# List of formulas for STK4011/9011 – Statistical Inference Theory (2022)

## Chapter 2: Transformations and Expectations

**Theorem 2.1.3** *Let $X$ have cdf $F_X(x)$, let $Y = g(X)$, and let $\mathcal{X}$ and $\mathcal{Y}$ be defined as in (2.1.7).*

**a.** *If $g$ is an increasing function on $\mathcal{X}$, $F_Y(y) = F_X\left(g^{-1}(y)\right)$ for $y \in \mathcal{Y}$.*

**b.** *If $g$ is a decreasing function on $\mathcal{X}$ and $X$ is a continuous random variable, $F_Y(y) = 1 - F_X\left(g^{-1}(y)\right)$ for $y \in \mathcal{Y}$.*

$$(2.1.7) \quad \mathcal{X} = \{x \colon f_X(x) > 0\} \quad \text{and} \quad \mathcal{Y} = \{y \colon y = g(x) \text{ for some } x \in \mathcal{X}\}.$$

**Theorem 2.1.5** *Let $X$ have pdf $f_X(x)$ and let $Y = g(X)$, where $g$ is a monotone function. Let $\mathcal{X}$ and $\mathcal{Y}$ be defined by (2.1.7). Suppose that $f_X(x)$ is continuous on $\mathcal{X}$ and that $g^{-1}(y)$ has a continuous derivative on $\mathcal{Y}$. Then the pdf of $Y$ is given by*

$$(2.1.10) \quad f_Y(y) = \begin{cases} f_X(g^{-1}(y)) \left| \dfrac{d}{dy} g^{-1}(y) \right| & y \in \mathcal{Y} \\ 0 & \text{otherwise.} \end{cases}$$

**Definition 2.2.1** *The expected value or mean of a random variable $g(X)$, denoted by $\mathrm{E}\, g(X)$, is*

$$\mathrm{E}\, g(X) = \begin{cases} \int_{-\infty}^{\infty} g(x) f_X(x)\, dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x) f_X(x) = \sum_{x \in \mathcal{X}} g(x) P(X = x) & \text{if } X \text{ is discrete,} \end{cases}$$

**Theorem 2.2.5** *Let $X$ be a random variable and let $a, b$, and $c$ be constants. Then for any functions $g_1(x)$ and $g_2(x)$ whose expectations exist,*

**a.** $\mathrm{E}(ag_1(X) + bg_2(X) + c) = a\mathrm{E}\, g_1(X) + b\mathrm{E}\, g_2(X) + c$.

**b.** *If $g_1(x) \geq 0$ for all $x$, then $\mathrm{E}\, g_1(X) \geq 0$.*

**c.** *If $g_1(x) \geq g_2(x)$ for all $x$, then $\mathrm{E}\, g_1(X) \geq \mathrm{E}\, g_2(X)$.*

**d.** *If $a \leq g_1(x) \leq b$ for all $x$, then $a \leq \mathrm{E}\, g_1(X) \leq b$.*

**Definition 2.3.1** For each integer $n$, the $n$th *moment of $X$* (or $F_X(x)$), $\mu'_n$, is

$$\mu'_n = \mathrm{E}\, X^n.$$

The $n$th *central moment of $X$*, $\mu_n$, is

$$\mu_n = \mathrm{E}(X - \mu)^n,$$

where $\mu = \mu'_1 = \mathrm{E}\, X$.

**Definition 2.3.2** The *variance* of a random variable $X$ is its second central moment, $\mathrm{Var}\, X = \mathrm{E}(X - \mathrm{E}\, X)^2$. The positive square root of $\mathrm{Var}\, X$ is the *standard deviation* of $X$.

**Theorem 2.3.4** *If $X$ is a random variable with finite variance, then for any constants $a$ and $b$,*

$$\mathrm{Var}(aX + b) = a^2\, \mathrm{Var}\, X.$$

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

$$M_X(t) = \mathrm{E}\, e^{tX},$$

**Theorem 2.3.11** *Let $F_X(x)$ and $F_Y(y)$ be two cdfs all of whose moments exist.*

**a.** *If $X$ and $Y$ have bounded support, then $F_X(u) = F_Y(u)$ for all $u$ if and only if $\mathrm{E}\, X^r = \mathrm{E}\, Y^r$ for all integers $r = 0, 1, 2, \dots$.*

**b.** *If the moment generating functions exist and $M_X(t) = M_Y(t)$ for all $t$ in some neighborhood of 0, then $F_X(u) = F_Y(u)$ for all $u$.*

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

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

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

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

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

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

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

# Chapter 3: Common Families of Distributions

A family of pdfs or pmfs is called an *exponential family* if it can be expressed as

$$(3.4.1) \qquad f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp\left(\sum_{i=1}^{k} w_i(\boldsymbol{\theta})t_i(x)\right).$$

Here $h(x) \geq 0$ and $t_1(x), \ldots, t_k(x)$ are real-valued functions of the observation $x$ (they cannot depend on $\boldsymbol{\theta}$), and $c(\boldsymbol{\theta}) \geq 0$ and $w_1(\boldsymbol{\theta}), \ldots, w_k(\boldsymbol{\theta})$ are real-valued functions of the possibly vector-valued parameter $\boldsymbol{\theta}$ (they cannot depend on $x$).

