



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 INTRODUCTION	5
METHOD OF MOMENTS	12
MAXIMUM LIKELIHOOD ESTIMATES	13
METHODS OF EVALUATING ESTIMATORS	15
SUFFICIENCY AND UNBIASEDNESS	18
CONSISTENCY	18
LARGE-SAMPLE COMPARISONS OF ESTIMATORS	20
EXPECTATION MAXIMIZATION ALGORITHM	22
FORMULATION OF THE EM ALGORITHM	22
HYPOTHESIS TESTING	25
METHODS OF FINDING TESTS	25
BAYESIAN INFERENCE	27
MONTE CARLO SAMPLING	27
SIMULTANEOUS CONFIDENCE BANDS	28
KERNEL DENSITY ESTIMATION	28
REFERENCES	30

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 $H(\mathbf{f})$	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.

(Combination)
$$\binom{n}{r} = \frac{n!}{r!(n-r)!}$$

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

or

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

(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.

$$(1) \quad \mathbb{E}(X) = \int_0^\infty (1 - F(X)) \, dx$$

$$(2) \quad \mathbb{E}(aX + b) = a\mathbb{E}X + b$$

$$(3) \quad \mathbb{E}(g(X) + h(X)) = \mathbb{E}g(X) + \mathbb{E}h(X)$$

$$(4) \quad \text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2$$

$$(5) \quad \text{Var}(aX + b) = a^2\text{Var}(X)$$

$$(6) \quad \text{Cov}(X, Y) = \mathbb{E}XY - \mathbb{E}X\mathbb{E}Y$$

$$(7) \quad \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

$$(8) \quad \mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]]$$

$$(9) \quad \text{Var}(Y) = \mathbb{E}[\text{Var}(Y|X)] + \text{Var}(\mathbb{E}[Y|X])$$

$$(10) \quad |\mathbb{E}(XY)|^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$$

$$(11) \quad |\text{Cov}(XY)|^2 \leq \text{Var}(X)\text{Var}(Y)$$

$$(12) \quad \mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$$

$$(\text{Bayes' Theorem}) \quad \mathbb{P}(A | B) = \frac{\mathbb{P}(B | A)\mathbb{P}(A)}{\mathbb{P}(B)}$$

$$(13) \quad \mathbb{P}(A_1, \dots, A_n) = \mathbb{P}(A_1) \mathbb{P}(A_2 | A_1) \mathbb{P}(A_3 | A_1, A_2) \cdots \mathbb{P}(A_n | A_1, A_2, \dots, A_{n-1})$$

$$(14)$$

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 $x \in \mathbb{R}$

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

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

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

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

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

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 $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ then

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

If $X \sim \chi_2^2$, then $X \sim \text{Exp}(1/2)$.

If $X \sim \text{U}(0, 1)$, then $-2 \ln(X) \sim \chi_2^2$.

If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Exp}(\lambda)$ then $\sum_i X_i \sim \text{Gam}(n, \lambda)$.

If $X \sim \text{Gam}(k, \lambda)$ then for any $c > 0$ we have $cX \sim \text{Gam}(k, c\lambda)$.

POINT ESTIMATION

Methods of Finding Estimates Introduction.

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 (Likelihood, Log-Likelihood and Score Function). *Let $f(\mathbf{x} \mid \theta)$ denote the joint pdf or pmf of the sample $\mathbf{X} = (X_1, \dots, X_n)$. Then, given that $\mathbf{X} = \mathbf{x}$ is observed, the function of θ defined by*

$$L(\theta \mid \mathbf{x}) = f(\mathbf{x} \mid \theta)$$

*is called the **likelihood function** [Cas01, page 290]. For a given outcome \mathbf{x} of \mathbf{X} , the **log-likelihood function**, denoted l , is the natural logarithm of the likelihood function*

$$l(\theta \mid \mathbf{x}) = \ln L(\theta \mid \mathbf{x}) = \ln f(\mathbf{x} \mid \theta).$$

*It's gradient with respect to θ , denoted S , is called the **score function***

$$S(\theta \mid \mathbf{x}) = \nabla_\theta l(\theta \mid \mathbf{x}) = \frac{\nabla_\theta f(\mathbf{x} \mid \theta)}{f(\mathbf{x} \mid \theta)}$$

[Kro13, page 165].

Theorem 10. *Under regularity conditions*

$$\mathbb{E}[S(\theta \mid \mathbf{x})] = 0$$

[Background Notes, page 10].

Proof. Since $L(\theta)$ is a density when viewed as a function of the observed data x_1, \dots, x_n we have the following identity in θ ,

$$\int \dots \int L(\theta) dx_1 \dots dx_n = 1.$$

On differentiating both sides of the above with respect to θ gives

$$\int \dots \int \left[\frac{\partial L(\theta)}{\partial \theta} \right] dx_1 \dots dx_n = 0.$$

Apply the chain rule to $\frac{\partial \ln L(\theta)}{\partial \theta}$ we find

$$\frac{\partial \ln L(\theta)}{\partial \theta} = \frac{\partial \ln L(\theta)}{\partial L(\theta)} \cdot \frac{\partial L(\theta)}{\partial \theta} = \frac{1}{L(\theta)} \frac{\partial L(\theta)}{\partial \theta}$$

meaning

$$\frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) = \frac{\partial L(\theta)}{\partial \theta}$$

so that

$$\begin{aligned} \int \cdots \int \left[\frac{\partial L(\theta)}{\partial \theta} \right] dx_1 \cdots dx_n &= 0 \\ \int \cdots \int \left[\frac{\partial \ln L(\theta)}{\partial \theta} \right] L(\theta) dx_1 \cdots dx_n &= 0 \\ \mathbb{E}[S(\theta)] &= 0 \end{aligned}$$

as wanted. \square

Definition 11 (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 12. 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 11, 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 13 (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 14. Let $f(\mathbf{x} \mid \boldsymbol{\theta})$ be the 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 $f(\mathbf{x} \mid \boldsymbol{\theta}) / f(\mathbf{y} \mid \boldsymbol{\theta})$ is constant as a function of $\boldsymbol{\theta}$ if and only if $T(\mathbf{x}) = T(\mathbf{y})$. Then $T(\mathbf{X})$ is a minimal sufficient statistic [Cas01, page 281].

