For part (b) you may assume that  $a > 0$  and establish the upper bound on  $P(S_n > a)$ . Use the definition  $M_X(0) = \lim_{t \rightarrow 0} M_X(t) = M_X(0)$  - do to show that there exists  $\delta > 0$  such that  $M_X(\delta) < 1$ . Use this result to define a suitable  $c$  and apply part (a). You may stop here since the result for  $a > 0$  is all that is needed to complete parts (c) and (d). The trick to getting the upper bound on  $P(S_n < a)$  is to assume that  $E(X) > 0$ . For parts (c) and (d), apply (b) with  $a = 0$ .

**CLT****Theorem 5.5.14 (Central Limit Theorem)**

Let  $X_1, X_2, \dots$  be a sequence of iid random variables whose mgfs exist in a neighborhood of 0 (that is,  $M_{X_i}(t)$  exists for  $|t| < h$ , for some positive  $h$ ). Let  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 > 0$ . (Both  $\mu$  and  $\sigma^2$  are finite since the mgf exists.) Define  $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$ . Let  $G_n(x)$  denote the cdf of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$ . Then, for any  $x, -\infty < x < \infty$ ,

$$\lim_{n \rightarrow \infty} G_n(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy;$$

that is,  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$  has a limiting standard normal distribution.

**Proof of CLT**

*Thought process: since to prove converge in probability, use MGF and combine it with Taylor series*

Proof of Theorem 5.5.14: We will show that, for  $|t| < h$ , the mgf of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$  converges to  $e^{t^2/2}$ , the mgf of a  $N(0, 1)$  random variable.

Define  $Y_i = (X_i - \mu)/\sigma$ , and let  $M_Y(t)$  denote the common mgf of the  $Y_i$  s, which exists for  $|t| < \sigma h$  and is given by Theorem 2.3.15. Since

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$$

we have, from the properties of mgfs (see Theorems 2.3.15 and 4.6.7),

$$\begin{aligned} M_{\sqrt{n}(\hat{X}_n - \mu)/\sigma}(t) &= M_{\Sigma_{i=1}^n} Y_i/\sqrt{n}(t) \\ &= M_{\Sigma_{i=1}^n} Y_i \left( \frac{t}{\sqrt{n}} \right) \\ &= \left( M_Y \left( \frac{t}{\sqrt{n}} \right) \right)^n. \end{aligned}$$

We now expand  $M_Y(t/\sqrt{n})$  in a Taylor series (power series) around 0. (See Definition 5.5.20.) We have

$$M_Y \left( \frac{t}{\sqrt{n}} \right) = \sum_{k=0}^{\infty} M_Y^{(k)}(0) \frac{(t/\sqrt{n})^k}{k!}$$

where  $M_Y^{(k)}(0) = (d^k/dt^k) M_Y(t)|_{t=0}$ . Since the mgfs exist for  $|t| < h$ , the power series expansion is valid if  $t < \sqrt{n}\sigma h$ .

Using the facts that  $M_Y^{(0)} = 1$ ,  $M_Y^{(1)} = 0$ , and  $M_Y^{(2)} = 1$  (by construction, the mean and variance of  $Y$  are 0 and 1), we have

$$M_Y \left( \frac{t}{\sqrt{n}} \right) = 1 + \frac{(t/\sqrt{n})^2}{2!} + R_Y \left( \frac{t}{\sqrt{n}} \right),$$

238 PROPERTIES OF A RANDOM SAMPLE Section 5.6 where  $R_Y$  is the remainder term in the Taylor expansion,

$$R_Y \left( \frac{t}{\sqrt{n}} \right) = \sum_{k=3}^{\infty} M_Y^{(k)}(0) \frac{(t/\sqrt{n})^k}{k!}$$

An application of Taylor's Theorem (Theorem 5.5.21) shows that, for fixed  $t \neq 0$ , we have

$$\lim_{n \rightarrow \infty} \frac{R_Y(t/\sqrt{n})}{(t/\sqrt{n})^2} = 0$$

Since  $t$  is fixed, we also have

$$\lim_{n \rightarrow \infty} \frac{R_Y(t/\sqrt{n})}{(1/\sqrt{n})^2} = \lim_{n \rightarrow \infty} n R_Y \left( \frac{t}{\sqrt{n}} \right) = 0,$$

and (5.5.5) is also true at  $t = 0$  since  $R_Y(0/\sqrt{n}) = 0$ . Thus, for any fixed  $t$ , we can write

$$\begin{aligned} \lim_{n \rightarrow \infty} \left( M_Y \left( \frac{t}{\sqrt{n}} \right) \right)^n &= \lim_{n \rightarrow \infty} \left[ 1 + \frac{(t/\sqrt{n})^2}{2!} + R_Y \left( \frac{t}{\sqrt{n}} \right) \right]^n \\ &= \lim_{n \rightarrow \infty} \left[ 1 + \frac{1}{n} \left( \frac{t^2}{2} + n R_Y \left( \frac{t}{\sqrt{n}} \right) \right) \right]^n \\ &= e^{t^2/2} \end{aligned}$$

by an application of Lemma 2.3.14, where we set  $a_n = (t^2/2) + n R_Y(t/\sqrt{n})$ . (Note that (5.5.5) implies that  $a_n \rightarrow t^2/2$  as  $n \rightarrow \infty$ .) Since  $e^{t^2/2}$  is the mgf of the  $n(0, 1)$  distribution, the theorem is proved.



## Delta Methods

- Application: how to use delta method to get the CI of non-linear regression
- 

## Taylor Series

Definition 5.5.20 If a function  $g(x)$  has derivatives of order  $r$ , that is,  $g^{(r)}(x) = \frac{d^r}{dx^r}g(x)$  exists, then for any constant  $a$ , the Taylor polynomial of order  $r$  about  $a$  is

$$T_r(x) = \sum_{i=0}^r \frac{g^{(i)}(a)}{i!} (x-a)^i.$$

Taylor's major theorem, which we will not prove here, is that the remainder from the approximation,  $g(x) - T_r(x)$ , always tends to 0 faster than the highest-order explicit term.

Theorem 5.5.21 (Taylor) If  $g^{(r)}(a) = \frac{d^r}{dx^r}g(x)|_{x=a}$  exists, then

$$\lim_{x \rightarrow a} \frac{g(x) - T_r(x)}{(x-a)^r} = 0.$$

In general, we will not be concerned with the explicit form of the remainder. Since we are interested in approximations, we are just going to ignore the remainder. There are, however, many explicit forms, one useful one being

$$g(x) - T_r(x) = \int_a^x \frac{g^{(r+1)}(t)}{r!} (x-t)^r dt$$

For the statistical application of Taylor's Theorem, we are most concerned with the first-order Taylor series, that is, an approximation using just the first derivative (taking  $r = 1$  in the above formulas). Furthermore, we will also find use for a multivariate Taylor series. Since the above detail is univariate, some of the following will have to be accepted on faith.

Let  $T_1, \dots, T_k$  be random variables with means  $\theta_1, \dots, \theta_k$ , and define  $\mathbf{T} = (T_1, \dots, T_k)$  and  $\theta = (\theta_1, \dots, \theta_k)$ . Suppose there is a differentiable function  $g(\mathbf{T})$  (an estimator of some parameter) for which we want an approximate estimate of variance. Define

$$g'_i(\theta) = \left. \frac{\partial}{\partial t_i} g(\mathbf{t}) \right|_{t_1=\theta_1, \dots, t_k=\theta_k}$$

The first-order Taylor series expansion of  $g$  about  $\theta$  is

$$g(\mathbf{t}) = g(\theta) + \sum_{i=1}^k g'_i(\theta) (t_i - \theta_i) + \text{Remainder}.$$

For our statistical approximation we forget about the remainder and write

$$g(t) \approx g(\theta) + \sum_{i=1}^k g'_i(\theta) (t_i - \theta_i)$$

Now, take expectations on both sides of (5.5.7) to get (5.5.8)

$$\begin{aligned} E_{\theta} g(\mathbf{T}) &\approx g(\theta) + \sum_{i=1}^k g'_i(\theta) E_{\theta} (T_i - \theta_i) \\ &= g(\theta). \end{aligned}$$

( $T_i$  has mean  $\theta_i$ )

We can now approximate the variance of  $g(\mathbf{T})$  by

$$\begin{aligned} \text{Var}_{\theta} g(\mathbf{T}) &\approx E_{\theta} ([g(\mathbf{T}) - g(\theta)]^2) \\ &\approx E_{\theta} \left( \left( \sum_{i=1}^k g'_i(\theta) (T_i - \theta_i) \right)^2 \right) \quad (\text{using (5.5.8)}) \\ &= \sum_{i=1}^k [g'_i(\theta)]^2 \text{Var}_{\theta} T_i + 2 \sum_{i>j} g'_i(\theta) g'_j(\theta) \text{Cov}_{\theta} (T_i, T_j), \end{aligned}$$

### Delta Method for Univariate

**Example 5.5.23** (Approximate mean and variance) Suppose  $X$  is a random variable with  $E_{\mu} X = \mu \neq 0$ . If we want to estimate a function  $g(\mu)$ , a first-order approximation would give us

$$g(X) \approx g(\mu) + g'(\mu)(X - \mu)$$

If we use  $g(X)$  as an estimator of  $g(\mu)$ , we can say that approximately

$$\begin{aligned} E_{\mu} g(X) &\approx g(\mu), \\ \text{Var}_{\mu} g(X) &\approx [g'(\mu)]^2 \text{Var}_{\mu} X. \end{aligned}$$

For a specific example, take  $g(\mu) = 1/\mu$ . We estimate  $1/\mu$  with  $1/X$ , and we can say

$$\begin{aligned} E_{\mu} \left( \frac{1}{X} \right) &\approx \frac{1}{\mu}, \\ \text{Var}_{\mu} \left( \frac{1}{X} \right) &\approx \left( \frac{1}{\mu} \right)^4 \text{Var}_{\mu} X. \end{aligned}$$

Using these Taylor series approximations for the mean and variance, we get the following useful generalization of the Central Limit Theorem, known as the Delta Method

#### Theorem 5.5.24 (Delta Method)

Let  $Y_n$  be a sequence of random variables that satisfies  $\sqrt{n} (Y_n - \theta) \rightarrow \mathcal{N}(0, \sigma^2)$  in distribution. For a given function  $g$  and a specific value of  $\theta$ , suppose that  $g'(\theta)$  exists and is not

0 . Then (5.5.10)  $\sqrt{n} [g(Y_n) - g(\theta)] \rightarrow n(0, \sigma^2 [g'(\theta)]^2)$  in distribution. Proof: The Taylor expansion of  $g(Y_n)$  around  $Y_n = \theta$  is

$$g(Y_n) = g(\theta) + g'(\theta)(Y_n - \theta) + \text{Remainder},$$

where the remainder  $\rightarrow 0$  as  $Y_n \rightarrow \theta$ . Since  $Y_n \rightarrow \theta$  in probability it follows that the remainder  $\rightarrow 0$  in probability. By applying Slutsky's Theorem (Theorem 5.5.17) to

$$\sqrt{n} [g(Y_n) - g(\theta)] = g'(\theta) \sqrt{n} (Y_n - \theta),$$

the result now follows. See Exercise 5.43 for details.

Example 5.5.25 (Continuation of Example 5.5.23) Suppose now that we have the mean of a random sample  $\bar{X}$ . For  $\mu \neq 0$ , we have

$$\sqrt{n} \left( \frac{1}{\bar{X}} - \frac{1}{\mu} \right) \rightarrow n \left( 0, \left( \frac{1}{\mu} \right)^4 \text{Var}_{\mu} X_1 \right)$$

in distribution. If we do not know the variance of  $X_1$ , to use the above approximation requires an estimate, say  $S^2$ . Moreover, there is the question of what to do with the  $1/\mu$  term, as we also do not know  $\mu$ . We can estimate everything, which gives us the approximate variance

$$\widehat{\text{Var}} \left( \frac{1}{\bar{X}} \right) \approx \left( \frac{1}{\bar{X}} \right)^4 S^2.$$

Furthermore, as both  $\bar{X}$  and  $S^2$  are consistent estimators, we can again apply Slutsky's Theorem to conclude that for  $\mu \neq 0$ ,

$$\frac{\sqrt{n} \left( \frac{1}{\bar{X}} - \frac{1}{\mu} \right)}{\left( \frac{1}{\bar{X}} \right)^2 S} \rightarrow n(0, 1)$$

### Theorem: Second Order Delta Method

In some cases  $\frac{d}{d\mu} g(\mu) = 0$ , ie. score function, we need to use second order delta method:

$$\begin{aligned} g(\hat{\theta})|_{\theta_0=\mu} &= g(\mu) + \underbrace{g'(\mu)}_{=0}(\hat{\theta} - \mu) + g''(\mu) \frac{(\hat{\theta} - \mu)^2}{2!} + \dots \\ &\Rightarrow g(\hat{\theta}) - g(\mu) \approx \frac{g''(\mu)}{2} (\hat{\theta} - \mu)^2 \end{aligned}$$

By applying Slutsky's theorem, where  $\sqrt{n}(\hat{\theta} - \mu) \rightarrow N(0, \sigma_{\theta}^2)$  in distribution

$$\frac{2}{g''(\mu)} \frac{[g(\hat{\theta}) - g(\mu)]^2}{\sigma^2/n} \rightarrow N(0, 1)^2 \sim \chi_1^2$$

## Delta Method for Multivariate

Given Taylor series:

$$h(B) \approx h(\beta) + \nabla h(\beta)^T \cdot (B - \beta)$$

We can have

$$\begin{aligned} \text{Var}(h(B)) &\approx \text{Var}[h(\beta) + \nabla h(\beta)^T \cdot (B - \beta)] \\ &= \text{Var}(h(\beta) + \nabla h(\beta)^T \cdot B - \nabla h(\beta)^T \cdot \beta) \\ &= \text{Var}(\nabla h(\beta)^T \cdot B) \\ &= \nabla h(\beta)^T \cdot \text{cov}(B) \cdot \nabla h(\beta) \\ &= \nabla h(\beta)^T \cdot \frac{\Sigma}{n} \cdot \nabla h(\beta) \end{aligned}$$

Where in this case the scale parameter is  $\alpha = n$ ,

So in general case we have:

$$\sqrt{\alpha} (h(B) - h(\beta)) \xrightarrow{D} N(0, \nabla h(\beta)^T \cdot \Sigma \cdot \nabla h(\beta))$$

In regression case, we have scale parameter  $\alpha = X^T X$  ( $\text{Var}(\hat{\beta}) = X^T X^{-1} \Sigma$ )

## Statistics & Estimator

### Sufficient Statistics

#### Definition

### Factorization Theorem: How to find Sufficient Statistics

## 2.3 Factorization Theorem

**Theorem 2.1** Let  $f_{\mathbf{X}}(\mathbf{x} \mid \boldsymbol{\theta})$  denote the joint pdf or pmf of a sample  $\mathbf{X}$ . A statistic  $T(\mathbf{X})$  is a sufficient statistic for  $\boldsymbol{\theta}$  if and only if there exist functions  $g(T(\mathbf{x}) \mid \boldsymbol{\theta})$  and  $h(\mathbf{x})$ , where  $h(\mathbf{x})$  does not depend on  $\boldsymbol{\theta}$ , such that

$$f_{\mathbf{X}}(\mathbf{x} \mid \boldsymbol{\theta}) = h(\mathbf{x})g(T(\mathbf{x}) \mid \boldsymbol{\theta})$$

 Factorization  
Theorem

**Proof:**

- Assume  $T(X)$  is a sufficient statistics for  $\theta$  (and prove factorization)

$$\begin{aligned} f(x|\theta) &= P_{\theta}(X = x) \\ &= P(X = x \text{ and } T(X) = T(x)) \\ &= P(X = x | T(X) = T(x)) * P_{\theta}(T(X) = T(x)) \\ &= g(T(x)|\theta)h(x) \end{aligned}$$

- Assume  $f_x(x|\theta) = h(x)g(T(x)|\theta)$  (and prove sufficient)

**Given the marginal distribution of  $T(X)$  q:**

$$q(T(x)|\theta) = \sum_{y=A_{T(x)}} g(T(y)|\theta)h(y)$$

**where**  $A_{T(x)} = \{y : T(y) = T(x)\}$

$$\begin{aligned} \frac{f(x|\theta)}{q(T(x)|\theta)} &= \frac{g(T(x|\theta)h(x))}{\sum_{y=A_{T(x)}} g(T(y)|\theta)h(y)} \\ &= \frac{g(T(x|\theta)h(x))}{g(T(x)|\theta) \sum_{y=A_{T(x)}} h(y)} \\ &= \frac{h(x))}{\sum_{y=A_{T(x)}} h(y)} \leftarrow \text{independent of } \theta \end{aligned}$$

■

**Example 2.1** Given data from  $N(\mu, \sigma^2)$ , what if we want to estimate  $\mu$  only (i.e given  $\sigma$ )

$$f_X(\mathbf{x} \mid \mu, \sigma) = \underbrace{(2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{(n-1)s_X^2}{2\sigma^2}\right)}_{\text{Does not depend on } \mu} \underbrace{\exp\left(-\frac{n}{2\sigma^2}(\bar{x} - \mu)^2\right)}_{\text{Depends only } \bar{x} \text{ and } \mu}$$

So,  $\bar{X}$  is a sufficient statistic for  $\mu$  when  $\sigma$  is known.

If both  $\mu'$  s and  $\sigma'$  s values are unknown then the model's parameter is  $\theta = (\mu, \sigma)$ .  $\bar{X}$  does not contain all of the information - in the joint distribution - about  $\theta$ . However,  $(\bar{X}, S_V^2)$  is jointly sufficient for  $\theta$ .

■

**Example 2.2** Given above, what if we only want to estimate  $\sigma$ ?

$$f_X(\mathbf{x} \mid \mu, \sigma) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

$\sum_{i=1}^n (x_i - \mu)^2$  is sufficient for  $\sigma$  ■

■

**Lemma 2.2** If  $T_1(\mathbf{X}), T_2(\mathbf{X}), \dots, T_k(\mathbf{X})$  are a set of jointly sufficient statistics for  $\theta$ , then any set of **one-to-one functions, or transformations**, of  $(T_1, \dots, T_k)$  is also jointly sufficient for  $\theta$ .

## 2.4 Exponential Family

**Theorem 2.3** Theorem: Let  $X_1, \dots, X_n$  be iid observations from a pdf or pmf  $f_X(x \mid \theta)$  that belongs to an exponential family given by

$$f_X(x \mid \theta) = h(x)c(\theta) \exp\left(\sum_{i=1}^k w_i(\theta)t_i(x)\right)$$

where  $\theta = (\theta_1, \dots, \theta_d), d \leq k$ . Then

$$T(\mathbf{X}) = \left(\sum_{j=1}^n t_1(X_j), \dots, \sum_{j=1}^n t_k(X_j)\right)$$

is a (jointly) sufficient statistic for  $\theta$ .

**Example 2.3**

$$f_X(\mathbf{x} \mid \mu, \sigma) = \left((2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right)\right)^n \exp\left(-\frac{1}{2\sigma^2} \left(\sum x_i^2 - 2\mu \sum x_i\right)\right)$$

- So  $h(x) = 1, c(\mu, \sigma) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right)$ ,

-  $w_1(\mu, \sigma) = -\frac{1}{2\sigma^2}, w_2(\mu, \sigma) = \frac{\mu}{\sigma^2}, T_1(\mathbf{x}) = \sum x_i^2$ , and  $T_2(\mathbf{x}) = \sum x_i$ .

- Thus  $(T_1(\mathbf{X}), T_2(\mathbf{X})) = (\sum X_i^2, \sum X_i)$  is sufficient for  $\theta = (\mu, \sigma)$ .

- By Lemma 2.2  $(T_1(\mathbf{X}), T_2(\mathbf{X})) = \left(\frac{n}{n-1} \sum (X_i - \bar{X})^2, \sum \bar{X}_i\right)$  is also sufficient for  $\theta = (\mu, \sigma)$ .

■

## Minimum Sufficient Statistics and Complete Statistics

### MSS

- The theorem only show MSS, but does not help finding MSS
- Complete statistics help finding MSS
- [Bob's Homework 1](#)
- NEED TO SHOW IF AND ONLY IF

### Example: Bob's final 1

$$Ratio = \left(\frac{\theta}{1-\theta}\right)^{x-y}(1-\theta)^{n-n'} * constant(x, y)$$

- proof of 'if':  
when  $n = n'$  and  $x = y$ , all the terms are equal to 1, so a constant as a function of  $\theta$
- proof of 'only if':  
Assuming the ratio is a constant  $c$  where  $c > 0$  then take the log of the ratio we have:

$$a(\log\theta - \log(1-\theta)) + b\log(1-\theta) = \log(c)$$

where  $a = x - y$ , and  $b = n - n'$ .

Take the derivative with respect to  $\theta$  (we want to create equation here (and also get rid of  $\theta$ ))

$$\frac{a}{\theta(1-\theta)} - \frac{b}{1-\theta} = 0 \Rightarrow a = b\theta$$

This apply to all  $\theta$  so we have  $a = 0$  and  $b = 0$ , making  $x = y$  and  $n = n'$  ■

### 3 Minimum Sufficient Statistics

**Definition 3.1** A sufficient statistic  $T(X)$  is called a minimal sufficient statistic if, for any other sufficient statistic  $T'(X)$ ,  $T(X)$  is a function of  $T'(X)$ .

Some notes:

- $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$  if and only if  $T'(\mathbf{x}) = T'(\mathbf{y}) \Rightarrow T(\mathbf{x}) = T(\mathbf{y})$ .
- **The MSS is not unique** (ie.  $(\sum X, \sum X^2)$  and  $(\bar{x}, s^2)$  both MMS for Normal). However, the minimal sufficient partition is unique! An MSS will provide a partition that is as coarse as any other sufficient statistic.
- By coarse we mean few partition elements.

show  
proof

**Theorem 3.1** Let  $f(\mathbf{x} | \theta)$  be the pmf or pdf of the 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  $\frac{f(\mathbf{x}|\theta)}{f(\mathbf{y}|\theta)}$  is a constant as a function of  $\theta$  **if and only if**  $T(\mathbf{x}) = T(\mathbf{y})$ . Then  $T(\mathbf{X})$  is a minimal sufficient statistic for  $\theta$ .

MMS the-  
orem

Proof:

To simplify the proof, we assume  $f(\mathbf{x} | \theta) > 0$  for all  $\mathbf{x} \in \mathcal{X}$  and  $\theta$ . First we show that  $T(\mathbf{X})$  is a sufficient statistic. Let  $\mathcal{T} = \{t : t = T(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathcal{X}\}$  be the image of  $\mathcal{X}$  under  $T(\mathbf{x})$ . Define the partition sets induced by  $T(\mathbf{x})$  as  $A_t = \{\mathbf{x} : T(\mathbf{x}) = t\}$ . For each  $A_t$ , choose and fix one element  $\mathbf{x}_t \in A_t$ . For any  $\mathbf{x} \in \mathcal{X}$ ,  $\mathbf{x}_{T(\mathbf{x})}$  is the fixed element that is in the same set,  $A_t$ , as  $\mathbf{x}$ . Since  $\mathbf{x}$  and  $\mathbf{x}_{T(\mathbf{x})}$  are in the same set  $A_t$ ,  $T(\mathbf{x}) = T(\mathbf{x}_{T(\mathbf{x})})$  and, hence,  $f(\mathbf{x} | \theta) / f(\mathbf{x}_{T(\mathbf{x})} | \theta)$  is constant as a function of  $\theta$ . Thus, we can define a function on  $\mathcal{X}$  by  $h(\mathbf{x}) = f(\mathbf{x} | \theta) / f(\mathbf{x}_{T(\mathbf{x})} | \theta)$  and  $h$  does not depend on  $\theta$ . Define a function on  $\mathcal{T}$  by  $g(t | \theta) = f(\mathbf{x}_t | \theta)$ . Then it can be seen that

$$f(\mathbf{x} | \theta) = \frac{f(\mathbf{x}_{T(\mathbf{x})} | \theta) f(\mathbf{x} | \theta)}{f(\mathbf{x}_{T(\mathbf{x})} | \theta)} = g(T(\mathbf{x}) | \theta) h(\mathbf{x})$$

and, by the Factorization Theorem,  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$ .

Now to show that  $T(\mathbf{X})$  is minimal, let  $T'(\mathbf{X})$  be any other sufficient statistic. By the Factorization Theorem, there exist functions  $g'$  and  $h'$  such that  $f(\mathbf{x} | \theta) = g'(T'(\mathbf{x}) | \theta) h'(\mathbf{x})$ . Let  $\mathbf{x}$  and  $\mathbf{y}$  be any two sample points with  $T'(\mathbf{x}) = T'(\mathbf{y})$ . Then

$$\frac{f(\mathbf{x} | \theta)}{f(\mathbf{y} | \theta)} = \frac{g'(T'(\mathbf{x}) | \theta) h'(\mathbf{x})}{g'(T'(\mathbf{y}) | \theta) h'(\mathbf{y})} = \frac{h'(\mathbf{x})}{h'(\mathbf{y})}$$

Since this ratio does not depend on  $\theta$ , the assumptions of the theorem imply that  $T(\mathbf{x}) = T(\mathbf{y})$ . Thus,  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$  and  $T(\mathbf{x})$  is minimal.

**Example 3.1** Example 6.2.14 (Normal minimal sufficient statistic)



Let  $X_1, \dots, X_n$  be iid  $n(\mu, \sigma^2)$ , both  $\mu$  and  $\sigma^2$  unknown. Let  $\mathbf{x}$  and  $\mathbf{y}$  denote two sample points, and let  $(\bar{x}, s_x^2)$  and  $(\bar{y}, s_y^2)$  be the sample means and variances corresponding to the  $\mathbf{x}$  and  $\mathbf{y}$  samples, respectively. Then, using (6.2.5), we see that the ratio of densities is

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

This ratio will be constant as a function of  $\mu$  and  $\sigma^2$  if and only if  $\bar{x} = \bar{y}$  and  $s_x^2 = s_y^2$ . Thus, by Theorem 6.2.13,  $(\bar{X}, S^2)$  is a minimal sufficient statistic for  $(\mu, \sigma^2)$ .

**Lemma 3.2** If  $f(\mathbf{x}|\theta)$  involve indicator function, we can use  $f(\mathbf{x}|\theta) = h(\mathbf{x}, \mathbf{y})f(\mathbf{y}|\theta)$  instead of ratio and show that *if and only if* some function of  $h(\mathbf{x}, \mathbf{y}) = \text{constant}$  existed.

**Example 3.2** why collapsing the sufficient statistics will not be a new sufficient statistics? ■

Where the fuck is this?

