

STAT 6720
Mathematical Statistics II
Spring Semester 2010

Dr. Jürgen Symanzik

Utah State University

Department of Mathematics and Statistics

3900 Old Main Hill

Logan, UT 84322-3900

Tel.: (435) 797-0696

FAX: (435) 797-1822

e-mail: symanzik@math.usu.edu

Contents

Acknowledgements	1
6 Limit Theorems	1
6.1 Modes of Convergence	2
6.2 Weak Laws of Large Numbers	16
6.3 Strong Laws of Large Numbers	20
6.4 Central Limit Theorems	29
7 Sample Moments	36
7.1 Random Sampling	36
7.2 Sample Moments and the Normal Distribution	39
8 The Theory of Point Estimation	44
8.1 The Problem of Point Estimation	44
8.2 Properties of Estimates	45
8.3 Sufficient Statistics	48
8.4 Unbiased Estimation	57
8.5 Lower Bounds for the Variance of an Estimate	66
8.6 The Method of Moments	74
8.7 Maximum Likelihood Estimation	76
8.8 Decision Theory — Bayes and Minimax Estimation	81
9 Hypothesis Testing	89
9.1 Fundamental Notions	89
9.2 The Neyman–Pearson Lemma	94
9.3 Monotone Likelihood Ratios	98
9.4 Unbiased and Invariant Tests	102
10 More on Hypothesis Testing	112
10.1 Likelihood Ratio Tests	112
10.2 Parametric Chi–Squared Tests	117
10.3 t –Tests and F –Tests	121
10.4 Bayes and Minimax Tests	125

11 Confidence Estimation	130
11.1 Fundamental Notions	130
11.2 Shortest-Length Confidence Intervals	134
11.3 Confidence Intervals and Hypothesis Tests	138
11.4 Bayes Confidence Intervals	143
12 Nonparametric Inference	145
12.1 Nonparametric Estimation	145
12.2 Single-Sample Hypothesis Tests	151
12.3 More on Order Statistics	158
13 Some Results from Sampling	162
13.1 Simple Random Samples	162
13.2 Stratified Random Samples	165
14 Some Results from Sequential Statistical Inference	169
14.1 Fundamentals of Sequential Sampling	169
14.2 Sequential Probability Ratio Tests	173

Acknowledgements

I would like to thank my students, Hanadi B. Eltahir, Rich Madsen, and Bill Morphet, who helped during the Fall 1999 and Spring 2000 semesters in typesetting these lecture notes using L^AT_EX and for their suggestions how to improve some of the material presented in class. Thanks are also due to about 40 students who took Stat 6710/20 with me since the Fall 2000 semester for their valuable comments that helped to improve and correct these lecture notes.

In addition, I particularly would like to thank Mike Minnotte and Dan Coster, who previously taught this course at Utah State University, for providing me with their lecture notes and other materials related to this course. Their lecture notes, combined with additional material from Casella/Berger (2002), Rohatgi (1976) and other sources listed below, form the basis of the script presented here.

The primary textbook required for this class is:

- Casella, G., and Berger, R. L. (2002): *Statistical Inference* (Second Edition), Duxbury Press/Thomson Learning, Pacific Grove, CA.

A Web page dedicated to this class is accessible at:

http://www.math.usu.edu/~symanzik/teaching/2010_stat6720/stat6720.html

This course closely follows Casella and Berger (2002) as described in the syllabus. Additional material originates from the lectures from Professors Hering, Trenkler, and Gather I have attended while studying at the Universität Dortmund, Germany, the collection of Masters and PhD Preliminary Exam questions from Iowa State University, Ames, Iowa, and the following textbooks:

- Bandelow, C. (1981): *Einführung in die Wahrscheinlichkeitstheorie*, Bibliographisches Institut, Mannheim, Germany.
- Büning, H., and Trenkler, G. (1978): *Nichtparametrische statistische Methoden*, Walter de Gruyter, Berlin, Germany.
- Casella, G., and Berger, R. L. (1990): *Statistical Inference*, Wadsworth & Brooks/Cole, Pacific Grove, CA.
- Fisz, M. (1989): *Wahrscheinlichkeitsrechnung und mathematische Statistik*, VEB Deutscher Verlag der Wissenschaften, Berlin, German Democratic Republic.
- Gibbons, J. D., and Chakraborti, S. (1992): *Nonparametric Statistical Inference* (Third Edition, Revised and Expanded), Dekker, New York, NY.

- Johnson, N. L., and Kotz, S., and Balakrishnan, N. (1994): *Continuous Univariate Distributions, Volume 1* (Second Edition), Wiley, New York, NY.
- Johnson, N. L., and Kotz, S., and Balakrishnan, N. (1995): *Continuous Univariate Distributions, Volume 2* (Second Edition), Wiley, New York, NY.
- Kelly, D. G. (1994): *Introduction to Probability*, Macmillan, New York, NY.
- Lehmann, E. L. (1983): *Theory of Point Estimation* (1991 Reprint), Wadsworth & Brooks/Cole, Pacific Grove, CA.
- Lehmann, E. L. (1986): *Testing Statistical Hypotheses* (Second Edition – 1994 Reprint), Chapman & Hall, New York, NY.
- Mood, A. M., and Graybill, F. A., and Boes, D. C. (1974): *Introduction to the Theory of Statistics* (Third Edition), McGraw-Hill, Singapore.
- Parzen, E. (1960): *Modern Probability Theory and Its Applications*, Wiley, New York, NY.
- Rohatgi, V. K. (1976): *An Introduction to Probability Theory and Mathematical Statistics*, John Wiley and Sons, New York, NY.
- Rohatgi, V. K., and Saleh, A. K. E. (2001): *An Introduction to Probability and Statistics* (Second Edition), John Wiley and Sons, New York, NY.
- Searle, S. R. (1971): *Linear Models*, Wiley, New York, NY.
- Tamhane, A. C., and Dunlop, D. D. (2000): *Statistics and Data Analysis – From Elementary to Intermediate*, Prentice Hall, Upper Saddle River, NJ.

Additional definitions, integrals, sums, etc. originate from the following formula collections:

- Bronstein, I. N. and Semendjajew, K. A. (1985): *Taschenbuch der Mathematik* (22. Auflage), Verlag Harri Deutsch, Thun, German Democratic Republic.
- Bronstein, I. N. and Semendjajew, K. A. (1986): *Ergänzende Kapitel zu Taschenbuch der Mathematik* (4. Auflage), Verlag Harri Deutsch, Thun, German Democratic Republic.
- Sieber, H. (1980): *Mathematische Formeln — Erweiterte Ausgabe E*, Ernst Klett, Stuttgart, Germany.

Jürgen Symanzik, January 10, 2010.

6 Limit Theorems

(Based on Rohatgi, Chapter 6, Rohatgi/Saleh, Chapter 6 & Casella/Berger, Section 5.5)

Motivation:

I found this slide from my Stat 250, Section 003, "Introductory Statistics" class (an undergraduate class I taught at George Mason University in Spring 1999):

Central Limit Theorem

If \bar{X} is the mean of a SRS of size n drawn from a population with mean μ and finite standard deviation σ , then

$$\bar{X} \underset{\text{approx}}{\sim} N\left(\mu, \left(\frac{\sigma}{\sqrt{n}}\right)^2\right) \text{ for large } n,$$

i.e., the sample mean \bar{X} is approximately Normal distributed with mean μ and standard deviation σ/\sqrt{n} if n is large.

Abt:

i, The central limit theorem is one of the most important theorems in statistics. It justifies working with a Normal distribution for large samples no matter what the underlying true distribution is, given it has a finite standard deviation. And this holds for most practical situations.

ii, A direct conclusion of the central limit theorem is that $\sum_{i=1}^n X_i$ where X_1, X_2, \dots, X_n is a SRS of size n drawn from a population with mean μ and finite standard deviation σ , it holds that

$$\sum_{i=1}^n X_i \underset{\text{approx}}{\sim} N\left(n \cdot \mu, \left(\sqrt{n} \sigma\right)^2\right),$$

i.e., the sum of random variables X_i (that all come from the same population with mean μ and standard deviation σ) is approximately Normal distributed with mean $n \cdot \mu$ and standard deviation $\sqrt{n} \sigma$ if n is large.

60,

What does this mean at a more theoretical level???

6.1 Modes of Convergence

Definition 6.1.1:

Let X_1, \dots, X_n be iid rv's with common cdf $F_X(x)$. Let $\underline{T} = \underline{T}(\underline{X})$ be any **statistic**, i.e., a Borel-measurable function of \underline{X} that does not involve the population parameter(s) ϑ , defined on the support \mathcal{X} of \underline{X} . The induced probability distribution of $\underline{T}(\underline{X})$ is called the **sampling distribution** of $\underline{T}(\underline{X})$. ■

Note:

(i) Commonly used statistics are:

$$\text{Sample Mean: } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

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

Sample Median, Order Statistics, Min, Max, etc.

(ii) Recall that if X_1, \dots, X_n are iid and if $E(X)$ and $Var(X)$ exist, then $E(\bar{X}_n) = \mu = E(X)$, $E(S_n^2) = \sigma^2 = Var(X)$, and $Var(\bar{X}_n) = \frac{\sigma^2}{n}$.

(iii) Recall that if X_1, \dots, X_n are iid and if X has mgf $M_X(t)$ or characteristic function $\Phi_X(t)$ then $M_{\bar{X}_n}(t) = (M_X(\frac{t}{n}))^n$ or $\Phi_{\bar{X}_n}(t) = (\Phi_X(\frac{t}{n}))^n$. ■

Note: Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's on some probability space (Ω, L, P) . Is there any meaning behind the expression $\lim_{n \rightarrow \infty} X_n = X$? Not immediately under the usual definitions of limits. We first need to define modes of convergence for rv's and probabilities. ■

Definition 6.1.2:

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's with cdf's $\{F_n\}_{n=1}^{\infty}$ and let X be a rv with cdf F . If $F_n(x) \rightarrow F(x)$ at all continuity points of F , we say that X_n **converges in distribution to** X ($X_n \xrightarrow{d} X$) or X_n **converges in law to** X ($X_n \xrightarrow{L} X$), or F_n **converges weakly to** F ($F_n \xrightarrow{w} F$). ■

Example 6.1.3:

Let $X_n \sim N(0, \frac{1}{n})$. Then

$$F_n(x) = \int_{-\infty}^x \frac{\exp\left(-\frac{1}{2}nt^2\right)}{\sqrt{\frac{2\pi}{n}}} dt$$

$$\begin{aligned}
&= \int_{-\infty}^{\sqrt{nx}} \frac{\exp(-\frac{1}{2}s^2)}{\sqrt{2\pi}} ds \\
&= \Phi(\sqrt{nx})
\end{aligned}$$

$$\implies F_n(x) \rightarrow$$

If $F_X(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$ the only point of discontinuity is at $x = 0$. Everywhere else, $\Phi(\sqrt{nx}) = F_n(x) \rightarrow F_X(x)$, where $\Phi(z) = P(Z \leq z)$ with $Z \sim N(0, 1)$.

So, $X_n \xrightarrow{d} X$, where $P(X = 0) = 1$, or $X_n \xrightarrow{d} 0$ since the limiting rv here is degenerate, i.e., it has a Dirac(0) distribution. ■

Example 6.1.4:

In this example, the sequence $\{F_n\}_{n=1}^{\infty}$ converges pointwise to something that is not a cdf:

Let $X_n \sim \text{Dirac}(n)$, i.e., $P(X_n = n) = 1$. Then,

$$F_n(x) = \begin{cases} 0, & x < n \\ 1, & x \geq n \end{cases}$$

It is $F_n(x) \rightarrow 0 \quad \forall x$ which is not a cdf. Thus, there is no rv X such that $X_n \xrightarrow{d} X$. ■

Example 6.1.5:

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's such that $P(X_n = 0) = 1 - \frac{1}{n}$ and $P(X_n = n) = \frac{1}{n}$ and let $X \sim \text{Dirac}(0)$, i.e., $P(X = 0) = 1$.

It is

$$F_n(x) = \begin{cases} 0, & x < 0 \\ 1 - \frac{1}{n}, & 0 \leq x < n \\ 1, & x \geq n \end{cases}$$

$$F_X(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

It holds that $F_n \xrightarrow{w} F_X$ but

$$E(X_n^k) = \quad \not\rightarrow \quad E(X^k) =$$

Thus, convergence in distribution does not imply convergence of moments/means. ■

Note:

Convergence in distribution does not say that the X_i 's are close to each other or to X . It only means that their cdf's are (eventually) close to some cdf F . The X_i 's do not even have to be defined on the same probability space. ■

Example 6.1.6:

Let X and $\{X_n\}_{n=1}^\infty$ be iid $N(0,1)$. Obviously, $X_n \xrightarrow{d} X$ but $\lim_{n \rightarrow \infty} X_n \neq X$. ■

Theorem 6.1.7:

Let X and $\{X_n\}_{n=1}^\infty$ be discrete rv's with support \mathcal{X} and $\{\mathcal{X}_n\}_{n=1}^\infty$, respectively. Define the countable set $A = \mathcal{X} \cup \bigcup_{n=1}^\infty \mathcal{X}_n = \{a_k : k = 1, 2, 3, \dots\}$. Let $p_k = P(X = a_k)$ and $p_{nk} = P(X_n = a_k)$. Then it holds that $p_{nk} \rightarrow p_k \forall k$ iff $X_n \xrightarrow{d} X$. ■

Theorem 6.1.8:

Let X and $\{X_n\}_{n=1}^\infty$ be continuous rv's with pdf's f and $\{f_n\}_{n=1}^\infty$, respectively. If $f_n(x) \rightarrow f(x)$ for almost all x as $n \rightarrow \infty$ then $X_n \xrightarrow{d} X$. ■

Theorem 6.1.9:

Let X and $\{X_n\}_{n=1}^\infty$ be rv's such that $X_n \xrightarrow{d} X$. Let $c \in \mathbb{R}$ be a constant. Then it holds:

(i) $X_n + c \xrightarrow{d} X + c$.

(ii) $cX_n \xrightarrow{d} cX$.

(iii) If $a_n \rightarrow a$ and $b_n \rightarrow b$, then $a_n X_n + b_n \xrightarrow{d} aX + b$.

Proof:

Part (iii):

Suppose that $a > 0, a_n > 0$. Let $Y_n = a_n X_n + b_n$ and $Y = aX + b$. It is

$$F_Y(y) = P(Y < y) = P(aX + b < y) = P(X < \frac{y-b}{a}) = F_X(\frac{y-b}{a}).$$

Likewise,

$$F_{Y_n}(y) = F_{X_n}(\frac{y-b_n}{a_n}).$$

If y is a continuity point of F_Y , $\frac{y-b}{a}$ is a continuity point of F_X . Since $a_n \rightarrow a, b_n \rightarrow b$ and $F_{X_n}(x) \rightarrow F_X(x)$, it follows that $F_{Y_n}(y) \rightarrow F_Y(y)$ for every continuity point y of F_Y . Thus, $a_n X_n + b_n \xrightarrow{d} aX + b$. ■

Definition 6.1.10:

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's defined on a probability space (Ω, L, P) . We say that X_n **converges in probability** to a rv X ($X_n \xrightarrow{p} X$, $P\text{-}\lim_{n \rightarrow \infty} X_n = X$) if

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0 \quad \forall \epsilon > 0.$$

■

Note:

The following are equivalent:

$$\begin{aligned} & \lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0 \\ \iff & \lim_{n \rightarrow \infty} P(|X_n - X| \leq \epsilon) = 1 \\ \iff & \lim_{n \rightarrow \infty} P(\{\omega : |X_n(\omega) - X(\omega)| > \epsilon\}) = 0 \end{aligned}$$

If X is degenerate, i.e., $P(X = c) = 1$, we say that X_n is **consistent** for c . For example, let X_n such that $P(X_n = 0) = 1 - \frac{1}{n}$ and $P(X_n = 1) = \frac{1}{n}$. Then

$$P(|X_n| > \epsilon) = \begin{cases} \frac{1}{n}, & 0 < \epsilon < 1 \\ 0, & \epsilon \geq 1 \end{cases}$$

Therefore, $\lim_{n \rightarrow \infty} P(|X_n| > \epsilon) = 0 \quad \forall \epsilon > 0$. So $X_n \xrightarrow{p} 0$, i.e., X_n is consistent for 0. ■

Theorem 6.1.11:

- (i) $X_n \xrightarrow{p} X \iff X_n - X \xrightarrow{p} 0$.
- (ii) $X_n \xrightarrow{p} X, X_n \xrightarrow{p} Y \implies P(X = Y) = 1$.
- (iii) $X_n \xrightarrow{p} X, X_m \xrightarrow{p} X \implies X_n - X_m \xrightarrow{p} 0$ as $n, m \rightarrow \infty$.
- (iv) $X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \implies X_n \pm Y_n \xrightarrow{p} X \pm Y$.
- (v) $X_n \xrightarrow{p} X, k \in \mathbb{R}$ a constant $\implies kX_n \xrightarrow{p} kX$.
- (vi) $X_n \xrightarrow{p} k, k \in \mathbb{R}$ a constant $\implies X_n^r \xrightarrow{p} k^r \quad \forall r \in \mathbb{N}$.
- (vii) $X_n \xrightarrow{p} a, Y_n \xrightarrow{p} b, a, b \in \mathbb{R} \implies X_n Y_n \xrightarrow{p} ab$.
- (viii) $X_n \xrightarrow{p} 1 \implies X_n^{-1} \xrightarrow{p} 1$.
- (ix) $X_n \xrightarrow{p} a, Y_n \xrightarrow{p} b, a \in \mathbb{R}, b \in \mathbb{R} - \{0\} \implies \frac{X_n}{Y_n} \xrightarrow{p} \frac{a}{b}$.
- (x) $X_n \xrightarrow{p} X, Y$ an arbitrary rv $\implies X_n Y \xrightarrow{p} XY$.

(xi) $X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \implies X_n Y_n \xrightarrow{p} XY$.

Proof:

See Rohatgi, page 244–245, and Rohatgi/Saleh, page 260–261 for partial proofs. ■

Theorem 6.1.12:

Let $X_n \xrightarrow{p} X$ and let g be a continuous function on \mathbb{R} . Then $g(X_n) \xrightarrow{p} g(X)$.

Proof:

Preconditions:

1.) X rv $\implies \forall \epsilon > 0 \exists k = k(\epsilon) : P(|X| > k) < \frac{\epsilon}{2}$

2.) g is continuous on \mathbb{R}

$\implies g$ is also uniformly continuous on $[-k, k]$ (see Definition of uniformly continuous in Theorem 3.3.3 (iii):

$\forall \epsilon > 0 \exists \delta > 0 \forall x_1, x_2 \in \mathbb{R} : |x_1 - x_2| < \delta \implies |g(x_1) - g(x_2)| < \epsilon$.)

$\implies \exists \delta = \delta(\epsilon, k) : |X| \leq k, |X_n - X| < \delta \implies |g(X_n) - g(X)| < \epsilon$

Let

$$A = \{|X| \leq k\} = \{\omega : |X(\omega)| \leq k\}$$

$$B = \{|X_n - X| < \delta\} = \{\omega : |X_n(\omega) - X(\omega)| < \delta\}$$

$$C = \{|g(X_n) - g(X)| < \epsilon\} = \{\omega : |g(X_n(\omega)) - g(X(\omega))| < \epsilon\}$$

■

Corollary 6.1.13:

- (i) Let $X_n \xrightarrow{p} c, c \in \mathbb{R}$ and let g be a continuous function on \mathbb{R} . Then $g(X_n) \xrightarrow{p} g(c)$.
- (ii) Let $X_n \xrightarrow{d} X$ and let g be a continuous function on \mathbb{R} . Then $g(X_n) \xrightarrow{d} g(X)$.
- (iii) Let $X_n \xrightarrow{d} c, c \in \mathbb{R}$ and let g be a continuous function on \mathbb{R} . Then $g(X_n) \xrightarrow{d} g(c)$.

■

Theorem 6.1.14:

$$X_n \xrightarrow{p} X \iff X_n \xrightarrow{d} X.$$

Proof:

$$X_n \xrightarrow{p} X \Leftrightarrow P(|X_n - X| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty \quad \forall \epsilon > 0$$

It holds:

■

Theorem 6.1.15:

Let $c \in \mathbb{R}$ be a constant. Then it holds:

$$X_n \xrightarrow{d} c \iff X_n \xrightarrow{p} c.$$

■

Example 6.1.16:

In this example, we will see that

$$X_n \xrightarrow{d} X \not\iff X_n \xrightarrow{p} X$$

for some rv X . Let X_n be identically distributed rv's and let (X_n, X) have the following joint distribution:

	X_n	0	1	
X				
0		0	$\frac{1}{2}$	$\frac{1}{2}$
1		$\frac{1}{2}$	0	$\frac{1}{2}$
		$\frac{1}{2}$	$\frac{1}{2}$	1

■

Theorem 6.1.17:

Let $\{X_n\}_{n=1}^\infty$ and $\{Y_n\}_{n=1}^\infty$ be sequences of rv's and X be a rv defined on a probability space (Ω, L, P) . Then it holds:

$$Y_n \xrightarrow{d} X, |X_n - Y_n| \xrightarrow{p} 0 \implies X_n \xrightarrow{d} X.$$

Proof:

Similar to the proof of Theorem 6.1.14. See also Rohatgi, page 253, Theorem 14, and Rohatgi/Saleh, page 269, Theorem 14. ■

Theorem 6.1.18: Slutsky's Theorem

Let $(X_n)_{n=1}^\infty$ and $(Y_n)_{n=1}^\infty$ be sequences of rv's and X be a rv defined on a probability space (Ω, L, P) . Let $c \in \mathbb{R}$ be a constant. Then it holds:

(i) $X_n \xrightarrow{d} X, Y_n \xrightarrow{p} c \implies X_n + Y_n \xrightarrow{d} X + c.$

$$(ii) X_n \xrightarrow{d} X, Y_n \xrightarrow{p} c \implies X_n Y_n \xrightarrow{d} cX.$$

If $c = 0$, then also $X_n Y_n \xrightarrow{p} 0$.

$$(iii) X_n \xrightarrow{d} X, Y_n \xrightarrow{p} c \implies \frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{c} \text{ if } c \neq 0.$$

Proof:

$$(i) Y_n \xrightarrow{p} c \xLeftrightarrow{\text{Th.6.1.11}(i)} Y_n - c \xrightarrow{p} 0$$

$$\implies Y_n - c = Y_n + (X_n - X_n) - c = (X_n + Y_n) - (X_n + c) \xrightarrow{p} 0 \quad (A)$$

$$X_n \xrightarrow{d} X \xRightarrow{\text{Th.6.1.9}(i)} X_n + c \xrightarrow{d} X + c \quad (B)$$

Combining (A) and (B), it follows from Theorem 6.1.17:

$$X_n + Y_n \xrightarrow{d} X + c$$

(ii) Case $c = 0$:

$\forall \epsilon > 0 \quad \forall k > 0$, it is

$$\begin{aligned} P(|X_n Y_n| > \epsilon) &= P(|X_n Y_n| > \epsilon, Y_n \leq \frac{\epsilon}{k}) + P(|X_n Y_n| > \epsilon, Y_n > \frac{\epsilon}{k}) \\ &\leq P(|X_n \frac{\epsilon}{k}| > \epsilon) + P(Y_n > \frac{\epsilon}{k}) \\ &\leq P(|X_n| > k) + P(|Y_n| > \frac{\epsilon}{k}) \end{aligned}$$

Since $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} 0$, it follows for any fixed $k > 0$

$$\overline{\lim}_{n \rightarrow \infty} P(|X_n Y_n| > \epsilon) \leq P(|X| > k).$$

As k is arbitrary, we can make $P(|X| > k)$ as small as we want by choosing k large.

Therefore, $X_n Y_n \xrightarrow{p} 0$.

Case $c \neq 0$:

Since $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} c$, it follows from (ii), case $c = 0$, that $X_n Y_n - c X_n = X_n(Y_n - c) \xrightarrow{p} 0$.

$$\implies X_n Y_n \xrightarrow{p} c X_n$$

$$\xrightarrow{\text{Th.6.1.14}} X_n Y_n \xrightarrow{d} c X_n$$

Since $c X_n \xrightarrow{d} c X$ by Theorem 6.1.9 (ii), it follows from Theorem 6.1.17:

$$X_n Y_n \xrightarrow{d} c X$$

(iii) Let $Z_n \xrightarrow{p} 1$ and let $Y_n = c Z_n$.

$$\xrightarrow{c \neq 0} \frac{1}{Y_n} = \frac{1}{Z_n} \cdot \frac{1}{c}$$

$$\xrightarrow{\text{Th.6.1.11(v,viii)}} \frac{1}{Y_n} \xrightarrow{p} \frac{1}{c}$$

With part (ii) above, it follows:

$$X_n \xrightarrow{d} X \text{ and } \frac{1}{Y_n} \xrightarrow{p} \frac{1}{c}$$

$$\implies \frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{c}$$

■

Definition 6.1.19:

Let $(X_n)_{n=1}^{\infty}$ be a sequence of rv's such that $E(|X_n|^r) < \infty$ for some $r > 0$. We say that X_n **converges in the r^{th} mean** to a rv X ($X_n \xrightarrow{r} X$) if $E(|X|^r) < \infty$ and

$$\lim_{n \rightarrow \infty} E(|X_n - X|^r) = 0.$$

■

Example 6.1.20:

Let $(X_n)_{n=1}^{\infty}$ be a sequence of rv's defined by $P(X_n = 0) = 1 - \frac{1}{n}$ and $P(X_n = 1) = \frac{1}{n}$.

It is $E(|X_n|^r) = \frac{1}{n} \forall r > 0$. Therefore, $X_n \xrightarrow{r} 0 \forall r > 0$.

■

Note:

The special cases $r = 1$ and $r = 2$ are called *convergence in absolute mean* for $r = 1$ ($X_n \xrightarrow{1} X$) and *convergence in mean square* for $r = 2$ ($X_n \xrightarrow{ms} X$ or $X_n \xrightarrow{2} X$).

■

Theorem 6.1.21:

Assume that $X_n \xrightarrow{r} X$ for some $r > 0$. Then $X_n \xrightarrow{p} X$.

Proof:

Using Markov's Inequality (Corollary 3.5.2), it holds for any $\epsilon > 0$:

■

Example 6.1.22:

Let $(X_n)_{n=1}^{\infty}$ be a sequence of rv's defined by $P(X_n = 0) = 1 - \frac{1}{n^r}$ and $P(X_n = n) = \frac{1}{n^r}$ for some $r > 0$.

For any $\epsilon > 0$, $P(|X_n| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$; so $X_n \xrightarrow{p} 0$.

For $0 < s < r$, $E(|X_n|^s) = \frac{1}{n^{r-s}} \rightarrow 0$ as $n \rightarrow \infty$; so $X_n \xrightarrow{s} 0$. But $E(|X_n|^r) = 1 \not\rightarrow 0$ as $n \rightarrow \infty$; so $X_n \not\xrightarrow{r} 0$. ■

Theorem 6.1.23:

If $X_n \xrightarrow{r} X$, then it holds:

- (i) $\lim_{n \rightarrow \infty} E(|X_n|^r) = E(|X|^r)$; and
- (ii) $X_n \xrightarrow{s} X$ for $0 < s < r$.

Proof:

- (i) For $0 < r \leq 1$, it holds:

For $r > 1$, it follows from Minkowski's Inequality (Theorem 4.8.3):

$$\begin{aligned}
[E(|X - X_n + X_n|^r)]^{\frac{1}{r}} &\leq [E(|X - X_n|^r)]^{\frac{1}{r}} + [E(|X_n|^r)]^{\frac{1}{r}} \\
&\implies [E(|X|^r)]^{\frac{1}{r}} - [E(|X_n|^r)]^{\frac{1}{r}} \leq [E(|X - X_n|^r)]^{\frac{1}{r}} \\
&\implies [E(|X|^r)]^{\frac{1}{r}} - \lim_{n \rightarrow \infty} [E(|X_n|^r)]^{\frac{1}{r}} \leq \lim_{n \rightarrow \infty} [E(|X_n - X|^r)]^{\frac{1}{r}} = 0 \text{ since } X_n \xrightarrow{r} X \\
&\implies [E(|X|^r)]^{\frac{1}{r}} \leq \lim_{n \rightarrow \infty} [E(|X_n|^r)]^{\frac{1}{r}} \quad (C)
\end{aligned}$$

Similarly,

$$\begin{aligned}
[E(|X_n - X + X|^r)]^{\frac{1}{r}} &\leq [E(|X_n - X|^r)]^{\frac{1}{r}} + [E(|X|^r)]^{\frac{1}{r}} \\
&\implies \lim_{n \rightarrow \infty} [E(|X_n|^r)]^{\frac{1}{r}} - \lim_{n \rightarrow \infty} [E(|X|^r)]^{\frac{1}{r}} \leq \lim_{n \rightarrow \infty} [E(|X_n - X|^r)]^{\frac{1}{r}} = 0 \text{ since } X_n \xrightarrow{r} X \\
&\implies \lim_{n \rightarrow \infty} [E(|X_n|^r)]^{\frac{1}{r}} \leq [E(|X|^r)]^{\frac{1}{r}} \quad (D)
\end{aligned}$$

Combining (C) and (D) gives

$$\begin{aligned}
\lim_{n \rightarrow \infty} [E(|X_n|^r)]^{\frac{1}{r}} &= [E(|X|^r)]^{\frac{1}{r}} \\
&\implies \lim_{n \rightarrow \infty} E(|X_n|^r) = E(|X|^r)
\end{aligned}$$

(ii) For $1 \leq s < r$, it follows from Lyapunov's Inequality (Theorem 3.5.4):

$$\begin{aligned}
[E(|X_n - X|^s)]^{\frac{1}{s}} &\leq [E(|X_n - X|^r)]^{\frac{1}{r}} \\
&\implies E(|X_n - X|^s) \leq [E(|X_n - X|^r)]^{\frac{s}{r}} \\
&\implies \lim_{n \rightarrow \infty} E(|X_n - X|^s) \leq \lim_{n \rightarrow \infty} [E(|X_n - X|^r)]^{\frac{s}{r}} = 0 \text{ since } X_n \xrightarrow{r} X
\end{aligned}$$

$$\implies X_n \xrightarrow{s} X$$

An additional proof is required for $0 < s < r < 1$. ■

Definition 6.1.24:

Let $\{X_n\}_{n=1}^\infty$ be a sequence of rv's on (Ω, L, P) . We say that X_n **converges almost surely** to a rv X ($X_n \xrightarrow{a.s.} X$) or X_n **converges with probability 1** to X ($X_n \xrightarrow{w.p.1} X$) or X_n converges strongly to X iff

$$P(\{\omega : X_n(\omega) \rightarrow X(\omega) \text{ as } n \rightarrow \infty\}) = 1.$$

Note:

An interesting characterization of convergence with probability 1 and convergence in probability can be found in Parzen (1960) "Modern Probability Theory and Its Applications" on page 416 (see Handout). ■

Example 6.1.25:

Let $\Omega = [0, 1]$ and P a uniform distribution on Ω . Let $X_n(\omega) = \omega + \omega^n$ and $X(\omega) = \omega$.

For $\omega \in [0, 1)$, $\omega^n \rightarrow 0$ as $n \rightarrow \infty$. So $X_n(\omega) \rightarrow X(\omega) \quad \forall \omega \in [0, 1)$.

However, for $\omega = 1$, $X_n(1) = 2 \neq 1 = X(1) \quad \forall n$, i.e., convergence fails at $\omega = 1$.

Anyway, since $P(\{\omega : X_n(\omega) \rightarrow X(\omega) \text{ as } n \rightarrow \infty\}) = P(\{\omega \in [0, 1)\}) = 1$, it is $X_n \xrightarrow{a.s.} X$. ■

Theorem 6.1.26:

$$X_n \xrightarrow{a.s.} X \implies X_n \xrightarrow{p} X.$$

Proof:

Choose $\epsilon > 0$ and $\delta > 0$. Find $n_0 = n_0(\epsilon, \delta)$ such that

$$P\left(\bigcap_{n=n_0}^{\infty} \{|X_n - X| \leq \epsilon\}\right) \geq 1 - \delta.$$

■

Example 6.1.27:

$$X_n \xrightarrow{p} X \not\Rightarrow X_n \xrightarrow{a.s.} X:$$

Let $\Omega = (0, 1]$ and P a uniform distribution on Ω .

Define A_n by

$$\begin{aligned} A_1 &= (0, \frac{1}{2}], A_2 = (\frac{1}{2}, 1] \\ A_3 &= (0, \frac{1}{4}], A_4 = (\frac{1}{4}, \frac{1}{2}], A_5 = (\frac{1}{2}, \frac{3}{4}], A_6 = (\frac{3}{4}, 1] \\ A_7 &= (0, \frac{1}{8}], A_8 = (\frac{1}{8}, \frac{1}{4}], \dots \end{aligned}$$

Let $X_n(\omega) = I_{A_n}(\omega)$.

It is $P(|X_n - 0| \geq \epsilon) \rightarrow 0 \quad \forall \epsilon > 0$ since X_n is 0 except on A_n and $P(A_n) \downarrow 0$. Thus $X_n \xrightarrow{p} 0$.

But $P(\{\omega : X_n(\omega) \rightarrow 0\}) = 0$ (and not 1) because any ω keeps being in some A_n beyond any n_0 , i.e., $X_n(\omega)$ looks like $0 \dots 010 \dots 010 \dots 010 \dots$, so $X_n \not\xrightarrow{q.s.} 0$. ■

Example 6.1.28:

$$X_n \xrightarrow{r} X \not\Rightarrow X_n \xrightarrow{a.s.} X:$$

Let X_n be independent rv's such that $P(X_n = 0) = 1 - \frac{1}{n}$ and $P(X_n = 1) = \frac{1}{n}$.

It is $E(|X_n - 0|^r) = E(|X_n|^r) = E(|X_n|) = \frac{1}{n} \rightarrow 0$ as $n \rightarrow \infty$, so $X_n \xrightarrow{r} 0 \quad \forall r > 0$ (and due to Theorem 6.1.21, also $X_n \xrightarrow{p} 0$).

But

$$P(X_n = 0 \quad \forall m \leq n \leq n_0) = \prod_{n=m}^{n_0} (1 - \frac{1}{n}) = (\frac{m-1}{m})(\frac{m}{m+1})(\frac{m+1}{m+2}) \dots (\frac{n_0-2}{n_0-1})(\frac{n_0-1}{n_0}) = \frac{m-1}{n_0}$$

As $n_0 \rightarrow \infty$, it is $P(X_n = 0 \quad \forall m \leq n \leq n_0) \rightarrow 0 \quad \forall m$, so $X_n \not\xrightarrow{q.s.} 0$. ■

Example 6.1.29:

$$X_n \xrightarrow{a.s.} X \not\Rightarrow X_n \xrightarrow{r} X:$$

Let $\Omega = [0, 1]$ and P a uniform distribution on Ω .

$$\text{Let } A_n = [0, \frac{1}{\ln n}].$$

$$\text{Let } X_n(\omega) = nI_{A_n}(\omega) \text{ and } X(\omega) = 0.$$

It holds that $\forall \omega > 0 \exists n_0 : \frac{1}{\ln n_0} < \omega \implies X_n(\omega) = 0 \forall n > n_0$ and $P(\omega = 0) = 0$. Thus, $X_n \xrightarrow{a.s.} 0$.

But $E(|X_n - 0|^r) = \frac{n^r}{\ln n} \rightarrow \infty \forall r > 0$, so $X_n \not\xrightarrow{r} X$. ■

6.2 Weak Laws of Large Numbers

Theorem 6.2.1: WLLN: Version I

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid rv's with mean $E(X_i) = \mu$ and variance $Var(X_i) = \sigma^2 < \infty$.

Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Then it holds

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| \geq \epsilon) = 0 \quad \forall \epsilon > 0,$$

i.e., $\bar{X}_n \xrightarrow{p} \mu$.

Proof:

■

Note:

For iid rv's with finite variance, \bar{X}_n is consistent for μ .

A more general way to derive a “WLLN” follows in the next Definition. ■

Definition 6.2.2:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of rv's. Let $T_n = \sum_{i=1}^n X_i$. We say that $\{X_i\}$ obeys the WLLN with respect to a sequence of **norming constants** $\{B_i\}_{i=1}^{\infty}$, $B_i > 0, B_i \uparrow \infty$, if there exists a sequence of **centering constants** $\{A_i\}_{i=1}^{\infty}$ such that

$$B_n^{-1}(T_n - A_n) \xrightarrow{p} 0.$$

■

Theorem 6.2.3:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of pairwise uncorrelated rv's with $E(X_i) = \mu_i$ and $Var(X_i) = \sigma_i^2$, $i \in \mathbb{N}$. If $\sum_{i=1}^n \sigma_i^2 \rightarrow \infty$ as $n \rightarrow \infty$, we can choose $A_n = \sum_{i=1}^n \mu_i$ and $B_n = \sum_{i=1}^n \sigma_i^2$ and get

$$\frac{\sum_{i=1}^n (X_i - \mu_i)}{\sum_{i=1}^n \sigma_i^2} \xrightarrow{p} 0.$$

Proof:

By Markov's Inequality (Corollary 3.5.2), it holds for all $\epsilon > 0$:

$$P\left(\left|\sum_{i=1}^n X_i - \sum_{i=1}^n \mu_i\right| > \epsilon \sum_{i=1}^n \sigma_i^2\right) \leq \frac{E\left(\left(\sum_{i=1}^n (X_i - \mu_i)\right)^2\right)}{\epsilon^2 \left(\sum_{i=1}^n \sigma_i^2\right)^2} = \frac{1}{\epsilon^2 \sum_{i=1}^n \sigma_i^2} \rightarrow 0 \text{ as } n \rightarrow \infty \quad \blacksquare$$

Note:

To obtain Theorem 6.2.1, we choose $A_n = n\mu$ and $B_n = n\sigma^2$. \blacksquare

Theorem 6.2.4:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of rv's. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. A necessary and sufficient condition for $\{X_i\}$ to obey the WLLN with respect to $B_n = n$ is that

$$E\left(\frac{\bar{X}_n^2}{1 + \bar{X}_n^2}\right) \rightarrow 0$$

as $n \rightarrow \infty$.