Example 15 (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

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

This ratio will be constant as a function of μ and σ^2 if and only if $\bar{x} = \bar{y}$ and $s_x^2 = s_y^2$. Thus by Theorem 14, (\bar{X}, S^2) is a minimal sufficient statistic for (μ, σ^2) .

Definition 16 (Ancillary Statistic). A statistic $S(\mathbf{X})$ whose distribution does not depend on the parameter θ is called an ancillary statistic [Cas01, page 282].

Definition 17 (Complete Distributions and Statistics). 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 $\mathbb{E}_\theta g(T) = 0$, for some function g , for all θ implies $\mathbb{P}(g(T) = 0) = 1$ for all θ . Equivalently, $T(\mathbf{X})$ is called a **complete statistic** [Cas01, page 285].

Theorem 18. Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be iid observations from a pdf or pmf $f(x | \theta)$ that belongs to an exponential family as seen in Definition 11, then the statistic

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

is complete as long as the parameter space is non-meager [Cas01, page 288].

Theorem 19. If a minimal sufficient statistic exists, then any complete statistic is also a minimal sufficient statistic [Cas01, page 289].

Theorem 20. A complete, sufficient statistic is always minimal [Background Notes, page 25].

Example 21 (Binomial Complete Statistic). Example taken from [Cas01, page 285]. Suppose that T has a $\text{Bin}(n, p)$ distribution, $0 < p < 1$. Let g be a function such that $\mathbb{E}_p g(T) = 0$. Then

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

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

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

for all, $0 < r < \infty$. But the last expression is a polynomial of degree n in r , where the coefficient of r^t is $g(t) \binom{n}{t}$. For the polynomial to be 0 for all r , each coefficient must be 0. Since none of the $\binom{n}{t}$ terms is 0, this implies that $g(t) = 0$ for $t = 0, 1, \dots, n$. Since T takes on the values $0, 1, \dots, n$ with probability 1, this means that $\mathbb{P}_p(g(T) = 0) = 1$ for all p , the desired conclusion. Hence, T is a complete statistic.

Example 22 (Sum of iid Bernoulli RVs). Example taken from [Tutorial Sheet 2, Q6]. Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$. The likelihood function for θ is given by

$$\begin{aligned} L(\theta) &= \prod_{j=1}^n \binom{n}{x_j} \theta^{x_j} (1-\theta)^{1-x_j} \\ &= \left[\prod_{j=1}^n \binom{n}{x_j} \right] \theta^t (1-\theta)^{n-t} \\ &= \left[\prod_{j=1}^n \binom{n}{x_j} \right] \exp[c(\theta)t] (1-\theta)^n \end{aligned}$$

$$= b(\mathbf{x}) \exp [c(\theta)t] / a(\theta)$$

where

$$\begin{aligned} t(\mathbf{X}) &= \sum_{i=1}^n X_i \\ c(\theta) &= \ln \frac{\theta}{1-\theta} \\ a(\theta) &= (1-\theta)^{-n} \\ b(\mathbf{x}) &= \prod_{j=1}^n \binom{n}{x_j}. \end{aligned}$$

Clearly, the likelihood belongs to the regular exponential family with canonical parameter $c(\theta)$ and complete sufficient statistic $T = t(\mathbf{X})$. Also, the score statistic (Definition 9) is given by

$$S(\theta) = \frac{\partial}{\partial \theta} \ln L(\theta) = \frac{n}{\theta(1-\theta)} \left(\frac{t}{n} - \theta \right)$$

showing that the estimator T attains the Cramer-Rao lower bound is estimating θ . Hence, it attains the MVB (Corollary 45) and is therefore also a UMVU estimator of θ . On the other hand, the estimator

$$V = (X_n, T_{n-1})^\top$$

where $T_{n-1} = \sum_{j=1}^{n-1} X_j$, while sufficient (with canonical parameter $c(\theta) = (\ln \frac{\theta}{1-\theta}, \ln \frac{\theta}{1-\theta})^\top$), is not complete. To demonstrate that V is not complete, we have that

$$\mathbb{E} \left[X_n - \frac{1}{n-1} T_{n-1} \right] = 0$$

however, consider

$$\mathbb{P} \left[X_n - \frac{1}{n-1} T_{n-1} = 0 \right].$$

Since, $X_n \sim \text{Ber}(\theta)$, $T_{n-1} \sim \text{Bin}(n-1, \theta)$ and X_i are iid

$$\begin{aligned} \mathbb{P} \left[X_n - \frac{1}{n-1} T_{n-1} = 0 \right] &= \mathbb{P} [T_{n-1} = 0 \mid X_n = 0] \cdot \mathbb{P} [X_n = 0] + \mathbb{P} [T_{n-1} = n-1 \mid X_n = 1] \cdot \mathbb{P} [X_n = 1] \\ &= (1-\theta)^n + \theta^n \neq 1 \end{aligned}$$

for $0 < \theta < 1$. So by Definition 17, V is not complete. Furthermore, as T is a complete, sufficient statistic, it is a minimal sufficient statistic (Theorem 20) for θ . It is a function of every other sufficient statistic (Definition 13) and here we can see it is a function of V with

$$T = (V)_1 + (V)_2 = X_n + T_{n-1}.$$

This also shows that V is not a (sufficient) minimal statistic (again by Definition 13). Now lets consider the variance between two estimators of θ , $T = \frac{1}{n} \sum_{i=1}^n X_i$ and $W(V) = \mathbb{E}[X_1 \mid V]$. We saw that T is UMVU and its variance attains MVB. Its variance can be computed as

$$\text{Var}(T) = \frac{1}{n^2} (n\theta(1-\theta)) = \frac{1}{n} \theta(1-\theta).$$