**Example 3.3** Show that  $S = (-1)^{X_1}T$  is also a sufficient statistics

If you combine or collapse sufficient statistics, it is not guaranteed that the new statistic will satisfy the factorization theorem, in which case it would not be a sufficient statistic. In some cases, combining sufficient statistics can lead to loss of information about the parameter of interest. ■

**Example 3.4** Location Exponential: Show that  $X_{(1)}$  is minimum sufficient ■

**Example 3.5** Two expressions of minimal sufficient statistic for Uniform with location shift:

*Example 6.2.15 (Uniform minimal sufficient statistic)* Suppose  $X_1, \dots, X_n$  are iid uniform observations on the interval  $(\theta, \theta + 1)$ ,  $-\infty < \theta < \infty$ . Then the joint pdf of  $\mathbf{X}$  is

$$f(\mathbf{x} | \theta) = \begin{cases} 1 & \theta < x_i < \theta + 1, i = 1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

which can be written as

$$f(\mathbf{x} | \theta) = \begin{cases} 1 & \max_i x_i - 1 < \theta < \min_i x_i \\ 0 & \text{otherwise.} \end{cases}$$

Thus, for two sample points  $\mathbf{x}$  and  $\mathbf{y}$ , the numerator and denominator of the ratio  $f(\mathbf{x} | \theta) / f(\mathbf{y} | \theta)$  will be positive for the same values of  $\theta$  if and only if  $\min_i x_i = \min_i y_i$  and  $\max_i x_i = \max_i y_i$ . And, if the minima and maxima are equal, then the ratio is constant and, in fact, equals 1. Thus, letting  $X_{(1)} = \min_i X_i$  and  $X_{(n)} = \max_i X_i$ , we have that  $T(\mathbf{X}) = (X_{(1)}, X_{(n)})$  is a minimal sufficient statistic. This is a case in which the dimension of a minimal sufficient statistic does not match the dimension of the parameter. ■

**Lemma 3.3** *If there is a 1-1 relationship, then the function of minimal sufficient is also a minimal sufficient, aka. minimal sufficient statistics is not unique (partition on the other hands is unique)*

### 3.0.1 MSS for a family of distributions

**Theorem 3.4** *Let  $\mathcal{P}_0$  be a finite family of pdf's or pmf's,  $f_0, f_1, \dots, f_k$ , all having the same support. Then, the statistic*

$$T(\mathbf{X}) = \left( \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})}, \frac{f_2(\mathbf{X})}{f_0(\mathbf{X})}, \dots, \frac{f_k(\mathbf{X})}{f_0(\mathbf{X})} \right)$$

*is minimal sufficient for  $\mathcal{P}_0$ .*

**Proof:**

Sufficiency

- Let  $g_i(T(\mathbf{X})) = T_i(\mathbf{X}) = f_i(\mathbf{X})/f_0(\mathbf{X})$  and  $g_0(T(\mathbf{X})) = 1$ .
- $f_i(\mathbf{x}) = g_i(T(\mathbf{x}))f_0(\mathbf{x})$  for  $i = 1, 2, \dots, k$
- By the factorization theorem,  $T$  is sufficient for  $\mathcal{P}_0$ .

Minimal Sufficiency

- Suppose  $S(\mathbf{X})$  is another sufficient statistic.
- $f_i(\mathbf{x}) = \tilde{g}_i(S(\mathbf{x}))h(\mathbf{x})$  for  $i = 1, 2, \dots, k$
- $T_i(\mathbf{x}) = f_i(\mathbf{x})/f_0(\mathbf{x}) = \tilde{g}_i(S(\mathbf{x}))/\tilde{g}_0(S(\mathbf{x}))$  for  $i = 1, 2, \dots, k$
- $T(\mathbf{x})$  is a function of  $S(\mathbf{x})$ , thus  $T(\mathbf{X})$  is minimal sufficient for  $\mathcal{P}_0$ .

**Lemma 3.5** *Lemma 5.2: If  $\mathcal{P}$  is a family of distributions with common support,  $T$  is minimal sufficient for  $\mathcal{P}_0 \subset \mathcal{P}$ , and  $T$  is sufficient for  $\mathcal{P}$ , then  $T$  is minimal sufficient for  $\mathcal{P}$ .*

**Proof:**

If  $U$  is sufficient for  $\mathcal{P}$ , it is also sufficient for  $\mathcal{P}_0$ , and hence  $T$  is a function of  $U$ . This implies  $T$  is minimal sufficient for  $\mathcal{P}$ .

### Lemma 3.6 MSS in exponential families

Let  $X_1, \dots, X_n$  be iid observations from an exponential family:  $f(x | \boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp \left( \sum_{j=1}^k w_j(\boldsymbol{\theta})t_j(x) \right)$  where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k) \in \Theta$ . The vector of canonical parameters is given by  $\mathbf{w}_{\boldsymbol{\theta}'} = (w_1(\boldsymbol{\theta}), \dots, w_k(\boldsymbol{\theta}))$ .

- Recall that  $T(\mathbf{X}) = (\sum_{i=1}^n t_1(X_i), \dots, \sum_{i=1}^n t_k(X_i))' \in \mathbb{R}^k$  is a sufficient for  $\mathcal{P} = \{f(x | \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$  and note that  $\sum_{j=1}^k w_j(\boldsymbol{\theta}) t_j(x) = \mathbf{w}'_{\boldsymbol{\theta}} T(\mathbf{x})$ .
- Let  $\Theta_0 = \{\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k\} \subset \Theta$  and  $\mathcal{P}_0 = \{f(x | \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta_0\}$ .
- By Theorem 5.1,  $S(\mathbf{X}) = \left( \frac{f(\mathbf{X}|\boldsymbol{\theta}_1)}{f(\mathbf{X}|\boldsymbol{\theta}_0)}, \dots, \frac{f(\mathbf{X}|\boldsymbol{\theta}_k)}{f(\mathbf{X}|\boldsymbol{\theta}_0)} \right)'$  is MSS for  $\mathcal{P}_0$ , where

$$\frac{f(\mathbf{X} | \boldsymbol{\theta}_i)}{f(\mathbf{X} | \boldsymbol{\theta}_0)} = \exp((\mathbf{w}'_{\boldsymbol{\theta}_i} - \mathbf{w}'_{\boldsymbol{\theta}_0}) T(\mathbf{X}))$$

- Define vectors  $\mathbf{d}_i = \mathbf{w}_{\boldsymbol{\theta}_i} - \mathbf{w}_{\boldsymbol{\theta}_0} \in \mathbb{R}^k, i = 1, 2, \dots, k$  and the matrix  $\mathbf{D} = [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_k] \in \mathbb{R}^k \times \mathbb{R}^k$
- Thus

$$\log S(\mathbf{X}) = \mathbf{D}T(\mathbf{X})$$

- If the matrix  $\mathbf{D}$  is full rank, then there is a 1-1 relationship between  $S(\mathbf{X})$  and  $T(\mathbf{X})$  which means that  $T(\mathbf{X})$  is also MSS for  $\mathcal{P}_0$ .
- The full rank property means that the  $\mathbf{d}_i$  's are linearly independent.
- Since  $T(\mathbf{X})$  is a sufficient for  $\mathcal{P}$  and is MSS for  $\mathcal{P}_0$ , then by Lemma 5.2,  $T(\mathbf{X})$  is MSS for  $\mathcal{P}$ .

## **Ancillary Statistics**

## 4 Ancillary Statistics and Complete Statistics

### 4.1 Ancillary Statistic

**Definition 4.1** A statistic  $S(\mathbf{X})$  with distribution that does not depend on the parameter  $\theta$  is called an ancillary statistic.

More precisely, a statistic  $S(\mathbf{X})$  is ancillary for  $\Theta$  if its distribution is the same for all  $\theta \in \Theta$ . That is,  $P_{\theta}\{S(\mathbf{X}) \in A\}$  is constant for  $\theta \in \Theta$  for any set  $A$ .

**Example 4.1** Revisit to  $N(\mu, \sigma^2)$

- We have previously showed that  $(\bar{X}, S_X^2)$  is sufficient for estimating  $\theta = (\mu, \sigma)$ . (Actually, it is MSS.)

- Note that the distribution of  $S_X^2$  depends on  $\sigma$  but not on  $\mu$ .

-  $S_X^2$  is ancillary for  $\Theta_1 = \{(\mu, \sigma^2) : \sigma^2 = \sigma_0^2\}$ . - Here  $\bar{X}$  is MSS and need not be paired with  $S_X^2$  to be sufficient for  $\Theta_1$ .

-  $S_X^2$  is not ancillary for  $\Theta_2 = \{(\mu, \sigma^2) : \sigma^2 > 0\}$ . ■

Why

### **Complete Statistics**

- The binomial example is useful to understand and proof completeness
- Exponential family has this open set thing
- Use complete statistics to find MVUE by lehman-sheffey

## 4.2 Complete statistics and completeness

- - Suppose that  $X_1, X_2, \dots, X_n$  are iid from  $\mathcal{F}(\theta)$ , a family of distributions indexed by  $\theta$ .
- - Let  $T(\mathbf{X})$  be a statistic and  $u(T)$  a real-valued function of  $T$  so that  $E_\theta[u(T)] = \theta$ . That is,  $u(T)$  is unbiased for  $\theta$ . [We will cover unbiasedness in more detail later.]
- - Under what conditions is  $u(T)$  the only function of  $T$  which is unbiased?
- - Let  $u_1(T)$  and  $u_2(T)$  be unbiased for  $\theta$
- - Define  $g(T) = u_1(T) - u_2(T)$ . Then  $E_\theta[g(T)] = 0$  for all  $\theta$ .
- - If the only function  $g(T)$  that satisfies  $E_\theta[g(T)] = 0$  is  $g(T) = 0$ , then this implies  $u_1(T) = u_2(T)$  and  $u(T)$  is unique for all  $\theta$

**Definition 4.2** Let  $f(t | \theta)$  be a family of pdfs or pmfs for a statistic  $T(\mathbf{X})$ . The family of probability distributions is called complete if  $E_\theta[g(T)] = 0$  for all  $\theta$  implies  $P_\theta\{g(T) = 0\} = 1$ , that is,  $g(T) \equiv 0$ , for all  $\theta$ . Equivalently,  $T(\mathbf{X})$  is called a complete statistic.

**Example 4.2** 6.2.22 (Binomial complete sufficient statistic)

Suppose that  $T$  has a binomial( $n, p$ ) distribution,  $0 < p < 1$ . Let  $g$  be a function such that  $E[g(T)] = 0$ . Then

$$\begin{aligned} 0 = E_p g(T) &= \sum_{t=0}^n g(t) \binom{n}{t} p^t (1-p)^{n-t} \\ &= (1-p)^n \sum_{t=0}^n g(t) \binom{n}{t} \left(\frac{p}{1-p}\right)^t \end{aligned}$$

for all  $p, 0 < p < 1$ . The factor  $(1-p)^n$  is not 0 for any  $p$  in this range. Thus it must be that

$$0 = \sum_{t=0}^n g(t) \binom{n}{t} \left(\frac{p}{1-p}\right)^t = \sum_{t=0}^n g(t) \binom{n}{t} r^t$$

for all  $r, 0 < r < \infty$ . But the last expression is a polynomial of degree  $n$  in  $r$ , where the coefficient of  $r^t$  is  $g(t) \binom{n}{t}$ . For the polynomial to be 0 for all  $r$ , each coefficient must

be 0. Since none of the  $\binom{n}{t}$  terms is 0, this implies that  $g(t) = 0$  for  $t = 0, 1, \dots, n$ .

Since  $T$  takes on the values  $0, 1, \dots, n$  with probability 1, this yields that  $P_p(g(T) = 0) = 1$  for all  $p$ , the desired conclusion. Hence,  $T$  is a complete statistic.

■

make a polynomial for  $g(t)$

**Theorem 4.1** *If a minimal sufficient statistic exists, then any complete sufficient statistic is also a minimal sufficient statistic.*

**Proof**

- Let  $S$  be a minimal sufficient statistic and  $T$  be any complete sufficient statistic.
- Define  $g_1(S) = E[T \mid S]$ . This is actually a function of  $T$  since  $S$  is a MSS, say  $S = h(T)$ .
- Define  $g(T) = T - g_1(S) = T - g_1(h(T))$ .
- Clearly  $E[g(T)] = 0$  for all  $\theta$  since  $E[g_1(S)] = E[E[T \mid S]] = E[T]$ .
- Since  $T$  is complete, then  $g(T) = 0$  for all  $\theta$ .
- This implies that  $T = g_1(S)$ .
- Since  $S$  is minimally sufficient, then it is a function of every other sufficient statistic. That is, if  $S^*$  is a sufficient statistic, then  $S = g_*(S^*)$  for some function  $g_*(\cdot)$ .
- Hence  $T = g_1(g_*(S^*))$  is a function of  $S^*$ . Since  $S^*$  is any sufficient statistic, then  $T$  is minimally sufficient.

■

I don't understand this step

**Example 4.3** *Not all minimal sufficient statistics are complete*

Recall our earlier example where  $\mathbf{X}$  is a random sample from  $\text{Uniform}(\theta, \theta + 1)$ .

- $T(\mathbf{X}) = (X_{(1)}, X_{(n)})$  is minimally sufficient for  $\theta$ .
- $R(\mathbf{X}) = X_{(n)} - X_{(1)} \sim \text{Beta}(n - 1, 2)$  is ancillary.  $E[R] = \frac{(n-1)}{(n+1)}$ .
- Define  $g(T) = R - \frac{(n-1)}{(n+1)}$ .
- Then  $E[g(T)] = 0$  for all  $\theta$ , but  $g(T) \neq 0$  for all  $\theta$ . ■

A key feature of the above example is the existence of an ancillary statistic. Suppose  $A(\mathbf{X}) = g_1(T(\mathbf{X}))$  is ancillary and let  $E[A(\mathbf{X})] = a$ , a constant, independent of  $\theta$ . Define  $g(T) = g_1(T(\mathbf{X})) - a$ . Then  $E[g(T)] = 0$  for all  $\theta$ , but  $g(T) \neq 0$  for all  $\theta$ .

**Example 4.4** 6.2.23 (Uniform complete sufficient statistic) Let  $X_1, \dots, X_n$  be iid uniform  $(0, \theta)$  observations,  $0 < \theta < \infty$ . Using an argument similar to that in Example 6.2.8, we can see that  $T(\mathbf{X}) = \max_i X_i$  is a sufficient statistic and, by Theorem 5.4.4, the pdf of  $T(\mathbf{X})$  is

$$f(t \mid \theta) = \begin{cases} nt^{n-1}\theta^{-n} & 0 < t < \theta \\ 0 & \text{otherwise} \end{cases}$$



Suppose  $g(t)$  is a function satisfying  $E_\theta g(T) = 0$  for all  $\theta$ . Since  $E_\theta g(T)$  is constant as a function of  $\theta$ , its derivative with respect to  $\theta$  is 0. Thus we have that

$$\begin{aligned} 0 &= \frac{d}{d\theta} E_\theta g(T) = \frac{d}{d\theta} \int_0^\theta g(t) n t^{n-1} \theta^{-n} dt \\ &= (\theta^{-n}) \frac{d}{d\theta} \int_0^\theta n g(t) t^{n-1} dt + \left( \frac{d}{d\theta} \theta^{-n} \right) \int_0^\theta n g(t) t^{n-1} dt \\ &= \theta^{-n} n g(\theta) \theta^{n-1} + 0 \quad \left( \begin{array}{c} \text{applying the product} \\ \text{rule for differentiation} \end{array} \right) \\ &= \theta^{-1} n g(\theta) \end{aligned}$$

The first term in the next to last line is the result of an application of the Fundamental Theorem of Calculus. The second term is 0 because the integral is, except for a constant, equal to  $E_\theta g(T)$ , which is 0. Since  $\theta^{-1} n g(\theta) = 0$  and  $\theta^{-1} n \neq 0$ , it must be that  $g(\theta) = 0$ . This is true for every  $\theta > 0$ ; hence,  $T$  is a complete statistic. (On a somewhat pedantic note, realize that the Fundamental Theorem of Calculus does not apply to all functions, but only to functions that are Riemann-integrable. ■)

#### Theorem 4.2 Basu's theorem

If  $T(X)$  is a complete and minimal sufficient statistic, then  $T(X)$  is independent of every ancillary statistic.

Proof:

- Define  $g(t) = f_{S|T}(s | t) - f_S(s)$ .

$$\begin{aligned} E_\theta[g(T)] &= \int_{-\infty}^{\infty} (f_{S|T}(s | t) - f_S(s)) f_T(t) dt \\ &= \int_{-\infty}^{\infty} f_{S|T}(s | t) f_T(t) dt - \int_{-\infty}^{\infty} f_S(s) f_T(t) dt \\ &= f_S(s) - f_S(s) \\ &= 0 \end{aligned}$$

- Since  $T$  is complete, then  $g(t) = 0$  for all  $\theta$ .

- Hence  $f_{S|T}(s | t) = f_S(s)$  which implies that  $T$  and  $S$  are independent.

■

#### Theorem 4.3 (Complete statistics in the exponential family)

Let  $X_1, \dots, X_n$  be iid observations from an **exponential family** with pdf or pmf of the form

$$f(x | \boldsymbol{\theta}) = h(x) c(\boldsymbol{\theta}) \exp \left( \sum_{j=1}^k w(\theta_j) t_j(x) \right)$$

where  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ . Then the statistic

$$T(\mathbf{X}) = \left( \sum_{i=1}^n t_1(X_i), \sum_{i=1}^n t_2(X_i), \dots, \sum_{i=1}^n t_k(X_i) \right)$$

is complete as long as the parameter space  $\Theta$  contains an open set in  $\mathbb{R}^k$ .

The condition that the parameter space contain an open set is needed to avoid a situation like the following. The  $n(\theta, \theta^2)$  distribution can be written in the form (6.2.7); however, the parameter space  $(\theta, \theta^2)$  does not contain a two-dimensional open set, as it consists of only the points on a parabola. As a result, we can find a transformation of the statistic  $T(\mathbf{X})$  that is an unbiased estimator of 0 (see Exercise 6.15). (Recall that exponential families such as the  $n(\theta, \theta^2)$ , where the parameter space is a lower-dimensional curve, are called curved exponential families; see Section 3.4.) The relationship between sufficiency, completeness, and minimality in exponential families is an interesting one. For a brief introduction, see Miscellanea 6.6.3.

see HW 2

**Example 4.5** Suppose  $X_1, X_2$  are iid Exponential  $(\beta)$  and  $Y = X_1 + X_2$ .

- Exponential $(\beta)$  is an exponential family

$$f_X(x) = \frac{1}{\beta} e^{-\frac{x}{\beta}} \mathbf{1}_{(0, \infty)}(x)$$

where  $h(x) = \mathbf{1}_{(0, \infty)}(x)$ ,  $c(\beta) = \frac{1}{\beta}$ ,  $w(\beta) = -\frac{1}{\beta}$ , and  $t(x) = x$ .

- Note that  $w(\beta) = -\frac{1}{\beta}$ , for  $\beta > 0$ , contains an open set in  $\mathbb{R}$ .

- Hence  $Y = T(\mathbf{X}) = X_1 + X_2$  is a complete sufficient statistic.

Let  $g(\mathbf{X}) = \frac{X_2}{Y}$ . Find the value of  $E_\beta[g(\mathbf{X})]$ .

-  $g(\mathbf{X})$  is an ancillary statistic:  $\frac{X_i}{\beta} \sim \text{Exponential}(1)$ , Thus  $\frac{Y}{\beta} \sim \text{Gamma}(2, 1)$ , Neither distribution dependent on  $\beta$

the key is to find the distribution of the estimator

thus the distribution of

$$g(\mathbf{X}) = \frac{\left(\frac{X_2}{\beta}\right)}{\left(\frac{Y}{\beta}\right)} = \frac{X_2}{Y} \text{ does not depend on } \beta$$

- We know that  $\beta = E_\beta[X_2] = E_\beta[g(\mathbf{X})Y] = E_\beta[g(\mathbf{X})]E_\beta[Y] = 2\beta E_\beta[g(\mathbf{X})]$ .

- Thus  $E_\beta[g(\mathbf{X})] = \frac{1}{2}$

Basu's theorem:  $Y$  and  $g(\mathbf{X})$  are independent

■

## Estimator

### MLE

$$\hat{\theta} = \arg \max_{\theta} f(x|\theta)$$

so we solve for

$$\frac{d}{d\theta} l(\hat{\theta}|x) = 0$$

### Invariance Property of MLE:

**Theorem:** If  $\hat{\theta}$  is the MLE of  $\theta$ , then for any function  $\tau(\theta)$ , the MLE of  $\tau(\theta)$  is  $\tau(\hat{\theta})$

## Methods of Momment

Given Moments  $M_i$ :

$$\begin{aligned} M_1 &= \frac{1}{n} \sum X_i^1, \mu_1 = EX^1 \\ M_2 &= \frac{1}{n} \sum X_i^2, \mu_2 = EX^2 \\ &\dots \end{aligned}$$

### Example: Normal method of moments:

Given i.i.d  $X_i \sim N(\theta, \sigma^2)$ , we have  $m_1 = \bar{X}$ ,  $m_2 = \frac{1}{n} \sum X_i^2$ ;

the first moment is  $\mu_1 = \theta$  and  $\mu_2 = \sigma^2 + \theta^2$ , hence we must solve for:

$$\begin{aligned} \bar{X} &= \theta, \frac{1}{n} \sum X_i^2 = \sigma^2 + \theta^2 \\ \Rightarrow \hat{\theta} &= \bar{X}, \hat{\sigma}^2 = \frac{1}{n} \sum X_i^2 - \bar{X}^2 \blacksquare \end{aligned}$$

**Linear Regression****Bob's Lab MLE:**

For random sample  $Y = Y_1, Y_2, \dots, Y_i \sim n(\mu_i, \phi)$ , such that

$$l(\theta|y) = -n\log(2\pi) - \frac{n}{2}\log\phi - \frac{1}{2\phi}\sum (y_i - \mu_i)^2$$

We have  $n$  data and  $n+1$  parameters, so we reparameterize:

With the **restriction** of  $\mu_i = \beta X_i + \alpha$ ,

$$l(\theta|y) = -n\log(2\pi) - \frac{n}{2}\log\phi - \frac{1}{2\phi}\sum (y_i - \beta x_i - \alpha)^2$$

By MLE:

$$\hat{\alpha} = \frac{\sum y_i - \beta x_i}{\sum x_i}$$
$$\hat{\beta} = \frac{\sum x_i(y_i - \alpha)}{\sum x_i^2}$$

Solving we have  $\hat{\mu}_i = \hat{\alpha} + \hat{\beta}x_i$  and  $\hat{\phi}$

**UMVUE****CRLB-Definition**

- You always mess up how to take the expectation of fisher information
- Don't forget about the integral condition when taking expectation
- You mess up the derivative of the bottom and the top
- fisher information **ONLY** applied to exponential family

[Bob's Homework 3](#)

## 6.2 Defining UMVUE

**Example 6.2** - Consider the family of estimators  $T_{\theta_0}(\mathbf{X}) \equiv \theta_0$  for  $\theta_0 \in \Theta$ .

- $\text{MSE}_{\theta} [T_{\theta_0}] = (\theta_0 - \theta)^2$ , since  $\text{Var}_{\theta} [T_{\theta_0}] = 0$ .
- Hence  $\min_{\theta \in \Theta} (\text{MSE}_{\theta} [T_{\theta_0}]) = 0$ , since  $\text{MSE}_{\theta_0} [T_{\theta_0}] = 0$ . - If  $T^*(\mathbf{X})$  is to be the uniformly minimum MSE estimator, then
- $\text{MSE}_{\theta} [T^*] \leq \min_{\theta \in \Theta} (\text{MSE}_{\theta} [T_{\theta_0}]) = 0$
- That is,  $T^*(\mathbf{X})$  must have zero variance and zero bias for all  $\theta \in \Theta$ . This is not possible, unless  $\Theta = \{\theta_0\}$ .

So, if  $\Theta$  is not restricted, UBE may not exist ■

**Example 6.3** - Define a class of estimators,  $\mathcal{C}$ , and find the estimator (or estimators) within  $\mathcal{C}$  which has minimal MSE.

- Example: Suppose  $\mathbf{X}$  is a sample point of size  $n$  from a  $N(\mu, \sigma^2)$  family.
- Consider  $\mathcal{C}_{\sigma^2} = \left\{ T_k(\mathbf{X}) = \frac{\sum (X_i - \bar{X})^2}{k}, k > 0 \right\}$ , a class of estimators of  $\sigma^2$ .

$$\begin{aligned} T &= \frac{n-1}{k} S^2 \\ E[T] &= \frac{n-1}{k} \sigma^2 \\ \text{Var}[T] &= \frac{2(n-1)}{k^2} \sigma^4 \\ \text{MSE}_k[T] &= \frac{n-1-k^2}{k} \sigma^4 + \frac{2(n-1)}{k^2} \sigma^4 \end{aligned}$$

- Argue that  $\text{argmin} \{ \text{MSE}_{\sigma^2} (T_k) \} = n+1$ .

- That is,  $\tilde{\sigma}^2 = \frac{\sum_{i=1}^{n+1} (X_i - \bar{X})^2}{n}$  has the minimum MSE among estimators in  $\mathcal{C}_{\sigma^2}$ .

■

**Definition 6.3** - Suppose  $\mathcal{C}_{\theta} = \{T(\mathbf{X}) : E_{\theta}[T(\mathbf{X})] = \theta, \text{ for all } \theta \in \Theta\}$ , the class of all unbiased estimators of  $\theta$ .

- Definition: An estimator  $T^*$  is a best unbiased estimator of  $\tau(\theta)$  if it satisfies  $E_{\theta} [T^*] = \tau(\theta)$  for all  $\theta \in \Theta$  and, for any other estimator  $T$  with  $E_{\theta}[T] = \tau(\theta)$ , we have  $\text{Var}_{\theta} [T^*] \leq \text{Var}_{\theta} [T]$  for all  $\theta$ .  $T^*$  is also called a uniform minimum variance unbiased estimator (UMVUE) of  $\tau(\theta)$

### 6.2.1 CRLB

**Definition 6.4** - Let  $T = T(\mathbf{X})$  be a statistic with  $E_{\theta}[T] = g(\theta)$  and  $\text{Var}_{\theta}[T] < \infty$ .

- For any random variable  $W(\mathbf{X}, \theta)$  which has a finite second moment, the Cauchy-Schwarz inequality says that

$$(\text{Cov}_\theta[T, W])^2 \leq \text{Var}_\theta[T] \text{Var}_\theta[W]$$

or equivalently,

$$\text{Var}_\theta[T] \geq \frac{(\text{Cov}_\theta[T, W])^2}{\text{Var}_\theta[W]}$$

- This is a lower bound for the variance of  $T(\mathbf{X})$ .