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

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

**Theorem 3.6.1 (Chebychev's Inequality)** *Let $X$ be a random variable and let $g(x)$ be a nonnegative function. Then, for any $r > 0$,*

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

# Chapter 4: Multiple Random Variables

If $(X, Y)$ is a discrete bivariate random vector, then there is only a countable set of values for which the joint pmf of $(X, Y)$ is positive. Call this set $\mathcal{A}$. Define the set $\mathcal{B} = \{(u, v) : u = g_1(x, y) \text{ and } v = g_2(x, y) \text{ for some } (x, y) \in \mathcal{A}\}$. Then $\mathcal{B}$ is the countable set of possible values for the discrete random vector $(U, V)$. And if, for any $(u, v) \in \mathcal{B}$, $A_{uv}$ is defined to be $\{(x, y) \in \mathcal{A} : g_1(x, y) = u \text{ and } g_2(x, y) = v\}$, then the joint pmf of $(U, V)$, $f_{U,V}(u, v)$, can be computed from the joint pmf of $(X, Y)$ by

$$(4.3.1) \quad f_{U,V}(u, v) = P(U = u, V = v) = P((X, Y) \in A_{uv}) = \sum_{(x,y) \in A_{uv}} f_{X,Y}(x, y).$$

If $(X, Y)$ is a continuous random vector with joint pdf $f_{X,Y}(x, y)$, then the joint pdf of $(U, V)$ can be expressed in terms of $f_{X,Y}(x, y)$ in a manner analogous to (2.1.8). As before, $\mathcal{A} = \{(x, y) : f_{X,Y}(x, y) > 0\}$ and $\mathcal{B} = \{(u, v) : u = g_1(x, y) \text{ and } v = g_2(x, y) \text{ for some } (x, y) \in \mathcal{A}\}$. The joint pdf $f_{U,V}(u, v)$ will be positive on the set $\mathcal{B}$. For the simplest version of this result we assume that the transformation $u = g_1(x, y)$ and $v = g_2(x, y)$ defines a one-to-one transformation of $\mathcal{A}$ onto $\mathcal{B}$. The transformation is onto because of the definition of $\mathcal{B}$. We are assuming that for each $(u, v) \in \mathcal{B}$ there is only one $(x, y) \in \mathcal{A}$ such that $(u, v) = (g_1(x, y), g_2(x, y))$. For such a one-to-one, onto transformation, we can solve the equations $u = g_1(x, y)$ and $v = g_2(x, y)$ for $x$ and $y$ in terms of $u$ and $v$. We will denote this inverse transformation by $x = h_1(u, v)$ and $y = h_2(u, v)$. The role played by a derivative in the univariate case is now played by a quantity called the *Jacobian of the transformation*. This function of $(u, v)$, denoted by $J$, is the *determinant of a matrix* of partial derivatives. It is defined by

$$J = \begin{vmatrix} \dfrac{\partial x}{\partial u} & \dfrac{\partial x}{\partial v} \\ \dfrac{\partial y}{\partial u} & \dfrac{\partial y}{\partial v} \end{vmatrix} = \frac{\partial x}{\partial u}\frac{\partial y}{\partial v} - \frac{\partial y}{\partial u}\frac{\partial x}{\partial v},$$

where

$$\frac{\partial x}{\partial u} = \frac{\partial h_1(u, v)}{\partial u}, \quad \frac{\partial x}{\partial v} = \frac{\partial h_1(u, v)}{\partial v}, \quad \frac{\partial y}{\partial u} = \frac{\partial h_2(u, v)}{\partial u}, \quad \text{and} \quad \frac{\partial y}{\partial v} = \frac{\partial h_2(u, v)}{\partial v}.$$