Now let us try and find an explicit expression for $W(V(\mathbf{x}))$. We have

$$\begin{aligned}
W(V(\mathbf{x})) &= \mathbb{E} \left[X_1 \mid X_n = x_n, \sum_{i=1}^{n-1} X_i = t_{n-1} \right] \\
&= \sum_{x_1=0}^1 x_1 \cdot \mathbb{P} \left[X_1 = x_1 \mid X_n = x_n, \sum_{i=1}^{n-1} X_i = t_{n-1} \right] \\
&= \mathbb{P} \left[X_1 = 1 \mid X_n = x_n, \sum_{i=1}^{n-1} X_i = t_{n-1} \right] \\
&= \frac{\mathbb{P} \left[X_1 = 1, X_n = x_n, \sum_{i=1}^{n-1} X_i = t_{n-1} \right]}{\mathbb{P} \left[X_n = x_n, \sum_{i=1}^{n-1} X_i = t_{n-1} \right]} \\
&= \frac{\mathbb{P} \left[X_1 = 1, X_n = x_n, \sum_{i=2}^{n-1} X_i = t_{n-1} - 1 \right]}{\mathbb{P} \left[X_n = x_n, \sum_{i=1}^{n-1} X_i = t_{n-1} \right]} \\
&= \frac{\mathbb{P} \left[X_1 = 1 \right] \mathbb{P} \left[X_n = x_n \right] \mathbb{P} \left[\sum_{i=2}^{n-1} X_i = t_{n-1} - 1 \right]}{\mathbb{P} \left[X_n = x_n \right] \mathbb{P} \left[\sum_{i=1}^{n-1} X_i = t_{n-1} \right]}.
\end{aligned}$$

Since $X_1 \sim \text{Ber}(\theta)$, $\sum_{i=1}^{n-1} X_i \sim \text{Bin}(n-1, \theta)$ and $\sum_{i=2}^{n-1} X_i \sim \text{Bin}(n-2, \theta)$, we have

$$\begin{aligned}
W(V(\mathbf{x})) &= \frac{\theta \binom{n-2}{t_{n-1}-1} \theta^{t_{n-1}-1} (1-\theta)^{(n-2)-(t_{n-1}-1)}}{\binom{n-1}{t_{n-1}} \theta^{t_{n-1}} (1-\theta)^{(n-1)-t_{n-1}}} \\
&= t_{n-1} / (n-1)
\end{aligned}$$

where $t_{n-1} = \sum_{i=1}^{n-1} x_i$. This means $W(V(\mathbf{X})) = \frac{1}{n-1} \sum_{i=1}^{n-1} X_i$ and

$$\text{Var}(W(V)) = \frac{(n-1)}{(n-1)^2} \theta(1-\theta) = \frac{1}{(n-1)} \theta(1-\theta) < \frac{1}{n} \theta(1-\theta).$$

Definition 23 (Point Estimator). A **point estimator** is any function $W(X_1, \dots, X_n)$ of a sample; that is, any statistic (see Definition 2) is a point estimator [Cas01, page 311].

Definition 24 (Fisher Information Matrix). For the model $\mathbf{X} \sim f(\cdot; \theta)$, let $S(\theta)$ be the score function (see Definition 9) of θ . The covariance matrix of the random vector $S(\theta)$, denoted by $\mathcal{J}(\theta)$, is called the **Fisher Information Matrix** where

$$\mathcal{J}(\theta) = \mathbb{E}_{\theta} [S(\theta) S(\theta)^{\top}]$$

in the multivariate case and

$$\mathcal{J}(\theta) = \mathbb{E}_{\theta} \left(\frac{d}{d\theta} \ln f(\mathbf{X}; \theta) \right)^2$$

in the one-dimensional case. Note that under regularity conditions $\mathbb{E}[S(\theta)] = 0$ (see Theorem 10) so that

$$\begin{aligned}
\mathcal{J}(\theta) &= \mathbb{E}_{\theta} \left[\frac{d}{d\theta} \ln f(\mathbf{X}; \theta) \right]^2 \\
&= \text{Var}_{\theta} \left(\frac{d}{d\theta} \ln f(\mathbf{X}; \theta) \right) + \left(\mathbb{E}_{\theta} \left[\frac{d}{d\theta} \ln f(\mathbf{X}; \theta) \right] \right)^2 \\
&= \text{Var}_{\theta} (S(\theta)) + (\mathbb{E}_{\theta} [S(\theta)])^2
\end{aligned}$$

$$= \text{Var}_{\theta}(S(\theta))$$

[Kro13, page 168].

Definition 25 (Observed Information). *For the model $\mathbf{X} \sim f(\cdot; \boldsymbol{\theta})$, let $S(\boldsymbol{\theta})$ be the score function (see Definition 9) of $\boldsymbol{\theta}$. The negative of the Hessian of the random vector $S(\boldsymbol{\theta})$, denoted by $I(\boldsymbol{\theta})$, is called the **Observed Information** where*

$$I(\boldsymbol{\theta}) = -\nabla \nabla S(\boldsymbol{\theta})$$

in the multivariate case and

$$I(\boldsymbol{\theta}) = -\frac{\partial^2}{\partial \theta^2} \ln f(\mathbf{X}; \theta)$$

in the one-dimensional case [Background Notes, page 8].

Theorem 26. *Under regularity conditions, the following equality holds*

$$\mathcal{J}(\boldsymbol{\theta}) = \mathbb{E}[I(\boldsymbol{\theta})]$$

[Kro13, page 169].