- **CRLB in general only apply to distribution which it's space doesn't depend on the parameter**

**Corollary 6.0.1** The cleverness in this theorem follows from choosing  $X$  to be the estimator  $W(\mathbf{X})$  and  $Y$  to be the quantity  $\frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta)$  and applying the Cauchy-Schwarz Inequality. First note that

$$\begin{aligned} \frac{d}{d\theta} \mathbb{E}_\theta W(\mathbf{X}) &= \int_{\mathbf{X}} W(\mathbf{x}) \left[ \frac{\partial}{\partial \theta} f(\mathbf{x} | \theta) \right] d\mathbf{x} \\ &= \mathbb{E}_\theta \left[ W(\mathbf{X}) \frac{\frac{\partial}{\partial \theta} f(\mathbf{X} | \theta)}{f(\mathbf{X} | \theta)} \right] \quad (\text{multiply by } f(\mathbf{X} | \theta) / f(\mathbf{X} | \theta)) \\ &= \mathbb{E}_\theta \left[ W(\mathbf{X}) \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right] \quad (\text{property of logs}) \end{aligned}$$

which suggests a covariance between  $W(\mathbf{X})$  and  $\frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta)$ . For it to be a covariance, we need to subtract the product of the expected values, so we calculate  $\mathbb{E}_\theta \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)$ . But if we apply (7.3.7) with  $W(\mathbf{x}) = 1$ , we have

$$\mathbb{E}_\theta \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right) = \frac{d}{d\theta} \mathbb{E}_\theta[1] = 0.$$

- It can be shown that defining  $W(\mathbf{X}, \theta)$  as

$$W(\mathbf{X}, \theta) = \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) = \frac{\frac{\partial}{\partial \theta} f(\mathbf{X} | \theta)}{f(\mathbf{X} | \theta)}$$

leads to the greatest lower bound.

$$\mathbb{E}_\theta[W] = \int W(\mathbf{x}, \theta) f(\mathbf{x} | \theta) d\mathbf{x} = \int \frac{\partial}{\partial \theta} f(\mathbf{x} | \theta) d\mathbf{x} = \frac{\partial}{\partial \theta} \int f(\mathbf{x} | \theta) d\mathbf{x} = \frac{\partial}{\partial \theta} 1 = 0$$

Therefore  $\text{Cov}_\theta \left( W(\mathbf{X}), \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)$  is equal to the expectation of the product, and it follows from (7.3.7) and (7.3.8) that

$$(7.3.9) \text{Cov}_\theta \left( W(\mathbf{X}), \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right) = \text{E}_\theta \left( W(\mathbf{X}) \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right) = \frac{d}{d\theta} \text{E}_\theta W(\mathbf{X})$$

Also, since  $\text{E}_\theta \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right) = 0$  we have

$$\text{Var}_\theta \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right) = \text{E}_\theta \left( \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)^2 \right) - 0^2.$$

Using the Cauchy-Schwarz Inequality together with (7.3.9) and (7.3.10), we obtain

$$\text{Var}_\theta(W(\mathbf{X})) \geq \frac{\left( \frac{d}{d\theta} \text{E}_\theta W(\mathbf{X}) \right)^2}{\text{E}_\theta \left( \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)^2 \right)},$$

**Theorem 6.1** *Cramér-Rao Inequality*

Let  $X_1, X_2, \dots, X_n$  be a sample from  $f(x | \theta)$  and let  $T(\mathbf{X})$  be any estimator satisfying

$$\frac{\partial}{\partial \theta} \text{E}_\theta[T(\mathbf{X})] = \int T(\mathbf{x}) \frac{\partial}{\partial \theta} f(\mathbf{x} | \theta) d\mathbf{x}$$

and

$$\text{Var}_\theta[T(\mathbf{X})] < \infty$$

Then

$$\text{Var}_\theta[T] \geq \frac{\left( \frac{\partial}{\partial \theta} \text{E}_\theta[T(\mathbf{X})] \right)^2}{\text{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)^2 \right]}$$

**Corollary 6.1.1** *Corollary: If the conditions of the Cramér-Rao Inequality theorem are met and if  $X_1, X_2, \dots, X_n$  are iid with pdf  $f(x | \theta)$ , then*

$$\text{Var}_\theta[T] \geq \frac{\left( \frac{\partial}{\partial \theta} \text{E}_\theta[T(\mathbf{X})] \right)^2}{n \text{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(X | \theta) \right)^2 \right]}$$

Proof:

$$\begin{aligned}
E_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)^2 \right] &= E_\theta \left[ \left( \sum_i \frac{\partial}{\partial \theta} \log f(X_i | \theta) \right)^2 \right] \\
&= E_\theta \left[ \sum_i \sum_j \left( \frac{\partial}{\partial \theta} \log f(X_i | \theta) \right) \left( \frac{\partial}{\partial \theta} \log f(X_j | \theta) \right) \right] \\
&= \sum_i \sum_j E_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(X_i | \theta) \right) \left( \frac{\partial}{\partial \theta} \log f(X_j | \theta) \right) \right] \\
&= \sum_i E_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(X_i | \theta) \right)^2 \right] \\
&= n E_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(X | \theta) \right)^2 \right]
\end{aligned}$$

For  $i \neq j$  we have

$$\begin{aligned}
E_\theta \left( \frac{\partial}{\partial \theta} \log f(X_i | \theta) \frac{\partial}{\partial \theta} \log f(X_j | \theta) \right) \\
&= E_\theta \left( \frac{\partial}{\partial \theta} \log f(X_i | \theta) \right) E_\theta \left( \frac{\partial}{\partial \theta} \log f(X_j | \theta) \right) \quad (\text{independence}) \\
&= 0. \quad (\text{from (7.3.8)})
\end{aligned}$$

**Example 6.4** *Example 7.3.13 (Unbiased estimator for the scale uniform)*

let  $X_1, \dots, X_n$  be iid with pdf  $f(x | \theta) = 1/\theta, 0 < x < \theta$ . Since  $\log f(x | \theta) = -1/\theta$ , we have

$$E_\theta \left( \left( \frac{\partial}{\partial \theta} \log f(X | \theta) \right)^2 \right) = \frac{1}{\theta^2}.$$

The Cramer-Rao Theorem would seem to indicate that if  $W$  is any unbiased estimator of  $\theta_1$

$$\text{Var } w \geq \frac{\theta^2}{n}$$

We would now like to find an unbiased estimator with small variance. As a first guess, consider the sufficient statistic  $Y = \max(X_1, \dots, X_n)$ , the largest order statistic. The pdf of  $Y$  is  $f_Y(y | \theta) = ny^{n-1}/\theta^n, 0 < y < \theta$ , so

$$E_0 Y = \int_0^\theta \frac{ny^n}{\theta^n} dy = \frac{n}{n+1} \theta_1$$



showing that  $\frac{n+1}{n}Y$  is an unbiased estimator of  $\theta$ . We next calculate

$$\begin{aligned}\text{Var}_\theta \left( \frac{n+1}{n}Y \right) &= \left( \frac{n+1}{n} \right)^2 \text{Var}_\theta Y \\ &= \left( \frac{n+1}{n} \right)^2 \left[ \mathbb{E}_\varphi Y^2 - \left( \frac{n}{n+1}\theta \right)^2 \right] \\ &= \left( \frac{n+1}{n} \right)^2 \left[ \frac{n}{n+2}\theta^2 - \left( \frac{n}{n+1}\theta \right)^2 \right] \\ &= \frac{1}{n(n+2)}\theta^2\end{aligned}$$

which is uniformly smaller than  $\theta^2/n$ . This indicnten that the Cramer-Rao Theorem is not applicable to this pdf. To see that this is so, we can use Leibnitz's Rule (Section 2.4) to calculate

$$\begin{aligned}\frac{d}{d\theta} \int_0^\theta h(x)f(x|\theta)dx &= \frac{d}{d\theta} \int_0^\theta h(x)\frac{1}{\theta}dx \\ &= \frac{h(\theta)}{\theta} + \int_0^\theta h(x)\frac{\partial}{\partial\theta} \left( \frac{1}{\theta} \right) dx \\ &\neq \int_0^\theta h(x)\frac{\partial}{\partial\theta} f(x|\theta)dx\end{aligned}$$

unless  $h(\theta)/\theta = 0$  for all  $\theta$ . Hence, the Cramér-Rao Theorem does not apply. In general, if the range of the pdf depends on the parameter, the theorem will not be applicable. ■

**Example 6.5** Estimate  $\beta$  for Exponential ( $\beta$ ) family

- Assume we have an iid sample:  $f(x|\beta) = \frac{1}{\beta^n} \exp(-\sum x/\beta)$
- $\log f(x|\beta) = -n \log \beta - \sum x/\beta = -n \log \beta - n\bar{x}/\beta$
- Score statistic:  $\frac{\partial[\log f(\mathbf{X}|\beta)]}{\partial\beta} = -\frac{n}{\beta} + \frac{n\bar{X}}{\beta^2}$ . Is the expectation equal to 0 ?
- $\mathbb{E}_\beta \left[ \left( \frac{\partial}{\partial\beta} \log f(\mathbf{X}|\beta) \right)^2 \right] = \text{Var}_\beta \left[ -\frac{n}{\beta} + \frac{n\bar{X}}{\beta^2} \right] = \frac{n^2}{\beta^4} \text{Var}_\beta[\bar{X}] = \frac{n^2}{\beta^4} \frac{\beta^2}{n} = \frac{n}{\beta^2}$
- $\frac{\partial}{\partial\beta} \mathbb{E}_\beta[\bar{X}] = \frac{\partial}{\partial\beta} \beta = 1$
- So the CRLB is  $\beta^2/n$  which is equal to  $\text{Var}_\beta[\bar{X}]$ . So  $\bar{X}$  is UMVUE for  $\beta$ .

■

**Theorem 6.2** *Attainment*

Let  $X_1, \dots, X_n$  be iid  $f(x | \theta)$ , where  $f(x | \theta)$  satisfies the conditions of the Cramér-Rao Theorem. Let  $L(\theta | \mathbf{x}) = \prod_{i=1}^n f(x_i | \theta)$  denote the likelihood function. If  $W(\mathbf{X}) = W(X_1, \dots, X_n)$  is any unbiased estimator of  $\tau(\theta)$ , then  $W(\mathbf{X})$  attains the Cramér-Rao Lower Bound if and only if

$$a(\theta)[W(\mathbf{x}) - \tau(\theta)] = \frac{\partial}{\partial \theta} \log L(\theta | \mathbf{x})$$

for some function  $a(\theta)$ .

**Proof:**

The Cramér-Rao Inequality, as given in (7.3.6), can be written as

$$\left[ \text{Cov}_\theta \left( W(\mathbf{X}), \frac{\partial}{\partial \theta} \log \prod_{i=1}^n f(X_i | \theta) \right) \right]^2 \leq \text{Var}_\theta W(\mathbf{X}) \text{Var}_\theta \left( \frac{\partial}{\partial \theta} \log \prod_{i=1}^n f(X_i | \theta) \right),$$

and, recalling that

$$\mathbb{E}_\theta W = \tau(\theta)$$

,

$$\mathbb{E}_\theta \left( \frac{\partial}{\partial \theta} \log \prod_{i=1}^n f(X_i | \theta) \right) = 0$$

, and using the results of Theorem 4.5.7, we can have equality if and only if  $W(\mathbf{x}) - \tau(\theta)$  is proportional to  $\frac{\partial}{\partial \theta} \log \prod_{i=1}^n f(x_i | \theta)$ . That is exactly what is expressed in (7.3.12).

\*Theorem 4.5.7 For any random variables  $X$  and  $Y$ ,

a.  $-1 \leq \rho_{XY} \leq 1$ .

b.  $|\rho_{XY}| = 1$  if and only if there exist numbers  $a \neq 0$  and  $b$  such that  $P(Y = aX + b) = 1$ .

If  $\rho_{XY} = 1$ , then  $a > 0$ , and if  $\rho_{XY} = -1$ , then  $a < 0$ .

**Example 6.6** *Example: Estimate  $\beta$  for Exponential ( $\beta$ ) family*

- Recall-Score statistic:  $\frac{\partial [\log f(x|\beta)]}{\partial \beta} = -\frac{n}{\beta} + \frac{n\bar{x}}{\beta^2}$
- Define  $a(\beta) = \frac{n}{\beta^2}$
- Then  $\frac{\partial}{\partial \beta} \log f(x | \beta) = a(\beta)(\bar{X} - \beta)$  So  $\text{Var}_\beta[\bar{X}]$  attains the  $\text{CRLB} = \beta^2/n$  and, again,  $\bar{X}$  is UMVUE for  $\beta$

**Corollary 6.2.1** *if the range of the pdf depends on the parameter, the theorem will not be applicable.*

**Definition 6.5** *Fisher's Information*

The quantity

$$I(\theta) = \mathbb{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X} | \theta) \right)^2 \right] = \underbrace{-\mathbb{E}_\theta \left( \frac{\partial^2}{\partial \theta^2} \log f(X | \theta) \right)}_{\text{true for exponential family}}$$

is called the Fisher's Information.  $I(\theta)$  depends on the particular parameterization of the model family. Suppose that  $\theta = h(v)$ , where  $h(\cdot)$  is differentiable, then the information that  $\mathbf{X}$  contains about  $v$  is

$$I^*(v) = I(h(v)) \cdot [h'(v)]^2$$

**Example 6.7** *Example 7.3.14 (Normal variance bound)* Let  $X_1, \dots, X_n$  be iid  $\mathcal{N}(\mu, \sigma^2)$ , and consider estimation of  $\sigma^2$ , where  $\mu$  is unknown. The normal pdf satisfies the assumptions of the Cramér-Rao Theorem and Lemma 7.3.11, so we have

$$\frac{\partial^2}{\partial (\sigma^2)^2} \log \left( \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-(1/2)(x-\mu)^2/\sigma^2} \right) = \frac{1}{2\sigma^4} - \frac{(x-\mu)^2}{\sigma^6}$$

and

$$\begin{aligned} -\mathbb{E} \left( \frac{\partial^2}{\partial (\sigma^2)^2} \log f(X | \mu, \sigma^2) \mid \mu, \sigma^2 \right) &= -\mathbb{E} \left( \frac{1}{2\sigma^4} - \frac{(X-\mu)^2}{\sigma^6} \mid \mu, \sigma^2 \right) \\ &= -\frac{1}{2\sigma^4} + \frac{1}{\sigma^4} \\ &= \frac{1}{2\sigma^4}. \end{aligned}$$

Thus, any unbiased estimator,  $W$ , of  $\sigma^2$  must satisfy

$$\text{Var}(W \mid \mu, \sigma^2) \geq \frac{2\sigma^4}{n}$$

■

**Rao-Blackwell-Finding UMVUE(Must be Unbiased)**

- Find by condition
- Find by proving it is unique/same
- How to find the conditional expectation?
- 2021 Theory 5a-c is good practice

### 6.3 Finding UMVUE

#### 6.3.1 Rao-Blackwell theorem

##### Theorem 6.3 Rao-Blackwell

Let  $W$  be an estimator of  $\tau(\theta)$  where  $E_{\theta}[W] = \eta(\theta)$ . Also let  $T$  be a sufficient statistic for  $\theta$ , possibly vector valued. Define  $\phi(T) = E_{\theta}[W | T]$ . Then

- $\phi(T)$  is a statistic which is a function of  $T$ .
- $E_{\theta}[\phi(T)] = \eta(\theta)$ ,
- $MSE_{\tau(\theta)}[\phi(T)] \leq MSE_{\tau(\theta)}[W]$  for all  $\theta$ , and
- $MSE_{\tau(\theta)}[\phi(T)] < MSE_{\tau(\theta)}[W]$  for some  $\theta$  unless  $\phi(T) = W$  with probability 1.

Proof

1. Since  $T$  is sufficient then  $f(x | T, \theta)$  does not depend on  $\theta$ .  $W$  is a statistic that does not depend on  $\theta$ , so  $E_{\theta}[W | T]$  does not depend on  $\theta$ .
2.  $E_{\theta}[\phi] = E_{\theta}[E[W | T]] = E_{\theta}[W] = \eta(\theta)$   
 $MSE_{\tau(\theta)}[W] = Var_{\theta}[W] + (\eta(\theta) - \tau(\theta))^2$   
 $= Var_{\theta}[E[W | T]] + E_{\theta}[Var[W | T]] + (\eta(\theta) - \tau(\theta))^2$   
 $\geq Var_{\theta}[\phi] + (\eta(\theta) - \tau(\theta))^2$   
 $= MSE_{\tau(\theta)}[\phi]$

##### Theorem 6.4 Rao-Blackwell for Unbiased estimator(by proving equal)

- Suppose  $W$  is an unbiased estimator, that is,  $E_{\theta}[W] = \eta(\theta) = \tau(\theta)$ .
- So, we can start with an unbiased estimator,  $W$ , and get a new estimator,  $\phi(T) = E_{\theta}[W | T]$ , that has variance that is no larger than the variance of  $W$ , and possibly smaller!
- Actually,  $MSE_{\theta}[W] = MSE_{\theta}[\phi]$  for all  $\theta \in \Theta$  if and only if

$$E_{\theta}[Var[W | T]] = 0$$

as seen in the proof of the theorem.

- $E_{\theta}[Var[W | T]] = 0 \Rightarrow Var[W | T] = 0$ , since  $Var[W | T] \geq 0$
- Thus,  $E_{\theta}[(W - \phi)^2 | T] = 0 \Rightarrow W = \phi$ .

- If the Blackwellization of  $W$  yields an estimator different from  $W$ , then the new estimator has a smaller MSE.

**Example 6.8** Suppose  $X_1, X_2, \dots, X_n \sim \text{iid Poisson}(\theta)$ . Let  $\tau(\theta) = P\{X = 0\} = e^{-\theta}$ .

- Consider the statistic  $W = \mathbf{1}_{\{0\}}(X_1) \sim \text{Bernoulli}(\tau(\theta))$
- $E_\theta[W] = \tau(\theta)$ , that is,  $W$  is unbiased for  $\tau(\theta)$ .
- $\text{Var}_\theta[W] = \tau(\theta)(1 - \tau(\theta))$
- From previous work we know that  $T = \sum X_i$  is sufficient for  $\theta$ . Note that  $T \sim \text{Poisson}(n\theta)$ .
- Find a better unbiased estimator, in terms of the variance.
- Confirm this estimator is unbiased.
- Is this estimator UMVUE?

$$\begin{aligned}\Phi(s) &= E[W|T = s] = P(X_1 = 0|T = s) \\ &= \frac{P(X_1 = 0, \sum_1 X_i = s)}{P(T = s)} \\ &= \frac{P(X_1 = 0) * P(\sum_2 X_i = s)}{P(T = s)} \\ &= \frac{\underbrace{Poi((n-1)\theta)}}{P(T = s)} \\ &= \left(\frac{n-1}{n}\right)^s \\ \Phi(T) &= \left(\frac{n-1}{n}\right)^{\sum X_i}\end{aligned}$$

This is a better estimator by Rao-blackwell.

$$\begin{aligned}E[\Phi(T)] &= E\left[\left(\frac{n-1}{n}\right)^{\sum X_i}\right] = E\left[\exp\left(\log\left(\frac{n-1}{n}\right) \sum X_i\right)\right] \\ &\text{by moment generate function: } M_T(t) = \exp(n\theta(e^t - 1)) \\ &= M\left(\log\left(\frac{n-1}{n}\right)\right) = \exp\left[n\theta\left(\frac{n-1}{n} - 1\right)\right] \\ &= \exp(-\theta)\end{aligned}$$

This shows that  $\Phi$  is unbiased.

$$\begin{aligned}
 \text{Var}(\Phi) &= \underbrace{E[\Phi^2]}_{E[\frac{n-1}{n}2^T]=M(2\log(\frac{n-1}{n}))} - \underbrace{(E[\Phi])^2}_{\exp(-2\theta)} \\
 &= \exp(-\theta) * \exp(-\frac{n-1}{n}\theta) - \exp(-2\theta) \\
 &< \exp(-\theta) * \underbrace{(1 - \exp(-\theta))}_{\frac{n-1}{n}\theta > 0} \\
 &= \text{Var}(W)
 \end{aligned}$$

This shows Rao-blackwell is true for this case and we have a better unbiased estimator.  
Finding the lower bound:

$$CRLB = \frac{\theta \exp(-2\theta)}{n}$$

■

**Theorem 6.5** *UMVUE is unique*

If  $W$  is a UMVUE of  $\tau(\theta)$ , then  $W$  is unique.

Proof: (A review of the proof in C&B, p. 343-344.) Suppose that both  $W$  and  $W'$  are UMVUE of  $\tau(\theta)$ . Then  $W^* = \frac{(W+W')}{2}$  is unbiased for  $\tau(\theta)$ . By supposition  $\text{Var}_\theta W' = \text{Var}_\theta W$ . This implies

$$\begin{aligned}
 \text{Var}_\theta W^* &= \text{Var}_\theta \left[ \frac{1}{2}W + \frac{1}{2}W' \right] \\
 &= \frac{1}{4} \text{Var}_\theta W + \frac{1}{4} \text{Var}_\theta W' + \frac{1}{2} \text{Cov}_\theta [W, W']
 \end{aligned}$$

$$\text{Cov}_\theta [W, W'] \leq \sqrt{\text{Var}_\theta W \cdot \text{Var}_\theta W'} \longrightarrow \leq \text{Var}_\theta W$$

However,  $\text{Var}_\theta W$  is the smallest variance among unbiased estimators, thus  $\text{Var}_\theta W^* = \text{Var}_\theta W$ .

The last result implies  $\text{Cov}_\theta [W, W'] = \sqrt{\text{Var}_\theta W \cdot \text{Var}_\theta W'} = \text{Var}_\theta W$ . We noted before that this implies  $W' = aW + b, a \neq 0$ , where  $a$  and  $b$  may depend on the parameters, but not the data. -  $E_\theta [W'] = aE_\theta [W] + b = a\tau(\theta) + b$ . But  $E_\theta [W'] = \tau(\theta)$ . So

$$(a - 1)\tau(\theta) + b = 0$$

$$- \text{Var}_\theta [W'] = a^2 \text{Var}_\theta [W].$$

But  $\text{Var}_\theta [W'] = \text{Var}_\theta [W]$ , thus  $a = 1$  and from above,  $b = 0$ .

if  $a = -1, b = 2\tau(\theta)$

Thus  $W' = W$  and the UMVUE is unique.

**Theorem 6.6** *The unbiased estimator  $W$  is the UMVUE of  $\tau(\theta)$  if and only if  $W$  is uncorrelated with all unbiased estimators of 0.*

Proof: (This is an expansion of the proof in C&B, p. 344-345.)

( $\Rightarrow$ ) Suppose  $W$  is the UMVUE of  $\tau(\theta)$ . Let  $V = W + aU$  where  $E_{\theta}[U] = 0$ ,  $\text{Var}_{\theta}[U] = 1$ , and  $a$  is any non-zero real constant. Then

$$E_{\theta}[V] = \tau(\theta)$$

and

$$\text{Var}_{\theta}[V] = \text{Var}_{\theta}[W] + a^2 \text{Var}_{\theta}[U] + 2a \text{Cov}_{\theta}[W, U] \geq \text{Var}_{\theta}[W]$$

So  $a + 2 \text{Cov}_{\theta}[W, U] > 0$  which implies  $|\text{Cov}_{\theta}[W, U]| \leq \frac{|a|}{2}$  for all  $a \neq 0$ . This inequality holds as  $|a| \rightarrow 0$ , thus  $\text{Cov}_{\theta}[W, U] = 0$  and  $\text{Var}_{\theta}(V) \rightarrow \text{Var}_{\theta}(W)$ .

( $\Leftarrow$ ) Let  $V$  be an estimator such that  $E_{\theta}[V] = \tau(\theta)$ . Note that

$$V = W + (V - W) \text{ and } E_{\theta}[V - W] = 0$$

Thus, by supposition,

$$\text{Cov}_{\theta}[W, V - W] = 0$$

The variance of  $V$  is thus

$$\text{Var}_{\theta}[V] = \text{Var}_{\theta}[W] + \text{Var}_{\theta}[V - W] \geq \text{Var}_{\theta}[W]$$

Since this applies for any unbiased estimator  $V$ ,  $W$  is the UMVUE of  $\tau(\theta)$ .

**Theorem 6.7** *Suppose that the UMVUE exists, then it is unique and is a function of a sufficient statistic.*

Proof:

We already have that the UMVUE, if it exists, is unique. Suppose  $W$  is the UMVUE of  $\tau(\theta)$  and  $T$  is a sufficient statistic for  $\theta$ . Then by the Rao-Blackwell theorem,  $\phi(T) = E_{\theta}[W | T]$  is unbiased and  $\text{Var}_{\theta}[\phi(T)] \leq \text{Var}_{\theta}[W]$ . Since  $\text{Var}_{\theta}[W]$  is the minimal achievable variance among unbiased estimators of  $\tau(\theta)$ ,  $\text{Var}_{\theta}[\phi(T)] = \text{Var}_{\theta}[W]$ . This implies that  $\phi(T)$  is a UMVUE of  $\tau(\theta)$ . Due to uniqueness,

$$\phi(T) = W.$$

Thus, the UMVUE is a function of a sufficient statistic.



**Lehmann-scheffe-Finding unique UMVUE**

- Completeness
- More example see *Mathematics Statistics*

### 6.3.2 Lehmann-scheffe

#### Theorem 6.8 Lehmann-Scheffé Theorem

Let  $T$  be any **complete sufficient statistic** for the parameter  $\theta$ , and let  $\phi(T)$  be any estimator based only on  $T$ . Then  $\phi(T)$  is the unique UMVUE of  $E_\theta[\phi(T)]$  (CRLB is not garenteed)

#### Theorem 6.9 Theorem 3.1 (Lehmann-Scheffé theorem) From Mathematics Statistics.

Suppose that there exists a sufficient and complete statistic  $T(X)$  for  $P \in \mathcal{P}$ . If  $\vartheta$  is estimable, then there is a unique unbiased estimator of  $\vartheta$  that is of the form  $h(T)$  with a Borel function  $h$ . (Two estimators that are equal a.s.  $\mathcal{P}$  are treated as one estimator.) Furthermore,  $h(T)$  is the unique UMVUE of  $\vartheta$ .

Proof:

In the last theorem, we required  $T$  to be a sufficient statistic for  $\theta$ . Suppose that  $T$  is also complete. Now consider  $\phi(T)$  to be an unbiased estimator of  $\tau(\theta)$

- Let  $W(T)$  be any other unbiased estimator of  $\tau(\theta)$  which is a function of  $T$ .
- $E_\theta[\phi(T) - W(T)] = 0$  for all  $\theta \in \Theta$
- By the completeness of  $T$ ,  $\phi(T) - W(T) = 0$  for all  $\theta \in \Theta$
- Hence  $\phi(T) = W(T)$ ; that is,  $\phi(T)$  is unique.
- So, there is at most one unbiased estimator of  $\tau(\theta)$  which is a function of a complete sufficient statistic.
- This means for any unbiased estimator of  $\tau(\theta)$ , say  $W'$ , we have  $\phi(T) = E_\theta[W' | T]$  and  $\text{Var}_\theta[\phi(T)] \leq \text{Var}_\theta[W']$
- Thus,  $\phi(T)$  is the unique UMVUE!