We assume that $J$ is not identically 0 on $\mathcal{B}$. Then the joint pdf of $(U, V)$ is 0 outside the set $\mathcal{B}$ and on the set $\mathcal{B}$ is given by

$$(4.3.2) \qquad f_{U,V}(u, v) = f_{X,Y}(h_1(u, v), h_2(u, v))|J|,$$

**Lemma 4.2.7** *Let $(X, Y)$ be a bivariate random vector with joint pdf or pmf $f(x, y)$. Then $X$ and $Y$ are independent random variables if and only if there exist functions $g(x)$ and $h(y)$ such that, for every $x \in \Re$ and $y \in \Re$,*

$$f(x, y) = g(x)h(y).$$

**Theorem 4.2.12** *Let $X$ and $Y$ be independent random variables with moment generating functions $M_X(t)$ and $M_Y(t)$. Then the moment generating function of the random variable $Z = X + Y$ is given by*

$$M_Z(t) = M_X(t)M_Y(t).$$

**Theorem 4.3.5** *Let $X$ and $Y$ be independent random variables. Let $g(x)$ be a function only of $x$ and $h(y)$ be a function only of $y$. Then the random variables $U = g(X)$ and $V = h(Y)$ are independent.*

**Theorem 4.4.3** *If $X$ and $Y$ are any two random variables, then*

$$(4.4.1) \qquad \mathrm{E} X = \mathrm{E}\left(\mathrm{E}(X|Y)\right),$$

*provided that the expectations exist.*

**Theorem 4.4.7 (Conditional variance identity)** *For any two random variables $X$ and $Y$,*

$$(4.4.4) \qquad \mathrm{Var}\, X = \mathrm{E}\left(\mathrm{Var}(X|Y)\right) + \mathrm{Var}\left(\mathrm{E}(X|Y)\right),$$

*provided that the expectations exist.*

**Definition 4.5.1** The *covariance of $X$ and $Y$* is the number defined by

$$\mathrm{Cov}(X, Y) = \mathrm{E}\left((X - \mu_X)(Y - \mu_Y)\right).$$

**Theorem 4.5.5** *If $X$ and $Y$ are independent random variables, then $\mathrm{Cov}(X, Y) = 0$ and $\rho_{XY} = 0$.*

**Theorem 4.5.6** *If $X$ and $Y$ are any two random variables and $a$ and $b$ are any two constants, then*

$$\mathrm{Var}(aX + bY) = a^2 \mathrm{Var}\, X + b^2 \mathrm{Var}\, Y + 2ab\, \mathrm{Cov}(X, Y).$$

*If $X$ and $Y$ are independent random variables, then*

$$\mathrm{Var}(aX + bY) = a^2 \mathrm{Var}\, X + b^2 \mathrm{Var}\, Y.$$

**Theorem 4.7.7 (Jensen's Inequality)** *For any random variable $X$, if $g(x)$ is a convex function, then*

$$\mathrm{E} g(X) \geq g(\mathrm{E} X).$$

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

# Chapter 5: Multiple Random Variables

**Theorem 5.2.6** *Let $X_1, \ldots, X_n$ be a random sample from a population with mean $\mu$ and variance $\sigma^2 < \infty$. Then*

**a.** $\mathrm{E} \bar{X} = \mu$,

**b.** $\mathrm{Var}\, \bar{X} = \dfrac{\sigma^2}{n}$,

**c.** $\mathrm{E} S^2 = \sigma^2$.

**Theorem 5.2.7** *Let $X_1, \ldots, X_n$ be a random sample from a population with mgf $M_X(t)$. Then the mgf of the sample mean is*

$$M_{\bar{X}}(t) = [M_X(t/n)]^n.$$

**Theorem 5.2.9** *If $X$ and $Y$ are independent continuous random variables with pdfs $f_X(x)$ and $f_Y(y)$, then the pdf of $Z = X + Y$ is*

$$(5.2.3) \qquad f_Z(z) = \int_{-\infty}^{\infty} f_X(w) f_Y(z - w)\, dw.$$

**Theorem 5.4.4** *Let* $X_{(1)}, \ldots, X_{(n)}$ *denote the order statistics of a random sample,* $X_1, \ldots, X_n$, *from a continuous population with cdf* $F_X(x)$ *and pdf* $f_X(x)$. *Then the pdf of* $X_{(j)}$ *is*

(5.4.4)
$$f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f_X(x)[F_X(x)]^{j-1}[1-F_X(x)]^{n-j}.$$

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

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

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

**Definition 5.5.1** A sequence of random variables, $X_1, X_2, \ldots$, *converges in probability* to a random variable $X$ if, for every $\epsilon > 0$,

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

**Theorem 5.5.2 (Weak Law of Large Numbers)** *Let* $X_1, X_2, \ldots$ *be iid random variables with* $\mathrm{E}X_i = \mu$ *and* $\mathrm{Var}\, X_i = \sigma^2 < \infty$. *Define* $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. *Then,*

*for every* $\epsilon > 0$,

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

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

**Theorem 5.5.4** *Suppose that* $X_1, X_2, \ldots$ *converges in probability to a random variable* $X$ *and that* $h$ *is a continuous function. Then* $h(X_1), h(X_2), \ldots$ *converges in probability to* $h(X)$.

