



THE UNIVERSITY OF QUEENSLAND
A U S T R A L I A

COURSE NOTES FOR STAT3001
MATHEMATICAL STATISTICS

CONTRIBUTORS:

MICHAEL CICCOTOSTO-CAMP
NAME2

THE UNIVERSITY OF QUEENSLAND
SCHOOL OF MATHEMATICS AND PHYSICS

CONTENTS

SYMBOLS AND NOTATION	iii
REVIEW	1
USEFUL FORMULAE AND THEOREMS	1
COMMON DISTRIBUTIONS	2
COMMON PROBABILISTIC PROPERTIES AND IDENTITIES	3
PROBABILISTIC PROPERTIES	3
PROBABILISTIC IDENTITIES	4
POINT ESTIMATION	5
METHODS OF FINDING ESTIMATES	5
REFERENCES	8

SYMBOLS AND NOTATION

Matrices are capitalized bold face letters while vectors are lowercase bold face letters.

<i>Syntax</i>	<i>Meaning</i>
\triangleq	An equality which acts as a statement
$ \mathbf{A} $	The determinate of a matrix.
$\mathbf{x}^\top, \mathbf{X}^\top$	The transpose operator.
$\mathbf{x}^*, \mathbf{X}^*$	The hermitian operator.
$\mathbf{a}.*\mathbf{b}$ or $\mathbf{A}.*\mathbf{B}$	Element-wise vector (matrix) multiplication, similar to Matlab.
\propto	Proportional to.
∇ or $\nabla_{\mathbf{f}}$	The partial derivative (with respect to \mathbf{f}).
$\nabla\nabla$ or $H(f)$	The Hessian.
\sim	Distributed according to, example $X \sim \mathcal{N}(0, 1)$
$\overset{\text{iid}}{\sim}$	Identically and independently distributed according to, example $X_1, X_2, \dots, X_n \overset{\text{iid}}{\sim} \mathcal{N}(0, 1)$
$\mathbf{0}$ or $\mathbf{0}_n$ or $\mathbf{0}_{n \times m}$	The zero vector (matrix) of appropriate length (size) or the zero vector of length n or the zero matrix with dimensions $n \times m$.
$\mathbf{1}$ or $\mathbf{1}_n$ or $\mathbf{1}_{n \times m}$	The one vector (matrix) of appropriate length (size) or the one vector of length n or the one matrix with dimensions $n \times m$.
$\mathbb{1}_A(x)$	The indicator function. $\mathbb{1}_A(x) = 1$ if $x \in A$, 0 otherwise.

$\mathbf{A}_{(:,)}$	Index slicing to extract a submatrix from the elements of $\mathbf{A} \in \mathbb{R}^{n \times m}$, similar to indexing slicing from the python and Matlab programming languages. Each parameter can receive a single value or a 'slice' consisting of a start and an end value separated by a semicolon. The first and second parameter describe what row and columns should be selected, respectively. A single value means that only values from the single specified row/column should be selected. A slice tells us that all rows/columns between the provided range should be selected. Additionally if now start and end values are specified in the slice then all rows/columns should be selected. For example, the slice $\mathbf{A}_{(1:3,j:j')}$ is the submatrix $\mathbb{R}^{3 \times (j'-j+1)}$ matrix containing the first three rows of \mathbf{A} and columns j to j' . As another example, $\mathbf{A}_{(:,j)}$ is the j^{th} column of \mathbf{A} .
\mathbf{A}^\dagger	Denotes the unique psuedo inverse or Moore-Penore inverse of \mathbf{A} .
\mathbb{C}	The complex numbers.
$\text{diag}(\mathbf{w})$	Vector argument, a diagonal matrix containing the elements of vector \mathbf{w} .
$\text{diag}(\mathbf{W})$	Matrix argument, a vector containing the diagonal elements of the matrix \mathbf{W} .
\mathbb{E} or $\mathbb{E}_{q(x)}[z(x)]$	Expectation, or expectation of $z(x)$ where $x \sim q(x)$.
\mathbb{R}	The real numbers.
$\text{tr}(\mathbf{A})$	The trace of a matrix.
\mathbb{V} or $\mathbb{V}_{q(x)}[z(x)]$	Variance, the variance of $z(x)$ when $x \sim q(x)$.
\mathbb{Z}	The integers, $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$.
Ω	The sample space.

REVIEW

Theorems and definitions here are mostly concepts seen before from other courses.

Useful Formulae and Theorems.