This theorem is a consequence of Theorem 2.5(ii) (Rao-Blackwell theorem). One can easily extend this theorem to the case of the uniformly minimum risk unbiased estimator under any loss function  $L(P, a)$  that is strictly convex in  $a$ . The uniqueness of the UMVUE follows from the completeness of  $T(X)$

There are two typical ways to derive a UMVUE when a sufficient and complete statistic  $T$  is available. The first one is solving for  $h$  when the distribution of  $T$  is available. The following are two typical examples.

#### Example 6.9 Example 3.1.

Let  $X_1, \dots, X_n$  be i.i.d. from the uniform distribution on  $(0, \theta)$ ,  $\theta > 0$ . Let  $\vartheta = g(\theta)$ , where  $g$  is a differentiable function on  $(0, \infty)$ . Since the sufficient and complete statistic

$X_{(n)}$  has the Lebesgue p.d.f.  $n\theta^{-n}x^{n-1}I_{(0,\theta)}(x)$ , an unbiased estimator  $h(X_{(n)})$  of  $\vartheta$  must satisfy

$$\theta^n g(\theta) = n \int_0^\theta h(x)x^{n-1}dx \quad \text{for all } \theta > 0.$$

Differentiating both sides of the previous equation and applying the result of differentiation of an integral (Royden (1968, §5.3) ) lead to

$$n\theta^{n-1}g(\theta) + \theta^n g'(\theta) = nh(\theta)\theta^{n-1}.$$

Hence, the UMVUE of  $\vartheta$  is  $h(X_{(n)}) = g(X_{(n)}) + n^{-1}X_{(n)}g'(X_{(n)})$ . In particular, if  $\vartheta = \theta$ , then the UMVUE of  $\theta$  is  $(1 + n^{-1})X_{(n)}$

**Example 6.10** • For some  $\tau(\theta)$  there is no unbiased estimator.

- Let  $X \sim \text{Binomial}(m, \theta)$  and  $\tau(\theta) = \ln\left(\frac{\theta}{1-\theta}\right)$ , the log-odds.
- Suppose  $T$  is any estimator of  $\tau(\theta)$ .
- $E_\theta[T] = \sum_{k=0}^m T(k) \binom{m}{k} \theta^k (1-\theta)^{m-k}$  is an  $m^{\text{th}}$ -degree polynomial in  $\theta$ .
- $\tau(\theta) = \ln\left(\frac{\theta}{1-\theta}\right)$  cannot be expressed as a finite-degree polynomial. (by Taylor expansion)
- Thus,  $T$  cannot be unbiased.
- Since  $T$  was an arbitrary estimator, no unbiased estimator of  $\tau(\theta)$  exists.
- For some situations there is an unbiased estimator, but no UMVUE. ■

**Consistency, Asymptotic Variance, efficiency of MLE**

## 7.2 Large Sample Property (and CI)

### Definition 7.5 Consistency

A sequence of estimators  $W_n = W_n(\mathbf{X})$  is a consistent sequence of estimators of the parameter  $\theta$  if, for every  $\epsilon > 0$  and every  $\theta \in \Theta$ ,

$$\lim_{n \rightarrow \infty} P_{\theta} \{ |W_n - \theta| < \epsilon \} = 1$$

- That is,  $W_n \rightarrow \theta$  in probability.
- For large  $n$ , almost all possible values of  $W_n(\mathbf{X})$  are close to  $\theta$ .
- $W_n$  is said to be a consistent estimator of  $\theta$ .
- Equivalently,  $\lim_{n \rightarrow \infty} P_{\theta} \{ |W_n - \theta| \geq \epsilon \} = 0$ .

### Example 7.9 -

Consider  $X_1, \dots, X_n$  as an iid random sample from Bernoulli( $p$ ). Is  $\hat{p} = \bar{X}$  is a consistent estimator of  $p$ ?

We know that when  $n$  is large,  $\frac{\hat{p}-p}{\sqrt{p(1-p)/n}}$  has an approximate  $N(0, 1)$  distribution; the larger the sample size, the better the approximation. Consider

$$P_p \{ |\hat{p} - p| < \epsilon \} = P_p \left\{ \left| \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \right| < \frac{\epsilon}{\sqrt{p(1-p)/n}} \right\} \approx P\{|Z| < \sqrt{n}K\}$$

where  $K = \epsilon / \sqrt{p(1-p)}$  and  $Z$  is a standard normal variable. Thus

$$\lim_{n \rightarrow \infty} P_p \{ |\hat{p} - p| < \epsilon \} \approx \lim_{n \rightarrow \infty} P\{|Z| < \sqrt{n}K\} = P\{|Z| < \infty\} = 1$$

Thus  $\hat{p}$  is a consistent estimator of  $p$ . ■

**Theorem 7.3** If  $W_n$  is a sequence of estimators of a parameter  $\theta$  satisfying

$$\begin{aligned} \lim_{n \rightarrow \infty} \text{Var}_{\theta} [W_n] &= 0 \\ \lim_{n \rightarrow \infty} \text{Bias}_{\theta} [W_n] &= 0 \end{aligned}$$

for every  $\theta \in \Theta$ , then  $W_n$  is a consistent sequence of estimators of  $\theta$ .

- By **Chebychev's inequality**,  $P_{\theta} \{ |W_n - \theta| \geq \epsilon \} \leq \frac{E_{\theta}[(W_n - \theta)^2]}{\epsilon^2}$ . The result follows given the sufficient conditions.
- These conditions are sufficient but may not be necessary.

### Example 7.10 -

- We have shown that  $\hat{p} = \bar{X}$  is an unbiased estimator of  $p$ , thus the bias is zero.
- $\text{Var}_p[\hat{p}] = \frac{p(1-p)}{n} \rightarrow 0$  as  $n \rightarrow \infty$ .
- By the previous theorem,  $\hat{p}$  is a consistent estimator of  $p$ . ■

**Example 7.11** - In general, if  $X_1, \dots, X_n$  is an iid random sample from a distribution where the first two moments exist (thus the distribution has a finite variance) and  $\mu = E[X]$ , then  $\bar{X}$  is a consistent estimator of  $\mu$ . Why?

■

**Theorem 7.4** Let  $W_n$  be a consistent sequence of estimators of a parameter  $\theta$ . Let  $a_1, a_2, \dots$  and  $b_1, b_2, \dots$  be sequences of constants satisfying

$$\lim_{n \rightarrow \infty} a_n = 1 (\text{variance factor})$$

$$\lim_{n \rightarrow \infty} b_n = 0 (\text{bias})$$

Then the sequence  $U_n = a_n W_n + b_n$  is a consistent sequence of estimators of  $\theta$ .

**Example 7.12** Is  $S_n$ , the square root of the sample variance of a random sample of size  $n$  from a normal distribution, a consistent estimator of  $\sigma$ , the standard deviation of the distribution?

It can be shown that  $E[S_n] = \sqrt{\frac{2}{n-1}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \sigma = c_n \sigma$  and thus

$\text{Var}[S_n] = \sigma^2 - c_n^2 \sigma^2$  and  $\text{Bias}[S_n] = c_n \sigma - \sigma$ . For both the variance and the bias to go to 0 as  $n \rightarrow \infty$ , we need  $c_n \rightarrow 1$ . We'll make use of Stirling's approximation:

$\Gamma(z) \approx \sqrt{\frac{2\pi}{z}} \left(\frac{z}{e}\right)^z$  where  $z$  is defined as a complex number but used here as real.

$$c_n = \sqrt{\frac{2}{n-1}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \approx \frac{1}{\sqrt{e}} \left(\frac{n}{n-1}\right)^{\frac{n-2}{2}} = \frac{1}{\sqrt{e}} \left\{ \frac{n}{n-1} \right\}^{-1} * \left[ \left(1 - \frac{1}{n}\right)^n \right]^{-\frac{1}{2}} \xrightarrow{n \rightarrow \infty} 1$$

So  $S_1, S_2, \dots$  is a consistent sequence of estimators of  $\sigma$ . ■

**Example 7.13** The MLE of  $\sigma$  is a function of  $S_n$ , namely  $T_n = a_n S_n$  where  $a_n = \sqrt{\frac{n-1}{n}}$ . Thus  $a_1, a_2, \dots$  is a sequence of constants such that  $\lim_{n \rightarrow \infty} a_n = 1$ .

So  $T_1, T_2, \dots$  is a consistent sequence of estimators of  $\sigma$ . That is, the MLE,  $T_n$ , is a consistent estimator of  $\sigma$ .

**Theorem 7.5** Consistency of MLE

Let  $X_1, X_2, \dots$ , be an iid random sample from a family of distributions indexed by  $\theta \in \Theta$ . Let  $\hat{\theta}$  denote the MLE of  $\theta$  and let  $\tau(\theta)$  be a continuous function of  $\theta$ . Under certain regularity conditions, for every  $\epsilon > 0$  and every  $\theta \in \Theta$ ,

$$\lim_{n \rightarrow \infty} P_\theta\{|\tau(\hat{\theta}) - \tau(\theta)| \geq \epsilon\} = 0$$

That is,  $\tau(\hat{\theta})$  is a consistent estimator of  $\tau(\theta)$ .

**Definition 7.6** *Asymptotic variance*

- First an example: We know that the sample mean,  $\bar{X}_n$ , from a  $N(\mu, \sigma^2)$  distribution has variance  $\sigma^2/n$ . Thus,  $\sigma_n^2 = \text{Var}[\sqrt{n}\bar{X}_n] = \sigma^2$ .

- In general, if  $X_1, \dots, X_n$  is an iid random sample from a distribution where the second moment exists (thus the distribution has a finite variance), then the central limit theorem says that  $\frac{\bar{X}_n - E[\bar{X}_n]}{\sqrt{\text{Var}[X]/n}} \rightarrow N(0, 1)$  in distribution.

- Define  $\sigma_n^2 = \text{Var}[\sqrt{n}\bar{X}_n]$ . Then  $\sigma_n^2 \rightarrow \text{Var}[X] = \tau^2$ . The constant  $\tau^2$  is the limiting variance or the limit of the variances.

- Definition: For an estimator  $T_n$ , suppose that  $k_n(T_n - \tau(\theta)) \rightarrow N(0, \sigma^2)$  in distribution. The parameter  $\sigma^2$  is called the asymptotic variance or variance of the limit distribution of  $T_n$ .

This is  
general  
Delta  
methods?

**Definition 7.7** *Asymptotic Efficiency*

Definition: A sequence of estimators  $W_n$  is asymptotically efficient for a parameter  $\tau(\theta)$  if

$$\sqrt{n}[W_n - \tau(\theta)] \rightarrow N(0, v(\theta))$$

in distribution and

$$v(\theta) = \frac{[\tau'(\theta)]^2}{E_\theta \left( \left( \frac{\partial}{\partial \theta} \ln f(X | \theta) \right)^2 \right)} = \frac{[\tau'(\theta)]^2}{-I(X)}$$

that is, the asymptotic variance of  $W_n$  achieves the Cramér-Rao Lower Bound.

**Example 7.14 :**

If  $\hat{\theta}_{jn}$  satisfies tA.GT with asymptotic covariance matrix  $V_{jn}(\theta)$ ,  $j = 1, 2$ , and  $V_{1,n}(\theta) \leq V_{2,n}(\theta)$  (in the sense that  $V_{2,n}(\theta) - V_{1,n}(\theta)$  is nonnegative definite for all  $\theta \in \Theta$ ), then  $\hat{\theta}_{1n}$  is said to be asymptotically more efficient than  $\hat{\theta}_{2n}$ .

**Theorem 7.6** *Asymptotic Efficiency of MLE*

Let  $X_1, X_2, \dots$  be an iid random sample from a family of distributions indexed by  $\theta \in \Theta$ . Let  $\hat{\theta}$  denote the MLE of  $\theta$  and let  $\tau(\theta)$  be a continuous function of  $\theta$ . Under certain regularity conditions,

$$\sqrt{n}[\tau(\hat{\theta}) - \tau(\theta)] \rightarrow N(0, v(\theta))$$

where  $v(\theta)$  is the Cramér-Rao Lower Bound. That is,  $\tau(\hat{\theta})$  is a consistent and asymptotically efficient estimator of  $\tau(\theta)$ .

Thus, the bounds of an approximate  $100(1 - \alpha)\%$  confidence interval for  $\theta$  is given as

$$\tau(\hat{\theta}) \pm z_{\alpha/2} \sqrt{v(\hat{\theta})/n} \text{ or } \tau(\hat{\theta}) \pm z_{\alpha/2} \sqrt{\hat{v}(\hat{\theta})/n}$$

**Corollary 7.6.1** • The quantity  $v(\hat{\theta})$  may be found by applying the delta method to the transformation  $\tau(\cdot)$  evaluated at  $\theta = \hat{\theta}$  and using  $\sigma^2 = I_1^{-1}(\theta)$ ; that is,  $v(\hat{\theta}) = \left. \frac{[\tau'(\theta)]^2}{I_1(\theta)} \right|_{\theta=\hat{\theta}}$

- For a given  $n$ , the approximate (unobservable) variance of  $\tau(\hat{\theta})$  is  $v(\theta)$ . To use this in practice, we must estimate the approximate variance.
  - First estimator: MLE of  $v(\theta)$ , given by  $v(\hat{\theta})$ .
  - This uses  $I_1(\hat{\theta})$ , the expectation of  $-\frac{\partial^2 \ln(f(X_i|\theta))}{\partial \theta^2}$  evaluated at  $\theta = \hat{\theta}$
  - We call this the expected information number.
  - Alternatively, we may estimate  $I_1(\theta)$  with  $-\frac{1}{n} \sum \frac{\partial^2 \ln(f(X_i|\theta))}{\partial \theta^2}$  evaluated at  $\theta = \hat{\theta}$ .
  - We call this the observed information number, denoted as  $\hat{I}_1(\hat{\theta})$ .
  - The variance estimator using the observed information number is denoted as  $\hat{v}(\hat{\theta})$ .
  - It has been shown that this provides a better estimator. (Efron & Hinkley, Biometrika, 1978)
  - If you can differentiate and evaluate the log-likelihood, you can calculate  $\hat{v}(\hat{\theta})$ .

**Example 7.15** Suppose  $X_1, \dots, X_n$  is an iid random sample from Poisson ( $\lambda$ ). The MLE of  $\lambda$  is  $\hat{\lambda} = \bar{X}$ .

- - Construct a 95%CI for  $\lambda$ , where  $n = 100$  and  $\bar{x} = 0.5$ .
- Large sample via the MLE

$$\begin{aligned}
 L(\lambda | x) &= e^{-\lambda} \lambda^x / x! \\
 l(\lambda | x) &= -\lambda + x \log \lambda - \log(x!) \\
 \frac{\partial l}{\partial \lambda} &= -1 + \frac{x}{\lambda}; \quad \frac{\partial^2 l}{\partial \lambda^2} = -\frac{x}{\lambda^2} \\
 I_1(\lambda) &= -E \left[ -\frac{x}{\lambda^2} \right] = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda} \\
 I_1(\hat{\lambda}) &= \frac{1}{\bar{x}} \\
 v(\lambda) &= \bar{x}
 \end{aligned}$$

- Pivot the cdf
- Recall that  $F_T(T | \theta) \sim \text{Uniform}(0, 1)$  if  $F_T$  is continuous in  $T$ . That is,  $F_T(T | \theta)$  has the same distribution for all values of  $\theta$ . Thus  $F_T(T | \theta)$  is a pivotal quantity.
- Hence  $P \{ \alpha_1 \leq F_T(t | \theta) \leq 1 - \alpha_2 \} = 1 - \alpha_1 - \alpha_2$ .



- The following also works if  $F_T$  is not continuous in  $T$ .
  - If  $F_T(T | \theta)$  is non-increasing in  $\theta$  then for given  $t$   $\theta_L = \inf \{\theta : F_T(t | \theta) \leq 1 - \alpha_2\}$  and  $\theta_U = \sup \{\theta : F_T(t | \theta) \geq \alpha_1\}$ .
  - If  $F_T(T | \theta)$  is non-decreasing in  $\theta$  then for given  $t$   $\theta_L = \inf \{\theta : F_T(t | \theta) \geq \alpha_1\}$  and  $\theta_U = \sup \{\theta : F_T(t | \theta) \leq 1 - \alpha_2\}$ .
- Note  $\sum_1^n x_i \sim \text{Poisson}(n\lambda)$

$$CLT \rightarrow \frac{\frac{1}{n} \sum x_i - n\lambda}{\sqrt{n\lambda}} = \sqrt{n} \left( \frac{\bar{x} - \lambda}{\sqrt{\lambda}} \right) \xrightarrow{\text{mis}_{dist.}} N(0, 1)$$

$\Phi \left( \sqrt{n} \left( \frac{\bar{x} - \lambda}{\sqrt{\lambda}} \right) \right)$  is a pivotal quantity.

Redo this

- - Construct a 95%CI for  $\theta = P_\lambda\{X = 0\}$ . - Large sample via the MLE

## Statistical Inference

### Hypothesis Testing

- What about testing  $H_0 : g(\theta) = a$  instead of  $H_0 : \theta = a$
- 

### Power/Size of the Test

---

update  
pages

## 8.2 Power, Size

**Definition 8.4** *The power function of a hypothesis test with rejection region  $R$  is the function of  $\theta$  defined by*

$$\beta(\theta) = P_{\theta}\{\mathbf{X} \in R\}$$

- This is a function of  $\theta$
- Note: In some textbooks,  $\beta$  refers to the probability of a Type II error. While the power function can be used to measure this probability,  $\beta(\theta)$  has a different meaning. Know the context in which  $\beta$  is used.
- The power function will be helpful in measuring and controlling the hypothesis testing errors. But first, an example.

**Example 8.6 :**

- Suppose we have  $H_0 : \theta \in \Theta_0$  and  $H_1 : \theta \in \Theta_1$ .
- Consider  $\arg \max_{\theta \in \Theta_0} \beta(\theta) = \theta_0$ . Then we may set  $\alpha$  as an upper bound on the Type I error probability by constructing our test rule so that  $\beta(\theta_0) \leq \alpha$
- For the previous example, we see that  $\beta(\lambda)$  is maximal at  $\lambda = 1$ . So we find  $c$  so that  $\beta(1) = P_{\lambda=1}\{\bar{T} \geq c\} = \alpha$ . We find  $c$  from the quantile of the distribution of  $T$
- When  $\alpha = 0.05$  and  $n = 25$ , we have  $c = 1.35$ .
- See the figure on the next slide.

**Definition 8.5** *A test with power function  $\beta(\theta)$  is a size  $\alpha$  test if*

$$\sup_{\theta \in \Theta_0} \beta(\theta) = \alpha$$

for  $0 \leq \alpha \leq 1$

(the maximum type-I Error can actually achieve)

**Definition 8.6** *A test with power function  $\beta(\theta)$  is a level  $\alpha$  test if*

$$\sup_{\theta \in \Theta_0} \beta(\theta) \leq \alpha$$

for  $0 \leq \alpha \leq 1$ . We say such a test is an alpha-level test.

( $\alpha$  is the upper bound of the type-I error rate)

**Example 8.7 :**

- Redefine the hypotheses  $H_0 : \lambda = \theta$  and  $H_1 : \lambda \neq \theta$ , where  $H_0$  is a point
- We can show that the LRT has the rejection region  $R = \{\mathbf{t} : \bar{T} \leq c_1 \text{ or } \bar{T} \geq c_2\}$  where  $c_1 \leq 1$  and  $c_2 \geq 1$  are constants.

REVIEW  
THIS

- The power function is defined as  $\beta(\lambda) = P_\lambda \{\bar{T} \leq c_1\} + P_\lambda \{\bar{T} \geq c_2\}$ .
- To make this an  $\alpha$ -level test (and a size  $\alpha$  test), we may choose  $c_1, c_2, \alpha_1$ , and  $\alpha_2$  so that  $P_{\lambda=1} \{\bar{T} \leq c_1\} = \alpha_1, P_{\lambda=1} \{\bar{T} \geq c_2\} = \alpha_2$ , and  $\alpha_1 + \alpha_2 = \alpha$

**Definition 8.7** Suppose we want the Type II error probability to be less than  $1 - \beta_*$  when  $\beta_*$  is a real number and  $\theta = \theta_* \in \theta_1$ . We seek the smallest sample size,  $n$ , such that

what is this?

$$\beta_n(\theta_*) \geq \beta_*$$

**Example 8.8** We may use the inverse cdf of the Gamma  $(n, \frac{\lambda}{n})$  distribution evaluated at  $\lambda = \lambda_*$  and perform a binary search for  $n$  with, say, the bisection algorithm.

**Definition 8.8** Unbiased Test:

A test with power function  $\beta(\theta)$  is unbiased if

$$\beta(\theta') \geq \beta(\theta'')$$

for every  $\theta' \in \Theta_1$  and  $\theta'' \in \Theta_0$ .

### Most Powerful Test

- According to Bob's note/homework, a **LRT test is most powerful of it's size if it exist for some distribution, by NP lemma**
- In LRT chapter  $\lambda = \frac{\sup f(x|\theta_0)}{\sup f(x|\theta)}$ , where in NP lemma,  $\lambda' = \frac{\sup f(x|\theta_1)}{\sup f(x|\theta_0)}$ , mark the difference.
- The rejection region of LRT is  $R = \{x : \lambda(x) \leq c\}, c \in (0, 1)$  vs rejection rejoin in the NP lemma:  $R' : \frac{\sup f(x|\theta_1)}{\sup f(x|\theta_0)} > k, k > 0$ , inverting the ratio we will have the same rejection region
- See [Bob's HW5](#)
- Don't ignore Karlin-Rubin Theorem ( Direction of MLR depends on the inequality of alternative test:  $\theta > \theta_0 \rightarrow$  non-decreasing/increasing)
- UMP is unbiased(can't find a proof yet)

Need to  
review  
This part

### 8.3 Most Powerful Test

A test for a hypothesis is a statistic  $T(X)$  taking values in  $[0, 1]$ . When  $X = x$  is observed, we reject  $H_0$  with probability  $T(x)$  and accept  $H_0$  with probability  $1 - T(x)$ . If  $T(X) = 1$  or 0 a.s.  $\mathcal{P}$ , then  $T(X)$  is a nonrandomized test. Otherwise  $T(X)$  is a randomized test. For a given test  $T(X)$ , the power function of  $T(X)$  is defined to be

$$\beta_T(P) = E[T(X)], \quad P \in \mathcal{P},$$

which is the type I error probability of  $T(X)$  when  $P \in \mathcal{P}_0$  and one minus the type II error probability of  $T(X)$  when  $P \in \mathcal{P}_1$ .

As we discussed in §2.4.2, with a sample of a fixed size, we are not able to minimize two error probabilities simultaneously. Our approach involves maximizing the power  $\beta_T(P)$  over all  $P \in \mathcal{P}_1$  (i.e., minimizing the type II error probability) and over all tests  $T$  satisfying

$$\sup_{P \in \mathcal{P}_0} \beta_T(P) \leq \alpha$$

where  $\alpha \in [0, 1]$  is a given level of significance. Recall that the left-hand side of (6.2) is defined to be the size of  $T$ .

**Definition 8.9** *Uniformly Most Powerful*

Let  $\mathcal{C}$  be a class of tests for testing  $H_0 : \theta \in \Theta_0$  vs  $H_1 : \theta \in \Theta_1$ . A test in class  $\mathcal{C}$  with power function  $\beta(\theta)$ , is a uniformly most powerful (UMP) class  $\mathcal{C}$  test if  $\beta(\theta) \geq \beta'(\theta)$  for every  $\theta \in \Theta_1$  and every other  $\beta'(\theta)$  that is a power function of a test in class  $\mathcal{C}$ .

- Typically, we choose a value of  $\alpha$  and define  $\mathcal{C}$  to be the class of all  $\alpha$ -level tests.
- This may include many or perhaps an unbounded number of tests with test size being less than  $\alpha$ .
- Tests with test size close to  $\alpha$  may be more powerful. Why?
- UMP tests may not exist in some situations.

compared the definition with unbiased test:

**Definition 8.10** *Simple Hypothesis:*

- A simple hypothesis completely specifies all of the parameters that index the related known or assumed family of distributions.
- That is, a simple hypothesis specifies the population distribution completely.

**Definition 8.11** *Composite Hypothesis:*

- A composite hypothesis is the union of a collection of simple hypotheses and specifies a subset of the related known or assumed family of distributions.
- That is, a composite hypothesis specifies possible population distributions but does not completely specify a single distribution.

**Lemma 8.2** *Neyman-Person lemme for two samples test:*

Consider testing the simple hypotheses  $H_0 : \theta = \theta_0$  vs  $H_1 : \theta = \theta_1$ , where  $\theta_0 \neq \theta_1$  and the pdf or pmf corresponding to  $\theta_i$  is  $f(\mathbf{x} | \theta_i)$ ,  $i = 0, 1$ , using a test with rejection region  $R$  that satisfies

1.  $\mathbf{x} \in R$  if  $f(\mathbf{x} | \theta_1) > k f(\mathbf{x} | \theta_0)$  and  $\mathbf{x} \in R^c$  if  $f(\mathbf{x} | \theta_1) < k f(\mathbf{x} | \theta_0)$  for some  $k \geq 0$ , and

2.  $\alpha = P_{\theta_0}(X \in R) = \beta(\theta_0)$

Then

- *Sufficient Condition:* Any test that satisfies (1) and (2) is a UMP level-  $\alpha$  test.

- *Necessary Condition:* If there exists a test satisfying (1) and (2) with  $k > 0$ , then every UMP level-  $\alpha$  test is a size  $\alpha$  test (satisfies (2)) and every UMP level-  $\alpha$  test satisfies (1) except perhaps on a set  $A$  satisfying  $P_{\theta_0}(\mathbf{X} \in A) = P_{\theta_1}(\mathbf{X} \in A) = 0$ .