**Definition 5.5.6** A sequence of random variables, $X_1, X_2, \ldots$, *converges almost surely* to a random variable $X$ if, for every $\epsilon > 0$,

$$P(\lim_{n \to \infty} |X_n - X| < \epsilon) = 1.$$

**Theorem 5.5.9 (Strong Law of Large Numbers)** *Let* $X_1, X_2, \ldots$ *be iid random variables with* $\mathrm{E}X_i = \mu$ *and* $\mathrm{Var}\, X_i = \sigma^2 < \infty$, *and define* $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. *Then, for every* $\epsilon > 0$,

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

*that is,* $\bar{X}_n$ *converges almost surely to* $\mu$.

**Definition 5.5.10** A sequence of random variables, $X_1, X_2, \ldots$, *converges in distribution* to a random variable $X$ if

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

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

**Theorem 5.5.12** *If the sequence of random variables,* $X_1, X_2, \ldots$, *converges in probability to a random variable* $X$, *the sequence also converges in distribution to* $X$.

**Theorem 5.5.13** *The sequence of random variables,* $X_1, X_2, \ldots$, *converges in probability to a constant* $\mu$ *if and only if the sequence also converges in distribution to* $\mu$. *That is, the statement*

$$P(|X_n - \mu| > \varepsilon) \to 0 \text{ for every } \varepsilon > 0$$

*is equivalent to*

$$P(X_n \le x) \to \begin{cases} 0 & \text{if } x < \mu \\ 1 & \text{if } x > \mu. \end{cases}$$

**Theorem 5.5.15 (Stronger form of the Central Limit Theorem)** *Let* $X_1, X_2, \ldots$ *be a sequence of iid random variables with* $\mathrm{E}X_i = \mu$ *and* $0 < \mathrm{Var}\, X_i = \sigma^2 < \infty$. *Define* $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. *Let* $G_n(x)$ *denote the cdf of* $\sqrt{n}(\bar{X}_n - \mu)/\sigma$. *Then, for any* $x$, $-\infty < x < \infty$,

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

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

**Theorem 5.5.17 (Slutsky's Theorem)** *If* $X_n \to X$ *in distribution and* $Y_n \to a$, *a constant, in probability, then*

**a.** $Y_n X_n \to aX$ *in distribution.*

**b.** $X_n + Y_n \to X + a$ *in distribution.*

**Theorem 5.5.24 (Delta Method)** *Let $Y_n$ be a sequence of random variables that satisfies $\sqrt{n}(Y_n - \theta) \to \mathrm{n}(0, \sigma^2)$ in distribution. For a given function $g$ and a specific value of $\theta$, suppose that $g'(\theta)$ exists and is not 0. Then*

$$(5.5.10) \qquad \sqrt{n}[g(Y_n) - g(\theta)] \to \mathrm{n}(0, \sigma^2[g'(\theta)]^2) \text{ in distribution.}$$

**Theorem 5.5.26 (Second-order Delta Method)** *Let $Y_n$ be a sequence of random variables that satisfies $\sqrt{n}(Y_n - \theta) \to \mathrm{n}(0, \sigma^2)$ in distribution. For a given function $g$ and a specific value of $\theta$, suppose that $g'(\theta) = 0$ and $g''(\theta)$ exists and is not 0. Then*

$$(5.5.13) \qquad n[g(Y_n) - g(\theta)] \to \sigma^2 \frac{g''(\theta)}{2} \chi_1^2 \text{ in distribution.}$$

# Chapter 6: Principles of Data Reduction

**Definition 6.2.1** A statistic $T(\mathbf{X})$ is a *sufficient statistic for $\theta$* if the conditional distribution of the sample $\mathbf{X}$ given the value of $T(\mathbf{X})$ does not depend on $\theta$.