Theorem 27 (Fisher Information Matrix for iid Data). *Let $\mathbf{X} = (X_1, \dots, X_n) \stackrel{iid}{\sim} f(x; \boldsymbol{\theta})$, and let $\mathring{\mathcal{J}}(\boldsymbol{\theta})$ be the information matrix corresponding to $X \sim f(x; \boldsymbol{\theta})$. Then the information matrix for \mathbf{X} is given by*

$$\mathcal{J}(\boldsymbol{\theta}) = n\mathring{\mathcal{J}}(\boldsymbol{\theta})$$

[Kro13, page 170].

Theorem 28. *If the $L(\theta)$ belongs to the regular exponential family, then the likelihood equation*

$$\frac{d}{d\boldsymbol{\theta}} \ln L(\boldsymbol{\theta}) = \mathbf{0},$$

can be expressed as

$$T(\mathbf{x}_1, \dots, \mathbf{x}_n) = \mathbb{E}[T(\mathbf{X}_1, \dots, \mathbf{X}_n)]$$

[Lecture Notes 1, page 8].

Method of Moments.

Definition 29 (Method of Moments). Let X_1, \dots, X_n be a random sample of size n from a population with pf $f(x \mid \theta_1, \dots, \theta_k)$. Method of moments estimators are found by equating the first k sample moments to the corresponding k population moments, and solving the resulting system of simultaneous equations. More precisely, define

$$\begin{aligned} m_1 &= \frac{1}{n} \sum_{i=1}^n X_i^1, & \mu'_1 &= \mathbb{E}X^1 \\ m_2 &= \frac{1}{n} \sum_{i=1}^n X_i^2, & \mu'_2 &= \mathbb{E}X^2 \\ &\vdots \\ m_k &= \frac{1}{n} \sum_{i=1}^n X_i^k, & \mu'_k &= \mathbb{E}X^k. \end{aligned}$$

The population moment μ'_j will typically be a function of $\theta_1, \dots, \theta_k$, say $\mu'_j(\theta_1, \dots, \theta_k)$. The method of moments estimator $(\tilde{\theta}_1, \dots, \tilde{\theta}_k)$ of $(\theta_1, \dots, \theta_k)$ is obtained by solving the following system of equations for $(\theta_1, \dots, \theta_k)$ in terms of (m_1, \dots, m_k)

$$\begin{aligned} m_1 &= \mu'_1(\theta_1, \dots, \theta_k) \\ m_2 &= \mu'_2(\theta_1, \dots, \theta_k) \\ &\vdots \\ m_k &= \mu'_k(\theta_1, \dots, \theta_k) \end{aligned}$$

[Cas01, page 312].

Example 30 (Normal Methods of Moments). Example taken from [Cas01, page 313]. Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathbf{N}(\theta, \sigma^2)$. In the preceding notation, $\theta_1 = \theta$ and $\theta_2 = \sigma^2$. We have $m_1 = \bar{X}$, $m_s = (1/n) \sum X_i^2$, $\mu'_1 = \theta$, $\mu'_2 = \theta^2 + \sigma^2$, and hence we must solve

$$\bar{X} = \theta, \quad \frac{1}{n} \sum X_i^2 = \theta^2 + \sigma^2.$$

Solving for θ and σ^2 yields the methods of moments estimators

$$\tilde{\theta} = \bar{X} \quad \text{and} \quad \tilde{\sigma}^2 = \frac{1}{n} \sum X_i^2 - \bar{X}^2 = \frac{1}{n} \sum (X_i^2 - \bar{X}^2).$$

Maximum Likelihood Estimates.

Definition 31 (Maximum Likelihood Estimator). For each sample point \mathbf{x} , let $\hat{\theta}(\mathbf{x})$ be a parameter value at which $L(\theta | \mathbf{x})$ attains its maximum as a function of θ , with \mathbf{x} held fixed. A **maximum likelihood estimator (MLE)** of the parameter θ based on a sample \mathbf{X} is $\hat{\theta}(\mathbf{X})$ [Cas01, page 316].

Example 32 (Normal Likelihood). Example taken from [Cas01, page 316]. Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, 1)$, and let $L(\theta | \mathbf{x})$ denote the likelihood function. Then

$$L(\theta | \mathbf{x}) = \prod_{i=1}^n \frac{1}{(2\pi)^{1/2}} \exp\left(-(1/2)(x_i - \theta)^2\right) = \frac{1}{(2\pi)^{1/2}} \exp\left(-(1/2) \sum_{i=1}^n (x_i - \theta)^2\right).$$

The equation $(d/d\theta)L(\theta | \mathbf{x}) = 0$ reduces to

$$\sum_{i=1}^n (x_i - \theta) = 0,$$

which has the solution $\hat{\theta} = \bar{x}$. Hence, \bar{x} is a candidate for the MLE. To verify that \bar{x} is, in fact, a global maximum of the likelihood function, we can use the following argument. First, note that $\hat{\theta} = \bar{x}$ is the only solution to $\sum_{i=1}^n (x_i - \theta) = 0$; hence \bar{x} is the only zero of the first derivative. Second, verify that

$$\frac{d^2}{d\theta^2} L(\theta | \mathbf{x})|_{\theta=\bar{x}} < 0.$$

Thus, \bar{x} is the only extreme point in the interior and it is a maximum. To finally verify that \bar{x} is a global maximum, we must check the boundaries at $\pm\infty$. So $\hat{\theta} = \bar{x}$ is a global maximum and hence \bar{X} is the MLE.

Theorem 33. If $\hat{\theta}$ is the MLE of θ , then for any function $\tau(\theta)$, the MLE of $\tau(\theta)$ is $\tau(\hat{\theta})$ [Cas01, page 320].