UMP at  
this  $\alpha$   
level test

**Proof**

This proof is adapted from Wikipedia proof . Define the rejection region of the null hypothesis for the Neyman-Pearson (NP) test as

$$R_{NP} = \left\{ \mathbf{x} : \frac{f(\mathbf{x} | \theta_0)}{f(\mathbf{x} | \theta_1)} \leq \eta \right\} \quad (8.2.1)$$

$$R_{NP}^c = \left\{ \mathbf{x} : \frac{f(\mathbf{x} | \theta_1)}{f(\mathbf{x} | \theta_0)} \leq \eta \right\} \quad (8.2.2)$$

where  $\eta$  is chosen so that the power function  $\beta_{NP}(\theta_0) = \alpha$ . Any alternative test will have a different rejection region that we denote by  $R_A$ . The probability of the data falling within either region  $R = R_A$  or  $R = R_{NP}$  given parameter  $\theta$  is

$$\beta_R(\theta) = P\{R | \theta\} = \int_R f(x | \theta) dx$$

For the test with critical region  $R_A$  to have significance level  $\alpha$ , it must be true that the power function  $\beta_A(\theta_0) \leq \alpha$ , hence

$$\alpha = \beta_{NP}(\theta_0) \geq \beta_A(\theta_0).$$

It will be useful to break these down into integrals over distinct regions:

$$\begin{aligned} \beta_{NP}(\theta) &= P\{R_{NP} | \theta\} = P\{R_{NP} \cap R_A | \theta\} + P\{R_{NP} \cap R_A^c | \theta\} \\ \beta_A(\theta) &= P\{R_A | \theta\} = P\{R_{NP} \cap R_A | \theta\} + P\{R_{NP}^c \cap R_A | \theta\} \end{aligned}$$

where  $R^c \equiv \{x : x \notin R\}$  is the complement of region  $R$ . Setting  $\theta = \theta_0$ , these two expressions and the above inequality yield that

$$P\{R_{NP} \cap R_A^c | \theta_0\} \geq P\{R_{NP}^c \cap R_A | \theta_0\} \quad (8.2.3)$$

We would like to prove that

$$\beta_{NP}(\theta_1) \geq \beta_A(\theta_1)$$

As similarly shown above this is equivalent to

$$P\{R_{NP} \cap R_A^c \mid \theta_1\} \geq P\{R_{NP}^c \cap R_A \mid \theta_1\}$$

In what follows we show this inequality holds:

$$\begin{aligned} P\{R_{NP} \cap R_A^c \mid \theta_1\} &= \int_{R_{NP} \cap R_A^c} f(x \mid \theta_1) dx \\ &\geq \frac{1}{\eta} \int_{R_{NP} \cap R_A^c} f(x \mid \theta_0) dx \quad \text{by definition of } R_{NP} \text{ this is true for its subset (8.2.1)} \\ &= \frac{1}{\eta} P\{R_{NP} \cap R_A^c \mid \theta_0\} \\ &\geq \frac{1}{\eta} P\{R_{NP}^c \cap R_A \mid \theta_0\} \quad (8.2.3) \\ &= \frac{1}{\eta} \int_{R_{NP}^c \cap R_A} f(x \mid \theta_0) dx \\ &\geq \int_{R_{NP}^c \cap R_A} f(x \mid \theta_1) dx \quad (8.2.2) \\ &\quad \text{by definition of } R_{NP} \text{ this is true for its complement and complement subsets;} \\ &\quad \text{equality if } R_{NP}^c \cap R_A \text{ is empty} \\ &= P\{R_{NP}^c \cap R_A \mid \theta_1\} \end{aligned}$$

■

### Example 8.9 Proof from *Mathematics Statistics*

(i) (Existence of a UMP test). For every  $\alpha$ , there exists a UMP test of size  $\alpha$ , which is equal to

$$T_*(X) = \begin{cases} 1 & f_1(X) > cf_0(X) \\ \gamma & f_1(X) = cf_0(X) \\ 0 & f_1(X) < cf_0(X) \end{cases}$$

where  $\gamma \in [0, 1]$  and  $c \geq 0$  are some constants chosen so that  $E[T_*(X)] = \alpha$  when  $P = P_0$  ( $c = \infty$  is allowed).

(ii) (Uniqueness). If  $T_n$  is a UMP test of size  $\alpha$ , then

$$T_{**}(X) = \begin{cases} 1 & f_1(X) > cf_0(X) \\ 0 & f_1(X) < cf_0(X) \end{cases} \quad \text{a.s. } P.$$

*Proof.*



The proof for the case of  $\alpha = 0$  or  $1$  is left as an exercise. Assume now that  $0 < \alpha < 1$ .

(i) We first show that there exist  $\gamma$  and  $c$  such that  $E_0 [T_*(X)] = \alpha$ , where  $E_j$  is the expectation wrt.  $P_j$ . Let  $\gamma(t) = P_0 (f_1(X) > t f_0(X))$ . Then  $\gamma(t)$  is nonincreasing,  $\gamma(0) = 1$ , and  $\gamma(\infty) = 0$  (why?). Thus, there exists a  $c \in (0, \infty)$  such that  $\gamma(c) \leq \alpha \leq \gamma(c-)$ . Set

$$\gamma = \begin{cases} \frac{\alpha - \gamma(c)}{\gamma(c-) - \gamma(c)} & \gamma(c-) \neq \gamma(c) \\ 0 & \gamma(c-) = \gamma(c) \end{cases}$$

Note that  $\gamma(c-) - \gamma(c) = P(f_1(X) = c f_0(X))$ . Then

$$E_0 [T_*(X)] = P_0 (f_1(X) > c f_0(X)) + \gamma P_0 (f_1(X) = c f_0(X)) = \alpha$$

Next, we show that  $T_*$  in (6.3) is a UMP test. Suppose that  $T(X)$  is a test satisfying  $E_0 [T(X)] \leq \alpha$ . If  $T_*(x) - T(x) > 0$ , then  $T_*(x) > 0$  and, therefore,  $f_1(x) \geq c f_0(x)$ . If  $T_*(x) - T(x) < 0$ , then  $T_*(x) < 1$  and, therefore,  $f_1(x) \leq c f_0(x)$ . In any case,  $[T_*(x) - T(x)] [f_1(x) - c f_0(x)] \geq 0$  and, therefore,

$$\int [T_*(x) - T(x)] [f_1(x) - c f_0(x)] d\nu \geq 0$$

i.e.,

$$\int [T_*(x) - T(x)] f_1(x) d\nu \geq c \int [T_*(x) - T(x)] f_0(x) d\nu$$

The left-hand side of (6.5) is  $E_1 [T \cdot (X)] - E_1 [T(X)]$  and the right-hand side of (6.5) is  $c \{E_0 [T_+(X)] - E_0 [T(X)]\} = c \{\alpha - E_0 [T(X)]\} \geq 0$ . This proves the result in (i).

(ii) Let  $T_n(X)$  be a UMP test of size  $\alpha$ . Define

$$A = \{x : T_*(x) \neq T_n(x), \quad f_1(x) \neq c f_0(x)\}$$

Then  $\{T_*(x) - T_n(x)\} [f_1(x) - c f_0(x)] > 0$  when  $x \in A$  and  $= 0$  when  $x \in A^c$ , and

$$\int [T_*(x) - T_n(x)] [f_1(x) - c f_0(x)] d\nu = 0,$$

since both  $T_*$  and  $T_n$  are UMP tests of size  $\alpha$ . By Proposition 1.6(ii),  $v(A) = 0$ . This proves (6.4) ■

**Corollary 8.2.1** Base on theorem 8.1

Consider testing the simple hypotheses  $H_0 : \theta = \theta_0$  vs  $H_1 : \theta = \theta_1$ , where  $\theta_0 \neq \theta_1$ . Suppose  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  and  $g(t | \theta_i)$  is the pdf or pmf of  $T$  corresponding to  $\theta_i, i = 0, 1$ . Then any test based on  $T$  with rejection region  $R$  is a UMP level- $\alpha$  test if it satisfies

1.  $t \in R$  if  $g(t | \theta_1) > k g(t | \theta_0)$  and  $t \in R^c$  if  $g(t | \theta_1) < k g(t | \theta_0)$  for some  $k \geq 0$ , and

$$2. \alpha = P_{\theta_0}(T \in S) = \beta(\theta_0)$$

The proof is a result of the Factorization Theorem which gives us  $f(\mathbf{x} | \theta_i) = g(T(\mathbf{x}) | \theta_i) h(\mathbf{x})$ . Applying the N-P lemma, with  $h(\mathbf{x})$  dividing out on both sides of the inequalities, we have the result of the corollary.

so  $g(t)$  does not change the monotone of  $f(x)$ ? since all  $> 0$

**Example 8.10** Let  $X \sim \text{Binomial}(2, \theta)$ . We want to test  $H_0 : \theta = 1/2$  versus  $H_1 : \theta = 3/4$ . Calculating the ratios of the pmfs gives

$$\frac{f(0 | \theta = 3/4)}{f(0 | \theta = 1/2)} = \frac{1}{4}, \quad \frac{f(1 | \theta = 3/4)}{f(1 | \theta = 1/2)} = \frac{3}{4}, \quad \text{and} \quad \frac{f(2 | \theta = 3/4)}{f(2 | \theta = 1/2)} = \frac{9}{4}$$

- If we choose  $3/4 < k < 9/4$ , the Neyman-Pearson Lemma says that the test that rejects  $H_0$  if  $X = 2$  is the UMP level  $\alpha = P\{X = 2 | \theta = 1/2\} = 1/4$  test.
- If we choose  $1/4 < k < 3/4$ , the Neyman-Pearson Lemma says that the test that rejects  $H_0$  if  $X = 1$  or  $2$  is the UMP level  $\alpha = P\{X = 1 \text{ or } 2 | \theta = 1/2\} = 3/4$  test.
- Choosing  $k < 1/4$  or  $k > 9/4$  yields the UMP level  $\alpha = 1$  or level  $\alpha = 0$  test. ■

**Example 8.11** - Note that if  $k = 3/4$ , then the NP testing rule says we must reject  $H_0$  for the sample point  $x = 2$  and accept  $H_0$  for  $x = 0$  but leaves our action for  $x = 1$  undetermined. But if we accept  $H_0$  for  $x = 1$ , we get the UMP level  $\alpha = 1/4$  test as before. If we reject  $H_0$  for  $x = 1$ , we get the UMP level  $\alpha = 3/4$  test as before. Such is the case with discrete distributions. No problem with continuous ones.

What is your opinion of this?:

- Suppose we want a  $\alpha$ -size test where  $\alpha < 1/4$ .
- We accept  $H_0$  if  $x < 2$ .
- If  $x = 2$  then we generate a uniform value  $u$  and reject  $H_0$  if  $u \leq 4\alpha$ . Otherwise, we accept  $H_0$ .
- Is this a size-  $\alpha$  test? Do we still have an UMP test? ■

**Definition 8.12** Monotone Likelihood Ratio (MLR)

A family of pdfs or pmfs  $\{g(t | \theta) : \theta \in \Theta\}$  for a univariate random variable  $T$  with real-valued parameter  $\theta$  has a monotone likelihood ratio (MLR) if, for every  $\theta_2 > \theta_1$ ,  $\frac{g(t|\theta_2)}{g(t|\theta_1)}$  is a monotone (nonincreasing or nondecreasing) function of  $t$  on  $\{t : g(t | \theta_1) > 0 \text{ or } g(t | \theta_2) > 0\}$

**Theorem 8.3** Karlin-Rubin Theorem

Consider testing the composite hypotheses  $H_0 : \theta \leq \theta_0$  vs  $H_1 : \theta > \theta_0$ . Suppose that  $T$  is a sufficient statistic for  $\theta$  and the family of pdfs or pmfs  $\mathcal{F} = \{g(t | \theta) : \theta \in \Theta\}$  has a nondecreasing MLR. Then for any  $t_0$ , the test that rejects  $H_0$  if and only if  $T > t_0$  is a UMP level-  $\alpha$  test, where  $\alpha = P_{\theta_0}(T > t_0) = \beta(\theta_0)$ .

- The Karlin-Rubin Theorem is essentially the Neyman-Pearson Lemma for composite hypotheses.

- Since  $\mathcal{F}$  has a nondecreasing MLR, we can show the test's power function,  $\beta(\theta)$ , is nondecreasing. (Exercise) This leads to the theorem's proof (C&B, p. 391-392).

- Under conditions of the theorem, the test above is also UMP for  $H_0 : \theta = \theta_0$  vs  $H_1 : \theta > \theta_0$  (limited the parameter space into  $\theta \in (\theta_0, \infty)$ )

**Example 8.12** - How would you restate the theorem for  $H_0 : \theta \geq \theta_0$  vs  $H_1 : \theta < \theta_0$  ?

*Proof:*

Let  $\beta(\theta) = P_\theta(T > t_0)$  be the power function of the test. Fix  $\theta' > \theta_0$  and consider testing  $H'_0 : \theta = \theta_0$  versus  $H'_1 : \theta = \theta'$ . Since the family of pdfs or pmfs of  $T$  has an MLR,  $\beta(\theta)$  is nondecreasing (see Exercise 8.34), so

i.  $\sup_{\theta \leq \theta_0} \beta(\theta) = \beta(\theta_0) = \alpha$ , and this is a level  $\alpha$  test.

ii. If we define

$$k' = \inf_{t \in \mathcal{T}} \frac{g(t | \theta')}{g(t | \theta_0)},$$

where  $\mathcal{T} = \{t : t > t_0 \text{ and either } g(t | \theta') > 0 \text{ or } g(t | \theta_0) > 0\}$ , it follows that

$$k > 0$$

and

$$T > t_0 \Leftrightarrow \frac{g(t | \theta')}{g(t | \theta_0)} > k'.$$

Together with Corollary 8.3.13, (i) and (ii) imply that this is an UMP and follows by UMP is unbiased,  $\beta(\theta') \geq \beta^*(\theta')$ , where  $\beta^*(\theta)$  is the power function for any other level  $\alpha$  test of  $H'_0$ , that is, any test satisfying  $\beta(\theta_0) \leq \alpha$ . However, any level  $\alpha$  test of  $H_0$  satisfies  $\beta^*(\theta_0) \leq \sup_{\theta \in \theta_0} \beta^*(\theta) \leq \alpha$ . Thus,  $\beta(\theta') \geq \beta^*(\theta')$  for any level  $\alpha$  test of  $H_0$ . Since  $\theta'$  was arbitrary, the test is a UMP level  $\alpha$  test. ■

is it because of monotone or Unbiased?

**Example 8.13** Define

$$\varphi(t) = \begin{cases} 1 & t > t_0 \\ \gamma & t = t_0 \\ 0 & t < t_0 \end{cases}$$

-  $\varphi(T)$  represents the test where we reject  $H_0$  with probability 1 if  $T > t_0$ ; with probability 0 if  $T < t_0$  (we fail to reject); and with probability  $\gamma$  if  $T = t_0$ .

- The values  $t_0$  and  $\gamma \in [0, 1]$  are chosen so that

$$\beta(\theta_0) = E_{\theta_0}[\varphi(T)] = P_{\theta_0}\{T > t_0\} + \gamma P_{\theta_0}\{T = t_0\} = \alpha.$$

This will allow us to construct a size  $\alpha$  test.

- This use of randomization when  $T = t_0$  is controversial but is not needed if  $T$  is a continuous random variable. ■

## Likelihood Ratio Test

- Learn how to construct Likelihood Ratio and how to find the cut off points

Table 2: Related Problem Sets

Year	Questions	Outcome
2021-T	3.d	Bad

## 8 Hypothesis Testing

### 8.1 classic hypothesis testing

**Definition 8.1** *Statistical Hypothesis*

- A statistical hypothesis is an assertion or conjecture about the distribution of one or more random variables. Usually, we will know the distribution to be a member of a family of distributions indexed by  $\theta \in \Theta$ .

**Definition 8.2** *Two complementary hypotheses*

The two complementary hypotheses in a hypothesis testing problem are called the null hypothesis and the alternative hypothesis. They are denoted by  $H_0$  and  $H_1$ , respectively.

- Again, suppose we posit that the mean,  $\mu$ , of a population is positive. By default, we may assume the complement is true, that is, the mean is zero or negative. Typically, the default hypothesis is the null hypotheses. Here we have  $H_0 : \mu \leq 0$  and  $H_1 : \mu > 0$

- In general, we specify the null hypothesis as  $H_0 : \theta \in \Theta_0$  and the alternative hypothesis as  $H_1 : \theta \in \Theta_1$  where  $\Theta_0, \Theta_1 \subset \Theta$ ,  $\Theta_0 \cup \Theta_1 = \Theta$ , and  $\Theta_0 \cap \Theta_1 = \emptyset$ ; that is,  $\Theta_0$  and  $\Theta_1$  are complementary subsets of  $\Theta$ .

The subset of the sample space,  $\mathcal{X}$ , for which  $H_0$  will be rejected is called the rejection region,  $R$ , or critical region. The complement of the rejection region is called the acceptance region,  $A$ .

In general, we denote the rejection region as  $R \subset \mathcal{X}$  and we

- Reject  $H_0$  and accept  $H_1$  as true if  $\mathbf{x} \in R$
- Accept  $H_0$  as true if  $\mathbf{x} \notin R$
- We do not support or accept the null hypothesis! We merely support our research hypothesis or fail to support it.
- The null hypothesis is a construct used solely for hypothesis testing. Failing to support the research hypothesis must not be confused with supporting the null hypothesis.

#### 8.1.1 Likelihood Ratio Test

**Definition 8.3** *The likelihood ratio test statistic for testing  $H_0 : \theta \in \Theta_0$  versus  $H_1 : \theta \in \Theta_1$ , where  $\Theta_1 = \Theta_0^c$ , is*

$$\lambda(\mathbf{x}) = \frac{\sup_{\Theta_0} L(\theta | \mathbf{x})}{\sup_{\Theta} L(\theta | \mathbf{x})}$$

A likelihood ratio test (LRT) is any test that has a rejection region of the form  $R = \{\mathbf{x} : \lambda(\mathbf{x}) \leq c\}$ , where  $c$  is some known number satisfying  $0 \leq c \leq 1$

- The denominator of the likelihood ratio statistic is the supremum of the likelihood function over all values of the parameter space. The numerator is the supremum of the likelihood function restricted to the subset of the parameter space corresponding to the null hypothesis.
- Often, we are interested in the natural log of the LRT statistic.

**Example 8.1** Suppose  $\Theta = \{\theta_0, \theta_1\}$  where  $\theta_0 \neq \theta_1$ . Further suppose that  $H_0 : \theta = \theta_0$ . Then the LRT statistic is

$$\lambda(\mathbf{x}) = \frac{L(\theta_0 | \mathbf{x})}{\max(L(\theta_0 | \mathbf{x}), L(\theta_1 | \mathbf{x}))}$$

- Note that in most cases  $\sup_{\theta} L(\theta | \mathbf{x}) = L(\hat{\theta} | \mathbf{x})$ , where  $\hat{\theta}$  is the overall MLE, and  $\sup_{\theta_0} L(\theta | \mathbf{x}) = L(\hat{\theta}_0 | \mathbf{x})$ , where  $\hat{\theta}_0$  is the MLE over  $\Theta_0$ . In this case, the LRT statistic is given by

$$\lambda(\mathbf{x}) = \frac{L(\hat{\theta}_0 | \mathbf{x})}{L(\hat{\theta} | \mathbf{x})}$$

**Example 8.2** Suppose we have an iid random sample from  $f(x | \theta) = \frac{1}{\theta} I_{[0, \theta]}(x)$ , where  $\Theta = (0, \infty)$ . We wish to test  $H_0 : \theta \geq 1$  versus  $H_1 : \theta < 1$ . The likelihood function is  $L(\theta | x) = (\theta)^{-n} I_{[0, \theta]}(x_{(n)})$ .

- The overall observed MLE is  $\hat{\theta} = X_{(n)}$  and  $L(\hat{\theta} | x) = (x_{(n)})^{-n}$ .
- $\Theta_0 = [1, \infty)$ . If  $X_{(n)} \geq 1$  then  $\hat{\theta}_0 = \hat{\theta} = X_{(n)}$ . Otherwise, if  $X_{(n)} < 1$  then  $\hat{\theta}_0 = 1$  (why?) and  $L(1 | x) = 1$ .
- The LRT statistic is thus

$$L(\hat{\theta}_0 | \mathbf{x}) = \begin{cases} L(\hat{\theta} | \mathbf{x}) & \text{if } x_{(n)} \geq 1 \\ 1 & \text{if } x_{(n)} < 1 \end{cases} \quad \text{and } \lambda(\mathbf{x}) = \begin{cases} 1 & \text{if } x_{(n)} \geq 1 \\ (x_{(n)})^{-n} & \text{if } x_{(n)} < 1 \end{cases}$$

- The rejection region is given as  $\{\mathbf{x} : \lambda(\mathbf{x}) \leq c\}$ . If we choose  $c = 1$ , then we will always reject  $H_0$ . If we choose  $0 \leq c < 1$ , then  $\lambda(\mathbf{x}) \leq c$  only when  $x_{(n)} < 1$

- Hence the LRT is given as Reject  $H_0 : \theta \geq 1$  if  $(x_{(n)})^n \leq c$  which is equivalent to  $x_{(n)} \leq c^{1/n}$ .

**Theorem 8.1** If  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$  and  $\lambda^*(t)$  and  $\lambda(\mathbf{x})$  are the LRT statistics based on  $T$  and  $\mathbf{X}$ , respectively, then  $\lambda^*(T(\mathbf{x})) = \lambda(\mathbf{x})$  for every  $\mathbf{x}$  in the sample space,  $\mathcal{X}$ .

**Example 8.3** Survival times for advanced-stage colon cancer patients are Assumed to be distributed as

$$f(x | \alpha) = \alpha e^{-\alpha x} I_{(0, \infty)}(x); \alpha > 0.$$

Suppose we have two treatments which are proposed to affect survival times.

1. Let  $X_1, \dots, X_n$  iid  $\sim f(x | \alpha)$  denote a random sample of patients given treatment 1.
2. Let  $Y_1, \dots, Y_n$  iid  $\sim f(y | \beta)$  denote a random sample of patients given treatment 2.

prove by factorization and substitution

The two samples are combined into one sample. Assume the two samples are independent. Note that  $(\bar{X}, \bar{Y})$  are jointly sufficient.

*Likelihood Function*

$$\begin{aligned} L(\alpha, \beta \mid \bar{x}, \bar{y}) &= L(\alpha \mid \bar{x}) \cdot L(\beta \mid \bar{y}), \text{ due to independence.} \\ &= \alpha^n e^{-n\alpha\bar{x}} \beta^n e^{-n\beta\bar{y}} \\ l(\alpha, \beta \mid \bar{x}, \bar{y}) &= n \ln \alpha - n\alpha\bar{x} + n \ln \beta - n\beta\bar{y} \end{aligned}$$

*MLE(General case without any constrain from the hypothesis)*

$$\begin{aligned} \frac{\partial l}{\partial \alpha} &= \frac{n}{\alpha} - n\bar{x} = 0 \Rightarrow \hat{\alpha} = \frac{1}{\bar{X}} \\ \left( \frac{\partial^2 l}{\partial \alpha^2} &= -\frac{n}{\alpha^2} < 0, \text{ so } \hat{\alpha} \text{ provides a maximum} \right) \\ \text{Similarly, } \hat{\beta} &= \frac{1}{\bar{Y}} \end{aligned}$$

*LRT for  $H_0 : \alpha = \beta$  vs  $H_1 : \alpha \neq \beta$*

*Restricted MLE(based on  $H_0$ )*

$$\begin{aligned} \text{set } \alpha &= \beta \\ l(\alpha, \alpha \mid \bar{x}, \bar{y}) &= 2n \ln \alpha - n\alpha(\bar{x} + \bar{y}) \\ \frac{\partial l}{\partial \alpha} &= \frac{2n}{\alpha} - n(\bar{x} + \bar{y}) \\ \hat{\alpha}_0 &= \frac{2}{\bar{x} + \bar{y}} = \hat{\beta}_0 \end{aligned}$$

*LR test statistic*

$$\lambda = \frac{L(\hat{\alpha}_0, \hat{\beta}_0 \mid \bar{x}, \bar{y})}{L(\hat{\alpha}, \hat{\beta} \mid \bar{x}, \bar{y})} = \left[ \frac{4\bar{x}\bar{y}}{(\bar{x} + \bar{y})^2} \right]^n$$

Let  $U = \bar{X}/(\bar{X} + \bar{Y})$ . Then ( $U$  here is the  $T(\cdot)$ )

$$\lambda = [4u(1-u)]^n.$$

*LR test  $R = \{\mathbf{x}, \mathbf{y} : \lambda \leq c\}$  for  $c \in (0, 1)$ . Reject  $H_0$  (support  $H_1$ ) if  $(\mathbf{x}, \mathbf{y}) \in R$ . Note that*

$$\begin{aligned} \lambda \leq c &\Rightarrow u(1-u) \leq c^*, \text{ where } c^* = c^{1/n}/4 \\ &\Rightarrow u \leq \frac{1 - \sqrt{1 - 4c^*}}{2} \text{ or } u \geq \frac{1 + \sqrt{1 - 4c^*}}{2} \\ &\Rightarrow \left| u - \frac{1}{2} \right| \geq k, \text{ where } k = \sqrt{1 - 4c^*} \end{aligned}$$

So equivalently, we may say reject  $H_0$  (support  $H_1$ ) if  $|U - \frac{1}{2}| \geq k, k \in (0, 1)$

■

**Example 8.4 (CONTINUE)**

LRT for  $H_0 : \alpha \geq \beta$  vs  $H_1 : \alpha < \beta$  The full parameter space is  $\Theta = \{(\alpha, \beta) : \alpha > 0, \beta > 0\}$ . The restricted (null) parameter space is  $\Theta_0 = \{(\alpha, \beta) : 0 < \beta \leq \alpha\}$ . We consider two cases.

1.  $\bar{x} \leq \bar{y} \Rightarrow \frac{1}{\bar{x}} \geq \frac{1}{\bar{y}} : \text{In this case, } \hat{\alpha} = \frac{1}{\bar{X}} \geq \hat{\beta} = \frac{1}{\bar{Y}}. \text{ So } \hat{\theta} = \hat{\theta}_0 \in \Theta_0.$

2.  $\bar{x} > \bar{y} \Rightarrow \frac{1}{\bar{x}} < \frac{1}{\bar{y}} : \text{In this case, } \hat{\alpha} = \frac{1}{\bar{X}} < \hat{\beta} = \frac{1}{\bar{Y}}. \text{ So } \hat{\theta} \in \Theta_0^c. \text{ To maximize the likelihood function for } \theta \in \Theta_0 \text{ the restricted MLE } \hat{\theta}_0 \text{ should be close to } \hat{\theta} \text{ which puts it on the border of } \Theta_0 \text{ where } \alpha = \beta. \text{ We've already found that } \hat{\alpha}_0 = \hat{\beta}_0 = \frac{2}{\bar{x} + \bar{y}} \text{ provides the maximum likelihood when } \alpha = \beta.$

LR Test Statistic

$$\lambda = \begin{cases} 1 & \text{for } \bar{x} \leq \bar{y} \\ [4u(1-u)]^n & \text{for } \bar{x} > \bar{y} \end{cases}$$

LR Test

$$R = \{\mathbf{x}, \mathbf{y} : \lambda \leq c\} \text{ for } c \in (0, 1).$$

Reject  $H_0$  (support  $H_1$ ) if  $(\mathbf{x}, \mathbf{y}) \in R$ .