**Theorem 6.2.2** *If $p(\mathbf{x}|\theta)$ is the joint pdf or pmf of $\mathbf{X}$ and $q(t|\theta)$ is the pdf or pmf of $T(\mathbf{X})$, then $T(\mathbf{X})$ is a sufficient statistic for $\theta$ if, for every $\mathbf{x}$ in the sample space, the ratio $p(\mathbf{x}|\theta)/q(T(\mathbf{x})|\theta)$ is constant as a function of $\theta$.*

**Theorem 6.2.6 (Factorization Theorem)** *Let $f(\mathbf{x}|\theta)$ denote the joint pdf or pmf of a sample $\mathbf{X}$. A statistic $T(\mathbf{X})$ is a sufficient statistic for $\theta$ if and only if there exist functions $g(t|\theta)$ and $h(\mathbf{x})$ such that, for all sample points $\mathbf{x}$ and all parameter points $\theta$,*

$$(6.2.3) \qquad f(\mathbf{x}|\theta) = g(T(\mathbf{x})|\theta)h(\mathbf{x}).$$

**Theorem 6.2.10** *Let $X_1, \ldots, X_n$ be iid observations from a pdf or pmf $f(x|\boldsymbol{\theta})$ that belongs to an exponential family given by*

$$f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp\left( \sum_{i=1}^{k} w_i(\boldsymbol{\theta})t_i(x) \right),$$

*where $\boldsymbol{\theta} = (\theta_1, \theta_2, \ldots, \theta_d)$, $d \leq k$. Then*

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

*is a sufficient statistic for $\boldsymbol{\theta}$.*

**Definition 6.2.11** A sufficient statistic $T(\mathbf{X})$ is called a *minimal sufficient statistic* if, for any other sufficient statistic $T'(\mathbf{X})$, $T(\mathbf{x})$ is a function of $T'(\mathbf{x})$.

**Theorem 6.2.13** *Let $f(\mathbf{x}|\theta)$ be the pmf or pdf of a sample $\mathbf{X}$. Suppose there exists a function $T(\mathbf{x})$ such that, for every two sample points $\mathbf{x}$ and $\mathbf{y}$, the ratio $f(\mathbf{x}|\theta)/f(\mathbf{y}|\theta)$ is constant as a function of $\theta$ if and only if $T(\mathbf{x}) = T(\mathbf{y})$. Then $T(\mathbf{X})$ is a minimal sufficient statistic for $\theta$.*

**Definition 6.2.21** Let $f(t|\theta)$ be a family of pdfs or pmfs for a statistic $T(\mathbf{X})$. The family of probability distributions is called *complete* if $\mathrm{E}_\theta g(T) = 0$ for all $\theta$ implies $P_\theta(g(T) = 0) = 1$ for all $\theta$. Equivalently, $T(\mathbf{X})$ is called a *complete statistic*.

**Theorem 6.2.25 (Complete statistics in the exponential family)** *Let $X_1, \ldots, X_n$ be iid observations from an exponential family with pdf or pmf of the form*

$$(6.2.7) \qquad f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp\left( \sum_{j=1}^{k} w(\theta_j)t_j(x) \right),$$

*where $\boldsymbol{\theta} = (\theta_1, \theta_2, \ldots, \theta_k)$. Then the statistic*

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

*is complete as long as the parameter space $\Theta$ contains an open set in $\Re^k$.*

# Chapter 7: Point Estimation

**Theorem 7.2.10 (Invariance property of MLEs)** *If $\hat{\theta}$ is the MLE of $\theta$, then for any function $\tau(\theta)$, the MLE of $\tau(\theta)$ is $\tau(\hat{\theta})$.*

**Definition 7.3.7** An estimator $W^*$ is a *best unbiased estimator* of $\tau(\theta)$ if it satisfies $\mathrm{E}_\theta W^* = \tau(\theta)$ for all $\theta$ and, for any other estimator $W$ with $\mathrm{E}_\theta W = \tau(\theta)$, we have $\mathrm{Var}_\theta W^* \leq \mathrm{Var}_\theta W$ for all $\theta$. $W^*$ is also called a *uniform minimum variance unbiased estimator* (UMVUE) of $\tau(\theta)$.

**Theorem 7.3.9 (Cramér–Rao Inequality)** *Let* $X_1, \ldots, X_n$ *be a sample with pdf* $f(\mathbf{x}|\theta)$, *and let* $W(\mathbf{X}) = W(X_1, \ldots, X_n)$ *be any estimator satisfying*

$$\frac{d}{d\theta} \mathrm{E}_\theta W(\mathbf{X}) = \int_{\mathcal{X}} \frac{\partial}{\partial\theta} \left[ W(\mathbf{x}) f(\mathbf{x}|\theta) \right] \, d\mathbf{x}$$

(7.3.4)        *and*

$$\mathrm{Var}_\theta W(\mathbf{X}) < \infty.$$

*Then*

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

**Corollary 7.3.10 (Cramér–Rao Inequality, iid case)** *If the assumptions of Theorem 7.3.9 are satisfied and, additionally, if* $X_1, \ldots, X_n$ *are iid with pdf* $f(x|\theta)$, *then*