Proof:

Rohatgi, page 258, Theorem 2, and Rohatgi/Saleh, page 275, Theorem 2. \blacksquare

Example 6.2.5:

Let (X_1, \dots, X_n) be jointly Normal with $E(X_i) = 0$, $E(X_i^2) = 1$ for all i , and $Cov(X_i, X_j) = \rho$ if $|i - j| = 1$ and $Cov(X_i, X_j) = 0$ if $|i - j| > 1$. Let $T_n = \sum_{i=1}^n X_i$. Then, $T_n \sim N(0, n + 2(n - 1)\rho) = N(0, \sigma^2)$. It is

$$\begin{aligned} E\left(\frac{\bar{X}_n^2}{1 + \bar{X}_n^2}\right) &= E\left(\frac{T_n^2}{n^2 + T_n^2}\right) \\ &= \frac{2}{\sqrt{2\pi}\sigma} \int_0^{\infty} \frac{x^2}{n^2 + x^2} e^{-\frac{x^2}{2\sigma^2}} dx \quad \left| \quad y = \frac{x}{\sigma}, \quad dy = \frac{dx}{\sigma} \right. \\ &= \frac{2}{\sqrt{2\pi}} \int_0^{\infty} \frac{\sigma^2 y^2}{n^2 + \sigma^2 y^2} e^{-\frac{y^2}{2}} dy \\ &= \frac{2}{\sqrt{2\pi}} \int_0^{\infty} \frac{(n + 2(n - 1)\rho)y^2}{n^2 + (n + 2(n - 1)\rho)y^2} e^{-\frac{y^2}{2}} dy \\ &\leq \frac{n + 2(n - 1)\rho}{n^2} \underbrace{\int_0^{\infty} \frac{2}{\sqrt{2\pi}} y^2 e^{-\frac{y^2}{2}} dy}_{=1, \text{ since Var of } N(0,1) \text{ distribution}} \\ &\rightarrow 0 \quad \text{as } n \rightarrow \infty \end{aligned}$$

$$\implies \bar{X}_n \xrightarrow{p} 0$$

■

Note:

We would like to have a WLLN that just depends on means but does not depend on the existence of finite variances. To approach this, we consider the following:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of rv's. Let $T_n = \sum_{i=1}^n X_i$. We truncate each $|X_i|$ at $c > 0$ and get

$$X_i^c = \begin{cases} X_i, & |X_i| \leq c \\ 0, & \text{otherwise} \end{cases}$$

Let $T_n^c = \sum_{i=1}^n X_i^c$ and $m_n = \sum_{i=1}^n E(X_i^c)$.

■

Lemma 6.2.6:

For T_n, T_n^c and m_n as defined in the Note above, it holds:

$$P(|T_n - m_n| > \epsilon) \leq P(|T_n^c - m_n| > \epsilon) + \sum_{i=1}^n P(|X_i| > c) \quad \forall \epsilon > 0$$

Proof:

■

Note:

If the X_i 's are identically distributed, then

$$P(|T_n - m_n| > \epsilon) \leq P(|T_n^c - m_n| > \epsilon) + nP(|X_1| > c) \quad \forall \epsilon > 0.$$

If the X_i 's are iid, then

$$P(|T_n - m_n| > \epsilon) \leq \frac{nE((X_1^c)^2)}{\epsilon^2} + nP(|X_1| > c) \quad \forall \epsilon > 0 \quad (*).$$

Note that $P(|X_i| > c) = P(|X_1| > c) \quad \forall i \in \mathbb{N}$ if the X_i 's are identically distributed and that $E((X_i^c)^2) = E((X_1^c)^2) \quad \forall i \in \mathbb{N}$ if the X_i 's are iid.

■

Theorem 6.2.7: Khintchine's WLLN

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid rv's with finite mean $E(X_i) = \mu$. Then it holds:

$$\bar{X}_n = \frac{1}{n}T_n \xrightarrow{p} \mu$$

Proof:

If we take $c = n$ and replace ϵ by $n\epsilon$ in (*) in the Note above, we get

$$P\left(\frac{|T_n - m_n|}{n} > \epsilon\right) = P(|T_n - m_n| > n\epsilon) \leq \frac{E((X_1^n)^2)}{n\epsilon^2} + nP(|X_1| > n).$$

Since $E(|X_1|) < \infty$, it is $nP(|X_1| > n) \rightarrow 0$ as $n \rightarrow \infty$ by Theorem 3.1.9. From Corollary 3.1.12 we know that $E(|X|^\alpha) = \alpha \int_0^\infty x^{\alpha-1}P(|X| > x)dx$. Therefore,

■

Note:

Theorem 6.2.7 meets the previously stated goal of not having a finite variance requirement. ■

6.3 Strong Laws of Large Numbers

Definition 6.3.1:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of rv's. Let $T_n = \sum_{i=1}^n X_i$. We say that $\{X_i\}$ obeys the SLLN with respect to a sequence of **norming constants** $\{B_i\}_{i=1}^{\infty}$, $B_i > 0, B_i \uparrow \infty$, if there exists a sequence of **centering constants** $\{A_i\}_{i=1}^{\infty}$ such that

$$B_n^{-1}(T_n - A_n) \xrightarrow{a.s.} 0.$$

■

Note:

Unless otherwise specified, we will only use the case that $B_n = n$ in this section. ■

Theorem 6.3.2:

$$X_n \xrightarrow{a.s.} X \iff \lim_{n \rightarrow \infty} P(\sup_{m \geq n} |X_m - X| > \epsilon) = 0 \quad \forall \epsilon > 0.$$

Proof: (see also Rohatgi, page 249, Theorem 11)

WLOG, we can assume that $X = 0$ since $X_n \xrightarrow{a.s.} X$ implies $X_n - X \xrightarrow{a.s.} 0$. Thus, we have to prove:

$$X_n \xrightarrow{a.s.} 0 \iff \lim_{n \rightarrow \infty} P(\sup_{m \geq n} |X_m| > \epsilon) = 0 \quad \forall \epsilon > 0$$

Choose $\epsilon > 0$ and define

$$\begin{aligned} A_n(\epsilon) &= \{\sup_{m \geq n} |X_m| > \epsilon\} \\ C &= \{\lim_{n \rightarrow \infty} X_n = 0\} \end{aligned}$$

“ \implies ”:

Since $X_n \xrightarrow{a.s.} 0$, we know that $P(C) = 1$ and therefore $P(C^c) = 0$.

Let $B_n(\epsilon) = C \cap A_n(\epsilon)$. Note that $B_{n+1}(\epsilon) \subseteq B_n(\epsilon)$ and for the limit set $\bigcap_{n=1}^{\infty} B_n(\epsilon) = \emptyset$. It follows that

$$\lim_{n \rightarrow \infty} P(B_n(\epsilon)) = P\left(\bigcap_{n=1}^{\infty} B_n(\epsilon)\right) = 0.$$

We also have

$$\begin{aligned} P(B_n(\epsilon)) &= P(A_n \cap C) \\ &= 1 - P(C^c \cup A_n^c) \\ &= 1 - \underbrace{P(C^c)}_{=0} - P(A_n^c) + \underbrace{P(C^c \cap A_n^c)}_{=0} \\ &= P(A_n) \end{aligned}$$

$$\implies \lim_{n \rightarrow \infty} P(A_n(\epsilon)) = 0$$

“ \Leftarrow ”:

Assume that $\lim_{n \rightarrow \infty} P(A_n(\epsilon)) = 0 \quad \forall \epsilon > 0$ and define $D(\epsilon) = \{\overline{\lim}_{n \rightarrow \infty} |X_n| > \epsilon\}$.

Since $D(\epsilon) \subseteq A_n(\epsilon) \quad \forall n \in \mathbb{N}$, it follows that $P(D(\epsilon)) = 0 \quad \forall \epsilon > 0$. Also,

$$C^c = \{\lim_{n \rightarrow \infty} X_n \neq 0\} \subseteq \bigcup_{k=1}^{\infty} \{\overline{\lim}_{n \rightarrow \infty} |X_n| > \frac{1}{k}\}.$$

$$\implies 1 - P(C) \leq \sum_{k=1}^{\infty} P(D(\frac{1}{k})) = 0$$

$$\implies X_n \xrightarrow{a.s.} 0$$

■

Note:

(i) $X_n \xrightarrow{a.s.} 0$ implies that $\forall \epsilon > 0 \quad \forall \delta > 0 \quad \exists n_0 \in \mathbb{N} : P(\sup_{n \geq n_0} |X_n| > \epsilon) < \delta$.

(ii) Recall that for a given sequence of events $\{A_n\}_{n=1}^{\infty}$,

$$A = \overline{\lim}_{n \rightarrow \infty} A_n = \lim_{n \rightarrow \infty} \bigcup_{k=n}^{\infty} A_k = \bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} A_k$$

is the event that infinitely many of the A_n occur. We write $P(A) = P(A_n \text{ i.o.})$ where *i.o.* stands for “infinitely often”.

(iii) Using the terminology defined in (ii) above, we can rewrite Theorem 6.3.2 as

$$X_n \xrightarrow{a.s.} 0 \iff P(|X_n| > \epsilon \text{ i.o.}) = 0 \quad \forall \epsilon > 0.$$

■

Theorem 6.3.3: Borel–Cantelli Lemma

Let A be defined as in (ii) of the previous Note.

(i) **1st BC–Lemma:**

Let $\{A_n\}_{n=1}^\infty$ be a sequence of events such that $\sum_{n=1}^\infty P(A_n) < \infty$. Then $P(A) = 0$.

(ii) **2nd BC–Lemma:**

Let $\{A_n\}_{n=1}^\infty$ be a sequence of independent events such that $\sum_{n=1}^\infty P(A_n) = \infty$. Then $P(A) = 1$.

Proof:

(i):

$$\begin{aligned} P(A) &= P\left(\lim_{n \rightarrow \infty} \bigcup_{k=n}^\infty A_k\right) \\ &= \lim_{n \rightarrow \infty} P\left(\bigcup_{k=n}^\infty A_k\right) \\ &\leq \lim_{n \rightarrow \infty} \sum_{k=n}^\infty P(A_k) \\ &= \lim_{n \rightarrow \infty} \left(\sum_{k=1}^\infty P(A_k) - \sum_{k=1}^{n-1} P(A_k) \right) \\ &= 0 \end{aligned}$$

(ii): We have $A^c = \bigcup_{n=1}^\infty \bigcap_{k=n}^\infty A_k^c$. Therefore,

$$P(A^c) = P\left(\lim_{n \rightarrow \infty} \bigcap_{k=n}^\infty A_k^c\right) = \lim_{n \rightarrow \infty} P\left(\bigcap_{k=n}^\infty A_k^c\right).$$

If we choose $n_0 > n$, it holds that

$$\bigcap_{k=n}^\infty A_k^c \subseteq \bigcap_{k=n}^{n_0} A_k^c.$$

Therefore,

$$\begin{aligned} P\left(\bigcap_{k=n}^\infty A_k^c\right) &\leq \lim_{n_0 \rightarrow \infty} P\left(\bigcap_{k=n}^{n_0} A_k^c\right) \\ &= \lim_{n_0 \rightarrow \infty} \prod_{k=n}^{n_0} (1 - P(A_k)) \\ &\stackrel{indep.}{\leq} \lim_{n_0 \rightarrow \infty} \exp\left(-\sum_{k=n}^{n_0} P(A_k)\right) \\ &= 0 \end{aligned}$$

$\implies P(A) = 1$ ■

Example 6.3.4:

Independence is necessary for 2nd BC–Lemma:

Let $\Omega = (0, 1)$ and P a uniform distribution on Ω .

Let $A_n = I_{(0, \frac{1}{n})}(\omega)$. Therefore,

$$\sum_{n=1}^{\infty} P(A_n) = \sum_{n=1}^{\infty} \frac{1}{n} = \infty.$$

But for any $\omega \in \Omega$, A_n occurs only for $1, 2, \dots, \lfloor \frac{1}{\omega} \rfloor$, where $\lfloor \frac{1}{\omega} \rfloor$ denotes the largest integer (“floor”) that is $\leq \frac{1}{\omega}$. Therefore, $P(A) = P(A_n \text{ i.o.}) = 0$. ■

Lemma 6.3.5: Kolmogorov’s Inequality

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of independent rv’s with common mean 0 and variances σ_i^2 . Let

$T_n = \sum_{i=1}^n X_i$. Then it holds:

$$P\left(\max_{1 \leq k \leq n} |T_k| \geq \epsilon\right) \leq \frac{\sum_{i=1}^n \sigma_i^2}{\epsilon^2} \quad \forall \epsilon > 0$$

Proof:

See Rohatgi, page 268, Lemma 2, and Rohatgi/Saleh, page 284, Lemma 1. ■

Lemma 6.3.6: Kronecker’s Lemma

If $\sum_{i=1}^{\infty} X_i$ converges to $s < \infty$ and $B_n \uparrow \infty$, then it holds:

$$\frac{1}{B_n} \sum_{k=1}^n B_k X_k \rightarrow 0$$

Proof:

See Rohatgi, page 269, Lemma 3, and Rohatgi/Saleh, page 285, Lemma 2. ■

Theorem 6.3.7: Cauchy Criterion

$X_n \xrightarrow{a.s.} X \iff \lim_{n \rightarrow \infty} P(\sup_m |X_{n+m} - X_n| \leq \epsilon) = 1 \quad \forall \epsilon > 0.$

Proof:

See Rohatgi, page 270, Theorem 5. ■

Theorem 6.3.8:

If $\sum_{n=1}^{\infty} \text{Var}(X_n) < \infty$, then $\sum_{n=1}^{\infty} (X_n - E(X_n))$ converges almost surely.

Proof:

See Rohatgi, page 272, Theorem 6, and Rohatgi/Saleh, page 286, Theorem 4. ■

Corollary 6.3.9:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of independent rv's. Let $\{B_i\}_{i=1}^{\infty}$, $B_i > 0$, $B_i \uparrow \infty$, a sequence of norming constants. Let $T_n = \sum_{i=1}^n X_i$. If $\sum_{i=1}^{\infty} \frac{\text{Var}(X_i)}{B_i^2} < \infty$ then it holds:

$$\frac{T_n - E(T_n)}{B_n} \xrightarrow{a.s.} 0$$

Proof:

This Corollary follows directly from Theorem 6.3.8 and Lemma 6.3.6. ■

Lemma 6.3.10: Equivalence Lemma

Let $\{X_i\}_{i=1}^{\infty}$ and $\{X'_i\}_{i=1}^{\infty}$ be sequences of rv's. Let $T_n = \sum_{i=1}^n X_i$ and $T'_n = \sum_{i=1}^n X'_i$.

If the series $\sum_{i=1}^{\infty} P(X_i \neq X'_i) < \infty$, then the series $\{X_i\}$ and $\{X'_i\}$ are **tail-equivalent** and T_n and T'_n are **convergence-equivalent**, i.e., for $B_n \uparrow \infty$ the sequences $\frac{1}{B_n}T_n$ and $\frac{1}{B_n}T'_n$ converge on the same event and to the same limit, except for a null set.

Proof:

See Rohatgi, page 266, Lemma 1. ■

Lemma 6.3.11:

Let X be a rv with $E(|X|) < \infty$. Then it holds:

$$\sum_{n=1}^{\infty} P(|X| \geq n) \leq E(|X|) \leq 1 + \sum_{n=1}^{\infty} P(|X| \geq n)$$

Proof:

Continuous case only:

Let X have a pdf f . Then it holds:

$$E(|X|) = \int_{-\infty}^{\infty} |x| f(x) dx = \sum_{k=0}^{\infty} \int_{k \leq |x| \leq k+1} |x| f(x) dx$$

$$\implies \sum_{k=0}^{\infty} kP(k \leq |X| \leq k+1) \leq E(|X|) \leq \sum_{k=0}^{\infty} (k+1)P(k \leq |X| \leq k+1)$$

It is

$$\begin{aligned} \sum_{k=0}^{\infty} kP(k \leq |X| \leq k+1) &= \sum_{k=0}^{\infty} \sum_{n=1}^k P(k \leq |X| \leq k+1) \\ &= \sum_{n=1}^{\infty} \sum_{k=n}^{\infty} P(k \leq |X| \leq k+1) \\ &= \sum_{n=1}^{\infty} P(|X| \geq n) \end{aligned}$$

Similarly,

$$\begin{aligned} \sum_{k=0}^{\infty} (k+1)P(k \leq |X| \leq k+1) &= \sum_{n=1}^{\infty} P(|X| \geq n) + \sum_{k=0}^{\infty} P(k \leq |X| \leq k+1) \\ &= \sum_{n=1}^{\infty} P(|X| \geq n) + 1 \end{aligned}$$

■

Theorem 6.3.12:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid rv's. Then it holds:

$$X_n \xrightarrow{a.s.} 0 \iff \sum_{n=1}^{\infty} P(|X_n| > \epsilon) < \infty \quad \forall \epsilon > 0$$

Proof:

See Rohatgi, page 265, Theorem 3.

■

Theorem 6.3.13: Kolmogorov's SLLN

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid rv's. Let $T_n = \sum_{i=1}^n X_i$. Then it holds:

$$\frac{T_n}{n} = \bar{X}_n \xrightarrow{a.s.} \mu < \infty \iff E(|X|) < \infty \text{ (and then } \mu = E(X))$$

Proof:

“ \implies ”:

Suppose that $\bar{X}_n \xrightarrow{a.s.} \mu < \infty$. It is

“ \impliedby ”:

Let $E(|X|) < \infty$.

It is

$$\begin{aligned}
 \sum_{n=k}^{\infty} \frac{1}{n^2} &= \frac{1}{k^2} + \frac{1}{(k+1)^2} + \frac{1}{(k+2)^2} + \dots \\
 &\leq \frac{1}{k^2} + \frac{1}{k(k+1)} + \frac{1}{(k+1)(k+2)} + \dots \\
 &= \frac{1}{k^2} + \sum_{n=k+1}^{\infty} \frac{1}{n(n-1)}
 \end{aligned}$$

From Bronstein, page 30, # 7, we know that

$$\begin{aligned}
 1 &= \frac{1}{1 \cdot 2} + \frac{1}{2 \cdot 3} + \frac{1}{3 \cdot 4} + \dots + \frac{1}{n(n+1)} + \dots \\
 &= \frac{1}{1 \cdot 2} + \frac{1}{2 \cdot 3} + \frac{1}{3 \cdot 4} + \dots + \frac{1}{(k-1) \cdot k} + \sum_{n=k+1}^{\infty} \frac{1}{n(n-1)} \\
 \implies \sum_{n=k+1}^{\infty} \frac{1}{n(n-1)} &= 1 - \frac{1}{1 \cdot 2} - \frac{1}{2 \cdot 3} - \frac{1}{3 \cdot 4} - \dots - \frac{1}{(k-1) \cdot k} \\
 &= \frac{1}{2} - \frac{1}{2 \cdot 3} - \frac{1}{3 \cdot 4} - \dots - \frac{1}{(k-1) \cdot k} \\
 &= \frac{1}{3} - \frac{1}{3 \cdot 4} - \dots - \frac{1}{(k-1) \cdot k} \\
 &= \frac{1}{4} - \dots - \frac{1}{(k-1) \cdot k} \\
 &= \dots \\
 &= \frac{1}{k} \\
 \implies \sum_{n=k}^{\infty} \frac{1}{n^2} &\leq \frac{1}{k^2} + \sum_{n=k+1}^{\infty} \frac{1}{n(n-1)} \\
 &= \frac{1}{k^2} + \frac{1}{k} \\
 &\leq \frac{2}{k}
 \end{aligned}$$



6.4 Central Limit Theorems

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's with cdf's $\{F_n\}_{n=1}^{\infty}$. Suppose that the mgf $M_n(t)$ of X_n exists.

Questions: Does $M_n(t)$ converge? Does it converge to a mgf $M(t)$? If it does converge, does it hold that $X_n \xrightarrow{d} X$ for some rv X ?

Example 6.4.1:

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's such that $P(X_n = -n) = 1$. Then the mgf is $M_n(t) = E(e^{tX}) = e^{-tn}$. So

$$\lim_{n \rightarrow \infty} M_n(t) = \begin{cases} 0, & t > 0 \\ 1, & t = 0 \\ \infty, & t < 0 \end{cases}$$

So $M_n(t)$ does not converge to a mgf and $F_n(x) \rightarrow F(x) = 1 \quad \forall x$. But $F(x)$ is not a cdf. ■

Note:

Due to Example 6.4.1, the existence of mgf's $M_n(t)$ that converge to something is not enough to conclude convergence in distribution.

Conversely, suppose that X_n has mgf $M_n(t)$, X has mgf $M(t)$, and $X_n \xrightarrow{d} X$. Does it hold that

$$M_n(t) \rightarrow M(t)?$$

Not necessarily! See Rohatgi, page 277, Example 2, and Rohatgi/Saleh, page 289, Example 2, as a counter example. Thus, convergence in distribution of rv's that all have mgf's does not imply the convergence of mgf's.

However, we can make the following statement in the next Theorem: ■

Theorem 6.4.2: Continuity Theorem

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of rv's with cdf's $\{F_n\}_{n=1}^{\infty}$ and mgf's $\{M_n(t)\}_{n=1}^{\infty}$. Suppose that $M_n(t)$ exists for $|t| \leq t_0 \quad \forall n$. If there exists a rv X with cdf F and mgf $M(t)$ which exists for $|t| \leq t_1 < t_0$ such that $\lim_{n \rightarrow \infty} M_n(t) = M(t) \quad \forall t \in [-t_1, t_1]$, then $F_n \xrightarrow{w} F$, i.e., $X_n \xrightarrow{d} X$. ■

Example 6.4.3:

Let $X_n \sim \text{Bin}(n, \frac{\lambda}{n})$. Recall (e.g., from Theorem 3.3.12 and related Theorems) that for $X \sim \text{Bin}(n, p)$ the mgf is $M_X(t) = (1 - p + pe^t)^n$. Thus,

■

Note:

Recall Theorem 3.3.11: Suppose that $\{X_n\}_{n=1}^{\infty}$ is a sequence of rv's with characteristic functions $\{\Phi_n(t)\}_{n=1}^{\infty}$. Suppose that

$$\lim_{n \rightarrow \infty} \Phi_n(t) = \Phi(t) \quad \forall t \in (-h, h) \text{ for some } h > 0,$$

and $\Phi(t)$ is the characteristic function of a rv X . Then $X_n \xrightarrow{d} X$.

■

Theorem 6.4.4: Lindeberg–Lévy Central Limit Theorem

Let $\{X_n\}_{n=1}^{\infty}$ be a sequence of iid rv's with $E(X_i) = \mu$ and $0 < \text{Var}(X_i) = \sigma^2 < \infty$. Then it holds for $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ that

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} Z$$

where $Z \sim N(0, 1)$.

Proof:

Let $Z \sim N(0, 1)$. According to Theorem 3.3.12 (v), the characteristic function of Z is $\Phi_Z(t) = \exp(-\frac{1}{2}t^2)$.

Let $\Phi(t)$ be the characteristic function of X_i . We now determine the characteristic function $\Phi_n(t)$ of $\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$:

Here we make use of the **Landau symbol** “ o ”. In general, if we write $u(x) = o(v(x))$ for $x \rightarrow L$, this implies $\lim_{x \rightarrow L} \frac{u(x)}{v(x)} = 0$, i.e., $u(x)$ goes to 0 faster than $v(x)$ or $v(x)$ goes to ∞ faster than $u(x)$. We say that $u(x)$ *is of smaller order than* $v(x)$ as $x \rightarrow L$. Examples are $\frac{1}{x^3} = o(\frac{1}{x^2})$ and $x^2 = o(x^3)$ for $x \rightarrow \infty$. See Rohatgi, page 6, for more details on the **Landau symbols** “ O ” and “ o ”.

■

Definition 6.4.5:

Let X_1, X_2 be iid non-degenerate rv's with common cdf F . Let $a_1, a_2 > 0$. We say that F is **stable** if there exist constants A and B (depending on a_1 and a_2) such that $B^{-1}(a_1X_1 + a_2X_2 - A)$ also has cdf F . ■

Note:

When generalizing the previous definition to sequences of rv's, we have the following examples for stable distributions:

- X_i iid Cauchy. Then $\frac{1}{n} \sum_{i=1}^n X_i \sim \text{Cauchy}$ (here $B_n = n, A_n = 0$).
 - X_i iid $N(0, 1)$. Then $\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \sim N(0, 1)$ (here $B_n = \sqrt{n}, A_n = 0$).
-

Definition 6.4.6:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid rv's with common cdf F . Let $T_n = \sum_{i=1}^n X_i$. F belongs to the **domain of attraction** of a distribution V if there exist norming and centering constants $\{B_n\}_{n=1}^{\infty}, B_n > 0$, and $\{A_n\}_{n=1}^{\infty}$ such that

$$P(B_n^{-1}(T_n - A_n) \leq x) = F_{B_n^{-1}(T_n - A_n)}(x) \rightarrow V(x) \text{ as } n \rightarrow \infty$$

at all continuity points x of V . ■

Note:

A very general Theorem from Loève states that only stable distributions can have domains of attraction. From the practical point of view, a wide class of distributions F belong to the domain of attraction of the Normal distribution. ■

Theorem 6.4.7: Lindeberg Central Limit Theorem

Let $\{X_i\}_{i=1}^\infty$ be a sequence of independent non-degenerate rv's with cdf's $\{F_i\}_{i=1}^\infty$. Assume that $E(X_k) = \mu_k$ and $Var(X_k) = \sigma_k^2 < \infty$. Let $s_n^2 = \sum_{k=1}^n \sigma_k^2$.

If the F_k are absolutely continuous with pdf's $f_k = F'_k$, assume that it holds for all $\epsilon > 0$ that

$$(A) \quad \lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_{k=1}^n \int_{\{|x-\mu_k| > \epsilon s_n\}} (x - \mu_k)^2 F'_k(x) dx = 0.$$

If the X_k are discrete rv's with support $\{x_{kl}\}$ and probabilities $\{p_{kl}\}$, $l = 1, 2, \dots$, assume that it holds for all $\epsilon > 0$ that

$$(B) \quad \lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_{k=1}^n \sum_{|x_{kl} - \mu_k| > \epsilon s_n} (x_{kl} - \mu_k)^2 p_{kl} = 0.$$

The conditions (A) and (B) are called **Lindeberg Condition (LC)**. If either LC holds, then

$$\frac{\sum_{k=1}^n (X_k - \mu_k)}{s_n} \xrightarrow{d} Z$$

where $Z \sim N(0, 1)$.

Proof:

Similar to the proof of Theorem 6.4.4, we can use characteristic functions again. An alternative proof is given in Rohatgi, pages 282–288. ■

Note:

Feller shows that the LC is a necessary condition if $\frac{\sigma_k^2}{s_n^2} \rightarrow 0$ and $s_n^2 \rightarrow \infty$ as $n \rightarrow \infty$. ■

Corollary 6.4.8:

Let $\{X_i\}_{i=1}^\infty$ be a sequence of iid rv's such that $\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$ has the same distribution for all n . If $E(X_i) = 0$ and $Var(X_i) = 1$, then $X_i \sim N(0, 1)$.

Proof:

Let F be the common cdf of $\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$ for all n (including $n = 1$). By the CLT,

$$\lim_{n \rightarrow \infty} P\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \leq x\right) = \Phi(x),$$

where $\Phi(x)$ denotes $P(Z \leq x)$ for $Z \sim N(0, 1)$. Also, $P\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \leq x\right) = F(x)$ for each n .

Therefore, we must have $F(x) = \Phi(x)$. ■

Note:

In general, if X_1, X_2, \dots , are independent rv's such that there exists a constant A with $P(|X_n| \leq A) = 1 \quad \forall n$, then the LC is satisfied if $s_n^2 \rightarrow \infty$ as $n \rightarrow \infty$. Why??

Suppose that $s_n^2 \rightarrow \infty$ as $n \rightarrow \infty$. Since the $|X_k|$'s are uniformly bounded (by A), so are the rv's $(X_k - E(X_k))$. Thus, for every $\epsilon > 0$ there exists an N_ϵ such that if $n \geq N_\epsilon$ then

$$P(|X_k - E(X_k)| < \epsilon s_n, \quad k = 1, \dots, n) = 1.$$

This implies that the LC holds since we would integrate (or sum) over the empty set, i.e., the set $\{|x - \mu_k| > \epsilon s_n\} = \emptyset$.

The converse also holds. For a sequence of uniformly bounded independent rv's, a necessary and sufficient condition for the CLT to hold is that $s_n^2 \rightarrow \infty$ as $n \rightarrow \infty$. ■

Example 6.4.9:

Let $\{X_i\}_{i=1}^\infty$ be a sequence of independent rv's such that $E(X_k) = 0$, $\alpha_k = E(|X_k|^{2+\delta}) < \infty$ for some $\delta > 0$, and $\sum_{k=1}^n \alpha_k = o(s_n^{2+\delta})$.

Does the LC hold? It is:

$$\begin{aligned} \frac{1}{s_n^2} \sum_{k=1}^n \int_{\{|x| > \epsilon s_n\}} x^2 f_k(x) dx &\stackrel{(A)}{\leq} \frac{1}{s_n^2} \sum_{k=1}^n \int_{\{|x| > \epsilon s_n\}} \frac{|x|^{2+\delta}}{\epsilon^\delta s_n^\delta} f_k(x) dx \\ &\leq \frac{1}{s_n^2 \epsilon^\delta s_n^\delta} \sum_{k=1}^n \int_{-\infty}^{\infty} |x|^{2+\delta} f_k(x) dx \\ &= \frac{1}{s_n^2 \epsilon^\delta s_n^\delta} \sum_{k=1}^n \alpha_k \\ &= \frac{1}{\epsilon^\delta} \left(\frac{\sum_{k=1}^n \alpha_k}{s_n^{2+\delta}} \right) \\ &\stackrel{(B)}{\rightarrow} 0 \quad \text{as } n \rightarrow \infty \end{aligned}$$

(A) holds since for $|x| > \epsilon s_n$, it is $\frac{|x|^\delta}{\epsilon^\delta s_n^\delta} > 1$. (B) holds since $\sum_{k=1}^n \alpha_k = o(s_n^{2+\delta})$.

Thus, the LC is satisfied and the CLT holds. ■

Note:

- (i) In general, if there exists a $\delta > 0$ such that

$$\frac{\sum_{k=1}^n E(|X_k - \mu_k|^{2+\delta})}{s_n^{2+\delta}} \longrightarrow 0 \text{ as } n \rightarrow \infty,$$

then the LC holds.

- (ii) Both the CLT and the WLLN hold for a large class of sequences of rv's $\{X_i\}_{i=1}^n$. If the $\{X_i\}$'s are independent uniformly bounded rv's, i.e., if $P(|X_n| \leq M) = 1 \quad \forall n$, the WLLN (as formulated in Theorem 6.2.3) holds. The CLT holds provided that $s_n^2 \rightarrow \infty$ as $n \rightarrow \infty$.

If the rv's $\{X_i\}$ are iid, then the CLT is a stronger result than the WLLN since the CLT provides an estimate of the probability $P(\frac{1}{n} \sum_{i=1}^n X_i - n\mu \geq \epsilon) \approx 1 - P(|Z| \leq \frac{\epsilon}{\sigma} \sqrt{n})$, where $Z \sim N(0, 1)$, and the WLLN follows. However, note that the CLT requires the existence of a 2^{nd} moment while the WLLN does not.

- (iii) If the $\{X_i\}$ are independent (but not identically distributed) rv's, the CLT may apply while the WLLN does not.
- (iv) See Rohatgi, pages 289–293, and Rohatgi/Saleh, pages 299–303, for additional details and examples.

■

7 Sample Moments

7.1 Random Sampling

(Based on Casella/Berger, Section 5.1 & 5.2)

Definition 7.1.1:

Let X_1, \dots, X_n be iid rv's with common cdf F . We say that $\{X_1, \dots, X_n\}$ is a **(random) sample** of size n from the **population distribution** F . The vector of values $\{x_1, \dots, x_n\}$ is called a **realization** of the sample. A rv $g(X_1, \dots, X_n)$ which is a Borel-measurable function of X_1, \dots, X_n and does not depend on any unknown parameter is called a **(sample) statistic**. ■

Definition 7.1.2:

Let X_1, \dots, X_n be a sample of size n from a population with distribution F . Then

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

is called the **sample mean** and

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right)$$

is called the **sample variance**. ■

Definition 7.1.3:

Let X_1, \dots, X_n be a sample of size n from a population with distribution F . The function

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, x]}(X_i)$$

is called **empirical cumulative distribution function** (empirical cdf). ■

Note:

For any fixed $x \in \mathbb{R}$, $\hat{F}_n(x)$ is a rv. ■

Theorem 7.1.4:

The rv $\hat{F}_n(x)$ has pmf

$$P(\hat{F}_n(x) = \frac{j}{n}) = \binom{n}{j} (F(x))^j (1 - F(x))^{n-j}, \quad j \in \{0, 1, \dots, n\},$$

with $E(\hat{F}_n(x)) = F(x)$ and $Var(\hat{F}_n(x)) = \frac{F(x)(1-F(x))}{n}$.

Proof:

It is $I_{(-\infty, x]}(X_i) \sim Bin(1, F(x))$. Then $n\hat{F}_n(x) \sim Bin(n, F(x))$.

The results follow immediately. ■

Corollary 7.1.5:

By the WLLN, it follows that

$$\hat{F}_n(x) \xrightarrow{p} F(x).$$

■

Corollary 7.1.6:

By the CLT, it follows that

$$\frac{\sqrt{n}(\hat{F}_n(x) - F(x))}{\sqrt{F(x)(1-F(x))}} \xrightarrow{d} Z,$$

where $Z \sim N(0, 1)$. ■

Theorem 7.1.7: Glivenko–Cantelli Theorem

$\hat{F}_n(x)$ converges uniformly to $F(x)$, i.e., it holds for all $\epsilon > 0$ that

$$\lim_{n \rightarrow \infty} P\left(\sup_{-\infty < x < \infty} |\hat{F}_n(x) - F(x)| > \epsilon\right) = 0.$$

■

Definition 7.1.8:

Let X_1, \dots, X_n be a sample of size n from a population with distribution F . We call

$$a_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

the **sample moment of order k** and

$$b_k = \frac{1}{n} \sum_{i=1}^n (X_i - a_1)^k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k$$

the **sample central moment of order k** . ■

Note:

It is $b_1 = 0$ and $b_2 = \frac{n-1}{n} S^2$. ■

Theorem 7.1.9:

Let X_1, \dots, X_n be a sample of size n from a population with distribution F . Assume that $E(X) = \mu$, $Var(X) = \sigma^2$, and $E((X - \mu)^k) = \mu_k$ exist. Then it holds:

- (i) $E(a_1) = E(\bar{X}) = \mu$
- (ii) $Var(a_1) = Var(\bar{X}) = \frac{\sigma^2}{n}$
- (iii) $E(b_2) = \frac{n-1}{n}\sigma^2$
- (iv) $Var(b_2) = \frac{\mu_4 - \mu_2^2}{n} - \frac{2(\mu_4 - 2\mu_2^2)}{n^2} + \frac{\mu_4 - 3\mu_2^2}{n^3}$
- (v) $E(S^2) = \sigma^2$
- (vi) $Var(S^2) = \frac{\mu_4}{n} - \frac{n-3}{n(n-1)}\mu_2^2$

Proof:

- (i)

See Casella/Berger, Page 214, and Rohatgi, page 303–306, for the proof of parts (iv) through (vi) and results regarding the 3rd and 4th moments and covariances. ■

7.2 Sample Moments and the Normal Distribution

(Based on Casella/Berger, Section 5.3)

Theorem 7.2.1:

Let X_1, \dots, X_n be iid $N(\mu, \sigma^2)$ rv's. Then $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $(X_1 - \bar{X}, \dots, X_n - \bar{X})$ are independent.

Proof:

By computing the joint mgf of $(\bar{X}, X_1 - \bar{X}, X_2 - \bar{X}, \dots, X_n - \bar{X})$, we can use Theorem 4.6.3 (iv) to show independence. We will use the following two facts:

(1):

From (1) and (2), it follows:

■

Corollary 7.2.2:

\bar{X} and S^2 are independent.

Proof:

This can be seen since S^2 is a function of the vector $(X_1 - \bar{X}, \dots, X_n - \bar{X})$, and $(X_1 - \bar{X}, \dots, X_n - \bar{X})$ is independent of \bar{X} , as previously shown in Theorem 7.2.1. We can use Theorem 4.2.7 to formally complete this proof. ■

Corollary 7.2.3:

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Proof:

Recall the following facts:

- If $Z \sim N(0, 1)$ then $Z^2 \sim \chi_1^2$.
- If $Y_1, \dots, Y_n \sim \text{iid } \chi_1^2$, then $\sum_{i=1}^n Y_i \sim \chi_n^2$.
- For χ_n^2 , the mgf is $M(t) = (1 - 2t)^{-n/2}$.
- If $X_i \sim N(\mu, \sigma^2)$, then $\frac{X_i - \mu}{\sigma} \sim N(0, 1)$ and $\frac{(X_i - \mu)^2}{\sigma^2} \sim \chi_1^2$.

$$\text{Therefore, } \sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^2} \sim \chi_n^2 \text{ and } \frac{(\bar{X} - \mu)^2}{(\frac{\sigma}{\sqrt{n}})^2} = n \frac{(\bar{X} - \mu)^2}{\sigma^2} \sim \chi_1^2. \quad (*)$$

Now consider

■

Corollary 7.2.4:

$$\frac{\sqrt{n}(\bar{X} - \mu)}{S} \sim t_{n-1}.$$

Proof:

Recall the following facts:

- If $Z \sim N(0, 1)$, $Y \sim \chi_n^2$ and Z, Y independent, then $\frac{Z}{\sqrt{\frac{Y}{n}}} \sim t_n$.
- $Z_1 = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N(0, 1)$, $Y_{n-1} = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$, and Z_1, Y_{n-1} are independent.

Therefore,

$$\frac{\sqrt{n}(\bar{X} - \mu)}{S} = \frac{\frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}}}{\frac{S/\sqrt{n}}{\sigma/\sqrt{n}}} = \frac{\frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}}}{\sqrt{\frac{S^2(n-1)}{\sigma^2(n-1)}}} = \frac{Z_1}{\sqrt{\frac{Y_{n-1}}{(n-1)}}} \sim t_{n-1}.$$

■

Corollary 7.2.5:

Let $(X_1, \dots, X_m) \sim \text{iid } N(\mu_1, \sigma_1^2)$ and $(Y_1, \dots, Y_n) \sim \text{iid } N(\mu_2, \sigma_2^2)$. Let X_i, Y_j be independent $\forall i, j$.

Then it holds:

$$\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{[(m-1)S_1^2/\sigma_1^2] + [(n-1)S_2^2/\sigma_2^2]}} \cdot \sqrt{\frac{m+n-2}{\sigma_1^2/m + \sigma_2^2/n}} \sim t_{m+n-2}$$

In particular, if $\sigma_1 = \sigma_2$, then:

$$\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{(m-1)S_1^2 + (n-1)S_2^2}} \cdot \sqrt{\frac{mn(m+n-2)}{m+n}} \sim t_{m+n-2}$$

Proof:

Homework.

■

Corollary 7.2.6:

Let $(X_1, \dots, X_m) \sim \text{iid } N(\mu_1, \sigma_1^2)$ and $(Y_1, \dots, Y_n) \sim \text{iid } N(\mu_2, \sigma_2^2)$. Let X_i, Y_j be independent $\forall i, j$.

Then it holds:

$$\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \sim F_{m-1, n-1}$$

In particular, if $\sigma_1 = \sigma_2$, then:

$$\frac{S_1^2}{S_2^2} \sim F_{m-1, n-1}$$

Proof:

Recall that, if $Y_1 \sim \chi_m^2$ and $Y_2 \sim \chi_n^2$, then

$$F = \frac{Y_1/m}{Y_2/n} \sim F_{m,n}.$$

Now, $C_1 = \frac{(m-1)S_1^2}{\sigma_1^2} \sim \chi_{m-1}^2$ and $C_2 = \frac{(n-1)S_2^2}{\sigma_2^2} \sim \chi_{n-1}^2$. Therefore,

$$\frac{C_1/(m-1)}{C_2/(n-1)} = \frac{\frac{(m-1)S_1^2}{\sigma_1^2(m-1)}}{\frac{(n-1)S_2^2}{\sigma_2^2(n-1)}} = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \sim F_{m-1,n-1}.$$

If $\sigma_1 = \sigma_2$, then

$$\frac{S_1^2}{S_2^2} \sim F_{m-1,n-1}.$$

■

8 The Theory of Point Estimation

(Based on Casella/Berger, Chapters 6 & 7)

8.1 The Problem of Point Estimation

Let \underline{X} be a rv defined on a probability space (Ω, L, P) . Suppose that the cdf F of \underline{X} depends on some set of parameters and that the functional form of F is known except for a finite number of these parameters.