Note that  $\bar{x} > \bar{y}$  iff  $u > 1/2$ . Thus,

$$\lambda \leq c \text{ when } u > \frac{1}{2} + k$$

So equivalently, we may say reject  $H_0$  (support  $H_1$ ) if  $U > \frac{1}{2} + k$  ■

**Example 8.5** Suppose  $X_1, \dots, X_n$  are a random sample from a  $N(\mu, \sigma^2)$  distribution, and an experimenter is interested only in inferences about  $\mu$ , such as testing  $H_0 : \mu \leq \mu_0$  versus  $H_1 : \mu > \mu_0$ . Then the parameter  $\sigma^2$  is a nuisance parameter. The LRT statistic is

$$\lambda(\mathbf{x}) = \frac{\max_{\{\mu, \sigma^2 : \mu \leq \mu_0, \sigma^2 \geq 0\}} L(\mu, \sigma^2 | \mathbf{x})}{\max_{\{\mu, \sigma^2 : -\infty < \mu < \infty, \sigma^2 \geq 0\}} L(\mu, \sigma^2 | \mathbf{x})}$$

Show that the LRT can be based on the Student's  $t$  statistic.

$$\begin{aligned} \lambda(\mathbf{x}) &= \frac{\max_{\{\mu, \sigma^2 : \mu \leq \mu_0, \sigma^2 \geq 0\}} L(\mu, \sigma^2 | \mathbf{x})}{\max_{\{\mu, \sigma^2 : -\infty < \mu < \infty, \sigma^2 \geq 0\}} L(\mu, \sigma^2 | \mathbf{x})} \\ &= \frac{L(\hat{\mu}_0, \hat{\sigma}_0^2 | \mathbf{x})}{L(\hat{\mu}, \hat{\sigma}^2 | \mathbf{x})} \end{aligned}$$

where

$$L(\mu, \sigma^2 | \mathbf{x}) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}$$



The (unrestricted) observed MLEs of  $\mu$  and  $\sigma^2$  are  $\hat{\mu} = \bar{x}$ ,  $\hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$ . Plugging these in we get

$$L(\hat{\mu}, \hat{\sigma}^2 | \mathbf{x}) = \left( \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} \right)^n \exp \left\{ -\frac{n}{2} \right\}$$

Now consider the restricted MLE where  $\mu \leq \mu_0$ . If  $\bar{x} \leq \mu_0$  then this restricted (observed) MLE is the same as the unrestricted MLE. Suppose  $\bar{x} > \mu_0$ . Recall that  $\sum (x_i - \mu)^2 = \sum (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2$ . To maximize  $L$  we need to minimize  $(\bar{x} - \mu)^2$ . Thus we choose  $\mu$  to be as close to  $\bar{x}$  as we can get; that is,  $\hat{\mu}_0 = \mu_0$ . Thus

$$L(\mu_0, \sigma^2 | \mathbf{x}) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum (x_i - \mu_0)^2 \right\}$$

Maximizing this with respect to  $\sigma^2$  yields  $\hat{\sigma}_0^2 = \frac{1}{n} \sum (x_i - \mu_0)^2$  and

$$L(\mu_0, \hat{\sigma}_0^2 | \mathbf{x}) = \left( \frac{1}{\sqrt{2\pi\hat{\sigma}_0^2}} \right)^n \exp \left\{ -\frac{n}{2} \right\}$$

Pulling it together, the LRT statistic is

$$\lambda(\mathbf{x}) = \begin{cases} 1 & \text{if } \bar{x} \leq \mu_0 \\ \frac{L(\mu_0, \hat{\sigma}_0^2 | \mathbf{x})}{L(\hat{\mu}, \hat{\sigma}^2 | \mathbf{x})} = \left( \frac{\hat{\sigma}^2}{\hat{\sigma}_0^2} \right)^{\frac{n}{2}} & \text{if } \bar{x} > \mu_0 \end{cases}$$

The rejection region is

$$R = \left\{ \mathbf{x} : \left( \frac{\hat{\sigma}^2}{\hat{\sigma}_0^2} \right)^{\frac{n}{2}} \leq c \right\}$$

To simplify, note that

$$\frac{\hat{\sigma}^2}{\hat{\sigma}_0^2} = \frac{\sum (x_i - \bar{x})^2}{\sum (x_i - \mu_0)^2} = \frac{\sum (x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2 + n(\bar{x} - \mu_0)^2} = \frac{1}{1 + \frac{n(\bar{x} - \mu_0)^2}{\sum (x_i - \bar{x})^2}}$$

An equivalent representation of the rejection region is thus

$$R = \left\{ \mathbf{x} : \frac{n(\bar{x} - \mu_0)^2}{\sum (x_i - \bar{x})^2} \geq k^* \right\}$$

Recognizing that  $\frac{n(\bar{x} - \mu_0)^2}{\sum (x_i - \bar{x})^2} = \frac{1}{n-1} \left( \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \right)^2$  we have the hypothesis test rule

Reject  $H_0$  in favor of  $H_1$  if  $\frac{\bar{x} - \mu_0}{s/\sqrt{n}} > k$  for some  $k > \mu_0$ .

We know that  $\frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim \text{t-dist} (df = n - 1)$ . ■

**Example from c.b:****Example 8.2.2**

Let  $X_i \sim N(\theta, 1)$ , test for  $H_0 : \theta = \theta_0$

$$\begin{aligned}\lambda(\mathbf{x}) &= \frac{(2\pi)^{-n/2} \exp \left[ -\sum_{i=1}^n (x_i - \theta_0)^2 / 2 \right]}{(2\pi)^{-n/2} \exp \left[ -\sum_{i=1}^n (x_i - \bar{x})^2 / 2 \right]} \\ &= \exp \left[ \left( -\sum_{i=1}^n (x_i - \theta_0)^2 + \sum_{i=1}^n (x_i - \bar{x})^2 \right) / 2 \right].\end{aligned}$$

The expression for  $\lambda(\mathbf{x})$  can be simplified by noting that

$$\sum_{i=1}^n (x_i - \theta_0)^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \theta_0)^2$$

Thus the LRT statistic is

$$\lambda(\mathbf{x}) = \exp \left[ -n(\bar{x} - \theta_0)^2 / 2 \right]$$

An LRT is a test that rejects  $H_0$  for small values of  $\lambda(\mathbf{x})$ . From (8.2.2), the rejection region,  $\{\mathbf{x} : \lambda(\mathbf{x}) \leq c\}$ , can be written as

$$\{\mathbf{x} : |\bar{x} - \theta_0| \geq \sqrt{-2(\log c)/n}\}$$

As  $c$  ranges between 0 and 1,  $\sqrt{-2(\log c)/n}$  ranges between 0 and  $\infty$ . Thus the LRTs are just those tests that reject  $H_0 : \theta = \theta_0$  if the sample mean differs from the hypothesized value  $\theta_0$  by more than a specified amount.

The analysis in Example 8.2.2 is typical in that first the expression for  $\lambda(\mathbf{X})$  from Definition 8.2.1 is found, as we did in (8.2.1). Then the description of the rejection region is simplified, if possible, to an expression involving a simpler statistic,  $|\bar{X} - \theta_0|$  in the example.

**Example 8.2.3 (Exponential LRT)**

Let  $X_1, \dots, X_n$  be a random sample from an exponential population with pdf

$$f(x | \theta) = \begin{cases} e^{-(x-\theta)} & x \geq \theta \\ 0 & x < \theta \end{cases}$$

where  $-\infty < \theta < \infty$ . The likelihood function is

$$L(\theta | \mathbf{x}) = \begin{cases} e^{-\sum x_i + n\theta} & \theta \leq x_{(1)} \\ 0 & \theta > x_{(1)}. \end{cases} \quad (x_{(1)} = \min x_i)$$

Consider testing  $H_0 : \theta \leq \theta_0$  versus  $H_1 : \theta > \theta_0$ , where  $\theta_0$  is a value specified by the experimenter. Clearly  $L(\theta | \mathbf{x})$  is an increasing function of  $\theta$  on  $-\infty < \theta \leq x_{(1)}$ . Thus, the denominator of  $\lambda(\mathbf{x})$ , the unrestricted maximum of  $L(\theta | \mathbf{x})$ , is

$$L(x_{(1)} | \mathbf{x}) = e^{-\sum x_i + nx_{(1)}}.$$

If  $x_{(1)} \leq \theta_0$ , the numerator of  $\lambda(\mathbf{x})$  is also  $L(x_{(1)} | \mathbf{x})$ . But since we are maximizing  $L(\theta | \mathbf{x})$  over  $\theta \leq \theta_0$ , the numerator of  $\lambda(\mathbf{x})$  is  $L(\theta_0 | \mathbf{x})$  if  $x_{(1)} > \theta_0$ . Therefore, the likelihood ratio test statistic is

$$\lambda(\mathbf{x}) = \begin{cases} 1 & x_{(1)} \leq \theta_0 \\ e^{-n(x_{(1)} - \theta_0)} & x_{(1)} > \theta_0 \end{cases}$$

A graph of  $\lambda(\mathbf{x})$  is shown in Figure 8.2.1. An LRT, a test that rejects  $H_0$  if  $\lambda(\mathbf{X}) \leq c$ , is a test with rejection region  $\{\mathbf{x} : x_{(1)} \geq \theta_0 - \frac{\log c}{n}\}$ . Note that the rejection region depends on the sample only through the sufficient statistic  $X_{(1)}$ . That this is generally the case will be seen in Theorem 8.2.4.

Example 8.2.3 again illustrates the point, expressed in Section 7.2.2, that differentiation of the likelihood function is not the only method of finding an MLE. In Example 8.2.3,  $L(\theta | \mathbf{x})$  is not differentiable at  $\theta = x_{(1)}$ .

### Example 8.2.5 (LRT and sufficiency)

In Example 8.2.2, we can recognize that  $\bar{X}$  is a sufficient statistic for  $\theta$ . We could use the likelihood function associated with  $\bar{X}$  ( $\bar{X} \sim n(\theta, \frac{1}{n})$ ) to more easily reach the conclusion that a likelihood ratio test of  $H_0 : \theta = \theta_0$  versus  $H_1 : \theta \neq \theta_0$  rejects  $H_0$  for large values of  $|\bar{X} - \theta_0|$ .

Similarly in Example 8.2.3,  $X_{(1)} = \min X_i$  is a sufficient statistic for  $\theta$ . The likelihood function of  $X_{(1)}$  (the pdf of  $X_{(1)}$ ) is

$$L^*(\theta | x_{(1)}) = \begin{cases} ne^{-n(x_{(1)} - \theta)} & \theta \leq x_{(1)} \\ 0 & \theta > x_{(1)}. \end{cases}$$

This likelihood could also be used to derive the fact that a likelihood ratio test of  $H_0 : \theta \leq \theta_0$  versus  $H_1 : \theta > \theta_0$  rejects  $H_0$  for large values of  $X_{(1)}$ .

## Large Sample Test Method

## 8.4 Large Sample ML-based method

Set up:

For large sample size we may use approximate asymptotic methods for hypothesis testing. These are based on the maximum likelihood and thus depend on the assumed family of distributions

-We will assume that the family of distributions  $\mathcal{T}$  is indexed by a p-dimensiona parameter vector =  $\theta$

-Hypotheses:  $H_0 : \theta \in \Theta_0$  versus  $H_1 : \theta \in \Theta_1$  where  $\Theta_0 \in \Theta$  and  $\Theta_1 \in \Theta_0^c$  is the complement of  $\Theta_0$  With respect to  $\Theta$

Given the Fisher's Information Matrix:

$$-E_{\theta} \left[ \frac{\partial^2 l(\theta|x)}{\partial \theta_i \partial \theta_j} \right]$$

We have

·Let  $R_i(\theta) = 0, i = 1, 2, \dots, r$ , represent  $r(< p)$  independent restrictions placed on the parameter vector  $\theta$

·Consider  $\Theta_0 = \{\theta; R_i(\theta) = 0, i = 1, 2, \dots, r\}$ .

·The composite hypotheses are  $H_0: \theta \in \Theta_0$  versus  $H_1: \theta \in \Theta_1$

·For example, with  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)$ , consider the  $r = 3$  linearly independent linear restrictionis  $R_1(\theta) = (\theta_1 - \theta_2) = 0$ ,  $R_2(\theta) = (\theta_1 - \theta_3) = 0$  and  $R_3(\theta) = (\theta_1 - \theta_4) = 0$  Then, the null hypothesis  $H_0 : R(\theta) = 0 \text{ for } i = 1, 2, 3$ , is equivalent to the null hypothesis

$$H_0: \theta_1 = \theta_2 = \theta_3 = \theta_4$$

OR We could just start with a specified  $\mathfrak{G}_{\mathfrak{G}}$  and denote its dimension as  $d_d$ . The dimension of  $\Theta$  is p. Then  $r = p - d_d$ , the dimension of the space of restrictions

### 8.4.1 Likelihood Ratio Test

The likelihood ratio test statsitics is defined as:

$$\hat{\lambda} = \frac{\max_{\theta \in \Theta_0} L(\theta|x)}{\max_{\theta \in \Theta} L(\theta|x)} = \frac{L(\hat{\theta}_0|x)}{L(\hat{\theta}|x)}$$

Under certain regularity conditions, for large n and assuming  $H_0: \theta \in \Theta_0$ ,

$$-2 \log \hat{\lambda} = 2 \left( l(\hat{\theta}|x) - l(\hat{\theta}_0|x) \right) \rightarrow^d \chi_r^2;$$

#### Proof

First expand  $\log L(\theta | \mathbf{x}) = l(\theta | \mathbf{x})$  in a Taylor series around  $\hat{\theta}$ , giving

$$l(\theta | \mathbf{x}) = l(\hat{\theta} | \mathbf{x}) + l'(\hat{\theta} | \mathbf{x})(\theta - \hat{\theta}) + l''(\hat{\theta} | \mathbf{x}) \frac{(\theta - \hat{\theta})^2}{2!} + \dots$$

Now substitute the expansion for  $l(\theta_0 | \mathbf{x})$  in  $-2 \log \lambda(\mathbf{x}) = -2l(\theta_0 | \mathbf{x}) + 2l(\hat{\theta} | \mathbf{x})$ , and get

$$\begin{aligned} -2 \log \lambda(\mathbf{x}) &= -2l(\hat{\theta} | \mathbf{x}) + -2 \underbrace{l'(\hat{\theta} | \mathbf{x})}_{\text{score}=0} (\theta_0 - \hat{\theta}) + -2l''(\hat{\theta} | \mathbf{x}) \frac{(\theta_0 - \hat{\theta})^2}{2!} + 2l(\hat{\theta} | \mathbf{x}) \\ &= -2l''(\hat{\theta} | \mathbf{x}) (\theta_0 - \hat{\theta})^2 \\ &\approx \frac{(\theta - \hat{\theta})^2}{I_n(\mathbf{x})}, \end{aligned}$$

where we use the fact that  $l'(\hat{\theta} | \mathbf{x}) = 0$ . Since the denominator is the observed information  $\hat{I}_n(\hat{\theta})$  and  $\frac{1}{n} \hat{I}_n(\hat{\theta}) \rightarrow I(\theta_0)$  it follows from Theorem 10.1.12 and Slutsky's Theorem (Theorem 5.5.17) that  $-2 \log \lambda(\mathbf{X}) \rightarrow \chi_1^2$ .

### 8.4.2 Wald Test

- Define the  $1 \times r$  row vector  $R(\theta)$  as

$$R(\theta) = (R_1(\theta), R_2(\theta), \dots, R_r(\theta))$$

- Let the  $r \times p$  matrix  $T'(\Theta)$  have its  $(i, j)$ th element defined as

$$\frac{\partial R_i(\theta)}{\partial \theta_j}$$

for  $i = 1, 2, \dots, r$  and  $j = 1, 2, \dots, p$

- Let the  $r \times r$  matrix  $A(\Theta)$  have the structure

$$\Lambda(\theta) = T(\theta)J^{-1}(\theta)T'(\theta)$$

.The Wald test statistic  $\hat{W} > 0$  is defined as

$$\tilde{W} = R(\bar{\theta})\Lambda^{-1}(\bar{\theta})R'(\bar{\theta})$$

. Under certain regularity conditions, for large  $n$  and assuming  $H_0: \theta \in \Theta_0$

$$W \rightarrow^d X_r^2$$

### 8.4.3 Score Test

- Define the  $1 \times p$  row vector  $S(\theta)$  as

$$S(\theta) = \left( \frac{\partial l(\theta | \mathbf{x})}{\partial \theta_1}, \frac{\partial l(\theta | \mathbf{x})}{\partial \theta_2}, \dots, \frac{\partial l(\theta | \mathbf{x})}{\partial \theta_p} \right)$$

- The score statistic  $\hat{S}$  is defined as

$$\hat{S} = S \left( \hat{\theta}_0 \right) J^{-1} \left( \hat{\theta}_0 \right) S' \left( \hat{\theta}_0 \right)$$

- Under certain regularity conditions, for large  $n$  and assuming  $H_0 : \theta \in \Theta_0$ ,

$$\hat{S} \rightarrow \chi_r^2 \text{ in distribution}$$

#### 8.4.4 Large Sample Rejection Region

The large sample rejection region for an  $\alpha$ -level LR test is

$$\left\{ \mathbf{x} : -2 \log \hat{\lambda} \geq \chi_{r,\alpha}^2 \right\}$$

where  $\chi_{r,\alpha}^2$  is the upper  $\alpha$ -quantile of the chi-squared distribution with  $r$  degrees of freedom.

- Replace  $-2 \log \hat{\lambda}$  with  $\widehat{W}$  for the Wald test rejection region.
- Replace  $-2 \log \hat{\lambda}$  with  $\hat{S}$  for the Score test rejection region

**Example 8.16** - Let  $X_1, X_2, \dots, X_n$  constitute a random sample of size  $n$  from a  $N(\mu, \sigma^2)$  population.

- Consider testing the null hypothesis

$$R_l(\mu, \sigma^2) = \mu - \mu_0$$

$$H_0 : \mu = \mu_0, 0 < \sigma^2 < +\infty \text{ versus } H_1 : \mu \neq \mu_0, 0 < \sigma^2 < +\infty$$

Note that this test is typically called a test of  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$ .

- The MLEs of  $\mu$  and  $\sigma^2$  are  $\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2) = (\bar{X}, (\frac{n-1}{n}) S^2)$ .
- We have that

$$T_{n-1} = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim t_{n-1} \text{ under } H_0 : \mu = \mu_0$$

#### Likelihood Ratio Test

It can be shown that

$$-2 \log \hat{\lambda} = n \log \left( 1 + \frac{T_{n-1}^2}{n-1} \right)$$

which for large  $n$  has an approximate  $\chi_1^2$  distribution.

$$-2 \log \hat{\lambda} \geq k \Rightarrow |T_{n-1}| \geq k^*$$

#### Wald Test

$$r = 1$$

$$R_1 = \mu - \mu_0$$

It can be shown that

$$\hat{W} = \left( \frac{n}{n-1} \right) T_{n-1}^2$$

which for large  $n$  has an approximate  $\chi_1^2$  distribution.

**Score Test**

It can be shown that

$$\hat{S} = \left[ \frac{\bar{X} - \mu_0}{\hat{\sigma}_0/\sqrt{n}} \right]^2$$

Test stat.

where

$$\hat{\sigma}_0^2 = n^{-1} \sum_{i=1}^n (X_i - \mu_0)^2$$

$$\left| \frac{\bar{x} - \mu_0}{\hat{\sigma}_0/\sqrt{n}} \right|$$

is the estimator of  $\sigma^2$  under the null hypothesis  $H_0 : \mu = \mu_0$ .

For large  $n$ ,  $\hat{S}$  has an approximate  $\chi_1^2$  distribution.

- Both the LRT and the Wald test (for this example) have rejection regions that are asymptotically equivalent to the rejection region

$$\{|T_{n-1}| \geq z_{\alpha/2}\}$$

- The exact LRT has rejection region  $\{|T_{n-1}| \geq t_{r,\alpha/2}\}$  where  $t_{r,\alpha/2}$  is the upper  $(\alpha/2)$ -quantile of a t-distribution with  $r$  degrees of freedom.

- The score test is a modification of  $T_{n-1}$  that uses the MLE of  $\sigma^2$  restricted to  $H_0 : \mu = \mu_0, 0 < \sigma^2 < +\infty$ . Its rejection region is equivalent to

$$\left\{ \left| \frac{\bar{X} - \mu_0}{\hat{\sigma}_0/\sqrt{n}} \right| \geq z_{\alpha/2} \right\}$$

## Confidence Interval

### Confidence Interval: Definition, Coverage Probability, Coverage Coefficient and Length

---

- Check Bob's 2023 final, you are really bad at it
- Note that in a confidence set, which one is Random Variable and which one is not
- How to find the coverage probability

You suck  
at this  
part



## 7 Interval Estimator

### 7.1 Confidence Intervals

#### 7.1.1 Definition, Coverage Probability, coefficient and length

**Definition 7.1** *Confidence Interval:*

An interval estimate of a real-valued parameter  $\theta$  is any pair of functions,  $(L(\mathbf{x}), U(\mathbf{x}))$ , of a sample point  $\mathbf{x}$  that satisfies

$$L(\mathbf{x}) \leq U(\mathbf{x})$$

for all  $\mathbf{x} \in \mathbf{X}$ . The random interval  $\underbrace{(L(\mathbf{X}), U(\mathbf{X}))}_{\text{Random Variables}}$  is called an interval estimator.

Given the sample point  $\mathbf{x}$ , we infer that

$$L(\mathbf{x}) \leq \theta \leq U(\mathbf{x})$$

# Not necessary center, but we will get to why center interval is good later

Although we are mainly concerned with confidence intervals, we occasionally will work with more general sets. When working in general, and not being quite sure of the exact form of our sets, we will speak of confidence sets. A confidence set with confidence coefficient equal to some value, say  $1 - \alpha$ , is simply called a  $1 - \alpha$  confidence set.

**Example 7.1** *Interpretation:* In  $100(2\Phi(c) - 1)\%$  of all random samples of size  $n$  from  $N(\mu, \sigma^2)$ , the interval  $(\bar{X} - c\frac{\sigma}{\sqrt{n}}, \bar{X} + c\frac{\sigma}{\sqrt{n}})$  covers the true population mean  $\mu$ .

**Definition 7.2** *Coverage Probability*

For an interval estimator  $(L(\mathbf{X}), U(\mathbf{X}))$  of a parameter  $\theta$ , the coverage probability of  $(L(\mathbf{X}), U(\mathbf{X}))$  is the probability that the random interval  $(L(\mathbf{X}), U(\mathbf{X}))$  covers the true parameter  $\theta$ . The coverage probability is denoted by

$$P_\theta\{\theta \in (L(\mathbf{X}), U(\mathbf{X}))\}$$

or

$$P\{\theta \in (L(\mathbf{X}), U(\mathbf{X})) \mid \theta\}$$

or

$$P_\theta\{L(X) < \theta, U(X) > \theta\}$$

There are a number of things to be aware of in these definitions. One, it is important to keep in mind that the interval is the random quantity, not the parameter. Therefore, when we write probability statements such as  $P_\theta(\theta \in [L(X), U(X)])$ , these probability statements refer to  $X$ , not  $\theta$ . In other words, think of  $P_\theta(\theta \in [L(X), U(X)])$ , which might look like a statement about a random  $\theta$ , as the algebraically equivalent  $P_\theta[L(X) < \theta, U(X) > \theta]$ , a statement about a random  $X$ .

see Bob's  
Final 4(b)

**Definition 7.3** *Confidence Coefficient*

For an interval estimator  $(L(\mathbf{X}), U(\mathbf{X}))$  of a parameter  $\theta$ , the confidence coefficient of  $(L(\mathbf{X}), U(\mathbf{X}))$  is the infimum of the coverage probabilities over all values of  $\theta \in \Theta$ ,

$$\inf_{\theta \in \Theta} P_{\theta}\{\theta \in (L(\mathbf{X}), U(\mathbf{X}))\}$$

- The confidence coefficient is typically a function of the sample size,  $n$ .
- The coverage probability of an interval estimator is at least the corresponding confidence coefficient.

what is this?

**Example 7.2** Suppose  $X_1, \dots, X_n \sim \text{iid Exponential}(\lambda)$ . Recall that  $E[X] = \sqrt{\text{Var}[X]} = \lambda$ .

- Consider the interval  $(\bar{X} - c\frac{\lambda}{\sqrt{n}}, \bar{X} + c\frac{\lambda}{\sqrt{n}})$  for some constant  $c > 0$ . Is this an interval estimator of  $\lambda$ ?
- Try replacing  $\lambda$  with  $\bar{X}$ .
- $L(\mathbf{X}) = \left(1 - \frac{c}{\sqrt{n}}\right) \bar{X}$ . Since  $\lambda > 0$ , it makes sense to restrict  $L(\mathbf{X}) > 0$ . Thus  $1 - \frac{c}{\sqrt{n}} > 0 \Rightarrow 0 < c < \sqrt{n}$
- $U(\mathbf{X}) = \left(1 + \frac{c}{\sqrt{n}}\right) \bar{X}$
- For a given  $\lambda$ , the coverage probability is

Is this a standard practice?

$$P_{\lambda} \left\{ \left(1 - \frac{c}{\sqrt{n}}\right) \bar{X} \leq \lambda \leq \left(1 + \frac{c}{\sqrt{n}}\right) \bar{X} \right\} = P_{\lambda} \left\{ \frac{1}{1 + \frac{c}{\sqrt{n}}} \leq \frac{\bar{X}}{\lambda} \leq \frac{1}{1 - \frac{c}{\sqrt{n}}} \right\}$$

- Note that  $\frac{\bar{X}}{\lambda} \sim \text{Gamma}(n, \frac{1}{n})$ . Let  $G(\cdot)$  denote the corresponding cdf.
- Then the coverage probability is

$$G\left(\left(1 - \frac{c}{\sqrt{n}}\right)^{-1}\right) - G\left(\left(1 + \frac{c}{\sqrt{n}}\right)^{-1}\right)$$

- Since this does not depend on  $\lambda$ , the coverage probability = confidence coefficient.

**Corollary 7.0.1** *Expected Length:*

The length (or expected length) of an interval estimate is related to the precision of the interval estimate.

**Example 7.3** *For our examples:*

- "Normal with known variance"

$$\text{length} = 2z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

- Note this is a fixed length independent of  $\mu$ .
- "Exponential"

$$\text{expected length} = 2c \frac{\lambda}{\sqrt{n}}$$

- This is also fixed for a given  $\lambda$  but varies linearly as  $\lambda$  varies.
- The two examples are similar since  $\lambda = \text{standard deviation}$ .

■

**Example 7.4** *Final 4:*

4. Let  $X_1, \dots, X_n$  be a random sample from a Poisson population with parameter  $\lambda$  and define  $Y = \sum X_i$ . Suppose  $Y = y_0$  is the observed value of  $Y$  in a given sample of size  $n$ . Recall the identity, from C&B: Example 3.3.1, p. 100, that links the Poisson and Gamma families. Applying that identity, we can show that

$$\sum_{k=0}^{y_0} e^{-n\lambda} \frac{(n\lambda)^k}{k!} = P\{Y \leq y_0 \mid \lambda\} = P\{V_{2(y_0+1)} > 2n\lambda\}$$

and

$$\sum_{k=y_0}^{\infty} e^{-n\lambda} \frac{(n\lambda)^k}{k!} = P\{Y \geq y_0 \mid \lambda\} = P\{V_{2y_0} < 2n\lambda\}$$

where  $V_d$  is a chi squared random variable with  $d$  degrees of freedom.