(Geometric Series)
$$\sum_{k=0}^{n-1} r^k = \left(\frac{1 - r^n}{1 - r} \right)$$

or

$$\sum_{i=0}^{\infty} r^i = \frac{1}{1 - r} \quad \text{with} \quad |r| < 1$$

(Euler's formula)
$$e^{ix} = \cos x + i \sin x$$

(Newton's Binomial formula)
$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k$$

Theorem 1 (Young's inequality for products). *If $a \geq 0$ and $b \geq 0$ are nonnegative real numbers and if $p > 1$ and $q > 1$ are real numbers such that $\frac{1}{p} + \frac{1}{q} = 1$, then*

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}.$$

Equality holds iff $a^p = b^q$.

Common Distributions. Common distributions seen from prior courses. Notations mostly borrowed from STAT2003.

<i>Name</i>	<i>Notation</i>	<i>Support</i>	<i>pf</i>	<i>Expectation</i>	<i>Variance</i>
Bernoulli	$\text{Ber}(p)$	$\{0, 1\}$	$p^k(1-p)^{1-k}$	p	$p(1-p)$
Binomial	$\text{Bin}(n, p)$	$\{0, \dots, n\}$	$\binom{n}{k} p^k (1-p)^{n-k}$	np	$np(1-p)$
Negative-Binomial	$\text{NB}(r, p)$	\mathbb{N}_0	$\binom{x+r-1}{x} p^x (1-p)^r$	$\frac{rp}{1-p}$	$\frac{rp}{(1-p)^2}$
Geometric	$\text{Geo}(n, p)$	\mathbb{N}_0	$(1-p)^k p$	$\frac{1-p}{p}$	$\frac{1-p}{p^2}$
Poisson	$\text{Poi}(\lambda)$	\mathbb{N}_0	$\frac{\lambda^x}{x!} e^{-\lambda}$	λ	λ
Uniform	$\text{U}[a, b]$	$[a, b]$	$\frac{1}{b-a}$	$\frac{a+b}{2}$	$\frac{(a-b)^2}{12}$
Exponential	$\text{Exp}(\lambda)$	\mathbb{R}^+	$\lambda e^{-\lambda x}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
Normal	$\text{N}(\mu, \sigma^2)$	\mathbb{R}	$\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$	μ	σ^2
Gamma	$\text{Gam}(\alpha, \lambda)$	\mathbb{R}^+	$\frac{\lambda^\alpha x^{\alpha-1} \exp(-\lambda x)}{\Gamma(\alpha)}$	$\frac{\alpha}{\lambda}$	$\frac{\alpha}{\lambda^2}$
Chi-Squared	χ_n^2	\mathbb{R}^+	$\frac{x^{\frac{n}{2}-1} \exp(-\frac{1}{2}x)}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}$	n	$2n$
White-Noise	$\text{WN}(\mu, \sigma^2)$	NA	NA	μ	σ^2

Common Probabilistic Properties and Identities. Common probabilistic properties seen from prior courses.

Probabilistic Properties. For any random variables, the following hold.

- (1) $\mathbb{E}(X) = \int_0^\infty (1 - F(X)) \, dx$
- (2) $\mathbb{E}(aX + b) = a\mathbb{E}X + b$
- (3) $\mathbb{E}(g(X) + h(X)) = \mathbb{E}g(X) + \mathbb{E}h(X)$
- (4) $\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2$
- (5) $\text{Var}(aX + b) = a^2\text{Var}(X)$
- (6) $\text{Cov}(X, Y) = \mathbb{E}XY - \mathbb{E}X\mathbb{E}Y$
- (7) $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$
- (8) $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]]$
- (9) $\text{Var}(Y) = \mathbb{E}[\text{Var}(Y|X)] + \text{Var}(\mathbb{E}[Y|X])$
- (10) $|\mathbb{E}(XY)|^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$
- (11) $|\text{Cov}(XY)|^2 \leq \text{Var}(X)\text{Var}(Y)$
- (12) $\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$
- (Bayes' Theorem) $\mathbb{P}(A | B) = \frac{\mathbb{P}(B | A)\mathbb{P}(A)}{\mathbb{P}(B)}$
- (13) $\mathbb{P}(A_1, \dots, A_n) = \mathbb{P}(A_1) \mathbb{P}(A_2 | A_1) \mathbb{P}(A_3 | A_1, A_2) \cdots \mathbb{P}(A_n | A_1, A_2, \dots, A_{n-1})$
- (14)

Let $\Omega = \bigcup_{i=1}^n B_i$ (that is B_i partitions the sample space) then

$$\text{(TLoP)} \quad \mathbb{P}(A) = \sum_{i=1}^n \mathbb{P}(A | B_i)\mathbb{P}(B_i)$$

$$\text{(TLoE)} \quad \mathbb{E}(A) = \sum_{i=1}^n \mathbb{E}(A | B_i)\mathbb{P}(B_i)$$

which, when **TLoP** used in conjunction with Bayes' Rule gives

$$(15) \quad \mathbb{P}(B_i | A) = \frac{\mathbb{P}(A | B_i)\mathbb{P}(B_i)}{\sum_{j=1}^n \mathbb{P}(A | B_j)\mathbb{P}(B_j)}.$$