Definition 8.1.1:

The set of admissible values of θ is called the **parameter space** Θ . If F_θ is the cdf of \underline{X} when θ is the parameter, the set $\{F_\theta : \theta \in \Theta\}$ is the **family of cdf's**. Likewise, we speak of the **family of pdf's** if \underline{X} is continuous, and the **family of pmf's** if \underline{X} is discrete. ■

Example 8.1.2:

$X \sim Bin(n, p)$, p unknown. Then $\theta = p$ and $\Theta = \{p : 0 < p < 1\}$.

$X \sim N(\mu, \sigma^2)$, (μ, σ^2) unknown. Then $\theta = (\mu, \sigma^2)$ and $\Theta = \{(\mu, \sigma^2) : -\infty < \mu < \infty, \sigma^2 > 0\}$. ■

Definition 8.1.3:

Let \underline{X} be a sample from F_θ , $\theta \in \Theta \subseteq \mathbb{R}$. Let a statistic $T(\underline{X})$ map \mathbb{R}^n to Θ . We call $T(\underline{X})$ an **estimator** of θ and $T(\underline{x})$ for a realization \underline{x} of \underline{X} an **(point) estimate** of θ . In practice, the term *estimate* is used for both. ■

Example 8.1.4:

Let X_1, \dots, X_n be iid $Bin(1, p)$, p unknown. Estimates of p include:

$$T_1(\underline{X}) = \bar{X}, T_2(\underline{X}) = X_1, T_3(\underline{X}) = \frac{1}{2}, T_4(\underline{X}) = \frac{X_1 + X_2}{3}$$

Obviously, not all estimates are equally good. ■

8.2 Properties of Estimates

Definition 8.2.1:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid rv's with cdf F_{θ} , $\theta \in \Theta$. A sequence of point estimates $T_n(X_1, \dots, X_n) = T_n$ is called

- **(weakly) consistent** for θ if $T_n \xrightarrow{p} \theta$ as $n \rightarrow \infty \forall \theta \in \Theta$
- **strongly consistent** for θ if $T_n \xrightarrow{a.s.} \theta$ as $n \rightarrow \infty \forall \theta \in \Theta$
- **consistent in the r^{th} mean** for θ if $T_n \xrightarrow{r} \theta$ as $n \rightarrow \infty \forall \theta \in \Theta$

■

Example 8.2.2:

Let $\{X_i\}_{i=1}^{\infty}$ be a sequence of iid $Bin(1, p)$ rv's. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Since $E(X_i) = p$, it follows by the WLLN that $\bar{X}_n \xrightarrow{p} p$, i.e., consistency, and by the SLLN that $\bar{X}_n \xrightarrow{a.s.} p$, i.e., strong consistency.

However, a consistent estimate may not be unique. We may even have infinite many consistent estimates, e.g.,

$$\frac{\sum_{i=1}^n X_i + a}{n + b} \xrightarrow{p} p \quad \forall \text{ finite } a, b \in \mathbb{R}.$$

■

Theorem 8.2.3:

If T_n is a sequence of estimates such that $E(T_n) \rightarrow \theta$ and $Var(T_n) \rightarrow 0$ as $n \rightarrow \infty$, then T_n is consistent for θ .

Proof:

■

Definition 8.2.4:

Let \mathcal{G} be a group of Borel-measurable functions of \mathbb{R}^n onto itself which is closed under composition and inverse. A family of distributions $\{P_\theta : \theta \in \Theta\}$ is **invariant** under \mathcal{G} if for each $g \in \mathcal{G}$ and for all $\theta \in \Theta$, there exists a unique $\theta' = \bar{g}(\theta)$ such that the distribution of $g(\underline{X})$ is $P_{\theta'}$ whenever the distribution of \underline{X} is P_θ . We call \bar{g} the **induced function** on θ since $P_\theta(g(\underline{X}) \in A) = P_{\bar{g}(\theta)}(\underline{X} \in A)$. ■

Example 8.2.5:

Let (X_1, \dots, X_n) be iid $N(\mu, \sigma^2)$ with pdf

$$f(x_1, \dots, x_n) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right).$$

The group of linear transformations \mathcal{G} has elements

$$g(x_1, \dots, x_n) = (ax_1 + b, \dots, ax_n + b), \quad a > 0, \quad -\infty < b < \infty.$$

The pdf of $g(\underline{X})$ is

$$f^*(x_1^*, \dots, x_n^*) = \frac{1}{(\sqrt{2\pi}a\sigma)^n} \exp\left(-\frac{1}{2a^2\sigma^2} \sum_{i=1}^n (x_i^* - a\mu - b)^2\right), \quad x_i^* = ax_i + b, \quad i = 1, \dots, n.$$

So $\{f : -\infty < \mu < \infty, \sigma^2 > 0\}$ is invariant under this group \mathcal{G} , with $\bar{g}(\mu, \sigma^2) = (a\mu + b, a^2\sigma^2)$, where $-\infty < a\mu + b < \infty$ and $a^2\sigma^2 > 0$. ■

Definition 8.2.6:

Let \mathcal{G} be a group of transformations that leaves $\{F_\theta : \theta \in \Theta\}$ invariant. An estimate T is **invariant** under \mathcal{G} if

$$T(g(X_1), \dots, g(X_n)) = T(X_1, \dots, X_n) \quad \forall g \in \mathcal{G}.$$

Definition 8.2.7:

An estimate T is:

- **location invariant** if $T(X_1 + a, \dots, X_n + a) = T(X_1, \dots, X_n)$, $a \in \mathbb{R}$
- **scale invariant** if $T(cX_1, \dots, cX_n) = T(X_1, \dots, X_n)$, $c \in \mathbb{R} - \{0\}$
- **permutation invariant** if $T(X_{i_1}, \dots, X_{i_n}) = T(X_1, \dots, X_n) \quad \forall$ permutations (i_1, \dots, i_n) of $1, \dots, n$

Example 8.2.8: ■

Let $F_\theta \sim N(\mu, \sigma^2)$.

S^2 is location invariant.

\bar{X} and S^2 are both permutation invariant.

Neither \bar{X} nor S^2 is scale invariant. ■

Note:

Different sources make different use of the term *invariant*. Mood, Graybill & Boes (1974) for example define *location invariant* as $T(X_1 + a, \dots, X_n + a) = T(X_1, \dots, X_n) + a$ (page 332) and *scale invariant* as $T(cX_1, \dots, cX_n) = cT(X_1, \dots, X_n)$ (page 336). According to their definition, \bar{X} is location invariant and scale invariant. ■

8.3 Sufficient Statistics

(Based on Casella/Berger, Section 6.2)

Definition 8.3.1:

Let $\underline{X} = (X_1, \dots, X_n)$ be a sample from $\{F_\theta : \theta \in \Theta \subseteq \mathbb{R}^k\}$. A statistic $T = T(\underline{X})$ is **sufficient** for θ (or for the family of distributions $\{F_\theta : \theta \in \Theta\}$) iff the conditional distribution of \underline{X} given $T = t$ does not depend on θ (except possibly on a null set A where $P_\theta(T \in A) = 0 \ \forall \theta$). ■

Note:

- (i) The sample \underline{X} is always sufficient but this is not particularly interesting and usually is excluded from further considerations.
- (ii) Idea: Once we have “reduced” from \underline{X} to $T(\underline{X})$, we have captured all the information in \underline{X} about θ .
- (iii) Usually, there are several sufficient statistics for a given family of distributions. ■

Example 8.3.2:

Let $\underline{X} = (X_1, \dots, X_n)$ be iid $Bin(1, p)$ rv's. To estimate p , can we ignore the order and simply count the number of “successes”?

Let $T(\underline{X}) = \sum_{i=1}^n X_i$. It is ■

Example 8.3.3:

Let $\underline{X} = (X_1, \dots, X_n)$ be iid $\text{Poisson}(\lambda)$. Is $T = \sum_{i=1}^n X_i$ sufficient for λ ? It is

■

Example 8.3.4:

Let X_1, X_2 be iid $\text{Poisson}(\lambda)$. Is $T = X_1 + 2X_2$ sufficient for λ ? It is

■

Note:

Definition 8.3.1 can be difficult to check. In addition, it requires a candidate statistic. We need something constructive that helps in finding sufficient statistics without having to check Definition 8.3.1. The next Theorem helps in finding such statistics. ■

Theorem 8.3.5: Factorization Criterion

Let X_1, \dots, X_n be rv's with pdf (or pmf) $f(x_1, \dots, x_n | \theta)$, $\theta \in \Theta$. Then $T(X_1, \dots, X_n)$ is sufficient for θ iff we can write

$$f(x_1, \dots, x_n | \theta) = h(x_1, \dots, x_n) g(T(x_1, \dots, x_n) | \theta),$$

where h does not depend on θ and g does not depend on x_1, \dots, x_n except as a function of T .

Proof:

Discrete case only.

“ \implies ”:

Suppose $T(\underline{X})$ is sufficient for θ . Let

“ \impliedby ”:

Suppose the factorization holds. For fixed t_0 , it is

■

Note:

- (i) In the Theorem above, θ and T may be vectors.
- (ii) If T is sufficient for θ , then also any 1-to-1 mapping of T is sufficient for θ . However, this does not hold for arbitrary functions of T .

■

Example 8.3.6:

Let X_1, \dots, X_n be iid $Bin(1, p)$. It is

■

Example 8.3.7:

Let X_1, \dots, X_n be iid $Poisson(\lambda)$. It is

■

Example 8.3.8:

Let X_1, \dots, X_n be iid $N(\mu, \sigma^2)$ where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$ are both unknown. It is

■

Example 8.3.9:

Let X_1, \dots, X_n be iid $U(\theta, \theta + 1)$ where $-\infty < \theta < \infty$. It is

■

Definition 8.3.10:

Let $\{f_\theta(x) : \theta \in \Theta\}$ be a family of pdf's (or pmf's). We say the family is **complete** if

$$E_\theta(g(X)) = 0 \quad \forall \theta \in \Theta$$

implies that

$$P_\theta(g(X) = 0) = 1 \quad \forall \theta \in \Theta.$$

We say a statistic $T(X)$ is **complete** if the family of distributions of T is complete. ■

Example 8.3.11:

Let X_1, \dots, X_n be iid $Bin(1, p)$. We have seen in Example 8.3.6 that $T = \sum_{i=1}^n X_i$ is sufficient for p . Is it also complete?

We know that $T \sim Bin(n, p)$. Thus,

■

Example 8.3.12:

Let X_1, \dots, X_n be iid $N(\theta, \theta^2)$. We know from Example 8.3.8 that $T = (\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ is sufficient for θ . Is it also complete?

■

Note:

Recall from Section 5.2 what it means if we say the family of distributions $\{f_\theta : \theta \in \Theta\}$ is a one-parameter (or k -parameter) exponential family. ■

Theorem 8.3.13:

Let $\{f_\theta : \theta \in \Theta\}$ be a k -parameter exponential family. Let T_1, \dots, T_k be statistics. Then the family of distributions of $(T_1(\underline{X}), \dots, T_k(\underline{X}))$ is also a k -parameter exponential family given by

$$g_\theta(\underline{t}) = \exp\left(\sum_{i=1}^k t_i Q_i(\theta) + D(\theta) + S^*(\underline{t})\right)$$

for suitable $S^*(t)$.

Proof:

The proof follows from our Theorems regarding the transformation of rv's. ■

Theorem 8.3.14:

Let $\{f_\theta : \theta \in \Theta\}$ be a k -parameter exponential family with $k \leq n$ and let T_1, \dots, T_k be statistics as in Theorem 8.3.13. Suppose that the range of $\underline{Q} = (Q_1, \dots, Q_k)$ contains an open set in \mathbb{R}^k . Then $\underline{T} = (T_1(\underline{X}), \dots, T_k(\underline{X}))$ is a complete sufficient statistic.

Proof:

Discrete case and $k = 1$ only.

Write $Q(\theta) = \theta$ and let $(a, b) \subseteq \Theta$.

It follows from the Factorization Criterion (Theorem 8.3.5) that T is sufficient for θ . Thus, we only have to show that T is complete, i.e., that

$$\begin{aligned} E_\theta(g(T(\underline{X}))) &= \sum_t g(t) P_\theta(T(\underline{X}) = t) \\ &\stackrel{(A)}{=} \sum_t g(t) \exp(\theta t + D(\theta) + S^*(t)) = 0 \quad \forall \theta \quad (B) \end{aligned}$$

implies $g(t) = 0 \quad \forall t$. Note that in (A) we make use of a result established in Theorem 8.3.13.

We now define functions g^+ and g^- as:

$$\begin{aligned} g^+(t) &= \begin{cases} g(t), & \text{if } g(t) \geq 0 \\ 0, & \text{otherwise} \end{cases} \\ g^-(t) &= \begin{cases} -g(t), & \text{if } g(t) < 0 \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

It is $g(t) = g^+(t) - g^-(t)$ where both functions, g^+ and g^- , are non-negative functions. Using g^+ and g^- , it turns out that (B) is equivalent to

$$\sum_t g^+(t) \exp(\theta t + S^*(t)) = \sum_t g^-(t) \exp(\theta t + S^*(t)) \quad \forall \theta \quad (C)$$

where the term $\exp(D(\theta))$ in (A) drops out as a constant on both sides.

If we fix $\theta_0 \in (a, b)$ and define

$$p^+(t) = \frac{g^+(t) \exp(\theta_0 t + S^*(t))}{\sum_t g^+(t) \exp(\theta_0 t + S^*(t))}, \quad p^-(t) = \frac{g^-(t) \exp(\theta_0 t + S^*(t))}{\sum_t g^-(t) \exp(\theta_0 t + S^*(t))},$$

it is obvious that $p^+(t) \geq 0 \quad \forall t$ and $p^-(t) \geq 0 \quad \forall t$ and by construction $\sum_t p^+(t) = 1$ and $\sum_t p^-(t) = 1$. Hence, p^+ and p^- are both pmf's.

From (C), it follows for the mgf's M^+ and M^- of p^+ and p^- that

$$\begin{aligned} M^+(\delta) &= \sum_t e^{\delta t} p^+(t) \\ &= \frac{\sum_t e^{\delta t} g^+(t) \exp(\theta_0 t + S^*(t))}{\sum_t g^+(t) \exp(\theta_0 t + S^*(t))} \\ &= \frac{\sum_t g^+(t) \exp((\theta_0 + \delta)t + S^*(t))}{\sum_t g^+(t) \exp(\theta_0 t + S^*(t))} \\ &\stackrel{(C)}{=} \frac{\sum_t g^-(t) \exp((\theta_0 + \delta)t + S^*(t))}{\sum_t g^-(t) \exp(\theta_0 t + S^*(t))} \\ &= \frac{\sum_t e^{\delta t} g^-(t) \exp(\theta_0 t + S^*(t))}{\sum_t g^-(t) \exp(\theta_0 t + S^*(t))} \\ &= \sum_t e^{\delta t} p^-(t) \\ &= M^-(\delta) \quad \forall \delta \in \underbrace{(a - \theta_0)}_{<0}, \underbrace{(b - \theta_0)}_{>0}. \end{aligned}$$

By the uniqueness of mgf's it follows that $p^+(t) = p^-(t) \quad \forall t$.

$$\implies g^+(t) = g^-(t) \quad \forall t$$

$$\implies g(t) = 0 \quad \forall t$$

$\implies T$ is complete ■

Definition 8.3.15:

Let $\underline{X} = (X_1, \dots, X_n)$ be a sample from $\{F_\theta : \theta \in \Theta \subseteq \mathbb{R}^k\}$ and let $T = T(\underline{X})$ be a sufficient statistic for θ . $T = T(\underline{X})$ is called a **minimal sufficient** statistic for θ if, for any other sufficient statistic $T' = T'(\underline{X})$, $T(\underline{x})$ is a function of $T'(\underline{x})$. ■

Note:

- (i) A minimal sufficient statistic achieves the greatest possible data reduction for a sufficient statistic.
- (ii) If T is minimal sufficient for θ , then also any 1-to-1 mapping of T is minimal sufficient for θ . However, this does not hold for arbitrary functions of T . ■

Definition 8.3.16:

Let $\underline{X} = (X_1, \dots, X_n)$ be a sample from $\{F_\theta : \theta \in \Theta \subseteq \mathbb{R}^k\}$. A statistic $T = T(\underline{X})$ is called **ancillary** if its distribution does not depend on the parameter θ . ■

Example 8.3.17:

Let X_1, \dots, X_n be iid $U(\theta, \theta + 1)$ where $-\infty < \theta < \infty$. As shown in Example 8.3.9, $T = (X_{(1)}, X_{(n)})$ is sufficient for θ . Define

$$R_n = X_{(n)} - X_{(1)}.$$

Use the result from Stat 6710, Homework Assignment 5, Question (viii) (a) to obtain

$$f_{R_n}(r | \theta) = f_{R_n}(r) = n(n-1)r^{n-2}(1-r)I_{(0,1)}(r).$$

This means that $R_n \sim \text{Beta}(n-1, 2)$. Moreover, R_n does not depend on θ and, therefore, R_n is ancillary. ■

Theorem 8.3.18: Basu's Theorem

Let $\underline{X} = (X_1, \dots, X_n)$ be a sample from $\{F_\theta : \theta \in \Theta \subseteq \mathbb{R}^k\}$. If $T = T(\underline{X})$ is a complete and minimal sufficient statistic, then T is independent of any ancillary statistic. ■

Theorem 8.3.19:

Let $\underline{X} = (X_1, \dots, X_n)$ be a sample from $\{F_\theta : \theta \in \Theta \subseteq \mathbb{R}^k\}$. If any minimal sufficient statistic $T = T(\underline{X})$ exists for θ , then any complete statistic is also a minimal sufficient statistic. ■

Note:

- (i) Due to the last Theorem, Basu's Theorem often only is stated in terms of a complete sufficient statistic (which automatically is also a minimal sufficient statistic).
- (ii) As already shown in Corollary 7.2.2, \bar{X} and S^2 are independent when sampling from a $N(\mu, \sigma^2)$ population. As outlined in Casella/Berger, page 289, we could also use Basu's Theorem to obtain the same result.
- (iii) The converse of Basu's Theorem is false, i.e., if $T(\underline{X})$ is independent of any ancillary statistic, it does not necessarily follow that $T(\underline{X})$ is a complete, minimal sufficient statistic.
- (iv) As seen in Examples 8.3.8 and 8.3.12, $T = (\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ is sufficient for θ but it is not complete when X_1, \dots, X_n are iid $N(\theta, \theta^2)$. However, it can be shown that T is minimal sufficient. So, there may be distributions where a minimal sufficient statistic exists but a complete statistic does not exist.
- (v) As with invariance, there exist several different definitions of ancillarity within the literature — the one defined in this chapter being the most commonly used.

■

8.4 Unbiased Estimation

(Based on Casella/Berger, Section 7.3)

Definition 8.4.1:

Let $\{F_\theta : \theta \in \Theta\}$, $\Theta \subseteq \mathbb{R}$, be a nonempty set of cdf's. A Borel-measurable function T from \mathbb{R}^n to Θ is called **unbiased** for θ (or an unbiased estimate for θ) if

$$E_\theta(T) = \theta \quad \forall \theta \in \Theta.$$

Any function $d(\theta)$ for which an unbiased estimate T exists is called an **estimable function**.

If T is biased,

$$b(\theta, T) = E_\theta(T) - \theta$$

is called the **bias** of T . ■

Example 8.4.2:

If the k^{th} population moment exists, the k^{th} sample moment is an unbiased estimate. If $\text{Var}(X) = \sigma^2$, the sample variance S^2 is an unbiased estimate of σ^2 .

However, note that for X_1, \dots, X_n iid $N(\mu, \sigma^2)$ S is not an unbiased estimate of σ :

$$\begin{aligned} \frac{(n-1)S^2}{\sigma^2} &\sim \chi_{n-1}^2 = \text{Gamma}\left(\frac{n-1}{2}, 2\right) \\ \implies E\left(\sqrt{\frac{(n-1)S^2}{\sigma^2}}\right) &= \int_0^\infty \sqrt{x} \frac{x^{\frac{n-1}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n-1}{2}} \Gamma(\frac{n-1}{2})} dx \\ &= \frac{\sqrt{2}\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \int_0^\infty \frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} dx \\ &\stackrel{(*)}{=} \frac{\sqrt{2}\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \\ \implies E(S) &= \sigma \sqrt{\frac{2}{n-1} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})}} \end{aligned}$$

(*) holds since $\frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}$ is the pdf of a $\text{Gamma}(\frac{n}{2}, 2)$ distribution and thus the integral is 1.

So S is biased for σ and

$$b(\sigma, S) = \sigma \left(\sqrt{\frac{2}{n-1} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})}} - 1 \right). \quad \blacksquare$$

Note:

If T is unbiased for θ , $g(T)$ is not necessarily unbiased for $g(\theta)$ (unless g is a linear function). ■

Example 8.4.3:

Unbiased estimates may not exist (see Rohatgi, page 351, Example 2) or they may be absurd as in the following case:

Let $X \sim \text{Poisson}(\lambda)$ and let $d(\lambda) = e^{-2\lambda}$. Consider $T(X) = (-1)^X$ as an estimate. It is

■

Note:

If there exist 2 unbiased estimates T_1 and T_2 of θ , then any estimate of the form $\alpha T_1 + (1 - \alpha)T_2$ for $0 < \alpha < 1$ will also be an unbiased estimate of θ . Which one should we choose? ■

Definition 8.4.4:

The **mean square error** of an estimate T of σ is defined as

$$\begin{aligned}MSE(\theta, T) &= E_{\theta}((T - \theta)^2) \\ &= \text{Var}_{\theta}(T) + (b(\theta, T))^2.\end{aligned}$$

Let $\{T_i\}_{i=1}^{\infty}$ be a sequence of estimates of θ . If

$$\lim_{i \rightarrow \infty} MSE(\theta, T_i) = 0 \quad \forall \theta \in \Theta,$$

then $\{T_i\}$ is called a **mean-squared-error consistent (MSE-consistent)** sequence of estimates of θ . ■

Note:

- (i) If we allow all estimates and compare their MSE, generally it will depend on θ which estimate is better. For example $\hat{\theta} = 17$ is perfect if $\theta = 17$, but it is lousy otherwise.
- (ii) If we restrict ourselves to the class of unbiased estimates, then $MSE(\theta, T) = \text{Var}_{\theta}(T)$.

(iii) MSE-consistency means that both the bias and the variance of T_i approach 0 as $i \rightarrow \infty$. ■

Definition 8.4.5:

Let $\theta_0 \in \Theta$ and let $U(\theta_0)$ be the class of all unbiased estimates T of θ_0 such that $E_{\theta_0}(T^2) < \infty$. Then $T_0 \in U(\theta_0)$ is called a **locally minimum variance unbiased estimate (LMVUE)** at θ_0 if

$$E_{\theta_0}((T_0 - \theta_0)^2) \leq E_{\theta_0}((T - \theta_0)^2) \quad \forall T \in U(\theta_0).$$
 ■

Definition 8.4.6:

Let U be the class of all unbiased estimates T of $\theta \in \Theta$ such that $E_{\theta}(T^2) < \infty \quad \forall \theta \in \Theta$. Then $T_0 \in U$ is called a **uniformly minimum variance unbiased estimate (UMVUE)** of θ if

$$E_{\theta}((T_0 - \theta)^2) \leq E_{\theta}((T - \theta)^2) \quad \forall \theta \in \Theta \quad \forall T \in U.$$
 ■

An Excursion into Logic II

In our first “Excursion into Logic” in Stat 6710 Mathematical Statistics I, we have established the following results:

$A \Rightarrow B$ is equivalent to $\neg B \Rightarrow \neg A$ is equivalent to $\neg A \vee B$:

A	B	$A \Rightarrow B$	$\neg A$	$\neg B$	$\neg B \Rightarrow \neg A$	$\neg A \vee B$
1	1	1	0	0	1	1
1	0	0	0	1	0	0
0	1	1	1	0	1	1
0	0	1	1	1	1	1

When dealing with formal proofs, there exists one more technique to show $A \Rightarrow B$. Equivalently, we can show $(A \wedge \neg B) \Rightarrow 0$, a technique called **Proof by Contradiction**. This means, assuming that A and $\neg B$ hold, we show that this implies 0, i.e., something that is always false, i.e., a contradiction. And here is the corresponding truth table:

A	B	$A \Rightarrow B$	$\neg B$	$A \wedge \neg B$	$(A \wedge \neg B) \Rightarrow 0$
1	1				
1	0				
0	1				
0	0				

Note:

We make use of this proof technique in the Proof of the next Theorem. ■

Example:

Let A : $x = 5$ and B : $x^2 = 25$. Obviously $A \Rightarrow B$.

But we can also prove this in the following way:

A : $x = 5$ and $\neg B$: $x^2 \neq 25$

$$\implies x^2 = 25 \wedge x^2 \neq 25$$

This is impossible, i.e., a contradiction. Thus, $A \Rightarrow B$. ■

Theorem 8.4.7:

Let U be the class of all unbiased estimates T of $\theta \in \Theta$ with $E_{\theta}(T^2) < \infty \quad \forall \theta$, and suppose that U is non-empty. Let U_0 be the set of all unbiased estimates of 0, i.e.,

$$U_0 = \{\nu : E_{\theta}(\nu) = 0, E_{\theta}(\nu^2) < \infty \quad \forall \theta \in \Theta\}.$$

Then $T_0 \in U$ is UMVUE iff

$$E_{\theta}(\nu T_0) = 0 \quad \forall \theta \in \Theta \quad \forall \nu \in U_0.$$

Proof:

Note that $E_{\theta}(\nu T_0)$ always exists.

■

Theorem 8.4.8:

Let U be the non-empty class of unbiased estimates of $\theta \in \Theta$ as defined in Theorem 8.4.7. Then there exists at most one UMVUE $T \in U$ for θ .

Proof:

Suppose $T_0, T_1 \in U$ are both UMVUE.

Then $T_1 - T_0 \in U$, $Var_\theta(T_0) = Var_\theta(T_1)$, and $E_\theta(T_0(T_1 - T_0)) = 0 \quad \forall \theta \in \Theta$

$$\implies E_\theta(T_0^2) = E_\theta(T_0 T_1)$$

$$\implies Cov_\theta(T_0, T_1) = E_\theta(T_0 T_1) - E_\theta(T_0)E_\theta(T_1)$$

$$= E_\theta(T_0^2) - (E_\theta(T_0))^2$$

$$= Var_\theta(T_0)$$

$$= Var_\theta(T_1) \quad \forall \theta \in \Theta$$

$$\implies \rho_{T_0 T_1} = 1 \quad \forall \theta \in \Theta$$

$$\implies P_\theta(aT_0 + bT_1 = 0) = 1 \quad \text{for some } a, b \quad \forall \theta \in \Theta$$

$$\implies \theta = E_\theta(T_0) = E_\theta(-\frac{b}{a}T_1) = E_\theta(T_1) \quad \forall \theta \in \Theta$$

$$\implies -\frac{b}{a} = 1$$

$$\implies P_\theta(T_0 = T_1) = 1 \quad \forall \theta \in \Theta \quad \blacksquare$$

Theorem 8.4.9:

- (i) If an UMVUE T exists for a real function $d(\theta)$, then λT is the UMVUE for $\lambda d(\theta)$, $\lambda \in \mathbb{R}$.
- (ii) If UMVUE's T_1 and T_2 exist for real functions $d_1(\theta)$ and $d_2(\theta)$, respectively, then $T_1 + T_2$ is the UMVUE for $d_1(\theta) + d_2(\theta)$.

Proof:

Homework. \blacksquare

Theorem 8.4.10:

If a sample consists of n independent observations X_1, \dots, X_n from the same distribution, the UMVUE, if it exists, is permutation invariant.

Proof:

Homework. ■

Theorem 8.4.11: Rao–Blackwell

Let $\{F_\theta : \theta \in \Theta\}$ be a family of cdf's, and let h be any statistic in U , where U is the non-empty class of all unbiased estimates of θ with $E_\theta(h^2) < \infty$. Let T be a sufficient statistic for $\{F_\theta : \theta \in \Theta\}$. Then the conditional expectation $E_\theta(h | T)$ is independent of θ and it is an unbiased estimate of θ . Additionally,

$$E_\theta((E(h | T) - \theta)^2) \leq E_\theta((h - \theta)^2) \quad \forall \theta \in \Theta$$

with equality iff $h = E(h | T)$.

Proof:

By Theorem 4.7.3, $E_\theta(E(h | T)) = E(h) = \theta$.

■

Theorem 8.4.12: Lehmann–Scheffée

If T is a complete sufficient statistic and if there exists an unbiased estimate h of θ , then $E(h | T)$ is the (unique) UMVUE.

Proof:

■

Note:

We can use Theorem 8.4.12 to find the UMVUE in two ways if we have a complete sufficient statistic T :

- (i) If we can find an unbiased estimate $h(T)$, it will be the UMVUE since $E(h(T) | T) = h(T)$.
- (ii) If we have any unbiased estimate h and if we can calculate $E(h | T)$, then $E(h | T)$ will be the UMVUE. The process of determining the UMVUE this way often is called *Rao–Blackwellization*.
- (iii) Even if a complete sufficient statistic does not exist, the UMVUE may still exist (see Rohatgi, page 357–358, Example 10).

■

Example 8.4.13:

Let X_1, \dots, X_n be iid $Bin(1, p)$. Then $T = \sum_{i=1}^n X_i$ is a complete sufficient statistic as seen in Examples 8.3.6 and 8.3.11.

Since $E(X_1) = p$, X_1 is an unbiased estimate of p . However, due to part (i) of the Note above, since X_1 is not a function of T , X_1 is not the UMVUE.

We can use part (ii) of the Note above to construct the UMVUE. It is

If we are interested in the UMVUE for $d(p) = p(1 - p) = p - p^2 = \text{Var}(X)$, we can find it in the following way:

■

8.5 Lower Bounds for the Variance of an Estimate

(Based on Casella/Berger, Section 7.3)

Theorem 8.5.1: Cramér–Rao Lower Bound (CRLB)

Let Θ be an open interval of \mathbb{R} . Let $\{f_\theta : \theta \in \Theta\}$ be a family of pdf's or pmf's. Assume that the set $\{\underline{x} : f_\theta(\underline{x}) = 0\}$ is independent of θ .

Let $\psi(\theta)$ be defined on Θ and let it be differentiable for all $\theta \in \Theta$. Let T be an unbiased estimate of $\psi(\theta)$ such that $E_\theta(T^2) < \infty \quad \forall \theta \in \Theta$. Suppose that

(i) $\frac{\partial f_\theta(\underline{x})}{\partial \theta}$ is defined for all $\theta \in \Theta$,

(ii) for a pdf f_θ

$$\frac{\partial}{\partial \theta} \left(\int f_\theta(\underline{x}) d\underline{x} \right) = \int \frac{\partial f_\theta(\underline{x})}{\partial \theta} d\underline{x} = 0 \quad \forall \theta \in \Theta$$

or for a pmf f_θ

$$\frac{\partial}{\partial \theta} \left(\sum_{\underline{x}} f_\theta(\underline{x}) \right) = \sum_{\underline{x}} \frac{\partial f_\theta(\underline{x})}{\partial \theta} = 0 \quad \forall \theta \in \Theta,$$

(iii) for a pdf f_θ

$$\frac{\partial}{\partial \theta} \left(\int T(\underline{x}) f_\theta(\underline{x}) d\underline{x} \right) = \int T(\underline{x}) \frac{\partial f_\theta(\underline{x})}{\partial \theta} d\underline{x} \quad \forall \theta \in \Theta$$

or for a pmf f_θ

$$\frac{\partial}{\partial \theta} \left(\sum_{\underline{x}} T(\underline{x}) f_\theta(\underline{x}) \right) = \sum_{\underline{x}} T(\underline{x}) \frac{\partial f_\theta(\underline{x})}{\partial \theta} \quad \forall \theta \in \Theta.$$

Let $\chi : \Theta \rightarrow \mathbb{R}$ be any measurable function. Then it holds

$$(\psi'(\theta))^2 \leq E_\theta((T(\underline{X}) - \chi(\theta))^2) E_\theta \left(\left(\frac{\partial \log f_\theta(\underline{X})}{\partial \theta} \right)^2 \right) \quad \forall \theta \in \Theta \quad (A).$$

Further, for any $\theta_0 \in \Theta$, either $\psi'(\theta_0) = 0$ and equality holds in (A) for $\theta = \theta_0$, or we have

$$E_{\theta_0}((T(\underline{X}) - \chi(\theta_0))^2) \geq \frac{(\psi'(\theta_0))^2}{E_{\theta_0} \left(\left(\frac{\partial \log f_{\theta_0}(\underline{X})}{\partial \theta} \right)^2 \right)} \quad (B).$$

Finally, if equality holds in (B), then there exists a real number $K(\theta_0) \neq 0$ such that

$$T(\underline{X}) - \chi(\theta_0) = K(\theta_0) \left. \frac{\partial \log f_{\theta_0}(\underline{X})}{\partial \theta} \right|_{\theta=\theta_0} \quad (C)$$

with probability 1, provided that T is not a constant. ■

Note:

- (i) Conditions (i), (ii), and (iii) are called **regularity conditions**. Conditions under which they hold can be found in Rohatgi, page 11–13, Parts 12 and 13.
- (ii) The right hand side of inequality (B) is called *Cramér–Rao Lower Bound* of θ_0 , or, in symbols $CRLB(\theta_0)$.

■

Proof:

From (ii), we get

$$\begin{aligned} E_{\theta} \left(\frac{\partial}{\partial \theta} \log f_{\theta}(\underline{X}) \right) &= \int \left(\frac{\partial}{\partial \theta} \log f_{\theta}(\underline{x}) \right) f_{\theta}(\underline{x}) d\underline{x} \\ &= \int \left(\frac{\partial}{\partial \theta} f_{\theta}(\underline{x}) \right) \frac{1}{f_{\theta}(\underline{x})} f_{\theta}(\underline{x}) d\underline{x} \\ &= \int \left(\frac{\partial}{\partial \theta} f_{\theta}(\underline{x}) \right) d\underline{x} \\ &= 0 \\ \implies E_{\theta} \left(\chi(\theta) \frac{\partial}{\partial \theta} \log f_{\theta}(\underline{X}) \right) &= 0 \end{aligned}$$

From (iii), we get

$$\begin{aligned} E_{\theta} \left(T(\underline{X}) \frac{\partial}{\partial \theta} \log f_{\theta}(\underline{X}) \right) &= \int \left(T(\underline{x}) \frac{\partial}{\partial \theta} \log f_{\theta}(\underline{x}) \right) f_{\theta}(\underline{x}) d\underline{x} \\ &= \int \left(T(\underline{x}) \frac{\partial}{\partial \theta} f_{\theta}(\underline{x}) \right) \frac{1}{f_{\theta}(\underline{x})} f_{\theta}(\underline{x}) d\underline{x} \\ &= \int \left(T(\underline{x}) \frac{\partial}{\partial \theta} f_{\theta}(\underline{x}) \right) d\underline{x} \\ &\stackrel{(iii)}{=} \frac{\partial}{\partial \theta} \left(\int T(\underline{x}) f_{\theta}(\underline{x}) d\underline{x} \right) \\ &= \frac{\partial}{\partial \theta} E(T(\underline{X})) \\ &= \frac{\partial}{\partial \theta} \psi(\theta) \\ &= \psi'(\theta) \\ \implies E_{\theta} \left((T(\underline{X}) - \chi(\theta)) \frac{\partial}{\partial \theta} \log f_{\theta}(\underline{X}) \right) &= \psi'(\theta) \end{aligned}$$

$$\begin{aligned} \implies (\psi'(\theta))^2 &= \left(E_\theta \left((T(\underline{X}) - \chi(\theta)) \frac{\partial}{\partial \theta} \log f_\theta(\underline{X}) \right) \right)^2 \\ &\stackrel{(*)}{\leq} E_\theta \left((T(\underline{X}) - \chi(\theta))^2 \right) E_\theta \left(\left(\frac{\partial}{\partial \theta} \log f_\theta(\underline{X}) \right)^2 \right), \end{aligned}$$

i.e., (A) holds. (*) follows from the Cauchy–Schwarz–Inequality (Theorem 4.5.7 (ii)).

If $\psi'(\theta_0) \neq 0$, then the left–hand side of (A) is > 0 . Therefore, the right–hand side is > 0 . Thus,

$$E_{\theta_0} \left(\left(\frac{\partial}{\partial \theta} \log f_{\theta_0}(\underline{X}) \right)^2 \right) > 0,$$

and (B) follows directly from (A).

If $\psi'(\theta_0) = 0$, but equality does not hold in (A), then

$$E_{\theta_0} \left(\left(\frac{\partial}{\partial \theta} \log f_{\theta_0}(\underline{X}) \right)^2 \right) > 0,$$

and (B) follows directly from (A) again.

Finally, if equality holds in (B), then $\psi'(\theta_0) \neq 0$ (because T is not constant). Thus, $MSE(\chi(\theta_0), T(\underline{X})) > 0$. The Cauchy–Schwarz–Inequality (Theorem 4.5.7 (iii)) gives equality iff there exist constants $(\alpha, \beta) \in \mathbb{R}^2 - \{(0, 0)\}$ such that

$$P \left(\alpha(T(\underline{X}) - \chi(\theta_0)) + \beta \left(\frac{\partial}{\partial \theta} \log f_\theta(\underline{X}) \Big|_{\theta=\theta_0} \right) = 0 \right) = 1.$$

This implies $K(\theta_0) = -\frac{\beta}{\alpha}$ and (C) holds. Since T is not a constant, it also holds that $K(\theta_0) \neq 0$. ■

Example 8.5.2:

If we take $\chi(\theta) = \psi(\theta)$, we get from (B)

$$Var_\theta(T(\underline{X})) \geq \frac{(\psi'(\theta))^2}{E_\theta \left(\left(\frac{\partial \log f_\theta(\underline{X})}{\partial \theta} \right)^2 \right)} \quad (*).$$

If we have $\psi(\theta) = \theta$, the inequality (*) above reduces to

$$Var_\theta(T(\underline{X})) \geq \left(E_\theta \left(\left(\frac{\partial \log f_\theta(\underline{X})}{\partial \theta} \right)^2 \right) \right)^{-1}.$$

Finally, if $\underline{X} = (X_1, \dots, X_n)$ iid with identical $f_\theta(x)$, the inequality (*) reduces to

$$Var_\theta(T(\underline{X})) \geq \frac{(\psi'(\theta))^2}{n E_\theta \left(\left(\frac{\partial \log f_\theta(X_1)}{\partial \theta} \right)^2 \right)}.$$

■

Example 8.5.3:

Let X_1, \dots, X_n be iid $Bin(1, p)$. Let $X \sim Bin(n, p)$, $p \in \Theta = (0, 1) \subset \mathbb{R}$. Let

$$\psi(p) = E(T(X)) = \sum_{x=0}^n T(x) \binom{n}{x} p^x (1-p)^{n-x}.$$

$\psi(p)$ is differentiable with respect to p under the summation sign since it is a finite polynomial in p .

■

Example 8.5.4:

Let $X \sim U(0, \theta)$, $\theta \in \Theta = (0, \infty) \subset \mathbb{R}$.