(a) Given  $\alpha \in (0, 1)$ , confirm that

$$\left\{ \lambda : \frac{1}{2n} \chi_{2y_0, 1-\alpha/2}^2 \leq \lambda \leq \frac{1}{2n} \chi_{2(y_0+1), \alpha/2}^2 \right\}$$

is a  $100(1 - \alpha)\%$  confidence interval for  $\lambda$ , where  $\chi_{d,\delta}^2$  denotes the upper  $\delta$ -probability quantile of a chi-squared distribution with  $d$  degrees of freedom. Define  $\chi_{0,\delta}^2 = 0$ .

(b) Argue that the coverage probability of the confidence interval in part (a) is given by

$$\sum_{k=0}^{\infty} \mathbf{1}_{[L(k), U(k)]}(\lambda) \frac{e^{-n\lambda} (n\lambda)^k}{k!}$$

where  $L(k)$  and  $U(k)$  are the lower and upper bounds of the confidence interval for observed  $Y = k$

(c) Letting  $n = 5$  and  $\alpha = 0.10$ , sketch the coverage probabilities for  $\lambda \in [0, 3]$ . Explain why the graph of the coverage probability has jumps occurring at numerous values of  $\lambda$ . [Suggestions: Design your graph to compute the coverage for a dense number of  $\lambda$  values, say 1001 values between 0 and 3. Zooming in of a portion of the graph ( $\lambda \in$  a subinterval of  $(0, 3)$ ) may aid your discussion. Also, look at the coverage graph when  $n = 1$ .]

**Unbiased CI**

### 7.1.2 Unbiased CI

#### Theorem 7.1 Unbiasedness of a confidence interval

A confidence interval is unbiased if it covers the true parameter value with probability no less than the probability it covers any other parameter value.

Let  $\Theta$  be the parameter space over which the family of distributions  $\mathcal{F}$  is indexed. Let  $\theta_* \in \Theta$  be the true parameter related to the random sample  $\mathbf{X}$  (that is, the sample is drawn from the distribution  $f(x | \theta_*)$ ). Then the confidence interval  $\mathcal{J}(\mathbf{X})$  for  $\theta_*$  is unbiased if

$$\operatorname{argmax}_{\theta \in \Theta} P_{\theta_*} \{\theta \in \mathcal{J}(\mathbf{X})\} = \theta_*$$

**Example 7.5** Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_*, \sigma^2)$  where  $\sigma^2$  is known. Consider the interval estimator

$$\mathcal{J}(\mathbf{X}) = \left( \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right).$$

$$\begin{aligned} P_{\mu_*} \{\mu \in \mathcal{J}(\mathbf{X})\} &= P_{\mu_*} \left\{ \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu_* - (\mu_* - \mu) \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right\} \\ &= P_{\mu_*} \left\{ \frac{\mu - \mu_*}{\sigma/\sqrt{n}} - z_{\alpha/2} \leq \frac{\bar{X} - \mu_*}{\sigma/\sqrt{n}} \leq \frac{\mu - \mu_*}{\sigma/\sqrt{n}} + z_{\alpha/2} \right\} \\ &= \Phi \left( \frac{\mu - \mu_*}{\sigma/\sqrt{n}} + z_{\alpha/2} \right) - \Phi \left( \frac{\mu - \mu_*}{\sigma/\sqrt{n}} - z_{\alpha/2} \right) \end{aligned}$$

show by integral shit  $\rightarrow \leq 2\Phi(z_{\alpha/2}) - 1 = \delta$

Thus,  $\mathcal{J}(\mathbf{X})$  is unbiased for  $\mu_*$ .

#### Example 7.6 CONTINUE:

Alternatively, consider the interval "estimator" with targeted confidence coefficient  $\delta$

$$\mathcal{J}(\mathbf{X}) = \left( -\infty, \bar{X} - z_{\delta/2} \frac{\sigma}{\sqrt{n}} \right) \cup \left( \bar{X} + z_{\delta/2} \frac{\sigma}{\sqrt{n}}, \infty \right).$$

$$\begin{aligned} P_{\mu_*} \{\mu \in \mathcal{J}(\mathbf{X})\} &= 1 - P_{\mu_*} \left\{ \bar{X} - z_{\delta/2} \frac{\sigma}{\sqrt{n}} \leq \mu_* - (\mu_* - \mu) \leq \bar{X} + z_{\delta/2} \frac{\sigma}{\sqrt{n}} \right\} \\ &= 1 - P_{\mu_*} \left\{ \frac{\mu_* - \mu}{\sigma/\sqrt{n}} - z_{\delta/2} \leq \frac{\bar{X} - \mu_*}{\sigma/\sqrt{n}} \leq \frac{\mu_* - \mu}{\sigma/\sqrt{n}} + z_{\delta/2} \right\} \\ &= 1 - \Phi \left( \frac{\mu_* - \mu}{\sigma/\sqrt{n}} + z_{\delta/2} \right) + \Phi \left( \frac{\mu_* - \mu}{\sigma/\sqrt{n}} - z_{\delta/2} \right) \end{aligned}$$

this probability is bigger than it's alpha level  $\geq 1 - (2\Phi(z_{\delta/2}) - 1) = \delta$

Thus,  $\mathcal{J}(\mathbf{X})$  is not unbiased for  $\mu_*$ .

$\mu, \mu_*$   
which is  
the R.V  
here?

### Find CI by Pivoting CDF

- More example in [homework4](#), that shit hard
- Need more example
- Example 7.7 is good to learn how to find pivot function

## 7.1.3 Pivoting

**Definition 7.4** *Pivotal Quantity*

A random variable  $Q(\mathbf{X}, \theta)$  is a pivotal quantity if its distribution is independent of  $\theta$ . That is  $Q(\mathbf{X}, \theta)$  has the same distribution for all values of  $\theta$ .

**Lemma 7.2** Suppose we have a statistic  $T$  with pdf  $f_T(t | \theta)$  that may be written as

$$f_T(t | \theta) = g(Q(t, \theta)) \left| \frac{d}{dt} Q(t, \theta) \right|$$

where  $g(\cdot)$  does not depend on  $\theta$  and  $Q(t, \theta)$  is monotone in  $t$ . Let  $y = Q(t, \theta)$ , then  $t = Q^{-1}(y, \theta)$ . As long as  $\frac{d}{dt} Q(t, \theta) \neq 0$ , then  $\frac{d}{dy} Q^{-1}(y, \theta) = 1 / \frac{d}{dt} Q(t, \theta)$ . So, making a change of variables,

$$\begin{aligned} f_Y(y | \theta) &= f_T(Q^{-1}(y, \theta) | \theta) \left| \frac{d}{dy} Q^{-1}(y, \theta) \right| \\ &= g(Q(Q^{-1}(y, \theta), \theta)) \left| \frac{d}{dt} Q(t, \theta) \right| \left| 1 / \frac{d}{dt} Q(t, \theta) \right| \\ &= g(y) \end{aligned}$$

This essentially transformation of random variable

what is this???

**Example 7.7** •  $X \sim \text{Uniform}(\mu - \frac{1}{2}, \mu + \frac{1}{2}) \rightarrow f_X(x | \mu) = I_{(\mu - \frac{1}{2}, \mu + \frac{1}{2})}(x) = I_{(-\frac{1}{2}, \frac{1}{2})}(x - \mu) = f(x - \mu)$ , where  $f(x - \mu)$  does not depend on  $\mu$ . Thus  $f(x - \mu)$  is a location pdf and  $\bar{X} - \mu$  is a pivotal quantity.

- The previous Exponential example is of the form  $f(x/\lambda)$ , a scale pdf, and  $\frac{\bar{X}}{\lambda}$  is a pivotal quantity.
- $X \sim N(\mu, \sigma^2) \rightarrow f_X(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right)$  where  $f(t) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{t^2}{2}\right\}$  is the standard normal distribution. Thus  $f\left(\frac{\bar{X}-\mu}{\sigma}\right)$  is a location-scale pdf and  $\frac{\bar{X}-\mu}{\sigma}$  is a pivotal quantity.

**Example 7.8** In one of our previous examples, we had  $\bar{X}/\lambda \sim \text{Gamma}(n, \frac{1}{n})$ . While the quantity  $\bar{X}/\lambda$  depends on  $\lambda$ , its distribution does not. Using this distribution, we can find  $t_1$  and  $t_2$  such that

$$\delta = P\{t_1 \leq \bar{X}/\lambda \leq t_2\} = P\left\{\frac{\bar{X}}{t_2} \leq \lambda \leq \frac{\bar{X}}{t_1}\right\}$$

So  $\left(\frac{\bar{x}}{t_2}, \frac{\bar{x}}{t_1}\right)$  is a  $100\delta\%$  confidence interval for  $\lambda$ . The expected length is  $\left(\frac{1}{t_1} - \frac{1}{t_2}\right)\lambda$ . How does this compare to  $2c\frac{\lambda}{\sqrt{n}}$  from our previous example?

So  $\bar{X}/\lambda$  is a pivotal quantity.

The idea is we can find  $q_1$  and  $q_2$  such that

$$\delta = P\{q_1 \leq Q(\mathbf{X}, \theta) \leq q_2\}$$

and construct the confidence set

$$\{\theta : q_1 \leq Q(\mathbf{X}, \theta) \leq q_2\}$$

- This set is not always an interval, but it is when  $Q(\mathbf{X}, \theta)$  is monotonic in  $\theta$ .

$T = \sum_{i=1}^n X_i/n$  is sufficient for the Exponential  $\lambda$ . We can show that  $T \sim \text{Gamma}(n, \lambda/n)$  where the form of the distribution is

$$f(t | n, \lambda/n) = \frac{n^n}{\Gamma(n)\lambda} \left(\frac{t}{\lambda}\right)^{n-1} e^{-\frac{nt}{\lambda}}, 0 \leq t < \infty, \lambda > 0$$

Let  $Q(T, \lambda) = T/\lambda$  and thus  $\frac{d}{dt}Q(t, \lambda) = 1/\lambda$ . So if

$$g(y) = \frac{n^n}{\Gamma(n)} (y)^{n-1} e^{-ny}$$

then

$$f(t | n, \lambda/n) = g(Q(t, \lambda)) \left| \frac{d}{dt}Q(t, \lambda) \right|$$

■

#### 7.1.4 Wald CI

#### 7.1.5 Wilson's Score CI

#### 7.1.6 Exact CI- Clopper Person

read  
Mathe-  
matics  
Statistics

How to  
construct  
a pivot  
interval

I am so  
confused  
what's the  
goal of  
this proof



**Find CI by MLE asymptotic**

- See How in [Consistency, Asymptotic Variance, efficiency of MLE]
- Not a lot of example for this topics
- Not familiar with this topics

**Definition 7.6** Asymptotic variance

- First an example: We know that the sample mean,  $\bar{X}_n$ , from a  $N(\mu, \sigma^2)$  distribution has variance  $\sigma^2/n$ . Thus,  $\sigma_n^2 = \text{Var}[\sqrt{n}\bar{X}_n] = \sigma^2$ .

- In general, if  $X_1, \dots, X_n$  is an iid random sample from a distribution where the second moment exists (thus the distribution has a finite variance), then the central limit theorem says that  $\frac{\bar{X}_n - E[\bar{X}_n]}{\sqrt{\text{Var}[X]/n}} \rightarrow N(0, 1)$  in distribution.

- Define  $\sigma_n^2 = \text{Var}[\sqrt{n}\bar{X}_n]$ . Then  $\sigma_n^2 \rightarrow \text{Var}[X] = \tau^2$ . The constant  $\tau^2$  is the limiting variance or the limit of the variances.

- Definition: For an estimator  $T_n$ , suppose that  $k_n(T_n - \tau(\theta)) \rightarrow N(0, \sigma^2)$  in distribution. The parameter  $\sigma^2$  is called the asymptotic variance or variance of the limit distribution of  $T_n$ .

This is  
general  
Delta  
methods?

**Definition 7.7** Asymptotic Efficiency

Definition: A sequence of estimators  $W_n$  is asymptotically efficient for a parameter  $\tau(\theta)$  if

$$\sqrt{n}[W_n - \tau(\theta)] \rightarrow N(0, v(\theta))$$

in distribution and

$$v(\theta) = \frac{[\tau'(\theta)]^2}{E_\theta \left( \left( \frac{\partial}{\partial \theta} \ln f(X | \theta) \right)^2 \right)} = \frac{[\tau'(\theta)]^2}{-I(X)}$$

that is, the asymptotic variance of  $W_n$  achieves the Cramér-Rao Lower Bound.

**Example 7.14 :**

If  $\hat{\theta}_{jn}$  satisfies tA.GT with asymptotic covariance matrix  $V_{jn}(\theta)$ ,  $j = 1, 2$ , and  $V_{1,n}(\theta) \leq V_{2,n}(\theta)$  (in the sense that  $V_{2,n}(\theta) - V_{1,n}(\theta)$  is nonnegative definite for all  $\theta \in \Theta$ ), then  $\hat{\theta}_{1n}$  is said to be asymptotically more efficient than  $\hat{\theta}_{2n}$ .

**Theorem 7.6** Asymptotic Efficiency of MLE

Let  $X_1, X_2, \dots$  be an iid random sample from a family of distributions indexed by  $\theta \in \Theta$ . Let  $\hat{\theta}$  denote the MLE of  $\theta$  and let  $\tau(\theta)$  be a continuous function of  $\theta$ . Under certain regularity conditions,

$$\sqrt{n}[\tau(\hat{\theta}) - \tau(\theta)] \rightarrow N(0, v(\theta))$$

where  $v(\theta)$  is the Cramér-Rao Lower Bound. That is,  $\tau(\hat{\theta})$  is a consistent and asymptotically efficient estimator of  $\tau(\theta)$ .

Thus, the bounds of an approximate  $100(1 - \alpha)\%$  confidence interval for  $\theta$  is given as

$$\tau(\hat{\theta}) \pm z_{\alpha/2} \sqrt{v(\hat{\theta})/n} \text{ or } \tau(\hat{\theta}) \pm z_{\alpha/2} \sqrt{\hat{v}(\hat{\theta})/n}$$

**Corollary 7.6.1** • The quantity  $v(\hat{\theta})$  may be found by applying the delta method to the transformation  $\tau(\cdot)$  evaluated at  $\theta = \hat{\theta}$  and using  $\sigma^2 = I_1^{-1}(\theta)$ ; that is,  $v(\hat{\theta}) = \left. \frac{[\tau'(\theta)]^2}{I_1(\theta)} \right|_{\theta=\hat{\theta}}$

- For a given  $n$ , the approximate (unobservable) variance of  $\tau(\hat{\theta})$  is  $v(\theta)$ . To use this in practice, we must estimate the approximate variance.
  - First estimator: MLE of  $v(\theta)$ , given by  $v(\hat{\theta})$ .
  - This uses  $I_1(\hat{\theta})$ , the expectation of  $-\frac{\partial^2 \ln(f(X_i|\theta))}{\partial \theta^2}$  evaluated at  $\theta = \hat{\theta}$
  - We call this the expected information number.
  - Alternatively, we may estimate  $I_1(\theta)$  with  $-\frac{1}{n} \sum \frac{\partial^2 \ln(f(X_i|\theta))}{\partial \theta^2}$  evaluated at  $\theta = \hat{\theta}$ .
  - We call this the observed information number, denoted as  $\hat{I}_1(\hat{\theta})$ .
  - The variance estimator using the observed information number is denoted as  $\hat{v}(\hat{\theta})$ .
  - It has been shown that this provides a better estimator. (Efron & Hinkley, Biometrika, 1978)
  - If you can differentiate and evaluate the log-likelihood, you can calculate  $\hat{v}(\hat{\theta})$ .

**Example 7.15** Suppose  $X_1, \dots, X_n$  is an iid random sample from Poisson ( $\lambda$ ). The MLE of  $\lambda$  is  $\hat{\lambda} = \bar{X}$ .

- - Construct a 95%CI for  $\lambda$ , where  $n = 100$  and  $\bar{x} = 0.5$ .
- Large sample via the MLE

$$\begin{aligned}
 L(\lambda | x) &= e^{-\lambda} \lambda^x / x! \\
 l(\lambda | x) &= -\lambda + x \log \lambda - \log(x!) \\
 \frac{\partial l}{\partial \lambda} &= -1 + \frac{x}{\lambda}; \quad \frac{\partial^2 l}{\partial \lambda^2} = -\frac{x}{\lambda^2} \\
 I_1(\lambda) &= -E \left[ -\frac{x}{\lambda^2} \right] = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda} \\
 I_1(\hat{\lambda}) &= \frac{1}{\bar{x}} \\
 v(\lambda) &= \bar{x}
 \end{aligned}$$

- Pivot the cdf
- Recall that  $F_T(T | \theta) \sim \text{Uniform}(0, 1)$  if  $F_T$  is continuous in  $T$ . That is,  $F_T(T | \theta)$  has the same distribution for all values of  $\theta$ . Thus  $F_T(T | \theta)$  is a pivotal quantity.
- Hence  $P \{ \alpha_1 \leq F_T(t | \theta) \leq 1 - \alpha_2 \} = 1 - \alpha_1 - \alpha_2$ .

- The following also works if  $F_T$  is not continuous in  $T$ .
  - If  $F_T(T \mid \theta)$  is non-increasing in  $\theta$  then for given  $t$   $\theta_L = \inf \{\theta : F_T(t \mid \theta) \leq 1 - \alpha_2\}$  and  $\theta_U = \sup \{\theta : F_T(t \mid \theta) \geq \alpha_1\}$ .
  - If  $F_T(T \mid \theta)$  is non-decreasing in  $\theta$  then for given  $t$   $\theta_L = \inf \{\theta : F_T(t \mid \theta) \geq \alpha_1\}$  and  $\theta_U = \sup \{\theta : F_T(t \mid \theta) \leq 1 - \alpha_2\}$ .
- Note  $\sum_1^n x_i \sim \text{Poisson}(n\lambda)$

$$CLT \rightarrow \frac{\frac{1}{n} \sum x_i - n\lambda}{\sqrt{n\lambda}} = \sqrt{n} \left( \frac{\bar{x} - \lambda}{\sqrt{\lambda}} \right) \xrightarrow{\text{mis}_{dist.}} N(0, 1)$$

$\Phi \left( \sqrt{n} \left( \frac{\bar{x} - \lambda}{\sqrt{\lambda}} \right) \right)$  is a pivotal quantity.

Redo this

- - Construct a 95%CI for  $\theta = P_\lambda\{X = 0\}$ . - Large sample via the MLE

**Find CI by Inverting Test**

- More Example on Mathematics Statistics
- This is related to [Delta Method]
- Example 9.1 is good example of LRT

## 9 Hypothesis and CI Inverting

**Example 9.1** Let  $X_1, X_2, \dots, X_n \sim iid N(\mu, \sigma^2)$  with  $(\mu, \sigma^2)$  unknown.

- Find the level  $\alpha$  LRT for  $H_0 : \sigma^2 = \sigma_0^2$  versus  $H_1 : \sigma^2 \neq \sigma_0^2$ .
- The unrestricted MLE of  $\sigma^2$  is given by  $\hat{\sigma}^2 = \frac{n-1}{n} S^2$  while the restricted MLE is  $\hat{\sigma}_0^2 = \sigma_0^2$ . In both cases,  $\hat{\mu} = \bar{X}$ . The likelihood ratio is given by

$$\begin{aligned} \lambda(\mathbf{x}) &= \frac{L(\hat{\mu}, \sigma_0^2)}{L(\hat{\mu}, \hat{\sigma}^2)} = \frac{[2\pi]^{-n/2} [\sigma_0^2]^{-n/2} \exp\left\{-\frac{1}{2\sigma_0^2} \sum (x_i - \bar{x})^2\right\}}{[2\pi]^{-n/2} \left[\frac{n-1}{n} S^2\right]^{-n/2} \exp\left\{-\frac{n}{2}\right\}} \\ &= \left[\frac{n-1}{\sigma_0^2} S^2\right]^{n/2} \exp\left\{-\frac{1}{2} \frac{n-1}{\sigma_0^2} S^2\right\} \left[\frac{1}{n}\right]^{n/2} \exp\left\{\frac{n}{2}\right\} \end{aligned}$$

- Let  $T = T_{\sigma_0^2}(\mathbf{X}) = \frac{n-1}{\sigma_0^2} S^2$
- Then

$$\begin{aligned} \lambda(\mathbf{X}) &= T^{n/2} \exp\left\{-\frac{1}{2}T\right\} \left[\frac{1}{n}\right]^{n/2} \exp\left\{\frac{n}{2}\right\} \\ \log \lambda(\mathbf{X}) &= \frac{n}{2} \left(\log T - \frac{1}{n}T - \log n + 1\right) \end{aligned}$$

- We reject the null hypothesis and support the research hypothesis if  $\lambda(\mathbf{X}) < k$  or  $\log \lambda(\mathbf{X}) < k^*$ .
- Given  $\mathbf{x}$ , let  $g(t) = \log \lambda(\mathbf{x})$ . When  $g(t) < k^*$ ,  $T < t_1$  or  $T > t_2$ , as illustrated in the figure (here,  $n = 25$ ).
- How might we use this test to construct an interval estimator for  $\sigma^2$  ?

graph

**Theorem 9.1** Given the rejection region  $R$  of a hypotheses test, the acceptance region is  $A = R^c$ .

For each  $\theta_0 \in \Theta$ , let  $A(\theta_0)$  be the acceptance region of a level  $\alpha$  test of  $H_0 : \theta = \theta_0$ . For each  $\mathbf{x} \in \mathcal{X}$ , define a set  $\mathcal{C}(\mathbf{x}) \subset \Theta$  by

$$\mathcal{C}(\mathbf{x}) = \{\theta_0 : \mathbf{x} \in A(\theta_0)\}$$

Then the random set  $\mathcal{C}(\mathbf{X})$  is a  $1 - \alpha$  confidence set.

Conversely, let  $\mathcal{C}(\mathbf{X})$  be a  $1 - \alpha$  confidence set. For any  $\theta_0 \in \Theta$ , define  $A(\theta_0) = \{\mathbf{x} : \theta_0 \in \mathcal{C}(\mathbf{x})\}$ . Then  $A(\theta_0)$  is the acceptance region of a level  $\alpha$  test of  $H_0 : \theta = \theta_0$ .

**Proof.** We prove the first assertion only. The proof for the second assertion is similar. Under the given condition, assuming an arbitrary  $\theta_0 \in \Theta$

$$\sup_{\theta=\theta_0} P(X \notin A(\theta_0)) = \sup_{\theta=\theta_0} P(T_{\theta_0} = 1) \leq \alpha,$$

which is the same as

$$1 - \alpha \leq \inf_{\theta=\theta_0} P(X \in A(\theta_0)) = \inf_{\theta=\theta_0} P(\theta_0 \in C(X)).$$

Since this holds for all  $\theta_0$ , the result follows from

$$\inf_{P \in \mathcal{P}} P(\theta \in C(X)) = \inf_{\theta_0 \in \Theta} \inf_{\theta=\theta_0} P(\theta_0 \in C(X)) \geq 1 - \alpha. \quad .$$

The converse of Theorem 7.2 is partially true, which is stated in the next result whose proof is left as an exercise.

**Corollary 9.1.1** *Proposition 7.2. Let  $C(X)$  be a confidence set for  $\theta$  with significance level (or confidence coefficient)  $1 - \alpha$ . For any  $\theta_0 \in \Theta$ , define a region  $A(\theta_0) = \{x : \theta_0 \in C(x)\}$ . Then the test  $T(X) = 1 - I_{A(\theta_0)}(X)$  has significance level  $\alpha$  for testing  $H_0 : \theta = \theta_0$  versus some  $H_1$ .*

**Example 9.2** *CONTINUE:*

- The rejection region is

$$R(\sigma_0^2) = \left\{ \mathbf{x} : T_{\sigma_0^2}(\mathbf{x}) < t_1 \text{ or } T_{\sigma_0^2}(\mathbf{x}) > t_2 \right\}$$

thus, the acceptance region is

$$A(\sigma_0^2) = \left\{ \mathbf{x} : t_1 \leq T_{\sigma_0^2}(\mathbf{x}) \leq t_2 \right\}$$

- Define  $\mathcal{C}(\mathbf{x}) = \{\sigma_0^2 : \mathbf{x} \in A(\sigma_0^2), \sigma_0^2 > 0\}$ . Then  $\mathcal{C}(\mathbf{X})$  is a  $1 - \alpha$  confidence set.

- Since  $A(\sigma_0^2) = \left\{ \mathbf{x} : t_1 \leq \frac{n-1}{\sigma_0^2} S^2 \leq t_2 \right\}$ , then

$$\mathcal{C}(\mathbf{x}) = \left\{ \sigma_0^2 : \frac{n-1}{t_2} S^2 \leq \sigma_0^2 \leq \frac{n-1}{t_1} S^2 \right\}$$

- For proper choice of  $(t_1, t_2)$ , we have  $\left( \frac{n-1}{t_2} S^2, \frac{n-1}{t_1} S^2 \right)$  is a  $100(1 - \alpha)\%$  CI for  $\sigma^2$ .

**Example 9.3** *CONTINUE*

**We choose  $t_1$  and  $t_2$  so that the test is size  $\alpha$ .**

- Values of  $t_1$  and  $t_2$  must be such that  $g(t_1) = g(t_2) = k^*$  (refer to earlier figure).

Thus  $g(t_1) - g(t_2) = 0$ .

-  $T_{\sigma_0^2}(\mathbf{X}) \sim \chi_{n-1}^2$ , central chi-squared, when  $H_0 : \sigma^2 = \sigma_0^2$  is true.