Example 34 (Normal MLE, μ and σ unknown). Example taken from [Cas01, page 321]. Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, \sigma^2)$ with both μ and σ^2 unknown. Then

$$L(\theta | \mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-(1/2) \sum_{i=1}^n (x_i - \theta)^2 / \sigma^2\right)$$

and

$$\ln L(\theta | \mathbf{x}) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 / \sigma^2.$$

The partial derivatives, with respect to θ and σ^2 are

$$\frac{\partial}{\partial \theta} \ln L(\theta | \mathbf{x}) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \theta)$$

and

$$\frac{\partial}{\partial \sigma^2} \ln L(\sigma^2 | \mathbf{x}) = -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \theta)^2.$$

Setting the partial derivatives equal to 0 and solving for the solution $\hat{\theta} = \bar{x}$, $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2$. To verify that this solution is, in fact, a global maximum, recall first that if $\theta \neq \bar{x}$, then $\sum (x_i - \theta)^2 >$

$\sum (x_i - \bar{x})^2$. Hence, for any value of σ^2 ,

$$\frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left(-(1/2) \sum_{i=1}^n (x_i - \bar{x})^2 / \sigma^2 \right) \geq \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left(-(1/2) \sum_{i=1}^n (x_i - \theta)^2 / \sigma^2 \right).$$

Therefore, verifying that we have found the maximum likelihood estimators is reduced to a one-dimensional problem, verifying that $(\sigma^2)^{-n/2} \exp \left(-\frac{1}{2} \sum (x_i - \bar{x})^2 / \sigma^2 \right)$ achieves its global maximum at $\sigma^2 = n^{-1} \sum (x_i - \bar{x})^2$. This is straightforward to do using univariate calculus and, in fact, the estimators $(\bar{X}, n^{-1} \sum (X_i - \bar{X})^2)$ are the MLEs.

Definition 35 (Quantile of Order α). *The quantile of order α , q_α , is the value of the random variable X such that*

$$\mathbb{P}[X \leq q_\alpha] = \alpha.$$

Example 36. Let X_1, \dots, X_n denote a random sample from a $N(\mu, \sigma^2)$ distribution, where both μ and σ are unknown. Let us consider a way to find the maximum likelihood for quantile of order α . Take $X \sim N(\mu, \sigma^2)$. As $Z = (X - \mu)/\sigma$ has a standard normal distribution with the function $\Phi(z)$, we can express the left hand for the expression of the quantile of order α as

$$\mathbb{P}[X \leq q_\alpha] = \mathbb{P} \left[\frac{X - \mu}{\sigma} \leq \frac{q_\alpha - \mu}{\sigma} \right] = \mathbb{P} \left[Z \leq \frac{q_\alpha - \mu}{\sigma} \right].$$

This means

$$\begin{aligned} \frac{q_\alpha - \mu}{\sigma} &= \Phi^{-1}(\alpha) \\ q_\alpha &= \mu + \sigma \Phi^{-1}(\alpha). \end{aligned}$$

Hence, as a consequence of Theorem 33, $g(\hat{\theta}) = \hat{\mu} + \hat{\sigma} \Phi^{-1}(\alpha)$ is the maximum likelihood estimate of q_α . Note, however, that this is not an unbiased estimate of q_α by virtue of the fact that $\mathbb{E}[\hat{\sigma}] \neq \sigma$. If we adjusted this estimate by multiplying by some constant k_n , that is,

$$\mathbb{E}[\hat{\sigma}] = k_n \sigma$$

then it would be an unbiased estimator of q_α . To compute such a k_n , we have that $n\hat{\sigma}^2/\sigma^2 \sim \chi_{n-1}^2$ that is, $\frac{1}{2}n\hat{\sigma}^2/\sigma^2 \sim \gamma(m/2)$ where $m = (n-1)/2$. This is equivalent to saying $\hat{\sigma} = \sigma \sqrt{2/n} \sqrt{Y}$ where $Y \sim \gamma(m)$. Thus $\mathbb{E}[\hat{\sigma}] = k_n \sigma$, where $k_n = \sqrt{2/n} \mathbb{E}[\sqrt{Y}]$ and where

$$\begin{aligned} \mathbb{E}[\sqrt{Y}] &= \frac{\int_0^\infty y^{1/2} \exp(-y) y^{m-1} dy}{\Gamma(m)} \\ &= \frac{\int_0^\infty \exp(-y) y^{m+\frac{1}{2}-1} dy}{\Gamma(m)} \\ &= \frac{\Gamma(m + \frac{1}{2})}{\Gamma(m)} \\ &= \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \end{aligned}$$

Upon substituting the result for $\mathbb{E}[\sqrt{Y}]$ into the right-hand side of expression for k_n we obtain

$$k_n = \sqrt{\frac{2}{n}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})}.$$

Methods of Evaluating Estimators.

Definition 37 (Mean Square Error). The **mean square error** (MSE) of an estimator W of a parameter θ is the function θ defined by $\mathbb{E}_\theta(W - \theta)^2$ [Cas01, page 330].

Definition 38 (Bias). The **bias** of an estimator W of a parameter θ is the difference between the expected value of W and θ ; that is $\text{Bias}_\theta W = \mathbb{E}_\theta W - \theta$. An estimator whose bias is identically (in θ) equal to 0 is called an **unbiased estimator** and satisfies $\mathbb{E}_\theta W = \theta$ for all θ [Cas01, page 330].