■

Theorem 8.5.5: Chapman, Robbins, Kiefer Inequality (CRK Inequality)

Let $\Theta \subseteq \mathbb{R}$. Let $\{f_\theta : \theta \in \Theta\}$ be a family of pdf's or pmf's. Let $\psi(\theta)$ be defined on Θ . Let T be an unbiased estimate of $\psi(\theta)$ such that $E_\theta(T^2) < \infty \quad \forall \theta \in \Theta$.

If $\theta \neq \vartheta$, θ and $\vartheta \in \Theta$, assume that $f_\theta(x)$ and $f_\vartheta(x)$ are different. Also assume that there exists such a $\vartheta \in \Theta$ such that $\theta \neq \vartheta$ and

$$S(\theta) = \{\underline{x} : f_\theta(\underline{x}) > 0\} \supset S(\vartheta) = \{\underline{x} : f_\vartheta(\underline{x}) > 0\}.$$

Then it holds that

$$Var_\theta(T(\underline{X})) \geq \sup_{\{\vartheta : S(\vartheta) \subset S(\theta), \vartheta \neq \theta\}} \frac{(\psi(\vartheta) - \psi(\theta))^2}{Var_\theta\left(\frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})}\right)} \quad \forall \theta \in \Theta.$$

Proof:

Since T is unbiased, it follows

$$E_\vartheta(T(\underline{X})) = \psi(\vartheta) \quad \forall \vartheta \in \Theta.$$

For $\vartheta \neq \theta$ and $S(\vartheta) \subset S(\theta)$, it follows

$$\int_{S(\theta)} T(\underline{x}) \frac{f_\vartheta(\underline{x}) - f_\theta(\underline{x})}{f_\theta(\underline{x})} f_\theta(\underline{x}) d\underline{x} = E_\vartheta(T(\underline{X})) - E_\theta(T(\underline{X})) = \psi(\vartheta) - \psi(\theta)$$

and

$$0 = \int_{S(\theta)} \frac{f_\vartheta(\underline{x}) - f_\theta(\underline{x})}{f_\theta(\underline{x})} f_\theta(\underline{x}) d\underline{x} = E_\theta\left(\frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})} - 1\right).$$

Therefore

$$Cov_\theta\left(T(\underline{X}), \frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})} - 1\right) = \psi(\vartheta) - \psi(\theta).$$

It follows by the Cauchy–Schwarz–Inequality (Theorem 4.5.7 (ii)) that

$$\begin{aligned} (\psi(\vartheta) - \psi(\theta))^2 &= \left(Cov_\theta\left(T(\underline{X}), \frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})} - 1\right)\right)^2 \\ &\leq Var_\theta(T(\underline{X})) Var_\theta\left(\frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})} - 1\right) \\ &= Var_\theta(T(\underline{X})) Var_\theta\left(\frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})}\right). \end{aligned}$$

Thus,

$$Var_\theta(T(\underline{X})) \geq \frac{(\psi(\vartheta) - \psi(\theta))^2}{Var_\theta\left(\frac{f_\vartheta(\underline{X})}{f_\theta(\underline{X})}\right)}.$$

Finally, we take the supremum of the right–hand side with respect to $\{\vartheta : S(\vartheta) \subset S(\theta), \vartheta \neq \theta\}$, which completes the proof. ■

Note:

(i) The CRK inequality holds without the previous regularity conditions.

(ii) An alternative form of the CRK inequality is:

Let $\theta, \theta + \delta$, $\delta \neq 0$, be distinct with $S(\theta + \delta) \subset S(\theta)$. Let $\psi(\theta) = \theta$. Define

$$J = J(\theta, \delta) = \frac{1}{\delta^2} \left(\left(\frac{f_{\theta+\delta}(\underline{X})}{f_{\theta}(\underline{X})} \right)^2 - 1 \right).$$

Then the CRK inequality reads as

$$\text{Var}_{\theta}(T(\underline{X})) \geq \frac{1}{\inf_{\delta} E_{\theta}(J)}$$

with the infimum taken over $\delta \neq 0 : S(\theta + \delta) \subset S(\theta)$.

(iii) The CRK inequality works for discrete Θ , the CRLB does not work in such cases. ■

Example 8.5.6:

Let $X \sim U(0, \theta)$, $\theta > 0$. The required conditions for the CRLB are not met. Recall from Example 8.5.4 that $\frac{n+1}{n}X_{(n)}$ is UMVUE with $\text{Var}(\frac{n+1}{n}X_{(n)}) = \frac{\theta^2}{n(n+2)} < \frac{\theta^2}{n} = \text{CRLB}$.

Definition 8.5.7:

Let T_1, T_2 be unbiased estimates of θ with $E_{\theta}(T_1^2) < \infty$ and $E_{\theta}(T_2^2) < \infty \forall \theta \in \Theta$. We define the **efficiency** of T_1 relative to T_2 by

$$\text{eff}_{\theta}(T_1, T_2) = \frac{\text{Var}_{\theta}(T_2)}{\text{Var}_{\theta}(T_1)}$$

and say that T_1 is **more efficient** than T_2 if $\text{eff}_{\theta}(T_1, T_2) < 1$. ■

Definition 8.5.8:

Assume the regularity conditions of Theorem 8.5.1 are satisfied by a family of cdf's $\{F_\theta : \theta \in \Theta\}$. An unbiased estimate T for θ is **most efficient** for $\{F_\theta\}$ if

$$Var_\theta(T) = \left(E_\theta \left(\left(\frac{\partial \log f_\theta(\underline{X})}{\partial \theta} \right)^2 \right) \right)^{-1}$$

■

Definition 8.5.9:

Let T be the most efficient estimate for the family of cdf's $\{F_\theta : \theta \in \Theta\}$, $\Theta \subseteq \mathcal{R}$. Then the **efficiency** of any unbiased T_1 of θ is defined as

$$ef f_\theta(T_1) = ef f_\theta(T_1, T) = \frac{Var_\theta(T_1)}{Var_\theta(T)}.$$

■

Definition 8.5.10:

T_1 is **asymptotically (most) efficient** if T_1 is asymptotically unbiased, i.e., $\lim_{n \rightarrow \infty} E_\theta(T_1) = \theta$, and $\lim_{n \rightarrow \infty} ef f_\theta(T_1) = 1$, where n is the sample size.

■

Theorem 8.5.11:

A necessary and sufficient condition for an estimate T of θ to be most efficient is that T is sufficient and

$$\frac{1}{K(\theta)}(T(\underline{x}) - \theta) = \frac{\partial \log f_\theta(\underline{x})}{\partial \theta} \quad \forall \theta \in \Theta \quad (*),$$

where $K(\theta)$ is defined as in Theorem 8.5.1 and the regularity conditions for Theorem 8.5.1 hold.

Proof:

“ \implies .”

Theorem 8.5.1 says that if T is most efficient, then (*) holds.

Assume that $\Theta = \mathcal{R}$. We define

$$C(\theta_0) = \int_{-\infty}^{\theta_0} \frac{1}{K(\theta)} d\theta, \quad \psi(\theta_0) = \int_{-\infty}^{\theta_0} \frac{\theta}{K(\theta)} d\theta, \quad \text{and } \lambda(\underline{x}) = \lim_{\theta \rightarrow -\infty} \log f_\theta(\underline{x}) - c(\underline{x}).$$

Integrating (*) with respect to θ gives

$$\begin{aligned} \int_{-\infty}^{\theta_0} \frac{1}{K(\theta)} T(\underline{x}) d\theta - \int_{-\infty}^{\theta_0} \frac{\theta}{K(\theta)} d\theta &= \int_{-\infty}^{\theta_0} \frac{\partial \log f_\theta(\underline{x})}{\partial \theta} d\theta \\ \implies T(\underline{x})C(\theta_0) - \psi(\theta_0) &= \log f_\theta(\underline{x}) \Big|_{-\infty}^{\theta_0} + c(\underline{x}) \\ \implies T(\underline{x})C(\theta_0) - \psi(\theta_0) &= \log f_{\theta_0}(\underline{x}) - \lim_{\theta \rightarrow -\infty} \log f_\theta(\underline{x}) + c(\underline{x}) \\ \implies T(\underline{x})C(\theta_0) - \psi(\theta_0) &= \log f_{\theta_0}(\underline{x}) - \lambda(\underline{x}) \end{aligned}$$

Therefore,

$$f_{\theta_0}(\underline{x}) = \exp(T(\underline{x})C(\theta_0) - \psi(\theta_0) + \lambda(\underline{x}))$$

which belongs to an exponential family. Thus, T is sufficient.

“ \Leftarrow :”

From (*), we get

$$E_{\theta} \left(\left(\frac{\partial \log f_{\theta}(\underline{X})}{\partial \theta} \right)^2 \right) = \frac{1}{(K(\theta))^2} \text{Var}_{\theta}(T(\underline{X})).$$

Additionally, it holds

$$E_{\theta} \left((T(\underline{X}) - \theta) \frac{\partial \log f_{\theta}(\underline{X})}{\partial \theta} \right) = 1$$

as shown in the Proof of Theorem 8.5.1.

Using (*) in the line above, we get

$$K(\theta) E_{\theta} \left(\left(\frac{\partial \log f_{\theta}(\underline{X})}{\partial \theta} \right)^2 \right) = 1,$$

i.e.,

$$K(\theta) = \left(E_{\theta} \left(\left(\frac{\partial \log f_{\theta}(\underline{X})}{\partial \theta} \right)^2 \right) \right)^{-1}.$$

Therefore,

$$\text{Var}_{\theta}(T(\underline{X})) = \left(E_{\theta} \left(\left(\frac{\partial \log f_{\theta}(\underline{X})}{\partial \theta} \right)^2 \right) \right)^{-1},$$

i.e., T is most efficient for θ . ■

Note:

Instead of saying “a necessary and sufficient condition for an estimate T of θ to be most efficient ...” in the previous Theorem, we could say that “an estimate T of θ is most efficient iff ...”, i.e., “necessary and sufficient” means the same as “iff”.

A is necessary for B means: $B \Rightarrow A$ (because $\neg A \Rightarrow \neg B$)

A is sufficient for B means: $A \Rightarrow B$ ■

8.6 The Method of Moments

(Based on Casella/Berger, Section 7.2.1)

Definition 8.6.1:

Let X_1, \dots, X_n be iid with pdf (or pmf) f_θ , $\theta \in \Theta$. We assume that first k moments m_1, \dots, m_k of f_θ exist. If θ can be written as

$$\theta = h(m_1, \dots, m_k),$$

where $h : \mathbb{R}^k \rightarrow \mathbb{R}$ is a Borel-measurable function, the **method of moments estimate (mom)** of θ is

$$\hat{\theta}_{mom} = T(X_1, \dots, X_n) = h\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2, \dots, \frac{1}{n} \sum_{i=1}^n X_i^k\right).$$

■

Note:

- (i) The Definition above can also be used to estimate joint moments. For example, we use $\frac{1}{n} \sum_{i=1}^n X_i Y_i$ to estimate $E(XY)$.
- (ii) Since $E\left(\frac{1}{n} \sum_{i=1}^n X_i^j\right) = m_j$, method of moment estimates are unbiased for the population moments. The WLLN and the CLT say that these estimates are consistent and asymptotically Normal as well.
- (iii) If θ is not a linear function of the population moments, $\hat{\theta}_{mom}$ will, in general, not be unbiased. However, it will be consistent and (usually) asymptotically Normal.
- (iv) Method of moments estimates do not exist if the related moments do not exist.
- (v) Method of moments estimates may not be unique. If there exist multiple choices for the mom, one usually takes the estimate involving the lowest-order sample moment.
- (vi) Alternative method of moment estimates can be obtained from central moments (rather than from raw moments) or by using moments other than the first k moments.

■

Example 8.6.2:

Let X_1, \dots, X_n be iid $N(\mu, \sigma^2)$.

Since $\mu = m_1$, it is $\hat{\mu}_{mom} = \bar{X}$.

This is an unbiased, consistent and asymptotically Normal estimate.

Since $\sigma = \sqrt{m_2 - m_1^2}$, it is $\hat{\sigma}_{mom} = \sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2}$.

This is a consistent, asymptotically Normal estimate. However, it is not unbiased. ■

Example 8.6.3:

Let X_1, \dots, X_n be iid Poisson(λ).

We know that $E(X_1) = Var(X_1) = \lambda$.

Thus, \bar{X} and $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ are possible choices for the mom of λ . Due to part (v) of the

Note above, one uses $\hat{\lambda}_{mom} = \bar{X}$. ■

8.7 Maximum Likelihood Estimation

(Based on Casella/Berger, Section 7.2.2)

Definition 8.7.1:

Let (X_1, \dots, X_n) be an n -rv with pdf (or pmf) $f_\theta(x_1, \dots, x_n)$, $\theta \in \Theta$. We call the function of θ

$$L(\theta; x_1, \dots, x_n) = f_\theta(x_1, \dots, x_n)$$

the **likelihood function**. ■

Note:

(i) Often θ is a vector of parameters.

(ii) If (X_1, \dots, X_n) are iid with pdf (or pmf) $f_\theta(x)$, then $L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i)$. ■

Definition 8.7.2:

A **maximum likelihood estimate (MLE)** is a non-constant estimate $\hat{\theta}_{ML}$ such that

$$L(\hat{\theta}_{ML}; x_1, \dots, x_n) = \sup_{\theta \in \Theta} L(\theta; x_1, \dots, x_n).$$

Note:

It is often convenient to work with $\log L$ when determining the maximum likelihood estimate. Since the log is monotone, the maximum is the same. ■

Example 8.7.3:

Let X_1, \dots, X_n be iid $N(\mu, \sigma^2)$, where μ and σ^2 are unknown.

$$L(\mu, \sigma^2; x_1, \dots, x_n) = \frac{1}{\sigma^n (2\pi)^{\frac{n}{2}}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Formally, we still have to verify that we found the maximum (and not a minimum) and that there is no parameter θ at the edge of the parameter space Θ such that the likelihood function

does not take its absolute maximum which is not detectable by using our approach for local extrema. ■

Example 8.7.4:

Let X_1, \dots, X_n be iid $U(\theta - \frac{1}{2}, \theta + \frac{1}{2})$.

Example 8.7.5:

Let $X \sim \text{Bin}(1, p)$, $p \in [\frac{1}{4}, \frac{3}{4}]$.

$$L(p; x) = p^x(1-p)^{1-x} = \begin{cases} p, & \text{if } x = 1 \\ 1-p, & \text{if } x = 0 \end{cases}$$

Theorem 8.7.6:

Let T be a sufficient statistic for $f_\theta(\underline{x})$, $\theta \in \Theta$. If a unique MLE of θ exists, it is a function of T .

Proof:

Since T is sufficient, we can write

$$f_\theta(\underline{x}) = h(\underline{x})g_\theta(T(\underline{x}))$$

due to the Factorization Criterion (Theorem 8.3.5). Maximizing the likelihood function with respect to θ takes $h(\underline{x})$ as a constant and therefore is equivalent to maximizing $g_\theta(\underline{x})$ with respect to θ . But $g_\theta(\underline{x})$ involves \underline{x} only through T . ■

Note:

- (i) MLE's may not be unique (however they frequently are).
- (ii) MLE's are not necessarily unbiased.
- (iii) MLE's may not exist.
- (iv) If a unique MLE exists, it is a function of a sufficient statistic.
- (v) Often (but not always), the MLE will be a sufficient statistic itself.

Theorem 8.7.7:

Suppose the regularity conditions of Theorem 8.5.1 hold and θ belongs to an open interval in \mathcal{R} . If an estimate $\hat{\theta}$ of θ attains the CRLB, it is the unique MLE.

Proof:

If $\hat{\theta}$ attains the CRLB, it follows by Theorem 8.5.1 that

$$\frac{\partial \log f_\theta(\underline{X})}{\partial \theta} = \frac{1}{K(\theta)}(\hat{\theta}(\underline{X}) - \theta) \text{ w.p. } 1.$$

Thus, $\hat{\theta}$ satisfies the likelihood equations.

We define $A(\theta) = \frac{1}{K(\theta)}$. Then it follows

$$\frac{\partial^2 \log f_\theta(\underline{X})}{\partial \theta^2} = A'(\theta)(\hat{\theta}(\underline{X}) - \theta) - A(\theta).$$

The Proof of Theorem 8.5.11 gives us

$$A(\theta) = E_\theta \left(\left(\frac{\partial \log f_\theta(\underline{X})}{\partial \theta} \right)^2 \right) > 0.$$

So

$$\left. \frac{\partial^2 \log f_\theta(\underline{X})}{\partial \theta^2} \right|_{\theta=\hat{\theta}} = -A(\theta) < 0,$$

i.e., $\log f_\theta(\underline{X})$ has a maximum in $\hat{\theta}$. Thus, $\hat{\theta}$ is the MLE. ■

Note:

The previous Theorem does not imply that every MLE is most efficient. ■

Theorem 8.7.8:

Let $\{f_\theta : \theta \in \Theta\}$ be a family of pdf's (or pmf's) with $\Theta \subseteq \mathbb{R}^k$, $k \geq 1$. Let $h : \Theta \rightarrow \Delta$ be a mapping of Θ onto $\Delta \subseteq \mathbb{R}^p$, $1 \leq p \leq k$. If $\hat{\theta}$ is an MLE of θ , then $h(\hat{\theta})$ is an MLE of $h(\theta)$.

Proof:

For each $\delta \in \Delta$, we define

$$\Theta_\delta = \{\theta : \theta \in \Theta, h(\theta) = \delta\}$$

and

$$M(\delta; \underline{x}) = \sup_{\theta \in \Theta_\delta} L(\theta; \underline{x}),$$

the likelihood function induced by h .

Let $\hat{\theta}$ be an MLE and a member of $\Theta_{\hat{\delta}}$, where $\hat{\delta} = h(\hat{\theta})$. It holds

$$M(\hat{\delta}; \underline{x}) = \sup_{\theta \in \Theta_{\hat{\delta}}} L(\theta; \underline{x}) \geq L(\hat{\theta}; \underline{x}),$$

but also

$$M(\hat{\delta}; \underline{x}) \leq \sup_{\delta \in \Delta} M(\delta; \underline{x}) = \sup_{\delta \in \Delta} \left(\sup_{\theta \in \Theta_\delta} L(\theta; \underline{x}) \right) = \sup_{\theta \in \Theta} L(\theta; \underline{x}) = L(\hat{\theta}; \underline{x}).$$

Therefore,

$$M(\hat{\delta}; \underline{x}) = L(\hat{\theta}; \underline{x}) = \sup_{\delta \in \Delta} M(\delta; \underline{x}).$$

Thus, $\hat{\delta} = h(\hat{\theta})$ is an MLE. ■

Example 8.7.9:

Let X_1, \dots, X_n be iid $Bin(1, p)$. Let $h(p) = p(1 - p)$.

Since the MLE of p is $\hat{p} = \bar{X}$, the MLE of $h(p)$ is $h(\hat{p}) = \bar{X}(1 - \bar{X})$. ■

Theorem 8.7.10:

Consider the following conditions a pdf f_θ can fulfill:

(i) $\frac{\partial \log f_\theta}{\partial \theta}$, $\frac{\partial^2 \log f_\theta}{\partial \theta^2}$, $\frac{\partial^3 \log f_\theta}{\partial \theta^3}$ exist for all $\theta \in \Theta$ for all x . Also,

$$\int_{-\infty}^{\infty} \frac{\partial f_\theta(x)}{\partial \theta} dx = E_\theta \left(\frac{\partial \log f_\theta(X)}{\partial \theta} \right) = 0 \quad \forall \theta \in \Theta.$$

(ii) $\int_{-\infty}^{\infty} \frac{\partial^2 f_\theta(x)}{\partial \theta^2} dx = 0 \quad \forall \theta \in \Theta$.

(iii) $-\infty < \int_{-\infty}^{\infty} \frac{\partial^2 \log f_{\theta}(x)}{\partial \theta^2} f_{\theta}(x) dx < 0 \quad \forall \theta \in \Theta.$

(iv) There exists a function $H(x)$ such that for all $\theta \in \Theta$:

$$\left| \frac{\partial^3 \log f_{\theta}(x)}{\partial \theta^3} \right| < H(x) \text{ and } \int_{-\infty}^{\infty} H(x) f_{\theta}(x) dx = M(\theta) < \infty.$$

(v) There exists a function $g(\theta)$ that is positive and twice differentiable for every $\theta \in \Theta$ and there exists a function $H(x)$ such that for all $\theta \in \Theta$:

$$\left| \frac{\partial^2}{\partial \theta^2} \left[g(\theta) \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right] \right| < H(x) \text{ and } \int_{-\infty}^{\infty} H(x) f_{\theta}(x) dx = M(\theta) < \infty.$$

In case that multiple of these conditions are fulfilled, we can make the following statements:

(i) (Cramér) Conditions (i), (iii), and (iv) imply that, with probability approaching 1, as $n \rightarrow \infty$, the likelihood equation has a consistent solution.

(ii) (Cramér) Conditions (i), (ii), (iii), and (iv) imply that a consistent solution $\hat{\theta}_n$ of the likelihood equation is asymptotically Normal, i.e.,

$$\frac{\sqrt{n}}{\sigma} (\hat{\theta}_n - \theta) \xrightarrow{d} Z$$

where $Z \sim N(0, 1)$ and $\sigma^2 = \left(E_{\theta} \left(\left(\frac{\partial \log f_{\theta}(X)}{\partial \theta} \right)^2 \right) \right)^{-1}$.

(iii) (Kulldorf) Conditions (i), (iii), and (v) imply that, with probability approaching 1, as $n \rightarrow \infty$, the likelihood equation has a consistent solution.

(iv) (Kulldorf) Conditions (i), (ii), (iii), and (v) imply that a consistent solution $\hat{\theta}_n$ of the likelihood equation is asymptotically Normal.

■

Note:

In case of a pmf f_{θ} , we can define similar conditions as in Theorem 8.7.10.

■

8.8 Decision Theory — Bayes and Minimax Estimation

(Based on Casella/Berger, Section 7.2.3 & 7.3.4)

Let $\{f_\theta : \theta \in \Theta\}$ be a family of pdf's (or pmf's). Let X_1, \dots, X_n be a sample from f_θ . Let \mathcal{A} be the set of possible **actions** (or decisions) that are open to the statistician in a given situation, e.g.,

$\mathcal{A} = \{\text{reject } H_0, \text{ do not reject } H_0\}$ (Hypothesis testing, see Chapter 9)

$\mathcal{A} = \{\text{artefact found is of } \{\text{Greek, Roman}\} \text{ origin}\}$ (Classification)

$\mathcal{A} = \Theta$ (Estimation)

Definition 8.8.1:

A **decision function** d is a statistic, i.e., a Borel-measurable function, that maps \mathbb{R}^n into \mathcal{A} . If $\underline{X} = \underline{x}$ is observed, the statistician takes action $d(\underline{x}) \in \mathcal{A}$. ■

Note:

For the remainder of this Section, we are restricting ourselves to $\mathcal{A} = \Theta$, i.e., we are facing the problem of estimation. ■

Definition 8.8.2:

A non-negative function L that maps $\Theta \times \mathcal{A}$ into \mathbb{R} is called a **loss function**. The value $L(\theta, a)$ is the loss incurred to the statistician if he/she takes action a when θ is the true parameter value. ■

Definition 8.8.3:

Let \mathcal{D} be a class of decision functions that map \mathbb{R}^n into \mathcal{A} . Let L be a loss function on $\Theta \times \mathcal{A}$. The function R that maps $\Theta \times \mathcal{D}$ into \mathbb{R} is defined as

$$R(\theta, d) = E_\theta(L(\theta, d(\underline{X})))$$

and is called the **risk function** of d at θ . ■

Example 8.8.4:

Let $\mathcal{A} = \Theta \subseteq \mathbb{R}$. Let $L(\theta, a) = (\theta - a)^2$. Then it holds that

$$R(\theta, d) = E_\theta(L(\theta, d(\underline{X}))) = E_\theta((\theta - d(\underline{X}))^2) = E_\theta((\theta - \hat{\theta})^2).$$

Note that this is just the MSE. If $\hat{\theta}$ is unbiased, this would just be $Var(\hat{\theta})$. ■

Note:

The basic problem of decision theory is that we would like to find a decision function $d \in \mathcal{D}$ such that $R(\theta, d)$ is minimized for all $\theta \in \Theta$. Unfortunately, this is usually not possible. ■

Definition 8.8.5:

The **minimax principle** is to choose the decision function $d^* \in \mathcal{D}$ such that

$$\max_{\theta \in \Theta} R(\theta, d^*) \leq \max_{\theta \in \Theta} R(\theta, d) \quad \forall d \in \mathcal{D}.$$

■

Note:

If the problem of interest is an estimation problem, we call a d^* that satisfies the condition in Definition 8.8.5 a **minimax estimate** of θ . ■

Example 8.8.6:

Let $X \sim \text{Bin}(1, p)$, $p \in \Theta = \{\frac{1}{4}, \frac{3}{4}\} = \mathcal{A}$.

We consider the following loss function:

p	a	$L(p, a)$
$\frac{1}{4}$	$\frac{1}{4}$	0
$\frac{1}{4}$	$\frac{3}{4}$	2
$\frac{3}{4}$	$\frac{1}{4}$	5
$\frac{3}{4}$	$\frac{3}{4}$	0

The set of decision functions consists of the following four functions:

$$\begin{aligned} d_1(0) &= \frac{1}{4}, & d_1(1) &= \frac{1}{4} \\ d_2(0) &= \frac{1}{4}, & d_2(1) &= \frac{3}{4} \\ d_3(0) &= \frac{3}{4}, & d_3(1) &= \frac{1}{4} \\ d_4(0) &= \frac{3}{4}, & d_4(1) &= \frac{3}{4} \end{aligned}$$

First, we evaluate the loss function for these four decision functions:

$$\begin{aligned} L\left(\frac{1}{4}, d_1(0)\right) &= L\left(\frac{1}{4}, \frac{1}{4}\right) = \\ L\left(\frac{1}{4}, d_1(1)\right) &= L\left(\frac{1}{4}, \frac{1}{4}\right) = \\ L\left(\frac{3}{4}, d_1(0)\right) &= L\left(\frac{3}{4}, \frac{1}{4}\right) = \end{aligned}$$

$$\begin{aligned}
L\left(\frac{3}{4}, d_1(1)\right) &= L\left(\frac{3}{4}, \frac{1}{4}\right) = \\
L\left(\frac{1}{4}, d_2(0)\right) &= L\left(\frac{1}{4}, \frac{1}{4}\right) = \\
L\left(\frac{1}{4}, d_2(1)\right) &= L\left(\frac{1}{4}, \frac{3}{4}\right) = \\
L\left(\frac{3}{4}, d_2(0)\right) &= L\left(\frac{3}{4}, \frac{1}{4}\right) = \\
L\left(\frac{3}{4}, d_2(1)\right) &= L\left(\frac{3}{4}, \frac{3}{4}\right) = \\
L\left(\frac{1}{4}, d_3(0)\right) &= L\left(\frac{1}{4}, \frac{3}{4}\right) = \\
L\left(\frac{1}{4}, d_3(1)\right) &= L\left(\frac{1}{4}, \frac{1}{4}\right) = \\
L\left(\frac{3}{4}, d_3(0)\right) &= L\left(\frac{3}{4}, \frac{3}{4}\right) = \\
L\left(\frac{3}{4}, d_3(1)\right) &= L\left(\frac{3}{4}, \frac{1}{4}\right) = \\
L\left(\frac{1}{4}, d_4(0)\right) &= L\left(\frac{1}{4}, \frac{3}{4}\right) = \\
L\left(\frac{1}{4}, d_4(1)\right) &= L\left(\frac{1}{4}, \frac{3}{4}\right) = \\
L\left(\frac{3}{4}, d_4(0)\right) &= L\left(\frac{3}{4}, \frac{3}{4}\right) = \\
L\left(\frac{3}{4}, d_4(1)\right) &= L\left(\frac{3}{4}, \frac{3}{4}\right) =
\end{aligned}$$

Then, the risk function

$$R(p, d_i(X)) = E_p(L(p, d(X))) = L(p, d(0)) \cdot P_p(X = 0) + L(p, d(1)) \cdot P_p(X = 1)$$

takes the following values:

i	$p = \frac{1}{4}: R\left(\frac{1}{4}, d_i\right)$	$p = \frac{3}{4}: R\left(\frac{3}{4}, d_i\right)$	$\max_{p \in \{1/4, 3/4\}} R(p, d_i)$
1			
2			
3			
4			

Hence,

$$\min_{i \in \{1, 2, 3, 4\}} \max_{p \in \{1/4, 3/4\}} R(p, d_i) = .$$

Thus, _____ is the minimax estimate. ■

Note:

Minimax estimation does not require any unusual assumptions. However, it tends to be very

conservative. ■

Definition 8.8.7:

Suppose we consider θ to be a rv with pdf $\pi(\theta)$ on Θ . We call π the **a priori distribution** (or **prior distribution**). ■

Note:

$f(\underline{x} | \theta)$ is the conditional density of \underline{x} given a fixed θ . The joint density of \underline{x} and θ is

$$f(\underline{x}, \theta) = \pi(\theta)f(\underline{x} | \theta),$$

the marginal density of \underline{x} is

$$g(\underline{x}) = \int f(\underline{x}, \theta)d\theta,$$

and the **a posteriori distribution** (or **posterior distribution**), which gives the distribution of θ after sampling, has pdf (or pmf)

$$h(\theta | \underline{x}) = \frac{f(\underline{x}, \theta)}{g(\underline{x})}.$$

Definition 8.8.8:

The **Bayes risk** of a decision function d is defined as

$$R(\pi, d) = E_{\pi}(R(\theta, d)),$$

where π is the a priori distribution. ■

Note:

If θ is a continuous rv and \underline{X} is of continuous type, then

$$R(\pi, d) = E_{\pi}(R(\theta, d))$$

Similar expressions can be written if θ and/or \underline{X} are discrete. ■

Definition 8.8.9:

A decision function d^* is called a **Bayes rule** if d^* minimizes the Bayes risk, i.e., if

$$R(\pi, d^*) = \inf_{d \in \mathcal{D}} R(\pi, d).$$

■

Theorem 8.8.10:

Let $\mathcal{A} = \Theta \subseteq \mathbb{R}$. Let $L(\theta, d(\underline{x})) = (\theta - d(\underline{x}))^2$. In this case, a Bayes rule is

$$d(\underline{x}) = E(\theta \mid \underline{X} = \underline{x}).$$

Proof:

Minimizing

$$R(\pi, d) = \int g(\underline{x}) \left(\int (\theta - d(\underline{x}))^2 h(\theta \mid \underline{x}) d\theta \right) d\underline{x},$$

where g is the marginal pdf of \underline{X} and h is the conditional pdf of θ given \underline{x} , is the same as minimizing

$$\int (\theta - d(\underline{x}))^2 h(\theta \mid \underline{x}) d\theta.$$

However, this is minimized when $d(\underline{x}) = E(\theta \mid \underline{X} = \underline{x})$ as shown in Stat 6710, Homework 3, Question (ii), for the unconditional case. ■

Note:

Under the conditions of Theorem 8.8.10, $d(\underline{x}) = E(\theta \mid \underline{X} = \underline{x})$ is called the **Bayes estimate**. ■

Example 8.8.11:

Let $X \sim \text{Bin}(n, p)$. Let $L(p, d(x)) = (p - d(x))^2$.

Let $\pi(p) = 1 \quad \forall p \in (0, 1)$, i.e., $\pi \sim U(0, 1)$, be the a priori distribution of p .

Then it holds:

$$\begin{aligned} f(x, p) &= \binom{n}{x} p^x (1-p)^{n-x} \\ g(x) &= \int f(x, p) dp \\ &= \int_0^1 \binom{n}{x} p^x (1-p)^{n-x} dp \\ h(p \mid x) &= \frac{f(x, p)}{g(x)} \end{aligned}$$

$$\begin{aligned}
&= \frac{\binom{n}{x} p^x (1-p)^{n-x}}{\int_0^1 \binom{n}{x} p^x (1-p)^{n-x} dp} \\
&= \frac{p^x (1-p)^{n-x}}{\int_0^1 p^x (1-p)^{n-x} dp} \\
E(p | x) &=
\end{aligned}$$

Thus, by Theorem 8.8.10, the Bayes rule is

$$\hat{p}_{Bayes} =$$

The Bayes risk of $d^*(X)$ is

$$\begin{aligned}
R(\pi, d^*(X)) &= E_\pi(R(p, d^*(X))) \\
&= \dots
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{(n+2)^2} \int_0^1 (1 - 4p + np - np^2 + 4p^2) dp \\
&= \frac{1}{(n+2)^2} \int_0^1 (1 + (n-4)p + (4-n)p^2) dp \\
&= \frac{1}{(n+2)^2} \left(p + \frac{n-4}{2}p^2 + \frac{4-n}{3}p^3 \right) \Big|_0^1 \\
&= \frac{1}{(n+2)^2} \left(1 + \frac{n-4}{2} + \frac{4-n}{3} \right) \\
&= \frac{1}{(n+2)^2} \frac{6 + 3n - 12 + 8 - 2n}{6} \\
&= \frac{1}{(n+2)^2} \frac{n+2}{6} \\
&= \frac{1}{6(n+2)}
\end{aligned}$$

Now we compare the Bayes rule $d^*(X)$ with the MLE $\hat{p}_{ML} = \frac{X}{n}$. This estimate has Bayes risk

$$R\left(\pi, \frac{X}{n}\right) =$$

■

Theorem 8.8.12:

Let $\{f_\theta : \theta \in \Theta\}$ be a family of pdf's (or pmf's). Suppose that an estimate d^* of θ is a Bayes estimate corresponding to some prior distribution π on Θ . If the risk function $R(\theta, d^*)$ is constant on Θ , then d^* is a minimax estimate of θ . ■

Proof:

Homework. ■

Definition 8.8.13:

Let F denote the class of pdf's (or pmf's) $f_\theta(x)$. A class Π of prior distributions is a **conjugate family** for F if the posterior distribution is in the class Π for all $f \in F$, all priors in Π , and all $x \in \mathcal{X}$. ■

Note:

The beta family is conjugate for the binomial family. Thus, if we start with a beta prior, we will end up with a beta posterior. (See Homework.) ■

9 Hypothesis Testing

9.1 Fundamental Notions

(Based on Casella/Berger, Section 8.1 & 8.3)

We assume that $\underline{X} = (X_1, \dots, X_n)$ is a random sample from a population distribution F_θ , $\theta \in \Theta \subseteq \mathbb{R}^k$, where the functional form of F_θ is known, except for the parameter θ . We also assume that Θ contains at least two points.

Definition 9.1.1:

A **parametric hypothesis** is an assumption about the unknown parameter θ .

The **null hypothesis** is of the form

$$H_0 : \theta \in \Theta_0 \subset \Theta.$$

The **alternative hypothesis** is of the form

$$H_1 : \theta \in \Theta_1 = \Theta - \Theta_0.$$

■

Definition 9.1.2:

If Θ_0 (or Θ_1) contains only one point, we say that H_0 and Θ_0 (or H_1 and Θ_1) are **simple**. In this case, the distribution of \underline{X} is completely specified under the null (or alternative) hypothesis.

If Θ_0 (or Θ_1) contains more than one point, we say that H_0 and Θ_0 (or H_1 and Θ_1) are **composite**.

■

Example 9.1.3:

Let X_1, \dots, X_n be iid $Bin(1, p)$. Examples for hypotheses are $p = \frac{1}{2}$ (simple), $p \geq \frac{1}{2}$ (composite), $p \neq \frac{1}{4}$ (composite), etc.

■

Note:

The problem of testing a hypothesis can be described as follows: Given a sample point \underline{x} , find a decision rule that will lead to a decision to accept or reject the null hypothesis. This means, we partition the space \mathbb{R}^n into two disjoint sets C and C^c such that, if $\underline{x} \in C$, we reject $H_0 : \theta \in \Theta_0$ (and we accept H_1). Otherwise, if $\underline{x} \in C^c$, we accept H_0 that $\underline{X} \sim F_\theta$, $\theta \in \Theta_0$.

■

Definition 9.1.4:

Let $\underline{X} \sim F_\theta$, $\theta \in \Theta$. Let C be a subset of \mathbb{R}^n such that, if $\underline{x} \in C$, then H_0 is rejected (with probability 1), i.e.,

$$C = \{\underline{x} \in \mathbb{R}^n : H_0 \text{ is rejected for this } \underline{x}\}.$$

The set C is called the **critical region**. ■

Definition 9.1.5:

If we reject H_0 when it is true, we call this a **Type I error**. If we fail to reject H_0 when it is false, we call this a **Type II error**. Usually, H_0 and H_1 are chosen such that the Type I error is considered more serious. ■

Example 9.1.6:

We first consider a non–statistical example, in this case a jury trial. Our hypotheses are that the defendant is innocent or guilty. Our possible decisions are guilty or not guilty. Since it is considered worse to punish the innocent than to let the guilty go free, we make innocence the null hypothesis. Thus, we have

Truth (unknown)	Innocent (H_0)	Guilty (H_1)
Decision (known)		
Not Guilty (H_0)	Correct	Type II Error
Guilty (H_1)	Type I Error	Correct

The jury tries to make a decision “beyond a reasonable doubt”, i.e., it tries to make the probability of a Type I error small. ■

Definition 9.1.7:

If C is the critical region, then $P_\theta(C)$, $\theta \in \Theta_0$, is a probability of Type I error, and $P_\theta(C^c)$, $\theta \in \Theta_1$, is a probability of Type II error. ■

Note:

We would like both error probabilities to be 0, but this is usually not possible. We usually settle for fixing the probability of Type I error to be small, e.g., 0.05 or 0.01, and minimizing the Type II error. ■

Definition 9.1.8:

Every Borel–measurable mapping ϕ of $\mathbb{R}^n \rightarrow [0, 1]$ is called a **test function**. $\phi(\underline{x})$ is the probability of rejecting H_0 when \underline{x} is observed.

If ϕ is the indicator function of a subset $C \subseteq \mathbb{R}^n$, ϕ is called a **nonrandomized test** and C is the critical region of this test function.

Otherwise, if ϕ is not an indicator function of a subset $C \subseteq \mathbb{R}^n$, ϕ is called a **randomized test**. ■

Definition 9.1.9:

Let ϕ be a test function of the hypothesis $H_0 : \theta \in \Theta_0$ against the alternative $H_1 : \theta \in \Theta_1$. We say that ϕ has a **level of significance** of α (or ϕ is a **level- α -test** or ϕ is of **size α**) if

$$E_{\theta}(\phi(\underline{X})) = P_{\theta}(\text{reject } H_0) \leq \alpha \quad \forall \theta \in \Theta_0.$$

In short, we say that ϕ is a test for the problem $(\alpha, \Theta_0, \Theta_1)$. ■

Definition 9.1.10:

Let ϕ be a test for the problem $(\alpha, \Theta_0, \Theta_1)$. For every $\theta \in \Theta$, we define

$$\beta_{\phi}(\theta) = E_{\theta}(\phi(\underline{X})) = P_{\theta}(\text{reject } H_0).$$

We call $\beta_{\phi}(\theta)$ the **power function** of ϕ . For any $\theta \in \Theta_1$, $\beta_{\phi}(\theta)$ is called the **power** of ϕ against the alternative θ . ■

Definition 9.1.11:

Let Φ_{α} be the class of all tests for $(\alpha, \Theta_0, \Theta_1)$. A test $\phi_0 \in \Phi_{\alpha}$ is called a **most powerful (MP) test** against an alternative $\theta \in \Theta_1$ if

$$\beta_{\phi_0}(\theta) \geq \beta_{\phi}(\theta) \quad \forall \phi \in \Phi_{\alpha}.$$

■

Definition 9.1.12:

Let Φ_{α} be the class of all tests for $(\alpha, \Theta_0, \Theta_1)$. A test $\phi_0 \in \Phi_{\alpha}$ is called a **uniformly most powerful (UMP) test** if

$$\beta_{\phi_0}(\theta) \geq \beta_{\phi}(\theta) \quad \forall \phi \in \Phi_{\alpha} \quad \forall \theta \in \Theta_1.$$

■

Example 9.1.13:

Let X_1, \dots, X_n be iid $N(\mu, 1)$, $\mu \in \Theta = \{\mu_0, \mu_1\}$, $\mu_0 < \mu_1$.