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

**Lemma 7.3.11** *If* $f(x|\theta)$ *satisfies*

$$\frac{d}{d\theta} \mathrm{E}_\theta \left( \frac{\partial}{\partial\theta} \log f(X|\theta) \right) = \int \frac{\partial}{\partial\theta} \left[ \left( \frac{\partial}{\partial\theta} \log f(x|\theta) \right) f(x|\theta) \right] \, dx$$

*(true for an exponential family), then*

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

**Corollary 7.3.15 (Attainment)** *Let* $X_1, \ldots, X_n$ *be iid* $f(x|\theta)$, *where* $f(x|\theta)$ *satisfies the conditions of the Cramér–Rao Theorem. Let* $L(\theta|\mathbf{x}) = \prod_{i=1}^{n} f(x_i|\theta)$ *denote the likelihood function. If* $W(\mathbf{X}) = W(X_1, \ldots, X_n)$ *is any unbiased estimator of* $\tau(\theta)$, *then* $W(\mathbf{X})$ *attains the Cramér–Rao Lower Bound if and only if*

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

*for some function* $a(\theta)$.

**Theorem 7.3.17 (Rao–Blackwell)** *Let* $W$ *be any unbiased estimator of* $\tau(\theta)$, *and let* $T$ *be a sufficient statistic for* $\theta$. *Define* $\phi(T) = \mathrm{E}(W|T)$. *Then* $\mathrm{E}_\theta \phi(T) = \tau(\theta)$ *and* $\mathrm{Var}_\theta \phi(T) \leq \mathrm{Var}_\theta W$ *for all* $\theta$; *that is,* $\phi(T)$ *is a uniformly better unbiased estimator of* $\tau(\theta)$.

**Theorem 7.3.19** *If* $W$ *is a best unbiased estimator of* $\tau(\theta)$, *then* $W$ *is unique.*

**Theorem 7.3.20** *If* $\mathrm{E}_\theta W = \tau(\theta), W$ *is the best unbiased estimator of* $\tau(\theta)$ *if and only if* $W$ *is uncorrelated with all unbiased estimators of* $0$.

**Theorem 7.3.23** *Let* $T$ *be a complete sufficient statistic for a parameter* $\theta$, *and let* $\phi(T)$ *be any estimator based only on* $T$. *Then* $\phi(T)$ *is the unique best unbiased estimator of its expected value.*

# Chapter 8: Hypothesis Testing

**Definition 8.2.1** The *likelihood ratio test statistic* for testing $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_0^c$ is

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

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

**Theorem 8.2.4** *If* $T(\mathbf{X})$ *is a sufficient statistic for* $\theta$ *and* $\lambda^*(t)$ *and* $\lambda(\mathbf{x})$ *are the LRT statistics based on* $T$ *and* $\mathbf{X}$, *respectively, then* $\lambda^*(T(\mathbf{x})) = \lambda(\mathbf{x})$ *for every* $\mathbf{x}$ *in the sample space.*

**Definition 8.3.1** The *power function* of a hypothesis test with rejection region $R$ is the function of $\theta$ defined by $\beta(\theta) = P_\theta(\mathbf{X} \in R)$.

**Definition 8.3.5** For $0 \leq \alpha \leq 1$, a test with power function $\beta(\theta)$ is a *size $\alpha$ test* if $\sup_{\theta \in \Theta_0} \beta(\theta) = \alpha$.

**Definition 8.3.6** For $0 \leq \alpha \leq 1$, a test with power function $\beta(\theta)$ is a *level $\alpha$ test* if $\sup_{\theta \in \Theta_0} \beta(\theta) \leq \alpha$.

**Definition 8.3.9** A test with power function $\beta(\theta)$ is *unbiased* if $\beta(\theta') \geq \beta(\theta'')$ for every $\theta' \in \Theta_0^c$ and $\theta'' \in \Theta_0$.

**Definition 8.3.11** Let $\mathcal{C}$ be a class of tests for testing $H_0: \theta \in \Theta_0$ versus $H_1: \theta \in \Theta_0^c$. A test in class $\mathcal{C}$, with power function $\beta(\theta)$, is a *uniformly most powerful* (UMP) *class $\mathcal{C}$ test* if $\beta(\theta) \geq \beta'(\theta)$ for every $\theta \in \Theta_0^c$ and every $\beta'(\theta)$ that is a power function of a test in class $\mathcal{C}$.