It is important to note that

$$\mathbb{E}_\theta (W - \theta)^2 = \text{Var}_\theta W + (\mathbb{E}_\theta W - \theta)^2 = \text{Var}_\theta W + (\text{Bias}_\theta W)^2.$$

Example 39 (Normal MSE). Example taken from [Cas01, page 331]. Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$. The statistics \bar{X} and S^2 are both unbiased estimators since

$$\mathbb{E} \bar{X} = \mu, \quad \mathbb{E} S^2 = \sigma^2, \quad \text{for all } \mu \text{ and } \sigma^2.$$

The MSEs of these estimators are given by

$$\begin{aligned} \mathbb{E} (\bar{X} - \mu)^2 &= \text{Var} \bar{X} = \frac{\sigma^2}{n} \\ \mathbb{E} (S^2 - \sigma^2)^2 &= \text{Var} S^2 = \frac{2\sigma^4}{n-1}. \end{aligned}$$

The MSE of \bar{X} remains σ^2/n even if the normality assumption is dropped. However, the above expression for the MSE of S^2 does not remain the same if the normality assumption is relaxed. An alternative estimator for σ^2 is the MLE $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n-1}{n} S^2$. It is straightforward to calculate

$$\mathbb{E} \hat{\sigma}^2 = \mathbb{E} \left(\frac{n-1}{n} S^2 \right) = \frac{n-1}{n} \sigma^2,$$

so that $\hat{\sigma}^2$ is a biased estimator of σ^2 . The variance of $\hat{\sigma}^2$ can also be calculated as

$$\text{Var} \hat{\sigma}^2 = \text{Var} \left(\frac{n-1}{n} S^2 \right) = \left(\frac{n-1}{n} \right)^2 \text{Var} S^2 = \frac{2(n-1)\sigma^4}{n^2},$$

and hence, its MSE is given by

$$\mathbb{E} (\hat{\sigma}^2 - \sigma^2)^2 = \frac{2(n-1)\sigma^4}{n^2} + \left(\frac{n-1}{n} \sigma^2 - \sigma^2 \right)^2 = \left(\frac{2n-1}{n^2} \right) \sigma^4.$$

Thus we have

$$\mathbb{E} (\hat{\sigma}^2 - \sigma^2)^2 = \left(\frac{2n-1}{n^2} \right) \sigma^4 < \left(\frac{2}{n-1} \right) \sigma^4 = \mathbb{E} (S^2 - \sigma^2)^2,$$

showing that $\hat{\sigma}^2$ has a smaller MSE than S^2 . Thus, by trading off variance for bias, the MSE is improved.

Definition 40 (Best Unbiased Estimator). An estimator W^* is a **best unbiased estimator** of $\tau(\theta)$ if it satisfies $\mathbb{E}_\theta W^* = \tau(\theta)$ for all θ and, for any other estimator W with $\mathbb{E}_\theta W = \tau(\theta)$, we have $\text{Var}_\theta W^* \leq \text{Var}_\theta W$ for all θ . W^* is also called a **uniform minimum variance unbiased estimator** (UMVUE) of $\tau(\theta)$ [Cas01, page 334].

Theorem 41 (Cramer-Rao Inequality). Let X_1, \dots, X_n be a sample with pdf $f(\mathbf{x} \mid \theta)$, and let $W(\mathbf{X}) = W(X_1, \dots, X_n)$ be any estimator satisfying

$$\frac{d}{d\theta} \mathbb{E}_\theta W(\mathbf{X}) = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} [W(\mathbf{x}) f(\mathbf{x} \mid \theta)]$$

and

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

Then

$$\text{Var}_\theta(W(\mathbf{X})) \geq \frac{\left(\frac{d}{d\theta} \mathbb{E}_\theta W(\mathbf{X})\right)^2}{\mathbb{E}_\theta \left(\left(\frac{\partial}{\partial \theta} \ln f(\mathbf{X} | \theta)\right)^2\right)} = \frac{\left(\frac{d}{d\theta} \mathbb{E}_\theta W(\mathbf{X})\right)^2}{\mathcal{J}(\theta)}$$

which is commonly referred to as the **minimum variance bound** (MVB). If $W(\mathbf{X})$ attains the MVB (for all values of θ), it is said to be a MVB estimator [Cas01, page 335].

Corollary 42 (Cramer-Rao Inequality, iid Case). *If the assumptions of Theorem 41 are satisfied and, additionally, if X_1, \dots, X_n are iid with pdf $f(x | \theta)$, then*

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

[Cas01, page 337].

Lemma 43. *If $f(x | \theta)$ satisfies*

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

(true for the exponential family), then

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

[Cas01, page 338].

Example 44 (Poisson Unbiased Estimate). Example taken from [Cas01, page 338]. Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Poi}(\lambda)$, and let \bar{X} and S^2 be the sample mean and variance, respectively. Recall that for the Poisson pmf both the mean and variance are equal to λ . We have

$$\mathbb{E}_\lambda \bar{X} = \lambda, \quad \text{for all } \lambda,$$

$$\mathbb{E}_\lambda S^2 = \lambda, \quad \text{for all } \lambda,$$

so both \bar{X} and S^2 are unbiased estimators of λ . To determine the better estimator, \bar{X} or S^2 , we should now compare the variances. We have $\text{Var}_\lambda \bar{X} = \lambda/n$, but $\text{Var}_\lambda S^2$ is quite a lengthy calculation. Not only this, even if we can establish that \bar{X} is better than S^2 , consider the class of estimators

$$W_a(\bar{X}, S^2) = a\bar{X} + (1-a)S^2.$$

For every constant a , $\mathbb{E}_\lambda W_a = \lambda$, so now we have infinitely many unbiased estimators of λ . Instead, let us show that \bar{X} is the best estimator directly using the Cramer-Rao inequality. Here we are estimating $\tau(\lambda) = \lambda$, so that $\tau'(\lambda) = 1$. Also, since we have an exponential family, using Lemma 43 gives us

$$\begin{aligned} \mathbb{E}_\lambda \left(\left(\frac{\partial}{\partial \lambda} \ln f(X | \lambda) \right)^2 \right) &= -n \mathbb{E}_\lambda \left(\frac{\partial^2}{\partial \lambda^2} \ln f(X | \lambda) \right) \\ &= -n \mathbb{E}_\lambda \left(\frac{\partial^2}{\partial \lambda^2} \ln \left(\frac{e^{-\lambda} \lambda^X}{X!} \right) \right) \end{aligned}$$

$$\begin{aligned}
&= -n\mathbb{E}_\lambda \left(\frac{\partial^2}{\partial \lambda^2} (-\lambda + X \ln \lambda - \ln X!) \right) \\
&= -n\mathbb{E}_\lambda \left(-\frac{X}{\lambda^2} \right) \\
&= \frac{n}{\lambda}.
\end{aligned}$$