Let $H_0 : X_i \sim N(\mu_0, 1)$ vs. $H_1 : X_i \sim N(\mu_1, 1)$.

Intuitively, reject H_0 when \bar{X} is too large, i.e., if $\bar{X} \geq k$ for some k .

Under H_0 it holds that $\bar{X} \sim N(\mu_0, \frac{1}{n})$. For a given α , we can solve the following equation for k :

$$P_{\mu_0}(\bar{X} > k) = P\left(\frac{\bar{X} - \mu_0}{1/\sqrt{n}} > \frac{k - \mu_0}{1/\sqrt{n}}\right) = P(Z > z_\alpha) = \alpha$$

Here, $\frac{\bar{X} - \mu_0}{1/\sqrt{n}} = Z \sim N(0, 1)$ and z_α is defined in such a way that $P(Z > z_\alpha) = \alpha$, i.e., z_α is the upper α -quantile of the $N(0, 1)$ distribution. It follows that $\frac{k - \mu_0}{1/\sqrt{n}} = z_\alpha$ and therefore, $k = \mu_0 + \frac{z_\alpha}{\sqrt{n}}$.

Thus, we obtain the nonrandomized test

■

Example 9.1.14:

Let $X \sim \text{Bin}(6, p)$, $p \in \Theta = (0, 1)$.

$H_0 : p = \frac{1}{2}$, $H_1 : p \neq \frac{1}{2}$.

Desired level of significance: $\alpha = 0.05$.

Reasonable plan: Since $E_{p=\frac{1}{2}}(X) = 3$, reject H_0 when $|X - 3| \geq c$ for some constant c . But how should we select c ?

x	$c = x - 3 $	$P_{p=\frac{1}{2}}(X = x)$	$P_{p=\frac{1}{2}}(X - 3 \geq c)$
0, 6			
1, 5			
2, 4			
3			

■

9.2 The Neyman–Pearson Lemma

(Based on Casella/Berger, Section 8.3.2)

Let $\{f_\theta : \theta \in \Theta = \{\theta_0, \theta_1\}\}$ be a family of possible distributions of \underline{X} . f_θ represents the pdf (or pmf) of \underline{X} . For convenience, we write $f_0(\underline{x}) = f_{\theta_0}(\underline{x})$ and $f_1(\underline{x}) = f_{\theta_1}(\underline{x})$.

Theorem 9.2.1: Neyman–Pearson Lemma (NP Lemma)

Suppose we wish to test $H_0 : \underline{X} \sim f_0(\underline{x})$ vs. $H_1 : \underline{X} \sim f_1(\underline{x})$, where f_i is the pdf (or pmf) of \underline{X} under H_i , $i = 0, 1$, where both, H_0 and H_1 , are simple.

(i) Any test of the form

$$\phi(\underline{x}) = \begin{cases} 1, & \text{if } f_1(\underline{x}) > k f_0(\underline{x}) \\ \gamma(\underline{x}), & \text{if } f_1(\underline{x}) = k f_0(\underline{x}) \\ 0, & \text{if } f_1(\underline{x}) < k f_0(\underline{x}) \end{cases} \quad (*)$$

for some $k \geq 0$ and $0 \leq \gamma(\underline{x}) \leq 1$, is most powerful of its significance level for testing H_0 vs. H_1 .

If $k = \infty$, the test

$$\phi(\underline{x}) = \begin{cases} 1, & \text{if } f_0(\underline{x}) = 0 \\ 0, & \text{if } f_0(\underline{x}) > 0 \end{cases} \quad (**)$$

is most powerful of size (or significance level) 0 for testing H_0 vs. H_1 .

(ii) Given $0 \leq \alpha \leq 1$, there exists a test of the form (*) or (**) with $\gamma(\underline{x}) = \gamma$ (i.e., a constant) such that

$$E_{\theta_0}(\phi(\underline{X})) = \alpha.$$

Proof:

We prove the continuous case only.

(i):

■

Theorem 9.2.2:

If a sufficient statistic T exists for the family $\{f_\theta : \theta \in \Theta = \{\theta_0, \theta_1\}\}$, then the Neyman–Pearson most powerful test is a function of T .

Proof:

Homework ■

Example 9.2.3:

We want to test $H_0 : X \sim N(0, 1)$ vs. $H_1 : X \sim \text{Cauchy}(1, 0)$, based on a single observation. It is

$$\frac{f_1(x)}{f_0(x)} = \frac{\frac{1}{\pi} \frac{1}{1+x^2}}{\frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})} = \sqrt{\frac{2}{\pi}} \frac{\exp(\frac{x^2}{2})}{1+x^2}.$$

The MP test is

$$\phi(x) = \begin{cases} 1, & \text{if } \sqrt{\frac{2}{\pi}} \frac{\exp(\frac{x^2}{2})}{1+x^2} > k \\ 0, & \text{otherwise} \end{cases}$$

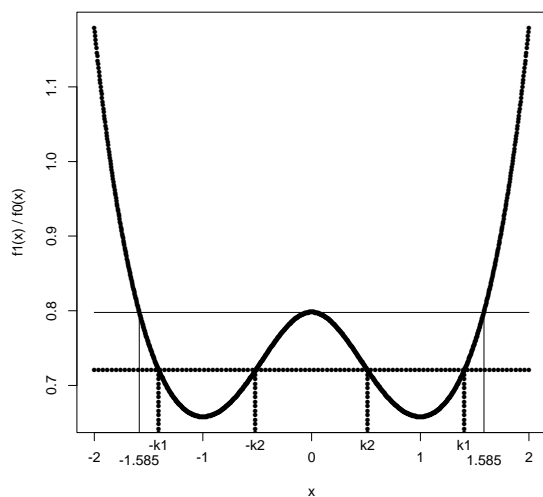
where k is determined such that $E_{H_0}(\phi(X)) = \alpha$.

If $\alpha < 0.113$, we reject H_0 if $|x| > z_{\frac{\alpha}{2}}$, where $z_{\frac{\alpha}{2}}$ is the upper $\frac{\alpha}{2}$ quantile of a $N(0, 1)$ distribution.

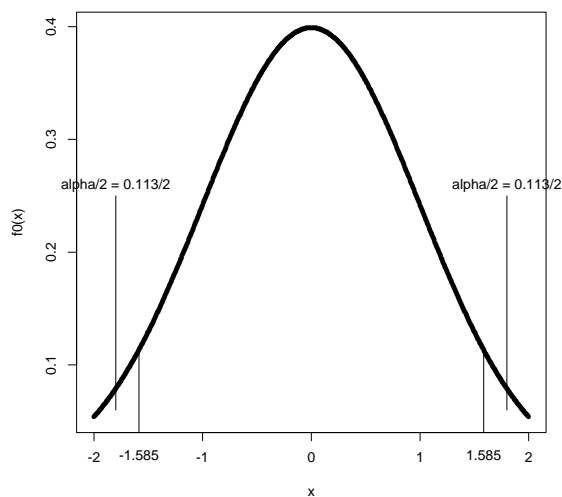
If $\alpha > 0.113$, we reject H_0 if $|x| > k_1$ or if $|x| < k_2$, where $k_1 > 0, k_2 > 0$, such that

$$\frac{\exp(\frac{k_1^2}{2})}{1+k_1^2} = \frac{\exp(\frac{k_2^2}{2})}{1+k_2^2} \quad \text{and} \quad \int_{k_2}^{k_1} \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2}) dx = \frac{1-\alpha}{2}.$$

Example 9.2.3a



Example 9.2.3b



Why is $\alpha = 0.113$ so interesting?

For $x = 0$, it is

$$\frac{f_1(x)}{f_0(x)} = \sqrt{\frac{2}{\pi}} \approx 0.7979.$$

Similarly, for $x \approx -1.585$ and $x \approx 1.585$, it is

$$\frac{f_1(x)}{f_0(x)} = \sqrt{\frac{2}{\pi}} \frac{\exp(\frac{(\pm 1.585)^2}{2})}{1 + (\pm 1.585)^2} \approx 0.7979 \approx \frac{f_1(0)}{f_0(0)}.$$

More importantly, $P_{H_0}(|X| > 1.585) = 0.113$.

■

9.3 Monotone Likelihood Ratios

(Based on Casella/Berger, Section 8.3.2)

Suppose we want to test $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$ for a family of pdf's $\{f_\theta : \theta \in \Theta \subseteq \mathbb{R}\}$. In general, it is not possible to find a UMP test. However, there exist conditions under which UMP tests exist.

Definition 9.3.1:

Let $\{f_\theta : \theta \in \Theta \subseteq \mathbb{R}\}$ be a family of pdf's (pmf's) on a one-dimensional parameter space. We say the family $\{f_\theta\}$ has a **monotone likelihood ratio (MLR)** in statistic $T(\underline{X})$ if for $\theta_1 < \theta_2$, whenever f_{θ_1} and f_{θ_2} are distinct, the ratio $\frac{f_{\theta_2}(\underline{x})}{f_{\theta_1}(\underline{x})}$ is a nondecreasing function of $T(\underline{x})$ for the set of values \underline{x} for which at least one of f_{θ_1} and f_{θ_2} is > 0 . ■

Note:

We can also define families of densities with nonincreasing MLR in $T(\underline{X})$, but such families can be treated by symmetry. ■

Example 9.3.2:

Let $X_1, \dots, X_n \sim U[0, \theta]$, $\theta > 0$. Then the joint pdf is

$$f_\theta(\underline{x}) = \begin{cases} \frac{1}{\theta^n}, & 0 \leq x_{(n)} \leq \theta \\ 0, & \text{otherwise} \end{cases} = \frac{1}{\theta^n} I_{[0, \theta]}(x_{(n)}),$$

where $x_{(n)} = x_{max} = \max_{i=1, \dots, n} x_i$.

Let $\theta_2 > \theta_1$, then

■

Theorem 9.3.3:

The one-parameter exponential family $f_\theta(\underline{x}) = \exp(Q(\theta)T(\underline{x}) + D(\theta) + S(\underline{x}))$, where $Q(\theta)$ is nondecreasing, has a MLR in $T(\underline{X})$.

Proof:

Homework. ■

Example 9.3.4:

Let $\underline{X} = (X_1, \dots, X_n)$ be a random sample from the Poisson family with parameter $\lambda > 0$. Then the joint pdf is

$$f_\lambda(\underline{x}) = \prod_{i=1}^n \left(e^{-\lambda} \lambda^{x_i} \frac{1}{x_i!} \right) = e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i} \prod_{i=1}^n \frac{1}{x_i!} = \exp \left(-n\lambda + \sum_{i=1}^n x_i \cdot \log(\lambda) - \sum_{i=1}^n \log(x_i!) \right),$$

which belongs to the one-parameter exponential family.

Since $Q(\lambda) = \log(\lambda)$ is a nondecreasing function of λ , it follows by Theorem 9.3.3 that the Poisson family with parameter $\lambda > 0$ has a MLR in $T(\underline{X}) = \sum_{i=1}^n X_i$.

We can verify this result by Definition 9.3.1:

$$\frac{f_{\lambda_2}(\underline{x})}{f_{\lambda_1}(\underline{x})} = \frac{\lambda_2^{\sum x_i} e^{-n\lambda_2}}{\lambda_1^{\sum x_i} e^{-n\lambda_1}} = \left(\frac{\lambda_2}{\lambda_1} \right)^{\sum x_i} e^{-n(\lambda_2 - \lambda_1)}.$$

If $\lambda_2 > \lambda_1$, then $\frac{\lambda_2}{\lambda_1} > 1$ and $\left(\frac{\lambda_2}{\lambda_1} \right)^{\sum x_i}$ is a nondecreasing function of $\sum x_i$.

Therefore, f_θ has a MLR in $T(\underline{X}) = \sum_{i=1}^n X_i$. ■

Theorem 9.3.5:

Let $X \sim f_\theta$, $\theta \in \Theta \subseteq \mathbb{R}$, where the family $\{f_\theta\}$ has a MLR in $T(\underline{X})$.

For testing $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$, $\theta_0 \in \Theta$, any test of the form

$$\phi(\underline{x}) = \begin{cases} 1, & \text{if } T(\underline{x}) > t_0 \\ \gamma, & \text{if } T(\underline{x}) = t_0 \\ 0, & \text{if } T(\underline{x}) < t_0 \end{cases} \quad (*)$$

has a nondecreasing power function and is UMP of its size $E_{\theta_0}(\phi(\underline{X})) = \alpha$, if the size is not 0.

Also, for every $0 \leq \alpha \leq 1$ and every $\theta_0 \in \Theta$, there exists a t_0 and a γ ($-\infty \leq t_0 \leq \infty$, $0 \leq \gamma \leq 1$), such that the test of form (*) is the UMP size α test of H_0 vs. H_1 .

Proof:

“ \implies ”:

“ \Leftarrow ”:

Use the Neyman–Pearson Lemma (Theorem 9.2.1). ■

Note:

By interchanging inequalities throughout Theorem 9.3.5 and its proof, we see that this Theorem also provides a solution of the dual problem $H'_0 : \theta \geq \theta_0$ vs. $H'_1 : \theta < \theta_0$. ■

Theorem: 9.3.6

For the one–parameter exponential family, there exists a UMP two–sided test of $H_0 : \theta \leq \theta_1$ or $\theta \geq \theta_2$, (where $\theta_1 < \theta_2$) vs. $H_1 : \theta_1 < \theta < \theta_2$ of the form

$$\phi(\underline{x}) = \begin{cases} 1, & \text{if } c_1 < T(\underline{x}) < c_2 \\ \gamma_i, & \text{if } T(\underline{x}) = c_i, \ i = 1, 2 \\ 0, & \text{if } T(\underline{x}) < c_1, \text{ or if } T(\underline{x}) > c_2 \end{cases}$$

■

Note:

UMP tests for $H_0 : \theta_1 \leq \theta \leq \theta_2$ and $H'_0 : \theta = \theta_0$ do not exist for one–parameter exponential families. ■

9.4 Unbiased and Invariant Tests

(Based on Rohatgi, Section 9.5, Rohatgi/Saleh, Section 9.5 & Casella/Berger, Section 8.3.2)

If we look at all size α tests in the class Φ_α , there exists no UMP test for many hypotheses. Can we find UMP tests if we reduce Φ_α by reasonable restrictions?

Definition 9.4.1:

A size α test ϕ of $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \in \Theta_1$ is **unbiased** if

$$E_\theta(\phi(\underline{X})) \geq \alpha \quad \forall \theta \in \Theta_1.$$

■

Note:

This condition means that $\beta_\phi(\theta) \leq \alpha \quad \forall \theta \in \Theta_0$ and $\beta_\phi(\theta) \geq \alpha \quad \forall \theta \in \Theta_1$. In other words, the power of this test is never less than α . ■

Definition 9.4.2:

Let U_α be the class of all unbiased size α tests of H_0 vs H_1 . If there exists a test $\phi \in U_\alpha$ that has maximal power for all $\theta \in \Theta_1$, we call ϕ a **UMP unbiased (UMPU)** size α test. ■

Note:

It holds that $U_\alpha \subseteq \Phi_\alpha$. A UMP test $\phi_\alpha \in \Phi_\alpha$ will have $\beta_{\phi_\alpha} \geq \alpha \quad \forall \theta \in \Theta_1$ since we must compare all tests ϕ_α with the trivial test $\phi(\underline{x}) = \alpha$. Thus, if a UMP test exists in Φ_α , it is also a UMPU test in U_α . ■

Example 9.4.3:

Let X_1, \dots, X_n be iid $N(\mu, \sigma^2)$, where $\sigma^2 > 0$ is known. Consider $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$. From the Neyman–Pearson Lemma, we know that for $\mu_1 > \mu_0$, the MP test is of the form

$$\phi_1(\underline{X}) = \begin{cases} 1, & \text{if } \bar{X} > \mu_0 + \frac{\sigma}{\sqrt{n}} z_\alpha \\ 0, & \text{otherwise} \end{cases}$$

and for $\mu_2 < \mu_0$, the MP test is of the form

$$\phi_2(\underline{X}) = \begin{cases} 1, & \text{if } \bar{X} < \mu_0 - \frac{\sigma}{\sqrt{n}} z_\alpha \\ 0, & \text{otherwise} \end{cases}$$

If a test is UMP, it must have the same rejection region as ϕ_1 and ϕ_2 . However, these 2 rejection regions are different (actually, their intersection is empty). Thus, there exists no

UMP test.

We next state a helpful Theorem and then continue with this example and see how we can find a UMPU test. ■

Theorem 9.4.4:

Let $c_1, \dots, c_n \in \mathbb{R}$ be constants and $f_1(\underline{x}), \dots, f_{n+1}(\underline{x})$ be real-valued functions. Let \mathcal{C} be the class of functions $\phi(\underline{x})$ satisfying $0 \leq \phi(\underline{x}) \leq 1$ and

$$\int_{-\infty}^{\infty} \phi(\underline{x}) f_i(\underline{x}) d\underline{x} = c_i \quad \forall i = 1, \dots, n.$$

If $\phi^* \in \mathcal{C}$ satisfies

$$\phi^*(\underline{x}) = \begin{cases} 1, & \text{if } f_{n+1}(\underline{x}) > \sum_{i=1}^n k_i f_i(\underline{x}) \\ 0, & \text{if } f_{n+1}(\underline{x}) < \sum_{i=1}^n k_i f_i(\underline{x}) \end{cases}$$

for some constants $k_1, \dots, k_n \in \mathbb{R}$, then ϕ^* maximizes $\int_{-\infty}^{\infty} \phi(\underline{x}) f_{n+1}(\underline{x}) d\underline{x}$ among all $\phi \in \mathcal{C}$.

Proof:

Let $\phi^*(\underline{x})$ be as above. Let $\phi(\underline{x})$ be any other function in \mathcal{C} . Since $0 \leq \phi(\underline{x}) \leq 1 \quad \forall \underline{x}$, it is

$$(\phi^*(\underline{x}) - \phi(\underline{x})) \left(f_{n+1}(\underline{x}) - \sum_{i=1}^n k_i f_i(\underline{x}) \right) \geq 0 \quad \forall \underline{x}.$$

This holds since if $\phi^*(\underline{x}) = 1$, the left factor is ≥ 0 and the right factor is ≥ 0 . If $\phi^*(\underline{x}) = 0$, the left factor is ≤ 0 and the right factor is ≤ 0 .

Therefore,

$$\begin{aligned} 0 &\leq \int (\phi^*(\underline{x}) - \phi(\underline{x})) \left(f_{n+1}(\underline{x}) - \sum_{i=1}^n k_i f_i(\underline{x}) \right) d\underline{x} \\ &= \int \phi^*(\underline{x}) f_{n+1}(\underline{x}) d\underline{x} - \int \phi(\underline{x}) f_{n+1}(\underline{x}) d\underline{x} - \sum_{i=1}^n k_i \underbrace{\left(\int \phi^*(\underline{x}) f_i(\underline{x}) d\underline{x} - \int \phi(\underline{x}) f_i(\underline{x}) d\underline{x} \right)}_{=c_i - c_i = 0} \end{aligned}$$

Thus,

$$\int \phi^*(\underline{x}) f_{n+1}(\underline{x}) d\underline{x} \geq \int \phi(\underline{x}) f_{n+1}(\underline{x}) d\underline{x}.$$

■

Note:

- (i) If f_{n+1} is a pdf, then ϕ^* maximizes the power.
- (ii) The Theorem above is the Neyman–Pearson Lemma if $n = 1$, $f_1 = f_{\theta_0}$, $f_2 = f_{\theta_1}$, and $c_1 = \alpha$.

■

Example 9.4.3: (continued)

So far, we have seen that there exists no UMP test for $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$.

We will show that

$$\phi_3(\underline{x}) = \begin{cases} 1, & \text{if } \bar{X} < \mu_0 - \frac{\sigma}{\sqrt{n}}z_{\alpha/2} \text{ or if } \bar{X} > \mu_0 + \frac{\sigma}{\sqrt{n}}z_{\alpha/2} \\ 0, & \text{otherwise} \end{cases}$$

is a UMPU size α test.

Due to Theorem 9.2.2, we only have to consider functions of sufficient statistics $T(\underline{X}) = \bar{X}$.

Let $\tau^2 = \frac{\sigma^2}{n}$.

To be unbiased and of size α , a test ϕ must have

$$(i) \int \phi(t)f_{\mu_0}(t)dt = \alpha, \text{ and}$$

$$(ii) \left. \frac{\partial}{\partial \mu} \int \phi(t)f_{\mu}(t)dt \right|_{\mu=\mu_0} = \int \phi(t) \left(\left. \frac{\partial}{\partial \mu} f_{\mu}(t) \right) \right|_{\mu=\mu_0} dt = 0, \text{ i.e., we have a minimum at } \mu_0.$$

We want to maximize $\int \phi(t)f_{\mu}(t)dt$, $\mu \neq \mu_0$ such that conditions (i) and (ii) hold.

We choose an arbitrary $\mu_1 \neq \mu_0$ and let

$$\begin{aligned} f_1(t) &= f_{\mu_0}(t) \\ f_2(t) &= \left. \frac{\partial}{\partial \mu} f_{\mu}(t) \right|_{\mu=\mu_0} \\ f_3(t) &= f_{\mu_1}(t) \end{aligned}$$

We now consider how the conditions on ϕ^* in Theorem 9.4.4 can be met:

$$\begin{aligned} f_3(t) &> k_1 f_1(t) + k_2 f_2(t) \\ \iff \frac{1}{\sqrt{2\pi\tau}} \exp\left(-\frac{1}{2\tau^2}(\bar{x} - \mu_1)^2\right) &> \frac{k_1}{\sqrt{2\pi\tau}} \exp\left(-\frac{1}{2\tau^2}(\bar{x} - \mu_0)^2\right) + \\ &\quad \frac{k_2}{\sqrt{2\pi\tau}} \exp\left(-\frac{1}{2\tau^2}(\bar{x} - \mu_0)^2\right) \left(\frac{\bar{x} - \mu_0}{\tau^2}\right) \\ \iff \exp\left(-\frac{1}{2\tau^2}(\bar{x} - \mu_1)^2\right) &> k_1 \exp\left(-\frac{1}{2\tau^2}(\bar{x} - \mu_0)^2\right) + k_2 \exp\left(-\frac{1}{2\tau^2}(\bar{x} - \mu_0)^2\right) \left(\frac{\bar{x} - \mu_0}{\tau^2}\right) \\ \iff \exp\left(\frac{1}{2\tau^2}((\bar{x} - \mu_0)^2 - (\bar{x} - \mu_1)^2)\right) &> k_1 + k_2 \left(\frac{\bar{x} - \mu_0}{\tau^2}\right) \\ \iff \exp\left(\frac{\bar{x}(\mu_1 - \mu_0)}{\tau^2} - \frac{\mu_1^2 - \mu_0^2}{2\tau^2}\right) &> k_1 + k_2 \left(\frac{\bar{x} - \mu_0}{\tau^2}\right) \end{aligned}$$

Note that the left hand side of this inequality is increasing in \bar{x} if $\mu_1 > \mu_0$ and decreasing in \bar{x} if $\mu_1 < \mu_0$. Either way, we can choose k_1 and k_2 such that the linear function in \bar{x} crosses the exponential function in \bar{x} at the two points

$$\mu_L = \mu_0 - \frac{\sigma}{\sqrt{n}}z_{\alpha/2}, \quad \mu_U = \mu_0 + \frac{\sigma}{\sqrt{n}}z_{\alpha/2}.$$

Obviously, ϕ_3 satisfies (i). We still need to check that ϕ_3 satisfies (ii) and that $\beta_{\phi_3}(\mu)$ has a minimum at μ_0 but omit this part from our proof here.

ϕ_3 is of the form ϕ^* in Theorem 9.4.4 and therefore ϕ_3 is UMP in \mathcal{C} . But the trivial test $\phi_t(\underline{x}) = \alpha$ also satisfies (i) and (ii) above. Therefore, $\beta_{\phi_3}(\mu) \geq \alpha \quad \forall \mu \neq \mu_0$. This means that ϕ_3 is unbiased.

Overall, ϕ_3 is a UMPU test of size α . ■

Definition 9.4.5:

A test ϕ is said to be α -**similar** on a subset Θ^* of Θ if

$$\beta_\phi(\theta) = E_\theta(\phi(\underline{X})) = \alpha \quad \forall \theta \in \Theta^*.$$

A test ϕ is said to be **similar** on $\Theta^* \subseteq \Theta$ if it is α -similar on Θ^* for some α , $0 \leq \alpha \leq 1$. ■

Note:

The trivial test $\phi(\underline{x}) = \alpha$ is α -similar on every $\Theta^* \subseteq \Theta$. ■

Theorem 9.4.6:

Let ϕ be an unbiased test of size α for $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \in \Theta_1$ such that $\beta_\phi(\theta)$ is a continuous function in θ . Then ϕ is α -similar on the boundary $\Lambda = \overline{\Theta_0} \cap \overline{\Theta_1}$, where $\overline{\Theta_0}$ and $\overline{\Theta_1}$ are the closures of Θ_0 and Θ_1 , respectively.

Proof:

Let $\theta \in \Lambda$. There exist sequences $\{\theta_n\}$ and $\{\theta'_n\}$ with $\theta_n \in \Theta_0$ and $\theta'_n \in \Theta_1$ such that $\lim_{n \rightarrow \infty} \theta_n = \theta$ and $\lim_{n \rightarrow \infty} \theta'_n = \theta$.

By continuity, $\beta_\phi(\theta_n) \rightarrow \beta_\phi(\theta)$ and $\beta_\phi(\theta'_n) \rightarrow \beta_\phi(\theta)$.

Since $\beta_\phi(\theta_n) \leq \alpha$ implies $\beta_\phi(\theta) \leq \alpha$ and since $\beta_\phi(\theta'_n) \geq \alpha$ implies $\beta_\phi(\theta) \geq \alpha$ it must hold that $\beta_\phi(\theta) = \alpha \quad \forall \theta \in \Lambda$. ■

Definition 9.4.7:

A test ϕ that is UMP among all α -similar tests on the boundary $\Lambda = \overline{\Theta_0} \cap \overline{\Theta_1}$ is called a **UMP α -similar test**. ■

Theorem 9.4.8:

Suppose $\beta_\phi(\theta)$ is continuous in θ for all tests ϕ of $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \in \Theta_1$. If a size α test of H_0 vs H_1 is UMP α -similar, then it is UMP unbiased.

Proof:

Let ϕ_0 be UMP α -similar and of size α . This means that $E_\theta(\phi(\underline{X})) \leq \alpha \quad \forall \theta \in \Theta_0$.

Since the trivial test $\phi(\underline{x}) = \alpha$ is α -similar, it must hold for ϕ_0 that $\beta_{\phi_0}(\theta) \geq \alpha \quad \forall \theta \in \Theta_1$ since ϕ_0 is UMP α -similar. This implies that ϕ_0 is unbiased.

Since $\beta_\phi(\theta)$ is continuous in θ , we see from Theorem 9.4.6 that the class of unbiased tests is a subclass of the class of α -similar tests. Since ϕ_0 is UMP in the larger class, it is also UMP in the subclass. Thus, ϕ_0 is UMPU. ■

Note:

The continuity of the power function $\beta_\phi(\theta)$ cannot always be checked easily. ■

Example 9.4.9:

Let $X_1, \dots, X_n \sim N(\mu, 1)$.

Let $H_0 : \mu \leq 0$ vs $H_1 : \mu > 0$.

Since the family of densities has a MLR in $\sum_{i=1}^n X_i$, we could use Theorem 9.3.5 to find a UMP test. However, we want to illustrate the use of Theorem 9.4.8 here.

It is $\Lambda = \{0\}$ and the power function

$$\beta_\phi(\mu) = \int_{\mathbb{R}^n} \phi(\underline{x}) \left(\frac{1}{\sqrt{2\pi}} \right)^n \exp \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 \right) d\underline{x}$$

of any test ϕ is continuous in μ . Thus, due to Theorem 9.4.6, any unbiased size α test of H_0 is α -similar on Λ .

We need a UMP test of $H'_0 : \mu = 0$ vs $H_1 : \mu > 0$.

By the NP Lemma, a MP test of $H''_0 : \mu = 0$ vs $H''_1 : \mu = \mu_1$, where $\mu_1 > 0$ is given by

$$\phi(\underline{x}) = \begin{cases} 1, & \text{if } \exp \left(\frac{\sum x_i^2}{2} - \frac{\sum (x_i - \mu)^2}{2} \right) > k' \\ 0, & \text{otherwise} \end{cases}$$

or equivalently, by Theorem 9.2.2,

$$\phi(\underline{x}) = \begin{cases} 1, & \text{if } T = \sum_{i=1}^n X_i > k \\ 0, & \text{otherwise} \end{cases}$$

Since under H_0 , $T \sim N(0, n)$, k is determined by $\alpha = P_{\mu=0}(T > k) = P\left(\frac{T}{\sqrt{n}} > \frac{k}{\sqrt{n}}\right)$, i.e., $k = \sqrt{n}z_\alpha$.

ϕ is independent of μ_1 for every $\mu_1 > 0$. So ϕ is UMP α -similar for H'_0 vs. H_1 .

Finally, ϕ is of size α , since for $\mu \leq 0$, it holds that

$$\begin{aligned} E_\mu(\phi(\underline{X})) &= P_\mu(T > \sqrt{n}z_\alpha) \\ &= P\left(\frac{T - n\mu}{\sqrt{n}} > z_\alpha - \sqrt{n}\mu\right) \\ &\stackrel{(*)}{\leq} P(Z > z_\alpha) \\ &= \alpha \end{aligned}$$

(*) holds since $\frac{T - n\mu}{\sqrt{n}} \sim N(0, 1)$ for $\mu \leq 0$ and $z_\alpha - \sqrt{n}\mu \geq z_\alpha$ for $\mu \leq 0$.

Thus all the requirements are met for Theorem 9.4.8, i.e., β_ϕ is continuous and ϕ is UMP α -similar and of size α , and thus ϕ is UMPU. ■

Note:

Rohatgi, page 428–430, lists Theorems (without proofs), stating that for Normal data, one- and two-tailed t -tests, one- and two-tailed χ^2 -tests, two-sample t -tests, and F -tests are all UMPU. ■

Note:

Recall from Definition 8.2.4 that a class of distributions is invariant under a group \mathcal{G} of transformations, if for each $g \in \mathcal{G}$ and for each $\theta \in \Theta$ there exists a unique $\theta' \in \Theta$ such that if $\underline{X} \sim P_\theta$, then $g(\underline{X}) \sim P_{\theta'}$. ■

Definition 9.4.10:

A group \mathcal{G} of transformations on \underline{X} leaves a hypothesis testing problem **invariant** if \mathcal{G} leaves both $\{P_\theta : \theta \in \Theta_0\}$ and $\{P_\theta : \theta \in \Theta_1\}$ invariant, i.e., if $\underline{y} = g(\underline{x}) \sim h_\theta(\underline{y})$, then $\{f_\theta(\underline{x}) : \theta \in \Theta_0\} \equiv \{h_\theta(\underline{y}) : \theta \in \Theta_0\}$ and $\{f_\theta(\underline{x}) : \theta \in \Theta_1\} \equiv \{h_\theta(\underline{y}) : \theta \in \Theta_1\}$. ■

Note:

We want two types of invariance for our tests:

Measurement Invariance: If $\underline{y} = g(\underline{x})$ is a 1-to-1 mapping, the decision based on \underline{y} should be the same as the decision based on \underline{x} . If $\phi(\underline{x})$ is the test based on \underline{x} and $\phi^*(\underline{y})$ is the test based on \underline{y} , then it must hold that $\phi(\underline{x}) = \phi^*(g(\underline{x})) = \phi^*(\underline{y})$.

Formal Invariance: If two tests have the same structure, i.e, the same Θ , the same pdf's (or pmf's), and the same hypotheses, then we should use the same test in both problems. So, if the transformed problem in terms of \underline{y} has the same formal structure as that of the problem in terms of \underline{x} , we must have that $\phi^*(\underline{y}) = \phi(\underline{x}) = \phi^*(g(\underline{x}))$.

We can combine these two requirements in the following definition: ■

Definition 9.4.11:

An **invariant test** with respect to a group \mathcal{G} of transformations is any test ϕ such that

$$\phi(\underline{x}) = \phi(g(\underline{x})) \quad \forall \underline{x} \quad \forall g \in \mathcal{G}.$$
■

Example 9.4.12:

Let $X \sim Bin(n, p)$. Let $H_0 : p = \frac{1}{2}$ vs. $H_1 : p \neq \frac{1}{2}$.

Let $\mathcal{G} = \{g_1, g_2\}$, where $g_1(x) = n - x$ and $g_2(x) = x$.

If ϕ is invariant, then $\phi(x) = \phi(n - x)$. Is the test problem invariant? For g_2 , the answer is obvious.

For g_1 , we get:

$$g_1(X) = n - X \sim Bin(n, 1 - p)$$

$$H_0 : p = \frac{1}{2} : \{f_p(x) : p = \frac{1}{2}\} = \{h_p(g_1(x)) : p = \frac{1}{2}\} = Bin(n, \frac{1}{2})$$

$$H_1 : p \neq \frac{1}{2} : \underbrace{\{f_p(x) : p \neq \frac{1}{2}\}}_{=Bin(n, p \neq \frac{1}{2})} = \underbrace{\{h_p(g_1(x)) : p \neq \frac{1}{2}\}}_{=Bin(n, p \neq \frac{1}{2})}$$

So all the requirements in Definition 9.4.10 are met. If, for example, $n = 10$, the test

$$\phi(x) = \begin{cases} 1, & \text{if } x = 0, 1, 2, 8, 9, 10 \\ 0, & \text{otherwise} \end{cases}$$

is invariant under \mathcal{G} . For example, $\phi(4) = 0 = \phi(10 - 4) = \phi(6)$, and, in general, $\phi(x) = \phi(10 - x) \quad \forall x \in \{0, 1, \dots, 9, 10\}$. ■

Example 9.4.13:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ where both μ and $\sigma^2 > 0$ are unknown. It is $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ and $\frac{n-1}{\sigma^2} S^2 \sim \chi_{n-1}^2$ and \bar{X} and S^2 independent.

Let $H_0 : \mu \leq 0$ vs. $H_1 : \mu > 0$.

Let \mathcal{G} be the group of scale changes:

$$\mathcal{G} = \{g_c(\bar{x}, s^2), c > 0 : g_c(\bar{x}, s^2) = (c\bar{x}, c^2 s^2)\}$$

The problem is invariant because, when $g_c(\bar{x}, s^2) = (c\bar{x}, c^2 s^2)$, then

- (i) $c\bar{X}$ and $c^2 S^2$ are independent.
- (ii) $c\bar{X} \sim N(c\mu, \frac{c^2 \sigma^2}{n}) \in \{N(\eta, \frac{\tau^2}{n})\}$.
- (iii) $\frac{n-1}{c^2 \sigma^2} c^2 S^2 \sim \chi_{n-1}^2$.

So, this is the same family of distributions and Definition 9.4.10 holds because $\mu \leq 0$ implies that $c\mu \leq 0$ (for $c > 0$).

An invariant test satisfies $\phi(\bar{x}, s^2) \equiv \phi(c\bar{x}, c^2 s^2)$, $c > 0, s^2 > 0, \bar{x} \in \mathbb{R}$.

Let $c = \frac{1}{s}$. Then $\phi(\bar{x}, s^2) \equiv \phi(\frac{\bar{x}}{s}, 1)$ so invariant tests depend on (\bar{x}, s^2) only through $\frac{\bar{x}}{s}$.

If $\frac{\bar{x}_1}{s_1} \neq \frac{\bar{x}_2}{s_2}$, then there exists no $c > 0$ such that $(\bar{x}_2, s_2^2) \equiv (c\bar{x}_1, c^2 s_1^2)$. So invariance places no restrictions on ϕ for different $\frac{\bar{x}_1}{s_1} = \frac{\bar{x}_2}{s_2}$. Thus, invariant tests are exactly those that depend only on $\frac{\bar{x}}{s}$, which are equivalent to tests that are based only on $t = \frac{\bar{x}}{s/\sqrt{n}}$. Since this mapping is 1-to-1, the invariant test will use $T = \frac{\bar{X}}{S/\sqrt{n}} \sim t_{n-1}$ if $\mu = 0$. Note that this test does not depend on the nuisance parameter σ^2 . Invariance often produces such results. ■

Definition 9.4.14:

Let \mathcal{G} be a group of transformations on the space of \underline{X} . We say a statistic $T(\underline{x})$ is **maximal invariant** under \mathcal{G} if

- (i) T is invariant, i.e., $T(\underline{x}) = T(g(\underline{x})) \quad \forall g \in \mathcal{G}$, and
- (ii) T is maximal, i.e., $T(\underline{x}_1) = T(\underline{x}_2)$ implies that $\underline{x}_1 = g(\underline{x}_2)$ for some $g \in \mathcal{G}$.

■

Example 9.4.15:

Let $\underline{x} = (x_1, \dots, x_n)$ and $g_c(\underline{x}) = (x_1 + c, \dots, x_n + c)$.

Consider $T(\underline{x}) = (x_n - x_1, x_n - x_2, \dots, x_n - x_{n-1})$.

It is $T(g_c(\underline{x})) = (x_n - x_1, x_n - x_2, \dots, x_n - x_{n-1}) = T(\underline{x})$, so T is invariant.

If $T(\underline{x}) = T(\underline{x}')$, then $x_n - x_i = x'_n - x'_i \quad \forall i = 1, 2, \dots, n-1$.

This implies that $x_i - x'_i = x_n - x'_n = c \quad \forall i = 1, 2, \dots, n-1$.

Thus, $g_c(\underline{x}') = (x'_1 + c, \dots, x'_n + c) = \underline{x}$.

Therefore, T is maximal invariant. ■

Definition 9.4.16:

Let I_α be the class of all invariant tests of size α of $H_0 : \theta \in \Theta_0$ vs. $H_1 : \theta \in \Theta_1$. If there exists a UMP member in I_α , it is called the **UMP invariant** test of H_0 vs H_1 . ■

Theorem 9.4.17:

Let $T(\underline{x})$ be maximal invariant with respect to \mathcal{G} . A test ϕ is invariant under \mathcal{G} iff ϕ is a function of T .

Proof:

“ \implies ”:

Let ϕ be invariant under \mathcal{G} . If $T(\underline{x}_1) = T(\underline{x}_2)$, then there exists a $g \in \mathcal{G}$ such that $\underline{x}_1 = g(\underline{x}_2)$. Thus, it follows from invariance that $\phi(\underline{x}_1) = \phi(g(\underline{x}_2)) = \phi(\underline{x}_2)$. Since ϕ is the same whenever $T(\underline{x}_1) = T(\underline{x}_2)$, ϕ must be a function of T .