**Theorem 8.3.12 (Neyman–Pearson Lemma)** *Consider testing $H_0: \theta = \theta_0$ versus $H_1: \theta = \theta_1$, where the pdf or pmf corresponding to $\theta_i$ is $f(\mathbf{x}|\theta_i), i = 0, 1$, using a test with rejection region $R$ that satisfies*

$$\mathbf{x} \in R \quad \text{if} \quad f(\mathbf{x}|\theta_1) > kf(\mathbf{x}|\theta_0)$$

(8.3.1) *and*

$$\mathbf{x} \in R^c \quad \text{if} \quad f(\mathbf{x}|\theta_1) < kf(\mathbf{x}|\theta_0),$$

*for some $k \geq 0$, and*

(8.3.2) $$\alpha = P_{\theta_0}(\mathbf{X} \in R).$$

*Then*

**a.** *(Sufficiency) Any test that satisfies (8.3.1) and (8.3.2) is a UMP level $\alpha$ test.*

**b.** *(Necessity) If there exists a test satisfying (8.3.1) and (8.3.2) with $k > 0$, then every UMP level $\alpha$ test is a size $\alpha$ test (satisfies (8.3.2)) and every UMP level $\alpha$ test satisfies (8.3.1) except perhaps on a set $A$ satisfying $P_{\theta_0}(\mathbf{X} \in A) = P_{\theta_1}(\mathbf{X} \in A) = 0$.*

**Corollary 8.3.13** *Consider the hypothesis problem posed in Theorem 8.3.12. Suppose $T(\mathbf{X})$ is a sufficient statistic for $\theta$ and $g(t|\theta_i)$ is the pdf or pmf of $T$ corresponding to $\theta_i$, $i = 0, 1$. Then any test based on $T$ with rejection region $S$ (a subset of the sample space of $T$) is a UMP level $\alpha$ test if it satisfies*

$$t \in S \quad \text{if} \quad g(t|\theta_1) > kg(t|\theta_0)$$

(8.3.4) *and*

$$t \in S^c \quad \text{if} \quad g(t|\theta_1) < kg(t|\theta_0),$$

*for some $k \geq 0$, where*

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

**Definition 8.3.16** A family of pdfs or pmfs $\{g(t|\theta): \theta \in \Theta\}$ for a univariate random variable $T$ with real-valued parameter $\theta$ has a *monotone likelihood ratio* (MLR) if, for every $\theta_2 > \theta_1$, $g(t|\theta_2)/g(t|\theta_1)$ is a monotone (nonincreasing or nondecreasing) function of $t$ on $\{t: g(t|\theta_1) > 0 \text{ or } g(t|\theta_2) > 0\}$. Note that $c/0$ is defined as $\infty$ if $0 < c$.

**Theorem 8.3.17 (Karlin–Rubin)** *Consider testing $H_0: \theta \leq \theta_0$ versus $H_1: \theta > \theta_0$. Suppose that $T$ is a sufficient statistic for $\theta$ and the family of pdfs or pmfs $\{g(t|\theta): \theta \in \Theta\}$ of $T$ has an MLR.* Then for any $t_0$, the test that rejects $H_0$ if and only if $T > t_0$ is a UMP level $\alpha$ test, where $\alpha = P_{\theta_0}(T > t_0)$. *Assumes nondecreasing LR.

# Chapter 9: Interval Estimation

**Definition 9.1.1** An *interval estimate* of a real-valued parameter $\theta$ is any pair of functions, $L(x_1, \ldots, x_n)$ and $U(x_1, \ldots, x_n)$, of a sample that satisfy $L(\mathbf{x}) \leq U(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$. If $\mathbf{X} = \mathbf{x}$ is observed, the inference $L(\mathbf{x}) \leq \theta \leq U(\mathbf{x})$ is made. The random interval $[L(\mathbf{X}), U(\mathbf{X})]$ is called an *interval estimator*.

**Definition 9.1.4** For an interval estimator $[L(\mathbf{X}), U(\mathbf{X})]$ of a parameter $\theta$, the *coverage probability* of $[L(\mathbf{X}), U(\mathbf{X})]$ is the probability that the random interval $[L(\mathbf{X}), U(\mathbf{X})]$ covers the true parameter, $\theta$. In symbols, it is denoted by either $P_\theta(\theta \in [L(\mathbf{X}), U(\mathbf{X})])$ or $P(\theta \in [L(\mathbf{X}), U(\mathbf{X})]|\theta)$.

**Definition 9.1.5** For an interval estimator $[L(\mathbf{X}), U(\mathbf{X})]$ of a parameter $\theta$, the *confidence coefficient* of $[L(\mathbf{X}), U(\mathbf{X})]$ is the infimum of the coverage probabilities, $\inf_\theta P_\theta(\theta \in [L(\mathbf{X}), U(\mathbf{X})])$.