- Let  $F$  denote the cdf of a  $\chi_{n-1}^2$  distribution. Then  $t_1$  and  $t_2$  must satisfy  $F(t_1) + 1 - F(t_2) = \alpha$

- Define

-  $h_1(t_1, t_2) = g(t_1) - g(t_2)$

- $h_2(t_1, t_2) = F(t_1) - F(t_2) + 1 - \alpha$
  - We seek  $t_1$  and  $t_2$  so that  $h_1(t_1, t_2) = 0$  and  $h_2(t_1, t_2) = 0$ .
  - Note that  $g'(t) = \frac{n}{2} \left[ \frac{1}{t} - \frac{1}{n} \right] = \frac{1}{2} \left[ \frac{n}{t} - 1 \right] = \frac{n-t}{2t}$ .
- We can find the value by Newton methods

**Example 9.4 CONTINUE****Alternate Cl: Equal-tail probabilities**

- Select  $t_1$  and  $t_2$  so that  $F(t_1) = 1 - F(t_2) = \alpha/2$ .
- $t_1 = F^{-1}(\alpha/2)$  and  $t_2 = F^{-1}(1 - \alpha/2)$
- $\left( \frac{n-1}{t_2} S^2, \frac{n-1}{t_1} S^2 \right)$  is a  $100(1 - \alpha)\%$  CI for  $\sigma^2$
- Use this to back-solve for the rejection region of a test where  $H_0 : \sigma^2 = \sigma_0^2$  :

$$R(\sigma_0^2) = \left\{ x : T_{\sigma_0^2}(\mathbf{x}) < t_1 \text{ or } T_{\sigma_0^2}(\mathbf{x}) > t_2 \right\}$$

**Example 9.5 CONTINUE****Alternate Cl: Minimum expected length**

- The expected length is  $\text{Length} = E_{\sigma^2} \left[ (n-1)S^2 \left( \frac{1}{t_1} - \frac{1}{t_2} \right) \right] = (n-1)\sigma^2 \left( \frac{1}{t_1} - \frac{1}{t_2} \right)$
- Since  $F(t_2) - F(t_1) = 1 - \alpha = \gamma$ , the related differential equation is  $f(t_2) dt_2 - f(t_1) dt_1 = 0$ . This implies that  $\frac{dt_2}{dt_1} = \frac{f(t_1)}{f(t_2)}$ .
- Thus,  $\frac{d \text{Length}}{dt_1} = (n-1)\sigma^2 \left( -\frac{1}{t_1^2} + \frac{1}{t_2^2} \frac{dt_2}{dt_1} \right)$ . Setting this equal to 0 yields

$$\frac{1}{t_2^2} \frac{f(t_1)}{f(t_2)} - \frac{1}{t_1^2} = 0 \rightarrow t_1^2 f(t_1) - t_2^2 f(t_2) = 0$$

- Define

$$\begin{aligned} h_1(t_1, t_2) &= t_1^2 f(t_1) - t_2^2 f(t_2) \\ h_2(t_1, t_2) &= F(t_1) - F(t_2) + 1 - \alpha \end{aligned}$$

- We seek  $t_1$  and  $t_2$  so that  $h_1(t_1, t_2) = 0$  and  $h_2(t_1, t_2) = 0$ .
- Note that  $\frac{d(f(x))}{dx} = -\frac{1}{2}f(x) + \left(\frac{d}{2} - 1\right) \frac{1}{x}f(x)$  where  $f(x) = \frac{1}{\Gamma(d/2)2^{d/2}} x^{(d/2)-1} \exp\{-x/2\}$  and  $d$  is the degrees of freedom.
- Thus

$$\frac{d(x^2 f(x))}{dx} = -\frac{1}{2}x^2 f(x) + \left(\frac{d}{2} + 1\right) x f(x) = \frac{1}{2}x f(x)(d + 2 - x)$$

- Hence

$$\frac{\partial h_1}{\partial t_1} = \frac{1}{2}t_1 f(t_1)(n + 1 - t_1) \text{ and } \frac{\partial h_1}{\partial t_2} = -\frac{1}{2}t_2 f(t_2)(n + 1 - t_2)$$

- The Hessian matrix,  $H$ , is

Differential  
Equation



P-value

---

**P-value**

## 10 P Value

### Definition 10.1 *p.value*

- A *p-value* is the minimum level of significance for a test to reject the null hypothesis with the observed data.

**Example 10.1** - Suppose the rejection region for  $H_0 : \sigma^2 = \sigma_0^2$  versus  $H_1 : \sigma^2 > \sigma_0^2$  is

$$R(\sigma_0^2) = \{ \mathbf{x} : T_{\sigma_0^2}(\mathbf{x}) > t_\alpha \}$$

where  $T_{\sigma_0^2}(\mathbf{X}) = \frac{n-1}{\sigma_0^2} S^2$  and  $t_\alpha$  is the critical value for the  $\alpha$ -level test.

- *p-value* =  $P \{ T_{\sigma_0^2}(\mathbf{X}) > T_{\sigma_0^2}(\mathbf{x}) \} = \inf_\alpha \{ \alpha : T_{\sigma_0^2}(\mathbf{x}) > t_\alpha \}$  where  $T_{\sigma_0^2}(\mathbf{x})$  is observed.
- This means that the *p-value* is a function of the data; that is, the *p* value is a statistic.

■

**Example 10.2** Suppose we have a random sample from a continuous distribution (can relax this to allow discrete distributions). We seek evidence to support our research hypothesis  $H_1 : \theta > \theta_0$  using the rejection region

$$R = \{ \mathbf{x} : T(\mathbf{x}) > c \}$$

Let's consider the function

$$h_\theta(c) = P_\theta \{ T(\mathbf{X}) > c \}$$

- Suppose  $X_1, X_2, \dots, X_n$  is an iid random sample from an Exponential ( $\lambda$ ) distribution.
- Hypotheses:  $H_0 : \lambda \leq 1$  versus  $H_1 : \lambda > 1$ .
- We'll use the rejection region  $R = \{ \mathbf{x} : T(\mathbf{x}) > c \}$  where  $T(\mathbf{X}) = \bar{X} \sim \sum \text{Exp} \sim \text{Gamma}$ .
- Then  $h_\lambda(c) = P_\lambda \{ T(\mathbf{X}) > c \} = 1 - G_\lambda(c)$  where  $G_\lambda(\cdot)$  is the  $\text{Gamma}(n, \lambda/n)$  cdf.
- For example,  $n = 10, \lambda = 1$ .
- Note that  $h_\lambda(c) = \beta_c(\lambda)$  where  $\beta_c(\cdot)$  is the power function for the test with critical value *c*. So  $h_{\lambda=1}(c) = \beta_c(1)$ . ■

**Example 10.3** - Define the random variable  $h_\lambda(T) = 1 - G_\lambda(T)$ . where  $\lambda$  is assumed to be the true parameter

by setting the  $\lambda$  as the edge, this test has the  $\alpha$  level size

$$\begin{aligned}
P_{\lambda} \{h_{\lambda}(T) \leq a\} &= P_{\lambda} \{1 - G_{\lambda}(T) \leq a\} \\
&= P_{\lambda} \{G_{\lambda}(T) \geq 1 - a\} \\
&= P_{\lambda} \{T \leq G_{\lambda}^{-1}(1 - a)\} \\
&= 1 - G_{\lambda}(G_{\lambda}^{-1}(1 - a)) \\
&= 1 - (1 - a) \\
&= a
\end{aligned}$$

- So,  $h_{\lambda}(T) \sim \text{Uniform}(0, 1)$ . ■

**Theorem 10.1** pdf for  $H_1$

- Let's define  $h_{\lambda}(T)$  for the null hypothesis, at the point in  $\lambda_0 \in \Theta_0$  where the maximum Type I error probability occurs (where the power function is maximum over  $\Theta_0$ ).

- For our example,  $h_{\lambda_0}(T)$ . The distribution (cdf) of this random variable for some  $\lambda_1 \in \Theta_1$  is

$$\begin{aligned}
F_{\lambda_1}(a) &= P_{\lambda_1} \{h_{\lambda_0}(T) \leq a\} = P_{\lambda_1} \{1 - G_{\lambda_0}(T) \leq a\} = P_{\lambda_1} \{G_{\lambda_0}(T) \geq 1 - a\} \\
&= P_{\lambda_1} \{T \geq G_{\lambda_0}^{-1}(1 - a)\} = 1 - G_{\lambda_1}(G_{\lambda_0}^{-1}(1 - a))
\end{aligned}$$

with density by taking the derivative  $\frac{d}{da}$ :

$$f_{\lambda_1}(a) = \frac{g_{\lambda_1}(G_{\lambda_0}^{-1}(1 - a))}{g_{\lambda_0}(G_{\lambda_0}^{-1}(1 - a))}$$

**Definition 10.2** p-value

A p-value  $p(\mathbf{X})$  is a test statistic satisfying  $p(\mathbf{x}) \in [0, 1]$  for every sample point  $\mathbf{x}$ . Small values of  $p(\mathbf{X})$  are supportive of  $H_1$ . A p-value is valid if, for every  $\theta \in \Theta_0$  and every  $\alpha \in [0, 1]$ ,

$$P_{\theta}\{p(\mathbf{X}) \leq \alpha\} \leq \alpha$$

- We may construct a test using the p-value by way of the rejection region

$$R = \{\mathbf{x} : p(\mathbf{x}) \leq \alpha\}$$

- That is, we reject  $H_0$  in support of  $H_1$  when the p-value is less than the size (or level) of the test.

- This assumes we set the size of our test prior to the derivation of the p-value.

## One Sample Testing

*Much like Two Sample Testing but just cannot relax the identical assumption*

## Two Sample Testing

### Assumption

Table 3: Hypothesis Testing

	Assumption	Normal	T-Distribution
A	1. Independent	<i>Approx</i> by CLT, Lindeberg-Feller for identity	NA
	2. Finite Variance		
B	1. Independent	<i>Approx</i> by CLT, Slutsky, and N large enough	<i>Approx</i> By CLT
	2. Finite Variance		
C	3. Identical	<i>Approx</i> by LLN	<i>equal variance:</i> <i>Exact</i> <i>unequal variance:</i> <i>Approx</i>
	1. Independent		
	2. Unknown Finite Variance		
	3. Identical		
	4. $X_i$ Normality		

Simon's Note Assumption Simon's Note hypothesis testing

## Two sample t-test

### Unequal variance

Given  $E[\bar{X}_1 - \bar{X}_2] = \mu_1 - \mu_2$

by independent

$$\begin{aligned} \text{Var}(\bar{X}_1 - \bar{X}_2) &= \frac{\sum \text{Var}(X_1)}{n_1^2} + \frac{\sum \text{Var}(X_1)}{n_2^2} \\ &= \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \end{aligned}$$

for unknown variance, use sample variance approx

$$\approx \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}$$

Test statistics:

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Under Null hypothesis  $\mu_1 = \mu_2$ ,  $T$  follow t-distribution with degrees of freedom:

$$\frac{(s_1^2/n_1 + s_2^2/n_2)^2}{(s_1^2/n_1)^2/(n_1-1) + (s_2^2/n_2)^2/(n_2-1)}$$

- *Because we need to estimate the degrees of freedom, this test required a larger sample size than equal variance test.*

### Equal variance

Given  $E[\bar{X}_1 - \bar{X}_2] = \mu_1 - \mu_2$

by independent

$$\begin{aligned} \text{Var}(\bar{X}_1 - \bar{X}_2) &= \frac{\sum \text{Var}(X_1)}{n_1^2} + \frac{\sum \text{Var}(X_1)}{n_2^2} \\ &= \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \end{aligned}$$

under null, use pooled sample variance approx

because sample size can be different, using sample-size adjusted weighted estimator for the pooled sample variance

$$s_{pool}^2 = \frac{n_1 - 1}{n_1 + n_2 - 2} s_1^2 + \frac{n_2 - 1}{n_1 + n_2 - 2} s_2^2$$

Test statistics:

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}}$$

where  $s_p$  is the pooled standard deviation:

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

with degrees of freedom:  $n_1 + n_2 - 2$

## Bayesian

### Posterior

Given a parameter  $\theta$ , which has a prior of  $P_{\Theta}(\theta)$ , observed data  $X_i$ , which has a likelihood function of  $P(X_i|\theta) = L(\theta|X_i)$ ,

The posterior probability is

$$\begin{aligned} P(\hat{\theta}|X) &= \frac{P(\theta, X)}{P(X)} \\ &= \frac{P(X|\theta)P(\theta)}{P(X)} \\ &= \frac{L(\theta|X)P(\theta)}{\int_{\Theta|X} L(\theta|X)P(\theta)d\theta} \\ &\propto L(\theta|X)P(\theta) \end{aligned}$$

### Credible Interval

They characterize the range of  $\theta$  that are most believable on the basis of the data (and the prior belief regarding  $\theta$ ).

In Bayesian world, we say there is a 95% probability that  $\theta$  lies between the endpoints of a 95% credible interval (but the word “chance” in place of “probability” is still not ideal).

Keep in mind: Talking about probability in the subjective sense, and not in the long-term frequency sense.

## Regression

### Linear Regression

#### Assumption

For inference  $\beta$  - i.i.d - We don't need *Linear* if N sufficient large and X random - We don't need *Normality* since we have CLT - We don't need *Homoscedecity aka. equal variance* since we have robest variance estimator and N sufficient large

For sub-group specific mean  $E[Y|X]$  - i.i.d - Linear - We don't need *normality* - We don't need *Homoscedacity aka. equal variance* since we have robust variance estimator and N sufficient large

For prediction  $E[\hat{Y}|X]$  - i.i.d - Linear - normality - Homoscedacity

### Algebra Form

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

-  $\hat{\beta}$  is unbiased

Assuming  $E[y|X] = X\beta$ :

$$\begin{aligned} E[\hat{\beta}] &= E[E[\hat{\beta}|X]] \\ &= E[(X^T X)^{-1} X^T y|X] \\ &= (X^T X)^{-1} X^T E[y|X] = (X^T X)^{-1} X^T X\beta \\ &= \beta \end{aligned}$$

- Variance of  $\hat{\beta}$

Assumption: 1. error terms are pairwise independent 2.  $Var[Y|X = x] = \sigma^2$  for all x (homoscedacity), this assumption can be eliminated with robust estimator

$$\begin{aligned} Var(\hat{\beta}) &= Var[E[\hat{\beta}|X]] + E[Var(\hat{\beta}|X)] && \text{total variance} \\ &= \cancel{Var(\beta)} + E[Var(\hat{\beta}|X)] && \text{constant} \\ &= E[Var((X^T X)^{-1} X^T y|X)] && \text{Linearity} \\ &= E[(X^T X)^{-1} X^T * Var(y|X) * ((X^T X)^{-1} X^T)^T] && \text{algebra} \\ &= E[(X^T X)^{-1} X^T * (\sigma^2 I) * ((X^T X)^{-1} X^T)^T] && \text{homoscedacity} \\ &= E[\sigma^2 (X^T X)^{-1}] \end{aligned}$$

Where  $\sigma^2$  is estimated with MSE under homoscedasticity. Under heteroscedasticity we use sandwich variance. If there is heteroscedasticity, then the MSE can instead be interpreted as estimating the average within-group variance (i.e., averaged over the values of X).

### Projection

$$H = X(X^T X)^{-1} X^T$$

$$H^T = H; H^T H = H$$

so that we have the projection:

$$\hat{y} = Hy = x\beta$$

## Interpretation

$\beta_0$ :

- $\beta_0$  corresponds to a theoretical subgroup that is not well represented by the data
- If the predictor of interest,  $X$ , were something like height, then the intercept ( $\beta_0$ ) would mark the mean HbA1c among those of height zero.
- Here,  $\beta_0$  does not even carry a real-world interpretation, let alone the fact that our data cannot reasonably be used to reliably estimate it.
- by removing the intercept from the model we've actually made a very huge, untestable, and very often untrue assumption that the line goes through the origin! That is, by the above model,  $E[Y|X = 0] = 0$  necessarily.

$\beta_1$ :

**For**  $E[Y|X] = \beta_0 + \beta_1 x$

$c\beta_1$  denotes difference in mean Y between subgroups differing in X by c units.

**For**  $E[\log(Y)|X] = \beta_0 + \beta_1(x)$

$\beta_1 \log(1 + q) \Rightarrow$  Difference in mean Y between subgroups differing in X by  $(100 \times q)\%$

**For**  $E[\log(Y)|X] = \beta_0 + \beta_1 x$

$e^{\beta_1}$  denotes Geometric mean ratio, comparing subgroups differing in X by one unit.

Right-skewness is again not a justifiable reason on its own to log-transform an outcome.

## Example:

Consider following example:

$$E[\log(Y)|X] = \beta_0 + \beta_1(\log(x) - 5)$$

- $\exp(\beta_0) \Rightarrow$  Geometric mean Y among subgroup  $\log(x) - 5 = 0 \rightarrow x = \exp(5)$ .
- $1.6^{\beta_1} \Rightarrow$  Geometric mean ratio of Y comparing subgroup differing in X by  $\exp(\log(1.6)\beta_1) \rightarrow \log(1.6)$  unit

## Multivariate Regression

### Confounding/Precision Variable

#### Confounding Variable

- we pre-specify variables we believe to be highly likely confounders, and then perhaps some other variables that we believe may be potential confounders. Then, you can look at the following three (pre-specified) models:



1. Model 1 (unadjusted).
  2. Model 2 (adjusting for likely confounders).
  3. Model 3 (adjusting for all potential confounders measured).
- If Z is caused by X, then we should not adjust for Z.
  - Z is a mediator, Adjusting for Z is not appropriate.

### Precision Variable

- Specifically in the case of **linear regression**(different in logistic), adjustment for a variable of this type reduces variance in estimating association between X and Y without changing its value.

$$\begin{aligned}
 X &\perp Z \\
 E[Y|X = x] &= \beta_0 + \beta_1 X \\
 E[Y|X = x, Z = z] &= \alpha_0 + \alpha_1 X + \alpha_2 Z \\
 &\Rightarrow \hat{\beta}_1 = \hat{\alpha}_1 \\
 &\Rightarrow \text{Var}(\hat{\beta}_1) \geq \text{Var}(\hat{\alpha})
 \end{aligned}$$

- This property, known as collapsibility, holds for linear regression models (but not for all models, as we will see).

### Spline

#### Degree's of Freedom and number of splines

Number of splines = total degree of freedom(knots+1 \*parameters) - number of constrains\* knots

#### Example: Cubic Splines with 3 knots

- (3+1) cubic functions, each has 4 parameters
- 3 constrains at each knots, ie. continuous, differentiable, and have continuous first derivative at each knot
- total degree of freedom is  $4 * 4 - 3 * 3 = 7$

#### Example: Natural Cubic Splines with 3 knots

- (3+1) cubic functions, each has 4 parameters
- for the head and tail cubic function, restricted to linear function so each free up 2 degrees of freedom  $\sum_{i=0}^3 \beta_i x^i \Rightarrow \beta_0 + \beta_1 x$
- 3 constrains at each knots, ie. continuous, differentiable, and have continuous first derivative at each knot
- total degree of freedom is  $4 * 4 - 3 * 3 - 2 * 2 = 3$

**Example: B5**

The most general form for  $f$  is as follows:

$$f(x) = \begin{cases} \alpha_0 & \text{if } 0 < x \leq c \\ \gamma_0 + \gamma_1 x + \gamma_2 x^2 & \text{if } x > c \end{cases}$$

for real numbers  $\alpha_0, \gamma_0, \gamma_1$ , and  $\gamma_2$ . First, let's impose the continuity constraint, which is that:

$$\lim_{x \rightarrow c^-} f(x) = \lim_{x \rightarrow c^+} f(x) \Rightarrow \alpha_0 = \gamma_0 + \gamma_1 c + \gamma_2 c^2.$$

Updating the function so far, it takes the following form:

$$f(x) = \begin{cases} \gamma_0 + \gamma_1 c + \gamma_2 c^2 & \text{if } 0 < x \leq c \\ \gamma_0 + \gamma_1 x + \gamma_2 x^2 & \text{if } x > c \end{cases}$$

for real numbers  $\gamma_0, \gamma_1$ , and  $\gamma_2$ . Next, let's impose the differentiability constraint, which is that:

$$\lim_{x \rightarrow c^-} f'(x) = \lim_{x \rightarrow c^+} f'(x) \Rightarrow 0 = \gamma_1 + 2\gamma_2 c \Rightarrow \gamma_1 = -2\gamma_2 c$$

Updating the function so far, it takes the following form:

$$f(x) = \begin{cases} \gamma_0 - 2\gamma_2 c^2 + \gamma_2 c^2 & \text{if } 0 < x \leq c \\ \gamma_0 - 2\gamma_2 c x + \gamma_2 x^2 & \text{if } x > c \end{cases} = \begin{cases} \gamma_0 + \gamma_2 (-c^2) & \text{if } 0 < x \leq c \\ \gamma_0 + \gamma_2 (x^2 - 2cx) & \text{if } x > c \end{cases}$$

for real numbers  $\gamma_0$  and  $\gamma_2$ . Re-expressing as a basis expansion,

$$f(x) = \beta_0 + \beta_1 \left[ (x^2 - 2cx) 1_{(c, \infty)}(x) - c^2 1_{(0, c]}(x) \right], \quad (\beta_0, \beta_1) \in \mathbb{R}^2$$

The basis function of interest is given by  $h_c(x) = (x^2 - 2cx) 1_{(c, \infty)}(x) - c^2 1_{(0, c]}(x)$ . Now,  $f'(x) = h'_c(x) = \beta_1(2x - 2c)1_{(c, \infty)}(x)$ , which is continuous; further,  $\lim_{x \rightarrow c^-} h''_c(x) = 0$  and  $\lim_{x \rightarrow c^+} h''_c(x) = 2\beta_1$ , so  $f''(x)$  clearly not defined at  $x = c$ .

The unconstrained model uses four degrees of freedom, which makes sense because it involves an unrestricted constant (one degree of freedom) and an unrestricted quadratic function (three degrees of freedom). One degree of freedom is taken away by imposing the continuity constraint, and another is taken away by imposing the differentiability constraint. Therefore, the total number of degrees of freedom is given by 2, matching the number of basis functions.

## Weighted Regression

### Gauss-Markov Theorem

**Gauss-Markov Theorem:** Correct specified weighted least square estimator has the minimum variance among all **unbiased linear estimator** of  $\beta$ , if normality, linearity, no collinearity and finite variance assumptions are satisfied.

aka. Best Linear Unbiased Estimator.

- The OLS estimator is the best (i.e., minimum variance) linear (in y) unbiased (for  $\beta$ ) estimator of  $\beta$  under the assumptions of the Gauss-Markov theorem if the errors are homoscedastic.
- The OLS estimator is unbiased even if homoscedasticity does not hold.
- The OLS estimator is less efficient than WLS (with  $W \propto V^{-1}$ ) if homoscedasticity does not hold.
- The sandwich variance is still valid even if  $W \not\propto V^{-1}$ , but WLS estimator is no longer BLUE.

## Saturated Model

**The idea of saturated model is that no information is borrowed across subgroup, such that we can estimate sub-group specific mean/odd/...**

- A model is saturated if the number of parameters and constraints add up to the number of groups
- *Example: Comps Applied 2022.3*

	Y = 1	Y = 2	Y = 3
X = 0	a	d	g
X = 1	b	e	h
X = 2	c	f	i

$$\text{logit}(P(Y \geq 1|X = x)) = \alpha_0 + \alpha_1 I(x \geq 1)$$

- This is a logistic model, and we have 4 groups ( $\{x = 0, x \geq 1\} \cap \{Y = 0, Y \geq 1\}$ ), and we have 2 parameters and 2 constraints ( $P(Y \geq 1|X = x) + P(Y = 0|X = x) = 1$ ) so the model is saturated.

$$\log\left(\frac{P(Y = k|x = x)}{P(Y = 0|X = x)}\right) = \beta_{0k} + \beta_{1k}I(x = 1) + \beta_{2k}I(x = 2), k = 1, 2$$

- This is a multinomial model, we have 9 groups ( $3 \times 3$ ), and the model has 6 parameters ( $\beta_{ij}$ ) and 3 constraints  $\sum_k P(Y = k|X = x) = 1$ , so the model is saturated.

$$\text{logit}(P(Y \leq k | X = x)) = \gamma_{0k} - \gamma_1 I(x = 1) - \gamma_2 I(x = 2), k = 0, 1$$

- This is a proportional odds model, this model cannot be saturated because it's borrowing information from the left-side of the equation.

### **Multivariate Regression**

### **Logistic Regression**

### **Ordinal Regression**

### **Poisson Regression**

### **Longitudinal Regression**

### **Survival Regression**

### **Integral**

### **Substitution**

### **Distribution**

### **Discrete**

### **Poisson**

Note:  $Y \sim Poi$ ,  $aY \sim Poi$  because  $Y$  is discrete.

### **Geometric**

Note:  $\underbrace{p(1-p)^{x-1}}_{n \text{ trials}}$  or  $\underbrace{p(1-p)^x}_{n+1 \text{ trials}}$

### **Simulation**

## COMPS Practice

Table 5: Comps Practice

Year	Question	Related field	Time spend	Turn out
2021-T	1a-c		short	good
		1. pdf, cdf		
		2. E, Var		
	1 d-f		long	don't know how to come up with the thing
		1. [Asymptotic Probability]		
	2		3h	Can't do Binormal pdf, Expectation and variance
		1. Conditional Probability		Can't do conditional expectation
		2. [Multivariate Transformation]		
		3. [Binormal Distribution]		
	3a	Transformation	short	good
	3b		Long	didn't finish
		1. Order Statistics		
		2.		
	3c-d	1. [Hypothesis Testing]	Long	Didn't finish
	4a-d	• CLT		
	4e	• Delta		Can't do it
		• Converge		Very easy if you think of second order Delta
		• Taylor Series		
		• Second Order Delta		
	5a-c	• MOM	Long	Fisher info can only applied to exponential family
	• UMVUE		CRLB only apply for family distribution that support doesn't depend on paramter	
	• Rao-blackwell			
5d-f	• Bayesian			

Year	Question	Related field	Time spend	Turn out
2021- A 2022- A	6	<ul style="list-style-type: none"> <li>Gal-Markov</li> <li>lagrange multiplier</li> <li>Order Statsitics</li> </ul>		Don't know how to get the order statistics distribution
	1c	[T-test]	Long	Didn't Finish, still don't understand equal variance vs. unequal variance
	1g	[T-test]	Long	Also don't understand the assumption to get a exact t-test, approx. t-test.
	2	<ol style="list-style-type: none"> <li>Log-transform</li> <li>Logistic</li> </ol>	medium	Not familiar with log transform value intepretation Not familiar with the idea of Risk and Odd
	3a	Saturated	Long	Not familiar with Logisitcs Don't understand how a model can be saturated Don't understand how constrains works on the saturated
	3b	<ol style="list-style-type: none"> <li>logistic</li> <li>Multiordinal</li> <li>Log-odd</li> </ol>	Long	Not familiar with the forms of the model Don't know the proportional odd assumption Don't know how to hand calculated logistic model Don't know how to hand calcualted Log-odd Model
	4	Baysian	Okayish	I don't know
	5	<ol style="list-style-type: none"> <li>Simulation</li> <li>Bootstrap CI</li> </ol>	Okayish	Not sure how to evaluate CI performance Not sure how bootstrap CI
	6	<ul style="list-style-type: none"> <li>Non-linear model</li> <li>Delta method</li> </ul>	Long	Not familiar with multivariate delta Methods
	1b	<ul style="list-style-type: none"> <li>Geometric</li> <li>MGF</li> </ul>		Don't know the MGF of Geometric Don't know the trick to find MGF
2022- T	1j	AUC		Don't understand the upper bound of AUC
	2c			

Year	Question	Related field	Time spend	Turn out
	2d	Transformation		Transformation with Jacobian apply to pdf Non-independent b/c of Indicator function
	2g	• Slutsky		good, but worth reviewing
	2h	• LLN • CLT		Confused at which variance should use in delta method
	3c	• Delta Method • UMVE		Completely forget about Lehmann-scheffe
	3d	• Complete		Don't know how to prove
	3e			Completely no idea
	4a-b	• MGF		Bad, not thinking of MGF and got stuck
	4c	Method of Moment		Still don't know how to get a MOM estimator
	4d	Asymptotic Distribution		How to determine asymptotic distribution? converge in distribution?
	5			
	6a	MLE	fast	Like bob's
	6b-c	bayesian		
	6d-f	Unbiased	Bad	F