“ \impliedby ”:

Let ϕ be a function of T , i.e., $\phi(\underline{x}) = h(T(\underline{x}))$. It follows that

$$\phi(g(\underline{x})) = h(T(g(\underline{x}))) \stackrel{(*)}{=} h(T(\underline{x})) = \phi(\underline{x}).$$

(*) holds since T is invariant.

This means that ϕ is invariant. ■

Example 9.4.18:

Consider the test problem

$$H_0 : \underline{X} \sim f_0(x_1 - \theta, \dots, x_n - \theta) \text{ vs. } H_1 : \underline{X} \sim f_1(x_1 - \theta, \dots, x_n - \theta),$$

where $\theta \in \mathbb{R}$.

Let \mathcal{G} be the group of transformations with

$$g_c(\underline{x}) = (x_1 + c, \dots, x_n + c),$$

where $c \in \mathbb{R}$ and $n \geq 2$.

As shown in Example 9.4.15, a maximal invariant statistic is $T(\underline{X}) = (X_1 - X_n, \dots, X_{n-1} - X_n) = (T_1, \dots, T_{n-1})$. Due to Theorem 9.4.17, an invariant test ϕ depends on \underline{X} only through T .

Since the transformation

$$\begin{pmatrix} T' \\ Z \end{pmatrix} = \begin{pmatrix} T_1 \\ \vdots \\ T_{n-1} \\ Z \end{pmatrix} = \begin{pmatrix} X_1 - X_n \\ \vdots \\ X_{n-1} - X_n \\ X_n \end{pmatrix}$$

is 1-to-1, there exists inverses $X_n = Z$ and $X_i = T_i + X_n = T_i + Z \quad \forall i = 1, \dots, n-1$. Applying Theorem 4.3.5 and integrating out the last component $Z (= X_n)$ gives us the joint pdf of $T = (T_1, \dots, T_{n-1})$.

Thus, under $H_i, i = 0, 1$, the joint pdf of T is given by $\int_{-\infty}^{\infty} f_i(t_1 + z, t_2 + z, \dots, t_{n-1} + z, z) dz$ which is independent of θ . The problem is thus reduced to testing a simple hypothesis against a simple alternative. By the NP Lemma (Theorem 9.2.1), the MP test is

$$\phi(t_1, \dots, t_{n-1}) = \begin{cases} 1, & \text{if } \lambda(\underline{t}) > c \\ 0, & \text{if } \lambda(\underline{t}) < c \end{cases}$$

where $\underline{t} = (t_1, \dots, t_{n-1})$ and $\lambda(\underline{t}) = \frac{\int_{-\infty}^{\infty} f_1(t_1 + z, t_2 + z, \dots, t_{n-1} + z, z) dz}{\int_{-\infty}^{\infty} f_0(t_1 + z, t_2 + z, \dots, t_{n-1} + z, z) dz}$.

In the homework assignment, we use this result to construct a UMP invariant test of

$$H_0 : \underline{X} \sim N(\theta, 1) \text{ vs. } H_1 : \underline{X} \sim \text{Cauchy}(1, \theta),$$

where a Cauchy(1, θ) distribution has pdf $f(x; \theta) = \frac{1}{\pi} \frac{1}{1 + (x - \theta)^2}$, where $\theta \in \mathbb{R}$. ■

10 More on Hypothesis Testing

10.1 Likelihood Ratio Tests

(Based on Casella/Berger, Section 8.2.1)

Definition 10.1.1:

The likelihood ratio test statistic for

$$H_0 : \theta \in \Theta_0 \text{ vs. } H_1 : \theta \in \Theta_1 = \Theta - \Theta_0$$

is

$$\lambda(\underline{x}) = \frac{\sup_{\theta \in \Theta_0} f_{\theta}(\underline{x})}{\sup_{\theta \in \Theta} f_{\theta}(\underline{x})}.$$

The likelihood ratio test (**LRT**) is the test function

$$\phi(\underline{x}) = I_{[0,c)}(\lambda(\underline{x})),$$

for some constant $c \in [0, 1]$, where c is usually chosen in such a way to make ϕ a test of size α . ■

Note:

- (i) We have to select c such that $0 \leq c \leq 1$ since $0 \leq \lambda(\underline{x}) \leq 1$.
- (ii) LRT's are strongly related to MLE's. If $\hat{\theta}$ is the unrestricted MLE of θ over Θ and $\hat{\theta}_0$ is the MLE of θ over Θ_0 , then $\lambda(\underline{x}) = \frac{f_{\hat{\theta}_0}(\underline{x})}{f_{\hat{\theta}}(\underline{x})}$. ■

Example 10.1.2:

Let X_1, \dots, X_n be a sample from $N(\mu, 1)$. We want to construct a LRT for

$$H_0 : \mu = \mu_0 \text{ vs. } H_1 : \mu \neq \mu_0.$$

It is $\hat{\mu}_0 = \mu_0$ and $\hat{\mu} = \bar{X}$. Thus,

$$\lambda(\underline{x}) = \frac{(2\pi)^{-n/2} \exp(-\frac{1}{2} \sum (x_i - \mu_0)^2)}{(2\pi)^{-n/2} \exp(-\frac{1}{2} \sum (x_i - \bar{x})^2)} = \exp(-\frac{n}{2} (\bar{x} - \mu_0)^2).$$

The LRT rejects H_0 if $\lambda(\underline{x}) \leq c$, or equivalently, $|\bar{x} - \mu_0| \geq \sqrt{-2 \frac{\log c}{n}}$. This means, the LRT rejects $H_0 : \mu = \mu_0$ if \bar{x} is *too far from* μ_0 . ■

Theorem 10.1.3:

If $T(\underline{X})$ is sufficient for θ and $\lambda^*(t)$ and $\lambda(\underline{x})$ are LRT statistics based on T and \underline{X} respectively, then

$$\lambda^*(T(\underline{x})) = \lambda(\underline{x}) \quad \forall \underline{x},$$

i.e., the LRT can be expressed as a function of every sufficient statistic for θ .

Proof:

Since T is sufficient, it follows from Theorem 8.3.5 that its pdf (or pdf) factorizes as $f_\theta(\underline{x}) = g_\theta(T)h(\underline{x})$. Therefore we get:

$$\begin{aligned} \lambda(\underline{x}) &= \frac{\sup_{\theta \in \Theta_0} f_\theta(\underline{x})}{\sup_{\theta \in \Theta} f_\theta(\underline{x})} \\ &= \frac{\sup_{\theta \in \Theta_0} g_\theta(T)h(\underline{x})}{\sup_{\theta \in \Theta} g_\theta(T)h(\underline{x})} \\ &= \frac{\sup_{\theta \in \Theta_0} g_\theta(T)}{\sup_{\theta \in \Theta} g_\theta(T)} \\ &= \lambda^*(T(\underline{x})) \end{aligned}$$

Thus, our simplified expression for $\lambda(\underline{x})$ indeed only depends on a sufficient statistic T . ■

Theorem 10.1.4:

If for a given α , $0 \leq \alpha \leq 1$, and for a simple hypothesis H_0 and a simple alternative H_1 a non-randomized test based on the NP Lemma and LRT's exist, then these tests are equivalent.

Proof:

See Homework. ■

Note:

Usually, LRT's perform well since they are often UMP or UMPU size α tests. However, this does not always hold. Rohatgi, Example 4, page 440–441, cites an example where the LRT is not unbiased and it is even worse than the trivial test $\phi(\underline{x}) = \alpha$. ■

Theorem 10.1.5:

Under some regularity conditions on $f_\theta(\underline{x})$, the rv $-2 \log \lambda(\underline{X})$ under H_0 has asymptotically a chi-squared distribution with ν degrees of freedom, where ν equals the difference between the number of independent parameters in Θ and Θ_0 , i.e.,

$$-2 \log \lambda(\underline{X}) \xrightarrow{d} \chi_\nu^2 \quad \text{under } H_0.$$

■

Note:

The regularity conditions required for Theorem 10.1.5 are basically the same as for Theorem 8.7.10. Under “independent” parameters we understand parameters that are unspecified, i.e., free to vary.

■

Example 10.1.6:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$ are both unknown.

Let $H_0: \mu = \mu_0$ vs. $H_1: \mu \neq \mu_0$.

We have $\theta = (\mu, \sigma^2)$, $\Theta = \{(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$ and $\Theta_0 = \{(\mu_0, \sigma^2) : \sigma^2 > 0\}$.

It is $\hat{\theta}_0 = (\mu_0, \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0)^2)$ and $\hat{\theta} = (\bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2)$.

Now, the LR test statistic $\lambda(\underline{x})$ can be determined:

Note that

$$f_{1,n-1}(f) = \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})\Gamma(\frac{1}{2})(n-1)^{\frac{1}{2}}} \frac{1}{\sqrt{f}} \left(1 + \frac{f}{n-1}\right)^{-\frac{n}{2}} \cdot I_{[0,\infty)}(f)$$

is the pdf of a $F_{1,n-1}$ distribution.

Let $y = \left(1 + \frac{f}{n-1}\right)^{-1}$, then $\frac{f}{n-1} = \frac{1-y}{y}$ and $df = -\frac{n-1}{y^2}dy$.

Thus,

As $n \rightarrow \infty$, we can apply Stirling's formula which states that

$$\Gamma(\alpha(n) + 1) \approx (\alpha(n))! \approx \sqrt{2\pi}(\alpha(n))^{\alpha(n)+\frac{1}{2}} \exp(-\alpha(n)).$$

So,

$$M_n(t) \approx \frac{\sqrt{2\pi}\left(\frac{n-2}{2}\right)^{\frac{n-1}{2}} \exp\left(-\frac{n-2}{2}\right) \sqrt{2\pi}\left(\frac{n(1-2t)-3}{2}\right)^{\frac{n(1-2t)-2}{2}} \exp\left(-\frac{n(1-2t)-3}{2}\right)}{\sqrt{2\pi}\left(\frac{n-3}{2}\right)^{\frac{n-2}{2}} \exp\left(-\frac{n-3}{2}\right) \sqrt{2\pi}\left(\frac{n(1-2t)-2}{2}\right)^{\frac{n(1-2t)-1}{2}} \exp\left(-\frac{n(1-2t)-2}{2}\right)}$$

10.2 Parametric Chi-Squared Tests

(Based on Rohatgi, Section 10.3 & Rohatgi/Saleh, Section 10.3)

Definition 10.2.1: Normal Variance Tests

Let X_1, \dots, X_n be a sample from a $N(\mu, \sigma^2)$ distribution where μ may be known or unknown and $\sigma^2 > 0$ is unknown. The following table summarizes the χ^2 tests that are typically being used:

			Reject H_0 at level α if	
	H_0	H_1	μ known	μ unknown
I	$\sigma \geq \sigma_0$	$\sigma < \sigma_0$	$\sum(x_i - \mu)^2 \leq \sigma_0^2 \chi_{n-1; 1-\alpha}^2$	$s^2 \leq \frac{\sigma_0^2}{n-1} \chi_{n-1; 1-\alpha}^2$
II	$\sigma \leq \sigma_0$	$\sigma > \sigma_0$	$\sum(x_i - \mu)^2 \geq \sigma_0^2 \chi_{n; \alpha}^2$	$s^2 \geq \frac{\sigma_0^2}{n-1} \chi_{n-1; \alpha}^2$
III	$\sigma = \sigma_0$	$\sigma \neq \sigma_0$	$\sum(x_i - \mu)^2 \leq \sigma_0^2 \chi_{n; 1-\alpha/2}^2$ or $\sum(x_i - \mu)^2 \geq \sigma_0^2 \chi_{n; \alpha/2}^2$	$s^2 \leq \frac{\sigma_0^2}{n-1} \chi_{n-1; 1-\alpha/2}^2$ or $s^2 \geq \frac{\sigma_0^2}{n-1} \chi_{n-1; \alpha/2}^2$

■

Note:

- (i) In Definition 10.2.1, σ_0 is any fixed positive constant.
- (ii) Tests I and II are UMPU if μ is unknown and UMP if μ is known.
- (iii) In test III, the constants have been chosen in such a way to give equal probability to each tail. This is the usual approach. However, this may result in a biased test.
- (iv) $\chi_{n; 1-\alpha}^2$ is the (lower) α quantile and $\chi_{n; \alpha}^2$ is the (upper) $1 - \alpha$ quantile, i.e., for $X \sim \chi_n^2$, it holds that $P(X \leq \chi_{n; 1-\alpha}^2) = \alpha$ and $P(X \leq \chi_{n; \alpha}^2) = 1 - \alpha$.
- (v) We can also use χ^2 tests to test for equality of binomial probabilities as shown in the next few Theorems.

■

Theorem 10.2.2:

Let X_1, \dots, X_k be independent rv's with $X_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, k$. Then it holds that

$$T = \sum_{i=1}^k \left(\frac{X_i - n_i p_i}{\sqrt{n_i p_i (1 - p_i)}} \right)^2 \xrightarrow{d} \chi_k^2$$

as $n_1, \dots, n_k \rightarrow \infty$.

Proof:

Homework ■

Corollary 10.2.3:

Let X_1, \dots, X_k be as in Theorem 10.2.2 above. We want to test the hypothesis that $H_0 : p_1 = p_2 = \dots = p_k = p$, where p is a known constant (vs. the alternative H_1 that at least one of the p_i 's is different from the other ones). An approximate level- α test rejects H_0 if

$$y = \sum_{i=1}^k \left(\frac{x_i - n_i p}{\sqrt{n_i p(1-p)}} \right)^2 \geq \chi_{k;\alpha}^2.$$
■

Theorem 10.2.4:

Let X_1, \dots, X_k be independent rv's with $X_i \sim \text{Bin}(n_i, p), i = 1, \dots, k$. Then the MLE of p is

$$\hat{p} = \frac{\sum_{i=1}^k x_i}{\sum_{i=1}^k n_i}.$$

Proof:

This can be shown by using the joint likelihood function or by the fact that $\sum X_i \sim \text{Bin}(\sum n_i, p)$ and for $X \sim \text{Bin}(n, p)$, the MLE is $\hat{p} = \frac{x}{n}$. ■

Theorem 10.2.5:

Let X_1, \dots, X_k be independent rv's with $X_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, k$. An approximate level- α test of $H_0 : p_1 = p_2 = \dots = p_k = p$, where p is unknown (vs. the alternative H_1 that at least one of the p_i 's is different from the other ones), rejects H_0 if

$$y = \sum_{i=1}^k \left(\frac{x_i - n_i \hat{p}}{\sqrt{n_i \hat{p}(1-\hat{p})}} \right)^2 \geq \chi_{k-1;\alpha}^2,$$

where $\hat{p} = \frac{\sum x_i}{\sum n_i}$. ■

Theorem 10.2.6:

Let (X_1, \dots, X_k) be a multinomial rv with parameters n, p_1, p_2, \dots, p_k where $\sum_{i=1}^k p_i = 1$ and

$\sum_{i=1}^k X_i = n$. Then it holds that

$$U_k = \sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i} \xrightarrow{d} \chi_{k-1}^2$$

as $n \rightarrow \infty$.

An approximate level- α test of $H_0 : p_1 = p'_1, p_2 = p'_2, \dots, p_k = p'_k$ rejects H_0 if

$$\sum_{i=1}^k \frac{(x_i - np'_i)^2}{np'_i} > \chi_{k-1; \alpha}^2.$$

Proof:

Case $k = 2$ only:

■

Theorem 10.2.7:

Let X_1, \dots, X_n be a sample from X . Let $H_0 : X \sim F$, where the functional form of F is known completely. We partition the real line into k disjoint Borel sets A_1, \dots, A_k and let $P(X \in A_i) = p_i$, where $p_i > 0 \quad \forall i = 1, \dots, k$.

Let $Y_j = \#X'_i \text{ s in } A_j = \sum_{i=1}^n I_{A_j}(X_i), \quad \forall j = 1, \dots, k$.

Then, (Y_1, \dots, Y_k) has multinomial distribution with parameters n, p_1, p_2, \dots, p_k . ■

Theorem 10.2.8:

Let X_1, \dots, X_n be a sample from X . Let $H_0 : X \sim F_{\underline{\theta}}$, where $\underline{\theta} = (\theta_1, \dots, \theta_r)$ is unknown. Let the MLE $\hat{\underline{\theta}}$ exist. We partition the real line into k disjoint Borel sets A_1, \dots, A_k and let $P_{\hat{\underline{\theta}}}(X \in A_i) = \hat{p}_i$, where $\hat{p}_i > 0 \quad \forall i = 1, \dots, k$.

Let $Y_j = \#X'_i \text{ s in } A_j = \sum_{i=1}^n I_{A_j}(X_i), \quad \forall j = 1, \dots, k$.

Then it holds that

$$V_k = \sum_{i=1}^k \frac{(Y_i - n\hat{p}_i)^2}{n\hat{p}_i} \xrightarrow{d} \chi_{k-r-1}^2.$$

An approximate level- α test of $H_0 : X \sim F_{\underline{\theta}}$ rejects H_0 if

$$\sum_{i=1}^k \frac{(y_i - n\hat{p}_i)^2}{n\hat{p}_i} > \chi_{k-r-1; \alpha}^2,$$

where r is the number of parameters in $\underline{\theta}$ that have to be estimated. ■

10.3 t -Tests and F -Tests

(Based on Rohatgi, Section 10.4 & 10.5 & Rohatgi/Saleh, Section 10.4 & 10.5)

Definition 10.3.1: One- and Two-Tailed t -Tests

Let X_1, \dots, X_n be a sample from a $N(\mu, \sigma^2)$ distribution where $\sigma^2 > 0$ may be known or unknown and μ is unknown. Let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$.

The following table summarizes the z - and t -tests that are typically being used:

			Reject H_0 at level α if	
	H_0	H_1	σ^2 known	σ^2 unknown
I	$\mu \leq \mu_0$	$\mu > \mu_0$	$\bar{x} \geq \mu_0 + \frac{\sigma}{\sqrt{n}} z_\alpha$	$\bar{x} \geq \mu_0 + \frac{s}{\sqrt{n}} t_{n-1; \alpha}$
II	$\mu \geq \mu_0$	$\mu < \mu_0$	$\bar{x} \leq \mu_0 + \frac{\sigma}{\sqrt{n}} z_{1-\alpha}$	$\bar{x} \leq \mu_0 + \frac{s}{\sqrt{n}} t_{n-1; 1-\alpha}$
III	$\mu = \mu_0$	$\mu \neq \mu_0$	$ \bar{x} - \mu_0 \geq \frac{\sigma}{\sqrt{n}} z_{\alpha/2}$	$ \bar{x} - \mu_0 \geq \frac{s}{\sqrt{n}} t_{n-1; \alpha/2}$

■

Note:

- (i) In Definition 10.3.1, μ_0 is any fixed constant.
- (ii) These tests are based on just one sample and are often called *one sample t -tests*.
- (iii) Tests I and II are UMP and test III is UMPU if σ^2 is known. Tests I, II, and III are UMPU and UMP invariant if σ^2 is unknown.
- (iv) For large n (≥ 30), we can use z -tables instead of t -tables. Also, for large n we can drop the Normality assumption due to the CLT. However, for small n , none of these simplifications is justified.

■

Definition 10.3.2: Two-Sample t -Tests

Let X_1, \dots, X_m be a sample from a $N(\mu_1, \sigma_1^2)$ distribution where $\sigma_1^2 > 0$ may be known or unknown and μ_1 is unknown. Let Y_1, \dots, Y_n be a sample from a $N(\mu_2, \sigma_2^2)$ distribution where $\sigma_2^2 > 0$ may be known or unknown and μ_2 is unknown.

$$\text{Let } \bar{X} = \frac{1}{m} \sum_{i=1}^m X_i \text{ and } S_1^2 = \frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})^2.$$

$$\text{Let } \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \text{ and } S_2^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

$$\text{Let } S_p^2 = \frac{(m-1)S_1^2 + (n-1)S_2^2}{m+n-2}.$$

The following table summarizes the z - and t -tests that are typically being used:

			Reject H_0 at level α if	
	H_0	H_1	σ_1^2, σ_2^2 known	σ_1^2, σ_2^2 unknown, $\sigma_1 = \sigma_2$
<i>I</i>	$\mu_1 - \mu_2 \leq \delta$	$\mu_1 - \mu_2 > \delta$	$\bar{x} - \bar{y} \geq \delta + z_\alpha \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}$	$\bar{x} - \bar{y} \geq \delta + t_{m+n-2; \alpha} s_p \sqrt{\frac{1}{m} + \frac{1}{n}}$
<i>II</i>	$\mu_1 - \mu_2 \geq \delta$	$\mu_1 - \mu_2 < \delta$	$\bar{x} - \bar{y} \leq \delta + z_{1-\alpha} \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}$	$\bar{x} - \bar{y} \leq \delta + t_{m+n-2; 1-\alpha} s_p \sqrt{\frac{1}{m} + \frac{1}{n}}$
<i>III</i>	$\mu_1 - \mu_2 = \delta$	$\mu_1 - \mu_2 \neq \delta$	$ \bar{x} - \bar{y} - \delta \geq z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}$	$ \bar{x} - \bar{y} - \delta \geq t_{m+n-2; \alpha/2} s_p \sqrt{\frac{1}{m} + \frac{1}{n}}$

■

Note:

- (i) In Definition 10.3.2, δ is any fixed constant.
- (ii) All tests are UMPU and UMP invariant.
- (iii) If $\sigma_1^2 = \sigma_2^2 = \sigma^2$ (which is unknown), then S_p^2 is an unbiased estimate of σ^2 . We should check that $\sigma_1^2 = \sigma_2^2$ with an F -test.
- (iv) For large $m + n$, we can use z -tables instead of t -tables. Also, for large m and large n we can drop the Normality assumption due to the CLT. However, for small m or small n , none of these simplifications is justified.

■

Definition 10.3.3: Paired t -Tests

Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be a sample from a bivariate $N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$ distribution where all 5 parameters are unknown.

Let $D_i = X_i - Y_i \sim N(\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2)$.

Let $\bar{D} = \frac{1}{n} \sum_{i=1}^n D_i$ and $S_d^2 = \frac{1}{n-1} \sum_{i=1}^n (D_i - \bar{D})^2$.

The following table summarizes the t -tests that are typically being used:

	H_0	H_1	Reject H_0 at level α if
<i>I</i>	$\mu_1 - \mu_2 \leq \delta$	$\mu_1 - \mu_2 > \delta$	$\bar{d} \geq \delta + \frac{s_d}{\sqrt{n}} t_{n-1; \alpha}$
<i>II</i>	$\mu_1 - \mu_2 \geq \delta$	$\mu_1 - \mu_2 < \delta$	$\bar{d} \leq \delta + \frac{s_d}{\sqrt{n}} t_{n-1; 1-\alpha}$
<i>III</i>	$\mu_1 - \mu_2 = \delta$	$\mu_1 - \mu_2 \neq \delta$	$ \bar{d} - \delta \geq \frac{s_d}{\sqrt{n}} t_{n-1; \alpha/2}$

■

Note:

- (i) In Definition 10.3.3, δ is any fixed constant.
- (ii) These tests are special cases of one-sample tests. All the properties stated in the Note following Definition 10.3.1 hold.
- (iii) We could do a test based on Normality assumptions if $\sigma^2 = \sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2$ were known, but that is a very unrealistic assumption.

■

Definition 10.3.4: F -Tests

Let X_1, \dots, X_m be a sample from a $N(\mu_1, \sigma_1^2)$ distribution where μ_1 may be known or unknown and σ_1^2 is unknown. Let Y_1, \dots, Y_n be a sample from a $N(\mu_2, \sigma_2^2)$ distribution where μ_2 may be known or unknown and σ_2^2 is unknown.

Recall that

$$\frac{\sum_{i=1}^m (X_i - \bar{X})^2}{\sigma_1^2} \sim \chi_{m-1}^2, \quad \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{\sigma_2^2} \sim \chi_{n-1}^2,$$

and

$$\frac{\frac{\sum_{i=1}^m (X_i - \bar{X})^2}{(m-1)\sigma_1^2}}{\frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{(n-1)\sigma_2^2}} = \frac{\sigma_2^2 S_1^2}{\sigma_1^2 S_2^2} \sim F_{m-1, n-1}.$$

The following table summarizes the F -tests that are typically being used:

	Reject H_0 at level α if			
	H_0	H_1	μ_1, μ_2 known	μ_1, μ_2 unknown
I	$\sigma_1^2 \leq \sigma_2^2$	$\sigma_1^2 > \sigma_2^2$	$\frac{\frac{1}{m} \sum (x_i - \mu_1)^2}{\frac{1}{n} \sum (y_i - \mu_2)^2} \geq F_{m, n; \alpha}$	$\frac{s_1^2}{s_2^2} \geq F_{m-1, n-1; \alpha}$
II	$\sigma_1^2 \geq \sigma_2^2$	$\sigma_1^2 < \sigma_2^2$	$\frac{\frac{1}{n} \sum (y_i - \mu_2)^2}{\frac{1}{m} \sum (x_i - \mu_1)^2} \geq F_{n, m; \alpha}$	$\frac{s_2^2}{s_1^2} \geq F_{n-1, m-1; \alpha}$
III	$\sigma_1^2 = \sigma_2^2$	$\sigma_1^2 \neq \sigma_2^2$	$\frac{\frac{1}{m} \sum (x_i - \mu_1)^2}{\frac{1}{n} \sum (y_i - \mu_2)^2} \geq F_{m, n; \alpha/2}$ or $\frac{\frac{1}{n} \sum (y_i - \mu_2)^2}{\frac{1}{m} \sum (x_i - \mu_1)^2} \geq F_{n, m; \alpha/2}$	$\frac{s_1^2}{s_2^2} \geq F_{m-1, n-1; \alpha/2}$ if $s_1^2 \geq s_2^2$ or $\frac{s_2^2}{s_1^2} \geq F_{n-1, m-1; \alpha/2}$ if $s_1^2 < s_2^2$

■

Note:

- (i) Tests I and II are UMPU and UMP invariant if μ_1 and μ_2 are unknown.
- (ii) Test III uses equal tails and therefore may not be unbiased.
- (iii) If an F -test (at level α_1) and a t -test (at level α_2) are both performed, the combined test has level $\alpha = 1 - (1 - \alpha_1)(1 - \alpha_2) \geq \max(\alpha_1, \alpha_2)$ ($\equiv \alpha_1 + \alpha_2$ if both are small).

■

10.4 Bayes and Minimax Tests

(Based on Rohatgi, Section 10.6 & Rohatgi/Saleh, Section 10.6)

Hypothesis testing may be conducted in a decision-theoretic framework. Here our action space \mathcal{A} consists of two options: $a_0 = \text{fail to reject } H_0$ and $a_1 = \text{reject } H_0$.

Usually, we assume no loss for a correct decision. Thus, our loss function looks like:

$$L(\theta, a_0) = \begin{cases} 0, & \text{if } \theta \in \Theta_0 \\ a(\theta), & \text{if } \theta \in \Theta_1 \end{cases}$$
$$L(\theta, a_1) = \begin{cases} b(\theta), & \text{if } \theta \in \Theta_0 \\ 0, & \text{if } \theta \in \Theta_1 \end{cases}$$

We consider the following special cases:

0–1 loss: $a(\theta) = b(\theta) = 1$, i.e., all errors are equally bad.

Generalized 0–1 loss: $a(\theta) = c_{II}$, $b(\theta) = c_I$, i.e., all Type I errors are equally bad and all Type II errors are equally bad and Type I errors are worse than Type II errors or vice versa.

Then, the risk function can be written as

$$R(\theta, d(\underline{X})) = L(\theta, a_0)P_\theta(d(\underline{X}) = a_0) + L(\theta, a_1)P_\theta(d(\underline{X}) = a_1)$$
$$= \begin{cases} a(\theta)P_\theta(d(\underline{X}) = a_0), & \text{if } \theta \in \Theta_1 \\ b(\theta)P_\theta(d(\underline{X}) = a_1), & \text{if } \theta \in \Theta_0 \end{cases}$$

The minimax rule minimizes

$$\max_{\theta} \{a(\theta)P_\theta(d(\underline{X}) = a_0), b(\theta)P_\theta(d(\underline{X}) = a_1)\}.$$

Theorem 10.4.1:

The minimax rule d for testing

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta = \theta_1$$

under the generalized 0–1 loss function rejects H_0 if

$$\frac{f_{\theta_1}(\underline{x})}{f_{\theta_0}(\underline{x})} \geq k,$$

where k is chosen such that

$$\begin{aligned} R(\theta_1, d(\underline{X})) &= R(\theta_0, d(\underline{X})) \\ \iff c_{II}P_{\theta_1}(d(\underline{X}) = a_0) &= c_I P_{\theta_0}(d(\underline{X}) = a_1) \\ \iff c_{II}P_{\theta_1}\left(\frac{f_{\theta_1}(\underline{X})}{f_{\theta_0}(\underline{X})} < k\right) &= c_I P_{\theta_0}\left(\frac{f_{\theta_1}(\underline{X})}{f_{\theta_0}(\underline{X})} \geq k\right). \end{aligned}$$

Proof:

Let d^* be any other rule.

- If $R(\theta_0, d) < R(\theta_0, d^*)$, then

- If $R(\theta_0, d) \geq R(\theta_0, d^*)$, then

■

Example 10.4.2:

Let X_1, \dots, X_n be iid $N(\mu, 1)$. Let $H_0 : \mu = \mu_0$ vs. $H_1 : \mu = \mu_1 > \mu_0$.

■

Note:

Now suppose we have a prior distribution $\pi(\theta)$ on Θ . Then the Bayes risk of a decision rule d (under the loss function introduced before) is

$$\begin{aligned} R(\pi, d) &= E_{\pi} R(\theta, d(\underline{X})) \\ &= \int_{\Theta} R(\theta, d) \pi(\theta) d\theta \\ &= \int_{\Theta_0} b(\theta) \pi(\theta) P_{\theta}(d(\underline{X}) = a_1) d\theta + \int_{\Theta_1} a(\theta) \pi(\theta) P_{\theta}(d(\underline{X}) = a_0) d\theta \end{aligned}$$

if π is a pdf.

The Bayes risk for a pmf π looks similar (see Rohatgi, page 461).

■

Theorem 10.4.3:

The Bayes rule for testing $H_0 : \theta = \theta_0$ vs. $H_1 : \theta = \theta_1$ under the prior $\pi(\theta_0) = \pi_0$ and $\pi(\theta_1) = \pi_1 = 1 - \pi_0$ and the generalized 0-1 loss function is to reject H_0 if

$$\frac{f_{\theta_1}(\underline{x})}{f_{\theta_0}(\underline{x})} \geq \frac{c_I \pi_0}{c_{II} \pi_1}.$$

Proof:

■

Note:

For minimax rules and Bayes rules, the significance level α is no longer predetermined.

■

Example 10.4.4:

Let X_1, \dots, X_n be iid $N(\mu, 1)$. Let $H_0 : \mu = \mu_0$ vs. $H_1 : \mu = \mu_1 > \mu_0$. Let $c_I = c_{II}$.

By Theorem 10.4.3, the Bayes rule d rejects H_0 if

■

Note:

We can generalize Theorem 10.4.3 to the case of classifying among k options $\theta_1, \dots, \theta_k$. If we use the 0–1 loss function

$$L(\theta_i, d) = \begin{cases} 1, & \text{if } d(\underline{X}) = \theta_j \quad \forall j \neq i \\ 0, & \text{if } d(\underline{X}) = \theta_i \end{cases},$$

then the Bayes rule is to pick θ_i if

$$\pi_i f_{\theta_i}(\underline{x}) \geq \pi_j f_{\theta_j}(\underline{x}) \quad \forall j \neq i.$$

■

Example 10.4.5:

Let X_1, \dots, X_n be iid $N(\mu, 1)$. Let $\mu_1 < \mu_2 < \mu_3$ and let $\pi_1 = \pi_2 = \pi_3$.

Choose $\mu = \mu_i$ if

$$\pi_i \exp\left(-\frac{\sum(x_k - \mu_i)^2}{2}\right) \geq \pi_j \exp\left(-\frac{\sum(x_k - \mu_j)^2}{2}\right), \quad j \neq i, j = 1, 2, 3.$$

Similar to Example 10.4.4, these conditions can be transformed as follows:

$$\bar{x}(\mu_i - \mu_j) \geq \frac{(\mu_i - \mu_j)(\mu_i + \mu_j)}{2}, \quad j \neq i, j = 1, 2, 3.$$

In our particular example, we get the following decision rules:

- (i) Choose μ_1 if $\bar{x} \leq \frac{\mu_1 + \mu_2}{2}$ (and $\bar{x} \leq \frac{\mu_1 + \mu_3}{2}$).
- (ii) Choose μ_2 if $\bar{x} \geq \frac{\mu_1 + \mu_2}{2}$ and $\bar{x} \leq \frac{\mu_2 + \mu_3}{2}$.
- (iii) Choose μ_3 if $\bar{x} \geq \frac{\mu_2 + \mu_3}{2}$ (and $\bar{x} \geq \frac{\mu_1 + \mu_3}{2}$).

Note that in (i) and (iii) the condition in parentheses automatically holds when the other condition holds.

If $\mu_1 = 0$, $\mu_2 = 2$, and $\mu_3 = 4$, we have the decision rules:

- (i) Choose μ_1 if $\bar{x} \leq 1$.
- (ii) Choose μ_2 if $1 \leq \bar{x} \leq 3$.
- (iii) Choose μ_3 if $\bar{x} \geq 3$.

We do not have to worry how to handle the boundary since the probability that the rv will realize on any of the two boundary points is 0. ■

11 Confidence Estimation

11.1 Fundamental Notions

(Based on Casella/Berger, Section 9.1 & 9.3.2)

Let X be a rv and a, b be fixed positive numbers, $a < b$. Then

$$P(a < X < b) =$$

The interval $I(X) = (\frac{aX}{b}, X)$ is an example of a **random interval**. $I(X)$ contains the value a with a certain fixed probability.

For example, if $X \sim U(0, 1)$, $a = \frac{1}{4}$, and $b = \frac{3}{4}$, then the interval $I(X) = (\frac{X}{3}, X)$ contains $\frac{1}{4}$ with probability $\frac{1}{2}$.

Definition 11.1.1:

Let $P_{\theta}, \theta \in \Theta \subseteq \mathbb{R}^k$, be a set of probability distributions of a rv \underline{X} . A family of subsets $S(\underline{x})$ of Θ , where $S(\underline{x})$ depends on \underline{x} but not on θ , is called a **family of random sets**. In particular, if $\theta \in \Theta \subseteq \mathbb{R}$ and $S(\underline{x})$ is an interval $(\underline{\theta}(\underline{x}), \bar{\theta}(\underline{x}))$ where $\underline{\theta}(\underline{x})$ and $\bar{\theta}(\underline{x})$ depend on \underline{x} but not on θ , we call $S(\underline{X})$ a **random interval**, with $\underline{\theta}(\underline{X})$ and $\bar{\theta}(\underline{X})$ as lower and upper bounds, respectively. $\underline{\theta}(\underline{X})$ may be $-\infty$ and $\bar{\theta}(\underline{X})$ may be $+\infty$. ■

Note:

Frequently in inference, we are not interested in estimating a parameter or testing a hypothesis about it. Instead, we are interested in establishing a lower or upper bound (or both) for one or multiple parameters. ■

Definition 11.1.2:

A family of subsets $S(\underline{x})$ of $\Theta \subseteq \mathbb{R}^k$ is called a **family of confidence sets at confidence level $1 - \alpha$** if

$$P_{\theta}(S(\underline{X}) \ni \theta) \geq 1 - \alpha \quad \forall \theta \in \Theta,$$

where $0 < \alpha < 1$ is usually small.

The quantity

$$\inf_{\theta} P_{\theta}(S(\underline{X}) \ni \theta) = 1 - \alpha$$

is called the **confidence coefficient** (i.e., the smallest probability of true coverage is $1 - \alpha$). ■

Definition 11.1.3:

For $k = 1$, we use the following names for some of the confidence sets defined in Definition 11.1.2:

- (i) If $S(\underline{x}) = (\underline{\theta}(\underline{x}), \infty)$, then $\underline{\theta}(\underline{x})$ is called a level $1 - \alpha$ **lower confidence bound**.
- (ii) If $S(\underline{x}) = (-\infty, \bar{\theta}(\underline{x}))$, then $\bar{\theta}(\underline{x})$ is called a level $1 - \alpha$ **upper confidence bound**.
- (iii) $S(\underline{x}) = (\underline{\theta}(\underline{x}), \bar{\theta}(\underline{x}))$ is called a level $1 - \alpha$ **confidence interval (CI)**.

■

Definition 11.1.4:

A family of $1 - \alpha$ level confidence sets $\{S(\underline{x})\}$ is called **uniformly most accurate (UMA)** if

$$P_{\underline{\theta}}(S(\underline{X}) \ni \underline{\theta}') \leq P_{\underline{\theta}'}(S(\underline{X}) \ni \underline{\theta}') \quad \forall \underline{\theta}, \underline{\theta}' \in \Theta, \quad \underline{\theta} \neq \underline{\theta}'$$

and for any $1 - \alpha$ level family of confidence sets $S'(\underline{X})$ (i.e., $S(\underline{x})$ minimizes the probability of false (or incorrect) coverage). ■

Theorem 11.1.5:

Let $X_1, \dots, X_n \sim F_{\theta}$, $\theta \in \Theta$, where Θ is an interval on \mathbb{R} . Let $T(\underline{X}, \theta)$ be a function on $\mathbb{R}^n \times \Theta$ such that for each θ , $T(\underline{X}, \theta)$ is a statistic, and as a function of θ , T is strictly monotone (either increasing or decreasing) in θ at every value of $\underline{x} \in \mathbb{R}^n$.

Let $\Lambda \subseteq \mathbb{R}$ be the range of T and let the equation $\lambda = T(\underline{x}, \theta)$ be solvable for θ for every $\lambda \in \Lambda$ and every $\underline{x} \in \mathbb{R}^n$.

If the distribution of $T(\underline{X}, \theta)$ is independent of θ , then we can construct a confidence interval for θ at any level.

Proof:

Choose α such that $0 < \alpha < 1$. Then we can choose $\lambda_1(\alpha) < \lambda_2(\alpha)$ (which may not necessarily be unique) such that

$$P_{\theta}(\lambda_1(\alpha) < T(\underline{X}, \theta) < \lambda_2(\alpha)) \geq 1 - \alpha \quad \forall \theta.$$

Since the distribution of $T(\underline{X}, \theta)$ is independent of θ , $\lambda_1(\alpha)$ and $\lambda_2(\alpha)$ also do not depend on θ .

If $T(\underline{X}, \theta)$ is increasing in θ , solve the equations $\lambda_1(\alpha) = T(\underline{X}, \theta)$ for $\underline{\theta}(\underline{X})$ and $\lambda_2(\alpha) = T(\underline{X}, \theta)$

for $\bar{\theta}(\underline{X})$.

If $T(\underline{X}, \theta)$ is decreasing in θ , solve the equations $\lambda_1(\alpha) = T(\underline{X}, \theta)$ for $\bar{\theta}(\underline{X})$ and $\lambda_2(\alpha) = T(\underline{X}, \theta)$ for $\underline{\theta}(\underline{X})$.

In either case, it holds that

$$P_{\theta}(\underline{\theta}(\underline{X}) < \theta < \bar{\theta}(\underline{X})) \geq 1 - \alpha \quad \forall \theta.$$

■

Note:

- (i) Solvability is guaranteed if T is continuous and strictly monotone as a function of θ .
- (ii) If T is not monotone, we can still use this Theorem to get confidence sets that may not be confidence intervals.