**Theorem 9.2.2** *For each $\theta_0 \in \Theta$, let $A(\theta_0)$ be the acceptance region of a level $\alpha$ test of $H_0: \theta = \theta_0$. For each $\mathbf{x} \in \mathcal{X}$, define a set $C(\mathbf{x})$ in the parameter space by*

(9.2.1) $$C(\mathbf{x}) = \{\theta_0: \mathbf{x} \in A(\theta_0)\}.$$

*Then the random set $C(\mathbf{X})$ is a $1 - \alpha$ confidence set. Conversely, let $C(\mathbf{X})$ be a $1 - \alpha$ confidence set. For any $\theta_0 \in \Theta$, define*

$$A(\theta_0) = \{\mathbf{x}: \theta_0 \in C(\mathbf{x})\}.$$

*Then $A(\theta_0)$ is the acceptance region of a level $\alpha$ test of $H_0: \theta = \theta_0$.*

# Chapter 10: Asymptotic Evaluations

**Definition 10.1.1**  A sequence of estimators $W_n = W_n(X_1, \ldots, X_n)$ is a *consistent sequence of estimators* of the parameter $\theta$ if, for every $\epsilon > 0$ and every $\theta \in \Theta$,

(10.1.1) $$\lim_{n \to \infty} P_\theta(|W_n - \theta| < \epsilon) = 1.$$

**Theorem 10.1.3**  *If $W_n$ is a sequence of estimators of a parameter $\theta$ satisfying*

  i. $\lim_{n \to \infty} \text{Var}_\theta W_n = 0$,

  ii. $\lim_{n \to \infty} \text{Bias}_\theta W_n = 0$,

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

**Definition 10.1.9**  For an estimator $T_n$, suppose that $k_n(T_n - \tau(\theta)) \to \text{n}(0, \sigma^2)$ in distribution. The parameter $\sigma^2$ is called the *asymptotic variance* or *variance of the limit distribution* of $T_n$.

**Definition 10.1.11**  A sequence of estimators $W_n$ is *asymptotically efficient* for a parameter $\tau(\theta)$ if $\sqrt{n}[W_n - \tau(\theta)] \to \text{n}[0, v(\theta)]$ in distribution and

$$v(\theta) = \frac{[\tau'(\theta)]^2}{\text{E}_\theta \left( \left( \frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right)};$$

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

**Theorem 10.1.12 (Asymptotic efficiency of MLEs)**  *Let $X_1, X_2, \ldots,$ be iid $f(x|\theta)$, let $\hat{\theta}$ denote the MLE of $\theta$, and let $\tau(\theta)$ be a continuous function of $\theta$. Under the regularity conditions in Miscellanea 10.6.2 on $f(x|\theta)$ and, hence, $L(\theta|\mathbf{x})$,*

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

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

**Definition 10.1.16**  If two estimators $W_n$ and $V_n$ satisfy

$$\sqrt{n}[W_n - \tau(\theta)] \to \text{n}[0, \sigma_W^2]$$
$$\sqrt{n}[V_n - \tau(\theta)] \to \text{n}[0, \sigma_V^2]$$

in distribution, the *asymptotic relative efficiency* (ARE) of $V_n$ with respect to $W_n$ is

$$\text{ARE}(V_n, W_n) = \frac{\sigma_W^2}{\sigma_V^2}.$$

**Theorem 10.3.1 (Asymptotic distribution of the LRT—simple $H_0$)**  *For testing $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$, suppose $X_1, \ldots, X_n$ are iid $f(x|\theta)$, $\hat{\theta}$ is the MLE of $\theta$, and $f(x|\theta)$ satisfies the regularity conditions in Miscellanea 10.6.2. Then under $H_0$, as $n \to \infty$,*

$$-2 \log \lambda(\mathbf{X}) \to \chi_1^2 \text{ in distribution,}$$

*where $\chi_1^2$ is a $\chi^2$ random variable with 1 degree of freedom.*

**Theorem 10.3.3**  *Let $X_1, \ldots, X_n$ be a random sample from a pdf or pmf $f(x|\theta)$. Under the regularity conditions in Miscellanea 10.6.2, if $\theta \in \Theta_0$, then the distribution of the statistic $-2 \log \lambda(\mathbf{X})$ converges to a chi squared distribution as the sample size $n \to \infty$. The degrees of freedom of the limiting distribution is the difference between the number of free parameters specified by $\theta \in \Theta_0$ and the number of free parameters specified by $\theta \in \Theta$.*