Hence for any unbiased estimator, W , of λ , from Corollary 42 we must have

$$\begin{aligned}
\text{Var}_\theta(W(\mathbf{X})) &\geq \frac{\left(\frac{d}{d\theta} \mathbb{E}_\theta W(\mathbf{X}) \right)^2}{n\mathbb{E}_\theta \left(\left(\frac{\partial}{\partial \theta} \ln f(X | \theta) \right)^2 \right)} \\
&= \frac{(1)^2}{\left(\frac{n}{\lambda} \right)} \\
&= \frac{\lambda}{n}.
\end{aligned}$$

Since $\text{Var}_\lambda \bar{X} = \lambda/n$, \bar{X} must be the best unbiased estimator.

Corollary 45 (Attainment). *Let X_1, \dots, X_n be a sample with pdf $f(x | \theta)$, where $f(x | \theta)$ satisfies the conditions of the Cramer-Rao Theorem. $L(\theta | \mathbf{x}) = \prod_{i=1}^n f(x_i | \theta)$ denote the likelihood function. If $W(\mathbf{X}) = W(X_1, \dots, X_n)$ is any unbiased estimator of $\tau(\theta)$, then $W(\mathbf{X})$ attains the Cramer-Rao Lower Bound if and only if*

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

for some function $a(\theta)$ [Cas01, page 341].

Example 46. Example taken from Tutorial Sheet 2 Q5. Let T be an estimator of the parameter θ , having bias $b(\theta)$. Assuming that the usual regularity conditions and using the Cramer-Rao lower bound (Theorem 41) for the variance of an unbiased estimator of θ , we can show that

$$\text{MSE}(T) \geq \left[1 + \frac{\partial}{\partial \theta} b(\theta) \right]^2 \cdot \mathcal{J}^{-1}(\theta) + [b(\theta)]^2$$

where $\mathcal{J}(\theta)$ is the Fisher information matrix (Definition 24). To start, since $b(\theta) = \mathbb{E}[\theta] - \theta$ we have

$$\mathbb{E}[T] = \theta + b(\theta) \triangleq g(\theta)$$

so that T is an unbiased estimate for $g(\theta)$. By the Cramer-Rao lower bound,

$$\text{Var}(T) \geq [g'(\theta)]^2 \cdot \mathcal{J}^{-1}(\theta) = \left[1 + \frac{\partial}{\partial \theta} b(\theta) \right]^2 \cdot \mathcal{J}^{-1}(\theta).$$

Now

$$\begin{aligned}
\text{MSE}(T) &= \text{Var}(T) + [\text{Bias}(T)]^2 \\
&= \text{Var}(T) + [b(\theta)]^2 \\
&\geq \left[1 + \frac{\partial}{\partial \theta} b(\theta) \right]^2 \cdot \mathcal{J}^{-1}(\theta) + [b(\theta)]^2.
\end{aligned}$$

Example 47 (Continuation of Example 34). Example taken from [Cas01, page 341]. Here we know

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

and hence

$$\frac{\partial}{\partial \sigma^2} \ln L(\mu, \sigma^2 \mid \mathbf{x}) = \frac{n}{2\sigma^4} \left(\sum_{i=1}^n \frac{(x_i - \mu)^2}{n} - \sigma^2 \right).$$

Thus, taking $a(\sigma^2) = n/(2\sigma^4)$ shows that the best unbiased estimator of σ^2 is $\frac{(x_i - \mu)^2}{n}$, which is calculable only if μ is known. If μ is not known, the bound *cannot* be attained.

Sufficiency and Unbiasedness.

Theorem 48 (Rao-Blackwell). *Let W be any unbiased estimator of $\tau(\theta)$, and let T be a sufficient statistic for θ . Define $\phi(T) = \mathbb{E}(W \mid T)$. Then $\mathbb{E}_\theta \phi(T) = \tau(\theta)$ and $\text{Var}_\theta \phi(T) \leq \text{Var}_\theta W$ for all θ ; that is, $\phi(T)$ is a uniformly better unbiased estimator of $\tau(\theta)$ [Cas01, page 342].*

Theorem 49. *If W is the best unbiased estimator of $\tau(\theta)$, then W is unique [Cas01, page 343].*

Theorem 50. *Let T be a complete sufficient statistic for a parameter θ , and let $\phi(T)$ be any estimator based only on T . Then $\phi(T)$ is the best unbiased estimator of its expected value [Cas01, page 347].*

Consistency.

Definition 51 (Consistency). *A sequence of estimators T_n of $g(\theta)$ is said to be **consistent** if for every $\theta \in \Omega$,*

$$T_n \xrightarrow{\mathbb{P}_\theta} g(\theta), \quad \text{as } n \rightarrow \infty$$

that is, given any $\varepsilon > 0$, then

$$\mathbb{P} [|T_n(\mathbf{X}_1, \dots, \mathbf{X}_n) - g(\theta)| \geq \varepsilon] \rightarrow 0, \quad \text{as } n \rightarrow \infty$$

[Kro13, page 176] [Background Notes, page 44].