■

Example 11.1.6:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where μ and $\sigma^2 > 0$ are both unknown. We seek a $1 - \alpha$ level confidence interval for μ .

■

Example 11.1.7:

Let $X_1, \dots, X_n \sim U(0, \theta)$.

We know that $\hat{\theta} = \max(X_i) = \text{Max}_n$ is the MLE for θ and sufficient for θ .

The pdf of Max_n is given by

$$f_n(y) = \frac{ny^{n-1}}{\theta^n} I_{(0, \theta)}(y).$$

Then the rv $T_n = \frac{Max_n}{\theta}$ has the pdf

$$h_n(t) = nt^{n-1}I_{(0,1)}(t),$$

which is independent of θ . T_n is monotone and decreasing in θ .

We now have to find numbers $\lambda_1(\alpha)$ and $\lambda_2(\alpha)$ such that

■

11.2 Shortest–Length Confidence Intervals

(Based on Casella/Berger, Section 9.2.2 & 9.3.1)

In practice, we usually want not only an interval with coverage probability $1 - \alpha$ for θ , but if possible the shortest (most precise) such interval.

Definition 11.2.1:

A rv $T(\underline{X}, \theta)$ whose distribution is independent of θ is called a **pivot**. ■

Note:

The methods we will discuss here can provide the shortest interval based on a given pivot. They will not guarantee that there is no other pivot with a shorter minimal interval. ■

Example 11.2.2:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where σ^2 is known. The obvious pivot for μ is

$$T_\mu(\underline{X}) = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

Suppose that (a, b) is an interval such that $P(a < Z < b) = 1 - \alpha$, where $Z \sim N(0, 1)$.

A $1 - \alpha$ level CI based on this pivot is found by

$$1 - \alpha = P\left(a < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < b\right) = P\left(\bar{X} - b\frac{\sigma}{\sqrt{n}} < \mu < \bar{X} - a\frac{\sigma}{\sqrt{n}}\right).$$

The length of the interval is $L = (b - a)\frac{\sigma}{\sqrt{n}}$.

To minimize L , we must choose a and b such that $b - a$ is minimal while

$$\Phi(b) - \Phi(a) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx = 1 - \alpha,$$

where $\Phi(z) = P(Z \leq z)$.

To find a minimum, we can differentiate these expressions with respect to a . However, b is not a constant but is an implicit function of a . Formally, we could write $\frac{db(a)}{da}$. However, this is usually shortened to $\frac{db}{da}$.

Here we get

$$\frac{d}{da}(\Phi(b) - \Phi(a)) = \phi(b)\frac{db}{da} - \phi(a) = 0$$

and

$$\frac{dL}{da} = \frac{\sigma}{\sqrt{n}}\left(\frac{db}{da} - 1\right) = \frac{\sigma}{\sqrt{n}}\left(\frac{\phi(a)}{\phi(b)} - 1\right).$$

The minimum occurs when $\phi(a) = \phi(b)$ which happens when $a = b$ or $a = -b$. If we select $a = b$, then $\Phi(b) - \Phi(a) = \Phi(a) - \Phi(a) = 0 \neq 1 - \alpha$. Thus, we must have that $b = -a = z_{\alpha/2}$. Thus, the shortest CI based on T_μ is

$$\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right).$$

■

Definition 11.2.3:

A pdf $f(x)$ is **unimodal** iff there exists a x^* such that $f(x)$ is nondecreasing for $x \leq x^*$ and $f(x)$ is nonincreasing for $x \geq x^*$. ■

Theorem 11.2.4:

Let $f(x)$ be a unimodal pdf. If the interval $[a, b]$ satisfies

- (i) $\int_a^b f(x)dx = 1 - \alpha$
- (ii) $f(a) = f(b) > 0$, and
- (iii) $a \leq x^* \leq b$, where x^* is a mode of $f(x)$,

then the interval $[a, b]$ is the shortest of all intervals which satisfy condition (i).

Proof:

Let $[a', b']$ be any interval with $b' - a' < b - a$. We will show that this implies $\int_{a'}^{b'} f(x)dx < 1 - \alpha$, i.e., a contradiction.

We assume that $a' \leq a$. The case $a < a'$ is similar.

- Suppose that $b' \leq a$. Then $a' \leq b' \leq a \leq x^*$. It follows

$$\begin{aligned} \int_{a'}^{b'} f(x)dx &\leq f(b')(b' - a') && |x \leq b' \leq x^* \Rightarrow f(x) \leq f(b') \\ &\leq f(a)(b' - a') && |b' \leq a \leq x^* \Rightarrow f(b') < f(a) \\ &< f(a)(b - a) && |b' - a' < b - a \text{ and } f(a) > 0 \\ &\leq \int_a^b f(x)dx && |f(x) \geq f(a) \text{ for } a \leq x \leq b \\ &= 1 - \alpha && | \text{by (i)} \end{aligned}$$

- Suppose $b' > a$. We can immediately exclude that $b' > b$ since then $b' - a' > b - a$, i.e., $b' - a'$ wouldn't be of shorter length than $b - a$. Thus, we have to consider the case that $a' \leq a < b' < b$. It holds that

$$\int_{a'}^{b'} f(x)dx = \int_a^b f(x)dx + \int_{a'}^a f(x)dx - \int_{b'}^b f(x)dx$$

Note that $\int_{a'}^a f(x)dx \leq f(a)(a - a')$ and $\int_{b'}^b f(x)dx \geq f(b)(b - b')$. Therefore, we get

$$\begin{aligned} \int_{a'}^a f(x)dx - \int_{b'}^b f(x)dx &\leq f(a)(a - a') - f(b)(b - b') \\ &= f(a)((a - a') - (b - b')) \quad \text{since } f(a) = f(b) \\ &= f(a)((b' - a') - (b - a)) \\ &< 0 \end{aligned}$$

Thus,

$$\int_{a'}^{b'} f(x)dx < a - \alpha.$$

■

Note:

Example 11.2.2 is a special case of Theorem 11.2.4. However, Theorem 11.2.4 is not immediately applicable in the following example since the length of that interval is proportional to $\frac{1}{a} - \frac{1}{b}$ (and not to $b - a$). ■

Example 11.2.5:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where μ is known. The obvious pivot for σ^2 is

$$T_{\sigma^2}(\underline{X}) = \frac{\sum(X_i - \mu)^2}{\sigma^2} \sim \chi_n^2.$$

So

$$\begin{aligned} P\left(a < \frac{\sum(X_i - \mu)^2}{\sigma^2} < b\right) &= 1 - \alpha \\ \iff P\left(\frac{\sum(X_i - \mu)^2}{b} < \sigma^2 < \frac{\sum(X_i - \mu)^2}{a}\right) &= 1 - \alpha \end{aligned}$$

We wish to minimize

$$L = \left(\frac{1}{a} - \frac{1}{b}\right) \sum(X_i - \mu)^2$$

such that $\int_a^b f_n(t)dt = 1 - \alpha$, where $f_n(t)$ is the pdf of a χ_n^2 distribution.

We get

$$f_n(b) \frac{db}{da} - f_n(a) = 0$$

and

$$\frac{dL}{da} = \left(-\frac{1}{a^2} + \frac{1}{b^2} \frac{db}{da}\right) \sum(X_i - \mu)^2 = \left(-\frac{1}{a^2} + \frac{1}{b^2} \frac{f_n(a)}{f_n(b)}\right) \sum(X_i - \mu)^2.$$

We obtain a minimum if $a^2 f_n(a) = b^2 f_n(b)$.

Note that in practice equal tails $\chi_{n;\alpha/2}^2$ and $\chi_{n;1-\alpha/2}^2$ are used, which do not result in shortest-length CI's. The reason for this selection is simple: When these tests were developed, computers did not exist that could solve these equations numerically. People in general had to rely on tabulated values. Manually solving the equation above for each case obviously wasn't a feasible solution. ■

Example 11.2.6:

Let $X_1, \dots, X_n \sim U(0, \theta)$. Let $Max_n = \max X_i = X_{(n)}$. Since $T_n = \frac{Max_n}{\theta}$ has pdf $nt^{n-1}I_{(0,1)}(t)$ which does not depend on θ , T_n can be selected as a our pivot. The density of T_n is strictly increasing for $n \geq 2$, so we cannot find constants a and b as in Example 11.2.5.

If $P(a < T_n < b) = 1 - \alpha$, then $P(\frac{Max_n}{b} < \theta < \frac{Max_n}{a}) = 1 - \alpha$.

We wish to minimize

$$L = Max_n \left(\frac{1}{a} - \frac{1}{b} \right)$$

such that $\int_a^b nt^{n-1}dt = b^n - a^n = 1 - \alpha$.

We get

$$nb^{n-1} - na^{n-1} \frac{da}{db} = 0 \implies \frac{da}{db} = \frac{b^{n-1}}{a^{n-1}}$$

and

$$\frac{dL}{db} = Max_n \left(-\frac{1}{a^2} \frac{da}{db} + \frac{1}{b^2} \right) = Max_n \left(-\frac{b^{n-1}}{a^{n+1}} + \frac{1}{b^2} \right) = Max_n \left(\frac{a^{n+1} - b^{n+1}}{b^2 a^{n+1}} \right) < 0 \quad \text{for } 0 \leq a < b \leq 1.$$

Thus, L does not have a local minimum. It is minimized when $b = 1$, i.e., when b is as large as possible. The corresponding a is selected as $a = \alpha^{1/n}$.

The shortest $1 - \alpha$ level CI based on T_n is $(Max_n, \alpha^{-1/n} Max_n)$. ■

11.3 Confidence Intervals and Hypothesis Tests

(Based on Casella/Berger, Section 9.2)

Example 11.3.1:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where $\sigma^2 > 0$ is known. In Example 11.2.2 we have shown that the interval

$$\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$$

is a $1 - \alpha$ level CI for μ .

Suppose we define a test ϕ of $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$ that rejects H_0 iff μ_0 does not fall in this interval.

Conversely, if $\phi(\underline{x}, \mu_0)$ is a family of size α tests of $H_0 : \mu = \mu_0$, the set $\{\mu_0 \mid \phi(\underline{x}, \mu_0) \text{ fails to reject } H_0\}$ is a level $1 - \alpha$ confidence set for μ_0 . ■

Theorem 11.3.2:

Denote $H_0(\theta_0)$ for $H_0 : \theta = \theta_0$, and $H_1(\theta_0)$ for the alternative. Let $A(\theta_0), \theta_0 \in \Theta$, denote the acceptance region of a level- α test of $H_0(\theta_0)$. For each possible observation \underline{x} , define

$$S(\underline{x}) = \{\theta : \underline{x} \in A(\theta), \theta \in \Theta\}.$$

Then $S(\underline{x})$ is a family of $1 - \alpha$ level confidence sets for θ .

If, moreover, $A(\theta_0)$ is UMP for $(\alpha, H_0(\theta_0), H_1(\theta_0))$, then $S(\underline{x})$ minimizes $P_{\theta'}(S(\underline{X}) \ni \theta') \forall \theta' \in H_1(\theta')$ among all $1 - \alpha$ level families of confidence sets, i.e., $S(\underline{x})$ is UMA.

Proof:

■

Example 11.3.3:

Let X be a rv that belongs to a one-parameter exponential family with pdf $f_{\theta}(x) = \exp(Q(\theta)T(x) + S'(x) + D(\theta))$, where $Q(\theta)$ is non-decreasing. We consider a test $H_0 : \theta = \theta_0$ vs. $H_1 : \theta < \theta_0$. The acceptance region of a UMP size α test of H_0 has the form $A(\theta_0) = \{x : T(x) > c(\theta_0)\}$.

■

Example 11.3.4:

Let $X \sim Exp(\theta)$ with $f_{\theta}(x) = \frac{1}{\theta}e^{-\frac{x}{\theta}}I_{(0,\infty)}(x)$, which belongs to a one-parameter exponential family. Then $Q(\theta) = -\frac{1}{\theta}$ is non-decreasing and $T(x) = x$.

We want to test $H_0 : \theta = \theta_0$ vs. $H_1 : \theta < \theta_0$.

■

Note:

Just as we frequently restrict the class of tests (when UMP tests don't exist), we can make the same sorts of restrictions on CI's. ■

Definition 11.3.5:

A family $S(\underline{x})$ of confidence sets for parameter θ is said to be **unbiased** at level $1 - \alpha$ if

$$P_{\theta}(S(\underline{X}) \ni \theta) \geq 1 - \alpha \text{ and } P_{\theta}(S(\underline{X}) \ni \theta') \leq 1 - \alpha \quad \forall \theta, \theta' \in \Theta, \quad \theta \neq \theta'.$$

If $S(\underline{x})$ is unbiased and minimizes $P_{\theta}(S(\underline{X}) \ni \theta')$ among all unbiased CI's at level $1 - \alpha$, it is called **uniformly most accurate unbiased (UMAUI)**. ■

Theorem 11.3.6:

Let $A(\theta_0)$ be the acceptance region of a UMPU size α test of $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$ (for all θ_0). Then $S(\underline{x}) = \{\theta : \underline{x} \in A(\theta)\}$ is a UMAUI family of confidence sets at level $1 - \alpha$.

Proof:

■

Theorem 11.3.7:

Let Θ be an interval on \mathbb{R} and f_θ be the pdf of \underline{X} . Let $S(\underline{X})$ be a family of $1 - \alpha$ level CI's, where $S(\underline{X}) = (\underline{\theta}(\underline{X}), \bar{\theta}(\underline{X}))$, $\underline{\theta}$ and $\bar{\theta}$ increasing functions of \underline{X} , and $\bar{\theta}(\underline{X}) - \underline{\theta}(\underline{X})$ is a finite rv.

Then it holds that

$$E_\theta(\bar{\theta}(\underline{X}) - \underline{\theta}(\underline{X})) = \int (\bar{\theta}(\underline{x}) - \underline{\theta}(\underline{x})) f_\theta(\underline{x}) d\underline{x} = \int_{\theta' \neq \theta} P_\theta(S(\underline{X}) \ni \theta') d\theta' \quad \forall \theta \in \Theta.$$

Proof:

It holds that $\bar{\theta} - \underline{\theta} = \int_{\underline{\theta}}^{\bar{\theta}} d\theta'$. Thus, for all $\theta \in \Theta$,

$$\begin{aligned} E_\theta(\bar{\theta}(\underline{X}) - \underline{\theta}(\underline{X})) &= \int_{\mathbb{R}^n} (\bar{\theta}(\underline{x}) - \underline{\theta}(\underline{x})) f_\theta(\underline{x}) d\underline{x} \\ &= \int_{\mathbb{R}^n} \left(\int_{\underline{\theta}(\underline{x})}^{\bar{\theta}(\underline{x})} d\theta' \right) f_\theta(\underline{x}) d\underline{x} \\ &= \int_{\mathbb{R}} \left(\int_{\substack{\underline{\theta}^{-1}(\theta') \\ \in \mathbb{R}^n}}^{\substack{\bar{\theta}^{-1}(\theta') \\ \in \mathbb{R}^n}} f_\theta(\underline{x}) d\underline{x} \right) d\theta' \\ &= \int_{\mathbb{R}} P_\theta(\underline{X} \in [\underline{\theta}^{-1}(\theta'), \bar{\theta}^{-1}(\theta')]) d\theta' \\ &= \int_{\mathbb{R}} P_\theta(S(\underline{X}) \ni \theta') d\theta' \\ &= \int_{\theta' \neq \theta} P_\theta(S(\underline{X}) \ni \theta') d\theta' \end{aligned} \quad \blacksquare$$

Note:

Theorem 11.3.7 says that the expected length of the CI is the probability that $S(\underline{X})$ includes the false θ' , averaged over all false values of θ' . \blacksquare

Corollary 11.3.8:

If $S(\underline{X})$ is UMAU, then $E_\theta(\bar{\theta}(\underline{X}) - \underline{\theta}(\underline{X}))$ is minimized among all unbiased families of CI's.

Proof:

In Theorem 11.3.7 we have shown that

$$E_\theta(\bar{\theta}(\underline{X}) - \underline{\theta}(\underline{X})) = \int_{\theta' \neq \theta} P_\theta(S(\underline{X}) \ni \theta') d\theta'.$$

Since a UMAU CI minimizes this probability for all θ' , the entire integral is minimized. \blacksquare

Example 11.3.9:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where $\sigma^2 > 0$ is known.

By Example 11.2.2, $(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$ is the shortest $1 - \alpha$ level CI for μ .

By Example 9.4.3, the equivalent test is UMPU. So by Theorem 11.3.6 this interval is UMAU and by Corollary 11.3.8 it has shortest expected length as well. ■

Example 11.3.10:

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where μ and $\sigma^2 > 0$ are both unknown.

Note that

$$T(\underline{X}, \sigma^2) = \frac{(n-1)S^2}{\sigma^2} = T_\sigma \sim \chi_{n-1}^2.$$

Thus,

Rohatgi, Theorem 4(b), page 428–429, states that the related test is UMPU. Therefore, by Theorem 11.3.6 and Corollary 11.3.8, our CI is UMAU with shortest expected length among all unbiased intervals.

Note that this CI is different from the equal-tail CI based on Definition 10.2.1, III, and from the shortest-length CI obtained in Example 11.2.5. ■

11.4 Bayes Confidence Intervals

(Based on Casella/Berger, Section 9.2.4)

Definition 11.4.1:

Given a posterior distribution $h(\theta | \underline{x})$, a level $1 - \alpha$ **credible set** (*Bayesian confidence set*) is any set A such that

$$P(\theta \in A | \underline{x}) = \int_A h(\theta | \underline{x}) d\theta = 1 - \alpha.$$

■

Example 11.4.2:

Let $X \sim \text{Bin}(n, p)$ and $\pi(p) \sim U(0, 1)$.

In Example 8.8.11, we have shown that

$$\begin{aligned} h(p | x) &= \frac{p^x(1-p)^{n-x}}{\int_0^1 p^x(1-p)^{n-x} dp} I_{(0,1)}(p) \\ &= B(x+1, n-x+1)^{-1} p^x(1-p)^{n-x} I_{(0,1)}(p) \\ &= \frac{\Gamma(n+2)}{\Gamma(x+1)\Gamma(n-x+1)} p^x(1-p)^{n-x} I_{(0,1)}(p) \\ \Rightarrow p | x &\sim \text{Beta}(x+1, n-x+1), \end{aligned}$$

where $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ is the beta function evaluated for a and b and $\text{Beta}(x+1, n-x+1)$ represents a Beta distribution with parameters $x+1$ and $n-x+1$.

Using the observed value for x and tables for incomplete beta integrals or a numerical approach, we can find λ_1 and λ_2 such that $P_{p|x}(\lambda_1 < p < \lambda_2) = 1 - \alpha$. So (λ_1, λ_2) is a credible interval for p . ■

Note:

- (i) The definitions and interpretations of credible intervals and confidence intervals are quite different. Therefore, very different intervals may result.
- (ii) We can often use Theorem 11.2.4 to find the shortest credible interval (if the preconditions hold).

■

Example 11.4.3:

Let X_1, \dots, X_n be iid $N(\mu, 1)$ and $\pi(\mu) \sim N(0, 1)$. We want to construct a Bayesian level $1 - \alpha$ CI for μ .

By Definition 8.8.7, the posterior distribution of μ given \underline{x} is

$$h(\mu | \underline{x}) = \frac{\pi(\mu)f(\underline{x} | \mu)}{g(\underline{x})}$$

where

$$g(\underline{x}) =$$

■

12 Nonparametric Inference

12.1 Nonparametric Estimation

Definition 12.1.1:

A statistical method which does not rely on assumptions about the distributional form of a rv (except, perhaps, that it is absolutely continuous, or purely discrete) is called a **nonparametric** or **distribution-free** method. ■

Note:

Unless otherwise specified, we make the following assumptions for the remainder of this chapter: Let X_1, \dots, X_n be iid $\sim F$, where F is unknown. Let \mathcal{P} be the class of all possible distributions of X . ■

Definition 12.1.2:

A statistic $T(\underline{X})$ is **sufficient** for a family of distributions \mathcal{P} if the conditional distribution of \underline{X} given $T = t$ is the same for all $F \in \mathcal{P}$. ■

Example 12.1.3:

Let X_1, \dots, X_n be absolutely continuous. Let $\underline{T} = (X_{(1)}, \dots, X_{(n)})$ be the order statistics.

It holds that

$$f(\underline{x} \mid \underline{T} = \underline{t}) = \frac{1}{n!},$$

so \underline{T} is sufficient for the family of absolutely continuous distributions on \mathbb{R} . ■

Definition 12.1.4:

A family of distributions \mathcal{P} is **complete** if the only unbiased estimate of 0 is the 0 itself, i.e.,

$$E_F(h(\underline{X})) = 0 \quad \forall F \in \mathcal{P} \quad \implies \quad h(\underline{x}) = 0 \quad \forall \underline{x}.$$

Definition 12.1.5:

A statistic $T(\underline{X})$ is **complete in relation to** \mathcal{P} if the class of induced distributions of T is complete. ■

Theorem 12.1.6:

The order statistic $(X_{(1)}, \dots, X_{(n)})$ is a complete sufficient statistic, provided that X_1, \dots, X_n are of either (pure) discrete or (pure) continuous type. ■

Definition 12.1.7:

A parameter $g(F)$ is called **estimable** if it has an unbiased estimate, i.e., if there exists a $T(\underline{X})$ such that

$$E_F(T(\underline{X})) = g(F) \quad \forall F \in \mathcal{P}.$$

■

Example 12.1.8:

Let \mathcal{P} be the class of distributions for which second moments exist. Then \bar{X} is unbiased for $\mu(F) = \int x dF(x)$. Thus, $\mu(F)$ is estimable.

■

Definition 12.1.9:

The **degree** m of an estimable parameter $g(F)$ is the smallest sample size for which an unbiased estimate exists for all $F \in \mathcal{P}$.

An unbiased estimate based on a sample of size m is called a **kernel**.

■

Lemma 12.1.10:

There exists a **symmetric kernel** for every estimable parameter.

Proof:

Let $T(X_1, \dots, X_m)$ be a kernel of $g(F)$. Define

$$T_s(X_1, \dots, X_m) = \frac{1}{m!} \sum_{\text{all permutations of } \{1, \dots, m\}} T(X_{i_1}, \dots, X_{i_m}).$$

where the summation is over all $m!$ permutations of $\{1, \dots, m\}$.

Clearly T_s is symmetric and $E(T_s) = g(F)$.

■

Example 12.1.11:

(i) $E(X_1) = \mu(F)$, so $\mu(F)$ has degree 1 with kernel X_1 .

(ii) $E(I_{(c, \infty)}(X_1)) = P_F(X > c)$, where c is a known constant. So $g(F) = P_F(X > c)$ has degree 1 with kernel $I_{(c, \infty)}(X_1)$.

(iii) There exists no $T(X_1)$ such that $E(T(X_1)) = \sigma^2(F) = \int (x - \mu(F))^2 dF(x)$.

But $E(T(X_1, X_2)) = E(X_1^2 - X_1 X_2) = \sigma^2(F)$. So $\sigma^2(F)$ has degree 2 with kernel $X_1^2 - X_1 X_2$. Note that $X_2^2 - X_2 X_1$ is another kernel.

(iv) A symmetric kernel for $\sigma^2(F)$ is

$$T_s(X_1, X_2) = \frac{1}{2}((X_1^2 - X_1 X_2) + (X_2^2 - X_1 X_2)) = \frac{1}{2}(X_1 - X_2)^2.$$

■

Definition 12.1.12:

Let $g(F)$ be an estimable parameter of degree m . Let X_1, \dots, X_n be a sample of size $n, n \geq m$. Given a kernel $T(X_{i_1}, \dots, X_{i_m})$ of $g(F)$, we define a U -statistic by

$$U(X_1, \dots, X_n) = \frac{1}{\binom{n}{m}} \sum_c T_s(X_{i_1}, \dots, X_{i_m}),$$

where T_s is defined as in Lemma 12.1.10 and the summation c is over all $\binom{n}{m}$ combinations of m integers (i_1, \dots, i_m) from $\{1, \dots, n\}$. $U(X_1, \dots, X_n)$ is symmetric in the X_i 's and $E_F(U(\underline{X})) = g(F)$ for all F . ■

Example 12.1.13:

For estimating $\mu(F)$ with degree m of $\mu(F) = 1$:

Symmetric kernel:

$$T_s(X_i) = X_i, \quad i = 1, \dots, n$$

U-statistic:

$$\begin{aligned} U_\mu(\underline{X}) &= \frac{1}{\binom{n}{1}} \sum_c X_i \\ &= \frac{1 \cdot (n-1)!}{n!} \sum_c X_i \\ &= \frac{1}{n} \sum_{i=1}^n X_i \\ &= \bar{X} \end{aligned}$$

For estimating $\sigma^2(F)$ with degree m of $\sigma^2(F) = 2$:

Symmetric kernel:

$$T_s(X_{i_1}, X_{i_2}) = \frac{1}{2}(X_{i_1} - X_{i_2})^2, \quad i_1, i_2 = 1, \dots, n, i_1 \neq i_2$$

U-statistic:

$$\begin{aligned} U_{\sigma^2}(\underline{X}) &= \frac{1}{\binom{n}{2}} \sum_{i_1 < i_2} \frac{1}{2}(X_{i_1} - X_{i_2})^2 \\ &= \frac{1}{\binom{n}{2}} \frac{1}{4} \sum_{i_1 \neq i_2} (X_{i_1} - X_{i_2})^2 \\ &= \frac{(n-2)! \cdot 2!}{n!} \frac{1}{4} \sum_{i_1 \neq i_2} (X_{i_1} - X_{i_2})^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2n(n-1)} \sum_{i_1} \sum_{i_2 \neq i_1} (X_{i_1}^2 - 2X_{i_1}X_{i_2} + X_{i_2}^2) \\
&= \frac{1}{2n(n-1)} \left[(n-1) \sum_{i_1=1}^n X_{i_1}^2 - 2 \left(\sum_{i_1=1}^n X_{i_1} \right) \left(\sum_{i_2=1}^n X_{i_2} \right) + 2 \sum_{i=1}^n X_i^2 + \right. \\
&\quad \left. (n-1) \sum_{i_2=1}^n X_{i_2}^2 \right] \\
&= \frac{1}{2n(n-1)} \left[n \sum_{i_1=1}^n X_{i_1}^2 - \sum_{i_1=1}^n X_{i_1}^2 - 2 \left(\sum_{i_1=1}^n X_{i_1} \right)^2 + 2 \sum_{i=1}^n X_i^2 + \right. \\
&\quad \left. n \sum_{i_2=1}^n X_{i_2}^2 - \sum_{i_2=1}^n X_{i_2}^2 \right] \\
&= \frac{1}{n(n-1)} \left[n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2 \right] \\
&= \frac{1}{n(n-1)} \left[n \sum_{i=1}^n \left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right)^2 \right] \\
&= \frac{1}{(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2 \\
&= S^2
\end{aligned}$$

■

Theorem 12.1.14:

Let \mathcal{P} be the class of all absolutely continuous or all purely discrete distribution functions on \mathbb{R} . Any estimable function $g(F)$, $F \in \mathcal{P}$, has a unique estimate that is unbiased and symmetric in the observations and has uniformly minimum variance among all unbiased estimates.

Proof:

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F \in \mathcal{P}$, with $T(X_1, \dots, X_n)$ an unbiased estimate of $g(F)$.

We define

$$T_i = T_i(X_1, \dots, X_n) = T(X_{i_1}, X_{i_2}, \dots, X_{i_n}), \quad i = 1, 2, \dots, n!,$$

over all possible permutations of $\{1, \dots, n\}$.

$$\text{Let } \bar{T} = \frac{1}{n!} \sum_{i=1}^{n!} T_i \text{ and } T = \sum_{i=1}^{n!} T_i.$$

Then

$$E_F(\bar{T}) = g(F)$$

and

$$\begin{aligned} \text{Var}(\bar{T}) &= E(\bar{T}^2) - (E(\bar{T}))^2 \\ &= E\left[\left(\frac{1}{n!} \sum_{i=1}^{n!} T_i\right)^2\right] - [g(F)]^2 \\ &= E\left[\left(\frac{1}{n!}\right)^2 \sum_{i=1}^{n!} \sum_{j=1}^{n!} T_i T_j\right] - [g(F)]^2 \\ &\leq E\left[\sum_{i=1}^{n!} \sum_{j=1}^{n!} T_i T_j\right] - [g(F)]^2 \\ &= E\left[\left(\sum_{i=1}^{n!} T_i\right) \left(\sum_{j=1}^{n!} T_j\right)\right] - [g(F)]^2 \\ &= E\left[\left(\sum_{i=1}^{n!} T_i\right)^2\right] - [g(F)]^2 \\ &= E(T^2) - [g(F)]^2 \\ &= \text{Var}(T) \end{aligned}$$

Equality holds iff $T_i = T_j \quad \forall i, j = 1, \dots, n!$

$\implies T$ is symmetric in (X_1, \dots, X_n) and $\bar{T} = T$

\implies by Rohatgi, Problem 4, page 538, T is a function of order statistics

\implies by Rohatgi, Theorem 1, page 535, T is a complete sufficient statistic

\implies by Note (i) following Theorem 8.4.12, T is UMVUE ■

Corollary 12.1.15:

If $T(X_1, \dots, X_n)$ is unbiased for $g(F)$, $F \in \mathcal{P}$, the corresponding U -statistic is an essentially unique UMVUE. ■

Definition 12.1.16:

Suppose we have independent samples $X_1, \dots, X_m \stackrel{\text{iid}}{\sim} F \in \mathcal{P}$, $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} G \in \mathcal{P}$ (G may or may not equal F .) Let $g(F, G)$ be an estimable function with unbiased estimator $T(X_1, \dots, X_k, Y_1, \dots, Y_l)$. Define

$$T_s(X_1, \dots, X_k, Y_1, \dots, Y_l) = \frac{1}{k!l!} \sum_{P_X} \sum_{P_Y} T(X_{i_1}, \dots, X_{i_k}, Y_{j_1}, \dots, Y_{j_l})$$

(where P_X and P_Y are permutations of X and Y) and

$$U(\underline{X}, \underline{Y}) = \frac{1}{\binom{m}{k} \binom{n}{l}} \sum_{C_X} \sum_{C_Y} T_s(X_{i_1}, \dots, X_{i_k}, Y_{j_1}, \dots, Y_{j_l})$$

(where C_X and C_Y are combinations of X and Y).

U is called a **generalized U -statistic**. ■

Example 12.1.17:

Let X_1, \dots, X_m and Y_1, \dots, Y_n be independent random samples from F and G , respectively, with $F, G \in \mathcal{P}$. We wish to estimate

$$g(F, G) = P_{F,G}(X \leq Y).$$

Let us define

$$Z_{ij} = \begin{cases} 1, & X_i \leq Y_j \\ 0, & X_i > Y_j \end{cases}$$

for each pair X_i, Y_j , $i = 1, 2, \dots, m$, $j = 1, 2, \dots, n$.

Then $\sum_{i=1}^m Z_{ij}$ is the number of X 's $\leq Y_j$, and $\sum_{j=1}^n Z_{ij}$ is the number of Y 's $> X_i$.

$$E(I(X_i \leq Y_j)) = g(F, G) = P_{F,G}(X \leq Y),$$

and degrees k and l are $= 1$, so we use

$$\begin{aligned} U(\underline{X}, \underline{Y}) &= \frac{1}{\binom{m}{1} \binom{n}{1}} \sum_{C_X} \sum_{C_Y} T_s(X_{i_1}, \dots, X_{i_k}, Y_{j_1}, \dots, Y_{j_l}) \\ &= \frac{(m-1)!(n-1)!}{m!n!} \sum_{C_X} \sum_{C_Y} \frac{1}{1!1!} \sum_{P_X} \sum_{P_Y} T(X_{i_1}, \dots, X_{i_k}, Y_{j_1}, \dots, Y_{j_l}) \\ &= \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n I(X_i \leq Y_j). \end{aligned}$$

This **Mann-Whitney estimator** (or **Wilcoxin 2-Sample estimator**) is unbiased and symmetric in the X 's and Y 's. It follows by Corollary 12.1.15 that it has minimum variance. ■

12.2 Single-Sample Hypothesis Tests

Let X_1, \dots, X_n be a sample from a distribution F . The **problem of fit** is to test the hypothesis that the sample X_1, \dots, X_n is from some specified distribution against the alternative that it is from some other distribution, i.e., $H_0 : F = F_0$ vs. $H_1 : F(x) \neq F_0(x)$ for some x .

Definition 12.2.1:

Let $X_1, \dots, X_n \stackrel{iid}{\sim} F$, and let the corresponding empirical cdf be

$$F_n^*(x) = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, x]}(X_i).$$

The statistic

$$D_n = \sup_x | F_n^*(x) - F(x) |$$

is called the **two-sided Kolmogorov–Smirnov statistic (K–S statistic)**.

The **one-sided K–S statistics** are

$$D_n^+ = \sup_x [F_n^*(x) - F(x)] \quad \text{and} \quad D_n^- = \sup_x [F(x) - F_n^*(x)].$$

■

Theorem 12.2.2:

For any continuous distribution F , the K–S statistics D_n, D_n^-, D_n^+ are distribution free.

Proof:

Let $X_{(1)}, \dots, X_{(n)}$ be the order statistics of X_1, \dots, X_n , i.e., $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$, and define $X_{(0)} = -\infty$ and $X_{(n+1)} = +\infty$.

Then,

$$F_n^*(x) = \frac{i}{n} \text{ for } X_{(i)} \leq x < X_{(i+1)}, \quad i = 0, \dots, n.$$

Therefore,

$$\begin{aligned} D_n^+ &= \max_{0 \leq i \leq n} \left\{ \sup_{X_{(i)} \leq x < X_{(i+1)}} \left[\frac{i}{n} - F(x) \right] \right\} \\ &= \max_{0 \leq i \leq n} \left\{ \frac{i}{n} - \left[\inf_{X_{(i)} \leq x < X_{(i+1)}} F(x) \right] \right\} \\ &\stackrel{(*)}{=} \max_{0 \leq i \leq n} \left\{ \frac{i}{n} - F(X_{(i)}) \right\} \\ &= \max \left\{ \max_{1 \leq i \leq n} \left\{ \frac{i}{n} - F(X_{(i)}) \right\}, 0 \right\} \end{aligned}$$

(*) holds since F is nondecreasing in $[X_{(i)}, X_{(i+1)})$.

Note that D_n^+ is a function of $F(X_{(i)})$. In order to make some inference about D_n^+ , the distribution of $F(X_{(i)})$ must be known. We know from the *Probability Integral Transformation* (see Rohatgi, page 203, Theorem 1) that for a rv X with continuous cdf F_X , it holds that $F_X(X) \sim U(0, 1)$.

Thus, $F(X_{(i)})$ is the i^{th} order statistic of a sample from $U(0, 1)$, independent from F . Therefore, the distribution of D_n^+ is independent of F .

Similarly, the distribution of

$$D_n^- = \max \left\{ \max_{1 \leq i \leq n} \left\{ F(X_{(i)}) - \frac{i-1}{n} \right\}, 0 \right\}$$

is independent of F .

Since

$$D_n = \sup_x |F_n^*(x) - F(x)| = \max \{D_n^+, D_n^-\},$$

the distribution of D_n is also independent of F . ■

Theorem 12.2.3:

If F is continuous, then

$$P(D_n \leq \nu + \frac{1}{2n}) = \begin{cases} 0, & \text{if } \nu \leq 0 \\ \int_{\frac{1}{2n}-\nu}^{\nu+\frac{1}{2n}} \int_{\frac{3}{2n}-\nu}^{\nu+\frac{3}{2n}} \cdots \int_{\frac{2n-1}{2n}-\nu}^{\nu+\frac{2n-1}{2n}} f(\underline{u}) d\underline{u}, & \text{if } 0 < \nu < \frac{2n-1}{2n} \\ 1, & \text{if } \nu \geq \frac{2n-1}{2n} \end{cases}$$

where

$$f(\underline{u}) = f(u_1, \dots, u_n) = \begin{cases} n!, & \text{if } 0 < u_1 < u_2 < \dots < u_n < 1 \\ 0, & \text{otherwise} \end{cases}$$

is the joint pdf of an order statistic of a sample of size n from $U(0, 1)$. ■

Note:

As Gibbons & Chakraborti (1992), page 108–109, point out, this result must be interpreted carefully. Consider the case $n = 2$.