If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{WN}(\mu, \sigma^2)$ and $S_n = \sum_{i=1}^n X_i$, then for all $\epsilon > 0$

$$\text{(Weak Law of Large Numbers)} \quad \mathbb{P}\left(\left|\frac{S_n}{n} - \mu\right| \geq \epsilon\right) = 0.$$

If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{WN}(\mu, \sigma^2)$ and $S_n = \sum_{i=1}^n X_i$, then for all $x \in \mathbb{R}$

(CLT)
$$\mathbb{P}\left(\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq x\right) = \Phi(x).$$

If X is a random variable and h is a convex function then

(Jensens Inequality)
$$h(\mathbb{E}(X)) \leq \mathbb{E}(h(X)).$$

Probabilistic Identities. If $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Ber}(p)$ then

(16)
$$\sum_{i=1}^n X_i \sim \text{Bin}(n, p).$$

If $X \sim \text{Bin}(n, p)$ and $Y \sim \text{Bin}(m, p)$, then $X + Y \sim \text{Bin}(n + m, p)$.

If $X \sim \text{N}(\mu_X, \sigma_X^2)$ and $Y \sim \text{N}(\mu_Y, \sigma_Y^2)$, then $X + Y \sim \text{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$.

If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma)$ then

(17)
$$\sum_{i=1}^n X_i^2 = \chi_n^2.$$

POINT ESTIMATION

Methods of Finding Estimates.

Definition 2 (Statistic). Let X_1, \dots, X_n be a random sample of size n from a population and let $T(x_1, \dots, x_n)$ be a real-valued or vector-valued function whose domain includes the sample space of (X_1, \dots, X_n) . The random variable or random vector $Y = T(X_1, \dots, X_n)$ is called a **statistic**. The probability distribution of a statistic Y is called the **sampling distribution** of Y [Cas01, page 211].

Definition 3 (Sample Mean). The **sample mean** is the arithmetic average of the values in a random sample. It is usually denoted by

$$(18) \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

[Cas01, page 212].

Definition 4 (Sample Variance and Standard Deviation). The **sample variance** is the statistic defined by

$$(19) \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

The **sample standard deviation** is the statistic defined by $S = \sqrt{S^2}$ [Cas01, page 212].

Definition 5 (Sufficient Statistic). A statistic $T(\mathbf{X})$ is a **sufficient statistic** for θ if the conditional distribution of the sample \mathbf{X} given the value of $T(\mathbf{X})$ does not depend on θ [Cas01, page 272].

Theorem 6. If $p(\mathbf{x} \mid \theta)$ is the joint pdf or pmf of \mathbf{X} and $q(\theta \mid \theta)$ is the pdf or pmf of $T(\mathbf{X})$, then $T(\mathbf{X})$ is a sufficient statistic for θ if, for every \mathbf{x} in the sample space, the ratio $p(\mathbf{x} \mid \theta)/q(T(\mathbf{x}) \mid \theta)$ is a constant function of θ [Cas01, page 274].

Theorem 7 (Factorization Theorem). Let $f(\mathbf{x} \mid \theta)$ denote the joint pdf or pmf of a sample \mathbf{X} . A statistic $T(\mathbf{X})$ is a sufficient statistic for θ , if and only if there exist function $g(t \mid \theta)$ and $h(\mathbf{x})$ such that, for all sample points \mathbf{x} and all parameter points θ ,

$$f(\mathbf{x} \mid \theta) = g(T(\mathbf{x}) \mid \theta)h(\mathbf{x})$$

[Cas01, page 276].

Example 8 (Uniform Sufficient Statistic). Example taken from [Cas01, page 277] and can also be found on tutorial sheet 3. Let X_1, \dots, X_n be iid observations from the discrete uniform distribution on $1, \dots, \theta$. That is, the unknown parameter, θ , is a positive integer and the pmf of X_i is

$$f(x \mid \theta) = \begin{cases} \frac{1}{\theta}, & x = 1, 2, \dots, \theta \\ 0, & \text{otherwise} \end{cases}.$$

The restriction $x_i \in \{1, \dots, \theta\}$ for $i = 1, \dots, n$ can be re-expressed as $x_i \in \{1, 2, \dots\}$ for $i = 1, \dots, n$ (note that there is no θ in this restriction) and $\max_i x_i \leq \theta$. If we define $T(\mathbf{x}) = \max_i x_i = x_{(n)}$,

$$h(\mathbf{x}) = \begin{cases} 1, & x_i \in \{1, \dots, \theta\} \text{ for } i = 1, \dots, n \\ 0, & \text{otherwise} \end{cases}$$

and

$$g(t \mid \theta) = \begin{cases} \theta^{-n} & t \leq \theta \\ 0, & \text{otherwise} \end{cases}.$$

It is easily verified that $f(\mathbf{x} \mid \theta) = g(T(\mathbf{x}) \mid \theta)$ for all \mathbf{x} and θ . Thus, according to Theorem 7, the largest order statistic, $T(\mathbf{X}) = X_{(n)}$, is a sufficient statistic in this problem. This type of analysis can sometimes be carried out more clearly and concisely using indicator function. Let \mathbb{N} be the set of natural numbers (discluding 0) and \mathbb{N}_θ be the natural numbers up to and including θ . Then the joint pmf of X_1, \dots, X_n is

$$f(\mathbf{x} \mid \theta) = \prod_{i=1}^n \theta^{-1} \mathbb{1}_{\mathbb{N}_\theta}(x_i) = \theta^{-n} \prod_{i=1}^n \mathbb{1}_{\mathbb{N}_\theta}(x_i).$$

Defining $T(\mathbf{x}) = x_{(n)}$, we see that

$$\prod_{i=1}^n \mathbb{1}_{\mathbb{N}_\theta}(x_i) = \left(\prod_{i=1}^n \mathbb{1}_{\mathbb{N}}(x_i) \right) \mathbb{1}_{\mathbb{N}_\theta}(T(\mathbf{x}))$$

thus providing the factorization

$$f(\mathbf{x} \mid \theta) = \theta^{-n} \mathbb{1}_{\mathbb{N}_\theta}(T(\mathbf{x})) \left(\prod_{i=1}^n \mathbb{1}_{\mathbb{N}}(x_i) \right).$$

The first factor depends on x_1, \dots, x_n only through the value of $T(\mathbf{x}) = x_{(n)}$, and the second factor does not depend on θ . Again, according to Theorem 7, $T(\mathbf{X}) = X_{(n)}$, is a sufficient statistic in this problem.

Definition 9 (Exponential Family). *In the case of p -dimensional observation $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{C}^p$, a d -dimensional parameter vector $\boldsymbol{\theta} \in \mathbb{C}^d$, and a q -dimensional sufficient statistic $T(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{C}^q$, the likelihood function $L(\boldsymbol{\theta})$ for the d -parameter vector $\boldsymbol{\theta}$ has the following form if it belongs to the d -parameter **exponential family***

$$L(\boldsymbol{\theta}) = b(\mathbf{x}_1, \dots, \mathbf{x}_n) \exp \{c(\boldsymbol{\theta})^\top T(\mathbf{x}_1, \dots, \mathbf{x}_n)\} / a(\boldsymbol{\theta})$$

where $c(\boldsymbol{\theta}) \in \mathbb{C}^q$ and $b(\mathbf{x}_1, \dots, \mathbf{x}_n)$ and $a(\boldsymbol{\theta})$ are scalar functions [Cas01, page 279].

Theorem 10. *Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be iid observations from a pdf or pmf $f(\mathbf{x} \mid \boldsymbol{\theta})$ that belongs to an exponential family as seen in Definition 9, then*

$$T(\mathbf{X}_1, \dots, \mathbf{X}_n) = \left(\sum_{j=1}^n t_1(\mathbf{X}_j), \dots, \sum_{j=1}^n t_k(\mathbf{X}_j) \right)$$

is a sufficient statistic for $\boldsymbol{\theta}$ [Cas01, page 279].

Definition 11 (Minimal Sufficient Statistic). *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 $T'(\mathbf{x})$ [Cas01, page 280].*

Theorem 12. *Let $f(\mathbf{x} \mid \theta)$ be the pd 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 $f(\mathbf{x} \mid \theta) / f(\mathbf{y} \mid \theta)$ is constant as a function of θ if and only if $T(\mathbf{x}) = T(\mathbf{y})$. Then $T(\mathbf{X})$ is a minimal sufficient statistic [Cas01, page 281].*

Example 13 (Normal Minimal Sufficient Statistic). Example taken from [Cas01, page 281]. Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, where both μ and σ^2 unknown. Let \mathbf{x} and \mathbf{y} denote two sample points, and let (\bar{x}, s_x^2) and (\bar{y}, s_y^2) be the

sample means and variances corresponding to the \mathbf{x} and \mathbf{y} samples, respectively. Then, the ratio of the densities becomes

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

This ratio will be constant as a function of μ and σ^2 if and only if $\bar{x} = \bar{y}$ and $s_{\mathbf{x}}^2 = s_{\mathbf{y}}^2$. Thus by Theorem 12, (\bar{X}, S^2) is a minimal sufficient statistic for (μ, σ^2) .

REFERENCES

[Cas01] George and Berger Casella Roger, *Statistical Inference*, Cengage, Mason, OH, 2001 (eng).