For $0 < \nu < \frac{3}{4}$, it holds that

$$P(D_2 \leq \nu + \frac{1}{4}) = \int_{\frac{1}{4}-\nu}^{\nu+\frac{1}{4}} \int_{0 < u_1 < u_2 < 1}^{\nu+\frac{3}{4}} 2! du_2 du_1.$$

Note that the integration limits overlap if

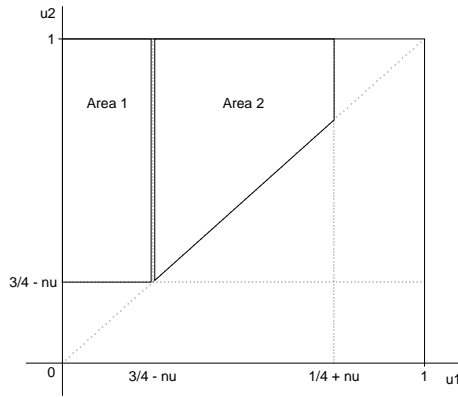
$$\begin{aligned} \nu + \frac{1}{4} &\geq -\nu + \frac{3}{4} \\ \iff \nu &\geq \frac{1}{4} \end{aligned}$$

When $0 < \nu < \frac{1}{4}$, it automatically holds that $0 < u_1 < u_2 < 1$. Thus, for $0 < \nu < \frac{1}{4}$, it holds that

$$\begin{aligned}
 P(D_2 \leq \nu + \frac{1}{4}) &= \int_{\frac{1}{4}-\nu}^{\nu+\frac{1}{4}} \int_{\frac{3}{4}-\nu}^{\nu+\frac{3}{4}} 2! \, du_2 \, du_1 \\
 &= 2! \int_{\frac{1}{4}-\nu}^{\nu+\frac{1}{4}} \left(u_2 \Big|_{\frac{3}{4}-\nu}^{\nu+\frac{3}{4}} \right) du_1 \\
 &= 2! \int_{\frac{1}{4}-\nu}^{\nu+\frac{1}{4}} 2\nu \, du_1 \\
 &= 2! (2\nu) u_1 \Big|_{\frac{1}{4}-\nu}^{\nu+\frac{1}{4}} \\
 &= 2! (2\nu)^2
 \end{aligned}$$

For $\frac{1}{4} \leq \nu < \frac{3}{4}$, the region of integration is as follows:

Note to Theorem 12.2.3



Thus, for $\frac{1}{4} \leq \nu < \frac{3}{4}$, it holds that

$$\begin{aligned}
 P(D_2 \leq \nu + \frac{1}{4}) &= \int_{\frac{1}{4}-\nu}^{\nu+\frac{1}{4}} \int_{0 < u_1 < u_2 < 1} \int_{\frac{3}{4}-\nu}^{\nu+\frac{3}{4}} 2! \, du_2 \, du_1 \\
 &= \int_{\frac{3}{4}-\nu}^{\nu+\frac{1}{4}} \int_{u_1}^1 2! \, du_2 \, du_1 + \int_0^{\frac{3}{4}-\nu} \int_{\frac{3}{4}-\nu}^1 2! \, du_2 \, du_1 \\
 &= 2 \left[\int_{\frac{3}{4}-\nu}^{\nu+\frac{1}{4}} \left(u_2 \Big|_{u_1}^1 \right) du_1 + \int_0^{\frac{3}{4}-\nu} \left(u_2 \Big|_{\frac{3}{4}-\nu}^1 \right) du_1 \right] \\
 &= 2 \left[\int_{\frac{3}{4}-\nu}^{\nu+\frac{1}{4}} (1 - u_1) \, du_1 + \int_0^{\frac{3}{4}-\nu} \left(1 - \frac{3}{4} + \nu \right) du_1 \right]
 \end{aligned}$$

$$\begin{aligned}
&= 2 \left[\left(u_1 - \frac{u_1^2}{2} \right) \Big|_{\frac{3}{4}-\nu}^{\nu+\frac{1}{4}} + \left(\frac{u_1}{4} + \nu u_1 \right) \Big|_0^{\frac{3}{4}-\nu} \right] \\
&= 2 \left[\left(\nu + \frac{1}{4} \right) - \frac{\left(\nu + \frac{1}{4} \right)^2}{2} - \left(-\nu + \frac{3}{4} \right) + \frac{\left(-\nu + \frac{3}{4} \right)^2}{2} + \frac{\left(-\nu + \frac{3}{4} \right)}{4} + \nu \left(-\nu + \frac{3}{4} \right) \right] \\
&= 2 \left[\nu + \frac{1}{4} - \frac{\nu^2}{2} - \frac{\nu}{4} - \frac{1}{32} + \nu - \frac{3}{4} + \frac{\nu^2}{2} - \frac{3\nu}{4} + \frac{9}{32} - \frac{\nu}{4} + \frac{3}{16} - \nu^2 + \frac{3}{4}\nu \right] \\
&= 2 \left[-\nu^2 + \frac{3}{2}\nu - \frac{1}{16} \right] \\
&= -2\nu^2 + 3\nu - \frac{1}{8}
\end{aligned}$$

Combining these results gives

$$P(D_2 \leq \nu + \frac{1}{4}) = \begin{cases} 0, & \text{if } \nu \leq 0 \\ 2! (2\nu)^2, & \text{if } 0 < \nu < \frac{1}{4} \\ -2\nu^2 + 3\nu - \frac{1}{8}, & \text{if } \frac{1}{4} \leq \nu < \frac{3}{4} \\ 1, & \text{if } \nu \geq \frac{3}{4} \end{cases}$$

■

Theorem 12.2.4:

Let F be a continuous cdf. Then it holds $\forall z \geq 0$:

$$\lim_{n \rightarrow \infty} P(D_n \leq \frac{z}{\sqrt{n}}) = L_1(z) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} \exp(-2i^2 z^2).$$

■

Theorem 12.2.5:

Let F be a continuous cdf. Then it holds:

$$P(D_n^+ \leq z) = P(D_n^- \leq z) = \begin{cases} 0, & \text{if } z \leq 0 \\ \int_{1-z}^1 \int_{\frac{n-1}{n}-z}^{u_n} \cdots \int_{\frac{2}{n}-z}^{u_3} \int_{\frac{1}{n}-z}^{u_2} f(\underline{u}) d\underline{u}, & \text{if } 0 < z < 1 \\ 1, & \text{if } z \geq 1 \end{cases}$$

where $f(\underline{u})$ is defined in Theorem 12.2.3.

■

Note:

It should be obvious that the statistics D_n^+ and D_n^- have the same distribution because of symmetry.

■

Theorem 12.2.6:

Let F be a continuous cdf. Then it holds $\forall z \geq 0$:

$$\lim_{n \rightarrow \infty} P(D_n^+ \leq \frac{z}{\sqrt{n}}) = \lim_{n \rightarrow \infty} P(D_n^- \leq \frac{z}{\sqrt{n}}) = L_2(z) = 1 - \exp(-2z^2)$$

■

Corollary 12.2.7:

Let $V_n = 4n(D_n^+)^2$. Then it holds $V_n \xrightarrow{d} \chi_2^2$, i.e., this transformation of D_n^+ has an asymptotic χ_2^2 distribution.

Proof:

Let $x \geq 0$. Then it follows:

$$\begin{aligned} \lim_{n \rightarrow \infty} P(V_n \leq x) &\stackrel{x=4z^2}{=} \lim_{n \rightarrow \infty} P(V_n \leq 4z^2) \\ &= \lim_{n \rightarrow \infty} P(4n(D_n^+)^2 \leq 4z^2) \\ &= \lim_{n \rightarrow \infty} P(\sqrt{n}D_n^+ \leq z) \\ &\stackrel{Th.12.2.6}{=} 1 - \exp(-2z^2) \\ &\stackrel{4z^2=x}{=} 1 - \exp(-x/2) \end{aligned}$$

Thus, $\lim_{n \rightarrow \infty} P(V_n \leq x) = 1 - \exp(-x/2)$ for $x \geq 0$. Note that this is the cdf of a χ_2^2 distribution. ■

Definition 12.2.8:

Let $D_{n;\alpha}$ be the smallest value such that $P(D_n > D_{n;\alpha}) \leq \alpha$. Likewise, let $D_{n;\alpha}^+$ be the smallest value such that $P(D_n^+ > D_{n;\alpha}^+) \leq \alpha$.

The **Kolmogorov–Smirnov test (K–S test)** rejects $H_0 : F(x) = F_0(x) \quad \forall x$ at level α if $D_n > D_{n;\alpha}$.

It rejects $H_0' : F(x) \geq F_0(x) \quad \forall x$ at level α if $D_n^- > D_{n;\alpha}^+$ and it rejects $H_0'' : F(x) \leq F_0(x) \quad \forall x$ at level α if $D_n^+ > D_{n;\alpha}^+$. ■

Note:

Rohatgi, Table 7, page 661, gives values of $D_{n;\alpha}$ and $D_{n;\alpha}^+$ for selected values of α and small n . Theorems 12.2.4 and 12.2.6 allow the approximation of $D_{n;\alpha}$ and $D_{n;\alpha}^+$ for large n . ■

Example 12.2.9:

Let $X_1, \dots, X_n \sim C(1, 0)$. We want to test whether $H_0 : X \sim N(0, 1)$.

The following data has been observed for $x_{(1)}, \dots, x_{(10)}$:

-1.42, -0.43, -0.19, 0.26, 0.30, 0.45, 0.64, 0.96, 1.97, and 4.68

The results for the K-S test have been obtained through the following S-Plus session, i.e., $D_{10}^+ = 0.02219616$, $D_{10}^- = 0.3025681$, and $D_{10} = 0.3025681$:

```
> x _ c(-1.42, -0.43, -0.19, 0.26, 0.30, 0.45, 0.64, 0.96, 1.97, 4.68)
> FX _ pnorm(x)
> FX
[1] 0.07780384 0.33359782 0.42465457 0.60256811 0.61791142 0.67364478
[7] 0.73891370 0.83147239 0.97558081 0.99999857
> Dp _ (1:10)/10 - FX
> Dp
[1] 2.219616e-02 -1.335978e-01 -1.246546e-01 -2.025681e-01 -1.179114e-01
[6] -7.364478e-02 -3.891370e-02 -3.147239e-02 -7.558081e-02 1.434375e-06
> Dm _ FX - (0:9)/10
> Dm
[1] 0.07780384 0.23359782 0.22465457 0.30256811 0.21791142 0.17364478
[7] 0.13891370 0.13147239 0.17558081 0.09999857
> max(Dp)
[1] 0.02219616
> max(Dm)
[1] 0.3025681
> max(max(Dp), max(Dm))
[1] 0.3025681
>
```

```
> ks.gof(x, alternative = "two.sided", mean = 0, sd = 1)
```

```
One-sample Kolmogorov-Smirnov Test
Hypothesized distribution = normal
```

```
data: x
ks = 0.3026, p-value = 0.2617
alternative hypothesis:
True cdf is not the normal distn. with the specified parameters
```

Using Rohatgi, Table 7, page 661, we have to use $D_{10;0.20} = 0.323$ for $\alpha = 0.20$. Since $D_{10} = 0.3026 < 0.323 = D_{10;0.20}$, it is $p > 0.20$. The K-S test does not reject H_0 at level $\alpha = 0.20$. As S-Plus shows, the precise p-value is even $p = 0.2617$. ■

Note:

Comparison between χ^2 and K-S goodness of fit tests:

- K-S uses all available data; χ^2 bins the data and loses information
- K-S works for all sample sizes; χ^2 requires large sample sizes
- it is more difficult to modify K-S for estimated parameters; χ^2 can be easily adapted for estimated parameters
- K-S is “conservative” for discrete data, i.e., it tends to accept H_0 for such data
- the order matters for K-S; χ^2 is better for unordered categorical data



12.3 More on Order Statistics

Definition 12.3.1:

Let F be a continuous cdf. A **tolerance interval** for F with **tolerance coefficient** γ is a random interval such that the probability is γ that this random interval covers at least a specified percentage $100p\%$ of the distribution. ■

Theorem 12.3.2:

If order statistics $X_{(r)} < X_{(s)}$ are used as the endpoints for a tolerance interval for a continuous cdf F , it holds that

$$\gamma = \sum_{i=0}^{s-r-1} \binom{n}{i} p^i (1-p)^{n-i}.$$

Proof:

According to Definition 12.3.1, it holds that

$$\gamma = P_{X_{(r)}, X_{(s)}} \left(P_X(X_{(r)} < X < X_{(s)}) \geq p \right).$$

Since F is continuous, it holds that $F_X(X) \sim U(0, 1)$. Therefore,

$$\begin{aligned} P_X(X_{(r)} < X < X_{(s)}) &= P(X < X_{(s)}) - P(X \leq X_{(r)}) \\ &= F(X_{(s)}) - F(X_{(r)}) \\ &= U_{(s)} - U_{(r)}, \end{aligned}$$

where $U_{(s)}$ and $U_{(r)}$ are the order statistics of a $U(0, 1)$ distribution.

Thus,

$$\gamma = P_{X_{(r)}, X_{(s)}} \left(P_X(X_{(r)} < X < X_{(s)}) \geq p \right) = P(U_{(s)} - U_{(r)} \geq p).$$

By Theorem 4.4.4, we can determine the joint distribution of order statistics and calculate γ as

$$\gamma = \int_p^1 \int_0^{y-p} \frac{n!}{(r-1)!(s-r-1)!(n-s)!} x^{r-1} (y-x)^{s-r-1} (1-y)^{n-s} dx dy.$$

Rather than solving this integral directly, we make the transformation

$$\begin{aligned} U &= U_{(s)} - U_{(r)} \\ V &= U_{(s)}. \end{aligned}$$

Then the joint pdf of U and V is

$$f_{U,V}(u, v) = \begin{cases} \frac{n!}{(r-1)!(s-r-1)!(n-s)!} (v-u)^{r-1} u^{s-r-1} (1-v)^{n-s}, & \text{if } 0 < u < v < 1 \\ 0, & \text{otherwise} \end{cases}$$

and the marginal pdf of U is

$$\begin{aligned}
f_U(u) &= \int_0^1 f_{U,V}(u,v) dv \\
&= \frac{n!}{(r-1)!(s-r-1)!(n-s)!} u^{s-r-1} I_{(0,1)}(u) \int_u^1 (v-u)^{r-1} (1-v)^{n-s} dv \\
&\stackrel{(A)}{=} \frac{n!}{(r-1)!(s-r-1)!(n-s)!} u^{s-r-1} (1-u)^{n-s+r} I_{(0,1)}(u) \underbrace{\int_0^1 t^{r-1} (1-t)^{n-s} dt}_{B(r,n-s+1)} \\
&= \frac{n!}{(r-1)!(s-r-1)!(n-s)!} u^{s-r-1} (1-u)^{n-s+r} \frac{(r-1)!(n-s)!}{(n-s+r)!} I_{(0,1)}(u) \\
&= \frac{n!}{(n-s+r)!(s-r-1)!} u^{s-r-1} (1-u)^{n-s+r} I_{(0,1)}(u) \\
&= n \binom{n-1}{s-r-1} u^{s-r-1} (1-u)^{n-s+r} I_{(0,1)}(u).
\end{aligned}$$

(A) is based on the transformation $t = \frac{v-u}{1-u}$, $v-u = (1-u)t$, $1-v = 1-u - (1-u)t = (1-u)(1-t)$ and $dv = (1-u)dt$.

It follows that

$$\begin{aligned}
\gamma &= P(U_{(s)} - U_{(r)} \geq p) \\
&= P(U \geq p) \\
&= \int_p^1 n \binom{n-1}{s-r-1} u^{s-r-1} (1-u)^{n-s+r} du \\
&\stackrel{(B)}{=} P(Y < s-r) \quad | \quad \text{where } Y \sim \text{Bin}(n, p) \\
&= \sum_{i=0}^{s-r-1} \binom{n}{i} p^i (1-p)^{n-i}.
\end{aligned}$$

(B) holds due to Rohatgi, Remark 3 after Theorem 5.3.18, page 216, since for $X \sim \text{Bin}(n, p)$, it holds that

$$P(X < k) = \int_p^1 n \binom{n-1}{k-1} x^{k-1} (1-x)^{n-k} dx.$$

■

Example 12.3.3:

Let $s = n$ and $r = 1$. Then,

$$\gamma = \sum_{i=0}^{n-2} \binom{n}{i} p^i (1-p)^{n-i} = 1 - p^n - np^{n-1}(1-p).$$

If $p = 0.8$ and $n = 10$, then

$$\gamma_{10} = 1 - (0.8)^{10} - 10 \cdot (0.8)^9 \cdot (0.2) = 0.624,$$

i.e., $(X_{(1)}, X_{(10)})$ defines a 62.4% tolerance interval for 80% probability.

If $p = 0.8$ and $n = 20$, then

$$\gamma_{20} = 1 - (0.8)^{20} - 20 \cdot (0.8)^{19} \cdot (0.2) = 0.931,$$

and if $p = 0.8$ and $n = 30$, then

$$\gamma_{30} = 1 - (0.8)^{30} - 30 \cdot (0.8)^{29} \cdot (0.2) = 0.989.$$

■

Theorem 12.3.4:

Let k_p be the p^{th} quantile of a continuous cdf F . Let $X_{(1)}, \dots, X_{(n)}$ be the order statistics of a sample of size n from F . Then it holds that

$$P(X_{(r)} \leq k_p \leq X_{(s)}) = \sum_{i=r}^{s-1} \binom{n}{i} p^i (1-p)^{n-i}.$$

Proof:

It holds that

$$\begin{aligned} P(X_{(r)} \leq k_p) &= P(\text{at least } r \text{ of the } X_i \text{'s are } \leq k_p) \\ &= \sum_{i=r}^n \binom{n}{i} p^i (1-p)^{n-i}. \end{aligned}$$

Therefore,

$$\begin{aligned} P(X_{(r)} \leq k_p \leq X_{(s)}) &= P(X_{(r)} \leq k_p) - P(X_{(s)} < k_p) \\ &= \sum_{i=r}^n \binom{n}{i} p^i (1-p)^{n-i} - \sum_{i=s}^n \binom{n}{i} p^i (1-p)^{n-i} \\ &= \sum_{i=r}^{s-1} \binom{n}{i} p^i (1-p)^{n-i}. \end{aligned}$$

■

Corollary 12.3.5:

$(X_{(r)}, X_{(s)})$ is a level $\sum_{i=r}^{s-1} \binom{n}{i} p^i (1-p)^{n-i}$ confidence interval for k_p . ■

Example 12.3.6:

Let $n = 10$. We want a 95% confidence interval for the median, i.e., k_p where $p = \frac{1}{2}$.

We get the following probabilities $p_{r,s} = \sum_{i=r}^{s-1} \binom{n}{i} p^i (1-p)^{n-i}$ that $(X_{(r)}, X_{(s)})$ covers $k_{0.5}$:

$p_{r,s}$		s								
		2	3	4	5	6	7	8	9	10
r	1	0.01	0.05	0.17	0.38	0.62	0.83	0.94	0.99	0.998
	2		0.04	0.16	0.37	0.61	0.82	0.93	0.98	0.99
	3			0.12	0.32	0.57	0.77	0.89	0.93	0.94
	4				0.21	0.45	0.66	0.77	0.82	0.83
	5					0.25	0.45	0.57	0.61	0.62
	6						0.21	0.32	0.37	0.38
	7							0.12	0.16	0.17
	8								0.04	0.05
	9									0.01

Only the random intervals $(X_{(1)}, X_{(9)})$, $(X_{(1)}, X_{(10)})$, $(X_{(2)}, X_{(9)})$, and $(X_{(2)}, X_{(10)})$ give the desired coverage probability. Therefore, we use the one that comes closest to 95%, i.e., $(X_{(2)}, X_{(9)})$, as the 95% confidence interval for the median. ■

13 Some Results from Sampling

13.1 Simple Random Samples

Definition 13.1.1:

Let Ω be a population of size N with mean μ and variance σ^2 . A sampling method (of size n) is called **simple** if the set S of possible samples contains all combinations of n elements of Ω (without repetition) and the probability for each sample $s \in S$ to become selected depends only on n , i.e., $p(s) = \frac{1}{\binom{N}{n}} \forall s \in S$. Then we call $s \in S$ a **simple random sample (SRS)** of size n . ■

Theorem 13.1.2:

Let Ω be a population of size N with mean μ and variance σ^2 . Let $Y : \Omega \rightarrow \mathbb{R}$ be a measurable function. Let n_i be the total number of times the parameter \tilde{y}_i occurs in the population and $p_i = \frac{n_i}{N}$ be the relative frequency the parameter \tilde{y}_i occurs in the population. Let (y_1, \dots, y_n) be a SRS of size n with respect to Y , where $P(Y = \tilde{y}_i) = p_i = \frac{n_i}{N}$.

Then the components $y_i, i = 1, \dots, n$, are identically distributed as Y and it holds for $i \neq j$:

$$P(y_i = \tilde{y}_k, y_j = \tilde{y}_l) = \frac{1}{N(N-1)} n_{kl}, \text{ where } n_{kl} = \begin{cases} n_k n_l, & k \neq l \\ n_k(n_k - 1), & k = l \end{cases}$$

Note:

- (i) In Sampling, many authors use capital letters to denote properties of the population and small letters to denote properties of the random sample. In particular, x_i 's and y_i 's are considered as random variables related to the sample. They are not seen as specific realizations.
- (ii) The following equalities hold in the scenario of Theorem 13.1.2:

$$\begin{aligned} N &= \sum_i n_i \\ \mu &= \frac{1}{N} \sum_i n_i \tilde{y}_i \\ \sigma^2 &= \frac{1}{N} \sum_i n_i (\tilde{y}_i - \mu)^2 \\ &= \frac{1}{N} \sum_i n_i \tilde{y}_i^2 - \mu^2 \end{aligned}$$

Theorem 13.1.3:

Let the same conditions hold as in Theorem 13.1.2. Let $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ be the sample mean of a SRS of size n . Then it holds:

(i) $E(\bar{y}) = \mu$, i.e., the sample mean is unbiased for the population mean μ .

(ii) $Var(\bar{y}) = \frac{1}{n} \frac{N-n}{N-1} \sigma^2 = \frac{1}{n} (1-f) \frac{N}{N-1} \sigma^2$, where $f = \frac{n}{N}$.

Proof:

(i)

$$E(\bar{y}) = \frac{1}{n} \sum_{i=1}^n E(y_i) = \mu, \quad \text{since } E(y_i) = \mu \quad \forall i.$$

(ii)

■

Theorem 13.1.4:

Let \bar{y}_n be the sample mean of a SRS of size n . Then it holds that

$$\sqrt{\frac{n}{1-f}} \frac{\bar{y}_n - \mu}{\sqrt{\frac{N}{N-1}} \sigma} \xrightarrow{d} N(0, 1),$$

where $N \rightarrow \infty$ and $f = \frac{n}{N}$ is a constant.

In particular, when the y_i 's are 0–1–distributed with $E(y_i) = P(y_i = 1) = p \quad \forall i$, then it holds that

$$\sqrt{\frac{n}{1-f}} \frac{\bar{y}_n - p}{\sqrt{\frac{N}{N-1} p(1-p)}} \xrightarrow{d} N(0, 1),$$

where $N \rightarrow \infty$ and $f = \frac{n}{N}$ is a constant.

■

13.2 Stratified Random Samples

Definition 13.2.1:

Let Ω be a population of size N , that is split into m disjoint sets Ω_j , called **strata**, of size N_j , $j = 1, \dots, m$, where $N = \sum_{j=1}^m N_j$. If we independently draw a random sample of size n_j in each strata, we speak of a **stratified random sample**. ■

Note:

- (i) The random samples in each strata are not always SRS's.
 - (ii) Stratified random samples are used in practice as a means to reduce the sample variance in the case that data in each strata is homogeneous and data among different strata is heterogeneous.
 - (iii) Frequently used strata in practice are gender, state (or county), income range, ethnic background, etc.
-

Definition 13.2.2:

Let $Y : \Omega \rightarrow \mathbb{R}$ be a measurable function. In case of a stratified random sample, we use the following notation:

Let Y_{jk} , $j = 1, \dots, m$, $k = 1, \dots, N_j$ be the elements in Ω_j . Then, we define

- (i) $Y_j = \sum_{k=1}^{N_j} Y_{jk}$ the total in the j^{th} strata,
- (ii) $\mu_j = \frac{1}{N_j} Y_j$ the mean in the j^{th} strata,
- (iii) $\mu = \frac{1}{N} \sum_{j=1}^m N_j \mu_j$ the expectation (or grand mean),
- (iv) $N\mu = \sum_{j=1}^m Y_j = \sum_{j=1}^m \sum_{k=1}^{N_j} Y_{jk}$ the total,
- (v) $\sigma_j^2 = \frac{1}{N_j} \sum_{k=1}^{N_j} (Y_{jk} - \mu_j)^2$ the variance in the j^{th} strata, and
- (vi) $\sigma^2 = \frac{1}{N} \sum_{j=1}^m \sum_{k=1}^{N_j} (Y_{jk} - \mu)^2$ the variance.

(vii) We denote an (ordered) sample in Ω_j of size n_j as $(y_{j1}, \dots, y_{jn_j})$ and $\bar{y}_j = \frac{1}{n_j} \sum_{k=1}^{n_j} y_{jk}$ the sample mean in the j^{th} strata. ■

Theorem 13.2.3:

Let the same conditions hold as in Definitions 13.2.1 and 13.2.2. Let $\hat{\mu}_j$ be an unbiased estimate of μ_j and $\widehat{Var}(\hat{\mu}_j)$ be an unbiased estimate of $Var(\hat{\mu}_j)$. Then it holds:

(i) $\hat{\mu} = \frac{1}{N} \sum_{j=1}^m N_j \hat{\mu}_j$ is unbiased for μ .

$$Var(\hat{\mu}) = \frac{1}{N^2} \sum_{j=1}^m N_j^2 Var(\hat{\mu}_j).$$

(ii) $\widehat{Var}(\hat{\mu}) = \frac{1}{N^2} \sum_{j=1}^m N_j^2 \widehat{Var}(\hat{\mu}_j)$ is unbiased for $Var(\hat{\mu})$.

Proof:

(i)

Theorem 13.2.4:

Let the same conditions hold as in Theorem 13.2.3. If we draw a SRS in each strata, then it holds:

(i) $\hat{\mu} = \frac{1}{N} \sum_{j=1}^m N_j \bar{y}_j$ is unbiased for μ , where $\bar{y}_j = \frac{1}{n_j} \sum_{k=1}^{n_j} y_{jk}$, $j = 1, \dots, m$.

$$Var(\hat{\mu}) = \frac{1}{N^2} \sum_{j=1}^m N_j^2 \frac{1}{n_j} (1 - f_j) \frac{N_j}{N_j - 1} \sigma_j^2, \text{ where } f_j = \frac{n_j}{N_j}.$$

(ii) $\widehat{Var}(\hat{\mu}) = \frac{1}{N^2} \sum_{j=1}^m N_j^2 \frac{1}{n_j} (1 - f_j) s_j^2$ is unbiased for $Var(\hat{\mu})$, where

$$s_j^2 = \frac{1}{n_j - 1} \sum_{k=1}^{n_j} (y_{jk} - \bar{y}_j)^2.$$

Proof:

■

Definition 13.2.5:

Let the same conditions hold as in Definitions 13.2.1 and 13.2.2. If the sample in each strata is of size $n_j = n \frac{N_j}{N}$, $j = 1, \dots, m$, where n is the total sample size, then we speak of **proportional selection**.

■

Note:

(i) In the case of proportional selection, it holds that $f_j = \frac{n_j}{N_j} = \frac{n}{N} = f$, $j = 1, \dots, m$.

(ii) Proportional strata cannot always be obtained for each combination of m , n , and N .

■

Theorem 13.2.6:

Let the same conditions hold as in Definition 13.2.5. If we draw a SRS in each strata, then it holds in case of proportional selection that

$$Var(\hat{\mu}) = \frac{1}{N^2} \frac{1-f}{f} \sum_{j=1}^m N_j \tilde{\sigma}_j^2,$$

where $\tilde{\sigma}_j^2 = \frac{N_j}{N_j-1} \sigma_j^2$.

Proof:

The proof follows directly from Theorem 13.2.4 (i).

■

Theorem 13.2.7:

If we draw (1) a stratified random sample that consists of SRS's of sizes n_j under proportional selection and (2) a SRS of size $n = \sum_{j=1}^m n_j$ from the same population, then it holds that

$$Var(\bar{y}) - Var(\hat{\mu}) = \frac{1}{n} \frac{N-n}{N(N-1)} \left(\sum_{j=1}^m N_j (\mu_j - \mu)^2 - \frac{1}{N} \sum_{j=1}^m (N - N_j) \tilde{\sigma}_j^2 \right).$$

Proof:

See Homework.



14 Some Results from Sequential Statistical Inference

14.1 Fundamentals of Sequential Sampling

Example 14.1.1:

A particular machine produces a large number of items every day. Each item can be either “defective” or “non-defective”. The unknown proportion of defective items in the production of a particular day is p .

Let (X_1, \dots, X_m) be a sample from the daily production where $x_i = 1$ when the item is defective and $x_i = 0$ when the item is non-defective. Obviously, $S_m = \sum_{i=1}^m X_i \sim \text{Bin}(m, p)$ denotes the total number of defective items in the sample (assuming that m is small compared to the daily production).

We might be interested to test $H_0 : p \leq p_0$ vs. $H_1 : p > p_0$ at a given significance level α and use this decision to trash the entire daily production and have the machine fixed if indeed $p > p_0$. A suitable test could be

$$\Phi_1(x_1, \dots, x_m) = \begin{cases} 1, & \text{if } s_m > c \\ 0, & \text{if } s_m \leq c \end{cases}$$

where c is chosen such that Φ_1 is a level- α test.

However, wouldn't it be more beneficial if we sequentially sample the items (e.g., take item # 57, 623, 1005, 1286, 2663, etc.) and stop the machine as soon as it becomes obvious that it produces too many bad items. (Alternatively, we could also finish the time consuming and expensive process to determine whether an item is defective or non-defective if it is impossible to surpass a certain proportion of defectives.) For example, if for some $j < m$ it already holds that $s_j > c$, then we could stop (and immediately call maintenance) and reject H_0 after only j observations.

More formally, let us define $T = \min\{j \mid S_j > c\}$ and $T' = \min\{T, m\}$. We can now consider a decision rule that stops with the sampling process at random time T' and rejects H_0 if $T \leq m$. Thus, if we consider $R_0 = \{(x_1, \dots, x_m) \mid t \leq m\}$ and $R_1 = \{(x_1, \dots, x_m) \mid s_m > c\}$ as critical regions of two tests Φ_0 and Φ_1 , then these two tests are equivalent. ■

Definition 14.1.2:

Let Θ be the parameter space and \mathcal{A} the set of actions the statistician can take. We assume that the rv's X_1, X_2, \dots are observed sequentially and iid with common pdf (or pmf) $f_\theta(x)$. A **sequential decision procedure** is defined as follows:

- (i) A **stopping rule** specifies whether an element of \mathcal{A} should be chosen without taking any further observation. If at least one observation is taken, this rule specifies for every set of observed values (x_1, x_2, \dots, x_n) , $n \geq 1$, whether to stop sampling and choose an action in \mathcal{A} or to take another observation x_{n+1} .
- (ii) A **decision rule** specifies the decision to be taken. If no observation has been taken, then we take action $d_0 \in \mathcal{A}$. If $n \geq 1$ observation have been taken, then we take action $d_n(x_1, \dots, x_n) \in \mathcal{A}$, where $d_n(x_1, \dots, x_n)$ specifies the action that has to be taken for the set (x_1, \dots, x_n) of observed values. Once an action has been taken, the sampling process is stopped.

■

Note:

In the remainder of this chapter, we assume that the statistician takes at least one observation.

■

Definition 14.1.3:

Let $R_n \subseteq \mathbb{R}^n$, $n = 1, 2, \dots$, be a sequence of Borel-measurable sets such that the sampling process is stopped after observing $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ if $(x_1, \dots, x_n) \in R_n$. If $(x_1, \dots, x_n) \notin R_n$, then another observation x_{n+1} is taken. The sets R_n , $n = 1, 2, \dots$ are called **stopping regions**.

■

Definition 14.1.4:

With every sequential stopping rule we associate a **stopping random variable** N which takes on the values $1, 2, 3, \dots$. Thus, N is a rv that indicates the total number of observations taken before the sampling is stopped.

■

Note:

We use the (sloppy) notation $\{N = n\}$ to denote the event that sampling is stopped after observing exactly n values x_1, \dots, x_n (i.e., sampling is not stopped before taking n samples). Then the following equalities hold:

$$\{N = 1\} = R_1$$

$$\begin{aligned}
\{N = n\} &= \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid \text{sampling is stopped after } n \text{ observations but not before}\} \\
&= (R_1 \cup R_2 \cup \dots \cup R_{n-1})^c \cap R_n \\
&= R_1^c \cap R_2^c \cap \dots \cap R_{n-1}^c \cap R_n
\end{aligned}$$

$$\{N \leq n\} = \bigcup_{k=1}^n \{N = k\}$$

Here we will only consider **closed** sequential sampling procedures, i.e., procedures where sampling eventually stops with probability 1, i.e.,

$$\begin{aligned}
P(N < \infty) &= 1, \\
P(N = \infty) &= 1 - P(N < \infty) = 0.
\end{aligned}$$

■

Theorem 14.1.5: Wald's Equation

Let X_1, X_2, \dots be iid rv's with $E(|X_1|) < \infty$. Let N be a stopping variable. Let $S_N = \sum_{k=1}^N X_k$.

If $E(N) < \infty$, then it holds

$$E(S_N) = E(X_1)E(N).$$

Proof:

Define a sequence of rv's Y_i , $i = 1, 2, \dots$, where

$$Y_i = \begin{cases} 1, & \text{if no decision is reached up to the } (i-1)^{\text{th}} \text{ stage, i.e., } N > (i-1) \\ 0, & \text{otherwise} \end{cases}$$

Then each Y_i is a function of X_1, X_2, \dots, X_{i-1} only and Y_i is independent of X_i .

Consider the rv

$$\sum_{n=1}^{\infty} X_n Y_n.$$

Obviously, it holds that

$$S_N = \sum_{n=1}^{\infty} X_n Y_n.$$

Thus, it follows that

$$E(S_N) = E\left(\sum_{n=1}^{\infty} X_n Y_n\right). \quad (*)$$

It holds that

$$\begin{aligned}
\sum_{n=1}^{\infty} E(|X_n Y_n|) &= \sum_{n=1}^{\infty} E(|X_n|)E(|Y_n|) \\
&= E(|X_1|) \sum_{n=1}^{\infty} P(N \geq n)
\end{aligned}$$

$$\begin{aligned}
&= E(|X_1|) \sum_{n=1}^{\infty} \sum_{k=n}^{\infty} P(N = k) \\
&\stackrel{(A)}{=} E(|X_1|) \sum_{n=1}^{\infty} nP(N = n) \\
&= E(|X_1|)E(N) \\
&< \infty
\end{aligned}$$

(A) holds due to the following rearrangement of indices:

n	k
1	1, 2, 3, ...
2	2, 3, ...
3	3, ...
\vdots	\vdots

We may therefore interchange the expectation and summation signs in (*) and get

$$\begin{aligned}
E(S_N) &= E\left(\sum_{n=1}^{\infty} X_n Y_n\right) \\
&= \sum_{n=1}^{\infty} E(X_n Y_n) \\
&= \sum_{n=1}^{\infty} E(X_n)E(Y_n) \\
&= E(X_1) \sum_{n=1}^{\infty} P(N \geq n) \\
&= E(X_1)E(N)
\end{aligned}$$

which completes the proof. ■

14.2 Sequential Probability Ratio Tests

Definition 14.2.1:

Let X_1, X_2, \dots be a sequence of iid rv's with common pdf (or pmf) $f_\theta(x)$. We want to test a simple hypothesis $H_0 : X \sim f_{\theta_0}$ vs. a simple alternative $H_1 : X \sim f_{\theta_1}$ when the observations are taken sequentially.

Let f_{0n} and f_{1n} denote the joint pdf's (or pmf's) of X_1, \dots, X_n under H_0 and H_1 respectively, i.e.,

$$f_{0n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{\theta_0}(x_i) \quad \text{and} \quad f_{1n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{\theta_1}(x_i).$$

Finally, let

$$\lambda_n(x_1, \dots, x_n) = \frac{f_{1n}(\underline{x})}{f_{0n}(\underline{x})},$$

where $\underline{x} = (x_1, \dots, x_n)$. Then a **sequential probability ratio test (SPRT)** for testing H_0 vs. H_1 is the following decision rule:

- (i) If at any stage of the sampling process it holds that

$$\lambda_n(\underline{x}) \geq A,$$

then stop and reject H_0 .

- (ii) If at any stage of the sampling process it holds that

$$\lambda_n(\underline{x}) \leq B,$$

then stop and accept H_0 , i.e., reject H_1 .

- (iii) If

$$B < \lambda_n(\underline{x}) < A,$$

then continue sampling by taking another observation x_{n+1} .

■

Note:

- (i) It is usually convenient to define

$$Z_i = \log \frac{f_{\theta_1}(X_i)}{f_{\theta_0}(X_i)},$$

where Z_1, Z_2, \dots are iid rv's. Then, we work with

$$\log \lambda_n(\underline{x}) = \sum_{i=1}^n z_i = \sum_{i=1}^n (\log f_{\theta_1}(x_i) - \log f_{\theta_0}(x_i))$$

instead of using $\lambda_n(\underline{x})$. Obviously, we now have to use constants $b = \log B$ and $a = \log A$ instead of the original constants B and A .

(ii) A and B (where $A > B$) are constants such that the SPRT will have strength (α, β) , where

$$\alpha = P(\text{Type I error}) = P(\text{Reject } H_0 \mid H_0)$$

and

$$\beta = P(\text{Type II error}) = P(\text{Accept } H_0 \mid H_1).$$

If N is the stopping rv, then

$$\alpha = P_{\theta_0}(\lambda_N(\underline{X}) \geq A) \quad \text{and} \quad \beta = P_{\theta_1}(\lambda_N(\underline{X}) \leq B).$$

■

Example 14.2.2:

Let X_1, X_2, \dots be iid $N(\mu, \sigma^2)$, where μ is unknown and $\sigma^2 > 0$ is known. We want to test $H_0 : \mu = \mu_0$ vs. $H_1 : \mu = \mu_1$, where $\mu_0 < \mu_1$.

If our data is sampled sequentially, we can construct a SPRT as follows:

$$\begin{aligned} \log \lambda_n(\underline{x}) &= \sum_{i=1}^n \left(-\frac{1}{2\sigma^2}(x_i - \mu_1)^2 - \left(-\frac{1}{2\sigma^2}(x_i - \mu_0)^2 \right) \right) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^n \left((x_i - \mu_0)^2 - (x_i - \mu_1)^2 \right) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^n \left(x_i^2 - 2x_i\mu_0 + \mu_0^2 - x_i^2 + 2x_i\mu_1 - \mu_1^2 \right) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^n \left(-2x_i\mu_0 + \mu_0^2 + 2x_i\mu_1 - \mu_1^2 \right) \\ &= \frac{1}{2\sigma^2} \left(\sum_{i=1}^n 2x_i(\mu_1 - \mu_0) + n(\mu_0^2 - \mu_1^2) \right) \\ &= \frac{\mu_1 - \mu_0}{\sigma^2} \left(\sum_{i=1}^n x_i - n \frac{\mu_0 + \mu_1}{2} \right) \end{aligned}$$

We decide for H_0 if

$$\begin{aligned} \log \lambda_n(\underline{x}) &\leq b \\ \iff \frac{\mu_1 - \mu_0}{\sigma^2} \left(\sum_{i=1}^n x_i - n \frac{\mu_0 + \mu_1}{2} \right) &\leq b \\ \iff \sum_{i=1}^n x_i &\leq n \frac{\mu_0 + \mu_1}{2} + b^*, \end{aligned}$$

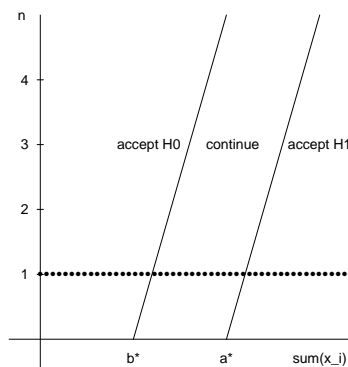
where $b^* = \frac{\sigma^2}{\mu_1 - \mu_0} b$.

We decide for H_1 if

$$\begin{aligned} \log \lambda_n(\underline{x}) &\geq a \\ \Leftrightarrow \frac{\mu_1 - \mu_0}{\sigma^2} \left(\sum_{i=1}^n x_i - n \frac{\mu_0 + \mu_1}{2} \right) &\geq a \\ \Leftrightarrow \sum_{i=1}^n x_i &\geq n \frac{\mu_0 + \mu_1}{2} + a^*, \end{aligned}$$

where $a^* = \frac{\sigma^2}{\mu_1 - \mu_0} a$.

Example 14.2.2



■

Theorem 14.2.3:

For a SPRT with stopping bounds A and B , $A > B$, and strength (α, β) , we have

$$A \leq \frac{1 - \beta}{\alpha} \quad \text{and} \quad B \geq \frac{\beta}{1 - \alpha},$$

where $0 < \alpha < 1$ and $0 < \beta < 1$.

■

Theorem 14.2.4:

Assume we select for given $\alpha, \beta \in (0, 1)$, where $\alpha + \beta \leq 1$, the stopping bounds

$$A' = \frac{1 - \beta}{\alpha} \quad \text{and} \quad B' = \frac{\beta}{1 - \alpha}.$$

Then it holds that the SPRT with stopping bounds A' and B' has strength (α', β') , where

$$\alpha' \leq \frac{\alpha}{1 - \beta}, \quad \beta' \leq \frac{\beta}{1 - \alpha}, \quad \text{and} \quad \alpha' + \beta' \leq \alpha + \beta.$$

■

Note:

- (i) The approximation $A' = \frac{1-\beta}{\alpha}$ and $B' = \frac{\beta}{1-\alpha}$ in Theorem 14.2.4 is called **Wald-Approximation** for the optimal stopping bounds of a SPRT.
- (ii) A' and B' are functions of α and β only and do not depend on the pdf's (or pmf's) f_{θ_0} and f_{θ_1} . Therefore, they can be computed once and for all f_{θ_i} 's, $i = 0, 1$.

■

THE END !!!