Theorem 52. *If $\text{Var}(T_n) \rightarrow 0$ and $\text{Bias}(T_n) \rightarrow 0$, as $n \rightarrow \infty$, then the sequence of estimates T_n is consistent for estimating $g(\theta)$ [Background Notes, page 44].*

Proof. Let $f(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta)$ denote the joint pdf of $\mathbf{X}_1, \dots, \mathbf{X}_n$. Then we have

$$\mathbb{P} [|T_n - g(\theta)| \geq \varepsilon] = \int \dots \int_{|T_n - g(\theta)| \geq \varepsilon} f(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta) d\mathbf{x}_1 \dots d\mathbf{x}_n.$$

On the region of integration in the above,

$$\begin{aligned} |T_n - g(\theta)| &\geq \varepsilon \\ (T_n - g(\theta))^2 &\geq \varepsilon^2 \\ \frac{(T_n - g(\theta))^2}{\varepsilon^2} &\geq 1, \end{aligned}$$

and since $f(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta)$ is non-negative,

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta) \leq f(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta) \frac{(T_n - g(\theta))^2}{\varepsilon^2}.$$

Thus

$$\begin{aligned}
& \mathbb{P} [|T_n - g(\boldsymbol{\theta})| \geq \varepsilon] \\
&= \int \dots \int_{|T_n - g(\boldsymbol{\theta})| \geq \varepsilon} f(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta}) d\mathbf{x}_1 \dots d\mathbf{x}_n \\
&\leq \int \dots \int_{|T_n - g(\boldsymbol{\theta})| \geq \varepsilon} f(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta}) \frac{(T_n - g(\boldsymbol{\theta}))^2}{\varepsilon^2} d\mathbf{x}_1 \dots d\mathbf{x}_n \\
&\leq \frac{1}{\varepsilon^2} \int \dots \int f(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta}) (T_n - g(\boldsymbol{\theta}))^2 d\mathbf{x}_1 \dots d\mathbf{x}_n \\
&= \frac{1}{\varepsilon^2} \text{MSE}(T_n) \\
&= \frac{1}{\varepsilon^2} [\text{Var}(T_n) + (\text{Bias}(T_n))^2]
\end{aligned}$$

which tends to 0 as $n \rightarrow \infty$, since $\text{Var}(T_n)$ and $\text{Bias}(T_n)$ both tend to 0. \square

Example 53 (Continuation of Example 39). Example taken from Background Notes, page 44. The estimators

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

are both consistent for σ^2 . This is easily seen with s^2 since

$$\mathbb{E}[s^2] = \sigma^2 \quad \text{and} \quad \text{Var}(s^2) \rightarrow 0$$

as $n \rightarrow \infty$. To show that $\hat{\sigma}^2$ is also a consistent estimator, note that

(see 17)
$$\frac{\sum_{j=1}^n (X_j - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2$$

meaning

$$\mathbb{E} \left[\frac{\sum_{j=1}^n (X_j - \bar{X})^2}{\sigma^2} \right] = n-1$$

and

$$\text{Var} \left[\frac{\sum_{j=1}^n (X_j - \bar{X})^2}{\sigma^2} \right] = 2(n-1)$$

so that

$$\mathbb{E} \left[\sum_{j=1}^n (X_j - \bar{X})^2 \right] = (n-1)\sigma^2$$

and

$$\text{Var} \left[\sum_{j=1}^n (X_j - \bar{X})^2 \right] = 2(n-1)\sigma^4.$$

Hence

$$\begin{aligned}
\mathbb{E}(\hat{\sigma}^2) &= \frac{n-1}{n} \sigma^2 = \sigma^2 - \frac{\sigma^2}{n} \\
\text{Var}(\hat{\sigma}^2) &= \frac{2(n-1)\sigma^4}{n^2}
\end{aligned}$$

meaning $\mathbb{E}(\hat{\sigma}^2) \rightarrow 0$ and $\text{Var}(\hat{\sigma}^2) \rightarrow 0$ as $n \rightarrow \infty$. Therefore, by Theorem 52, $\hat{\sigma}^2$ is a consistent estimator of σ^2 .

Large-Sample Comparisons of Estimators.

Theorem 54 (Asymptotic Distribution of the MLE). *Suppose that $\hat{\theta}_n$ is a sequence of consistent ML estimates for θ . Then $\sqrt{n}(\hat{\theta}_n - \theta)$ converges in distribution to a $\mathcal{N}(\mathbf{0}, \mathring{\mathcal{J}}^{-1}(\theta))$ distributed random vector as $n \rightarrow \infty$. In other words,*

$$\hat{\theta}_n \overset{\text{approx}}{\sim} \mathcal{N}(\theta, \mathring{\mathcal{J}}^{-1}(\theta)/n).$$

Theorem 55 (Multivariate Central Limit Theorem). *Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \overset{\text{iid}}{\sim} \text{WN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. For large n , the random vector $\sum_i \mathbf{X}_i$ approximately has a $\mathcal{N}(n\boldsymbol{\mu}, \boldsymbol{\Sigma})$ distribution.*

Theorem 56 (Delta Method). *Let $\mathbf{Z}_1, \mathbf{Z}_2, \dots$ be a sequence of random vectors such that $\sqrt{n}(\mathbf{Z}_n - \boldsymbol{\mu}) \rightarrow \mathbf{K} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ as $n \rightarrow \infty$. Then for any continuously differentiable function g of \mathbf{Z}_n ,*

$$\sqrt{n}(g(\mathbf{Z}_n) - g(\boldsymbol{\mu})) \rightarrow \mathbf{R} \sim \mathcal{N}(\mathbf{0}, \mathbf{J}\boldsymbol{\Sigma}\mathbf{J}^\top),$$

where $\mathbf{J} = \mathbf{J}(\boldsymbol{\mu}) = \left(\frac{\partial g_i(\boldsymbol{\mu})}{\partial x_j}\right)$ is the Jacobian of the matrix g evaluated at $\boldsymbol{\mu}$ [Kro13, page 92].

Definition 57 (Asymptotic Relative Efficiency). *Suppose that $\hat{\theta}_{n_1}$ and $\hat{\theta}_{n_2}$ are two single variable estimates such that*

$$\begin{aligned} \hat{\theta}_{n_1} &\overset{\text{approx}}{\sim} \mathcal{N}(\theta, \tau_1^2/n) \\ \hat{\theta}_{n_2} &\overset{\text{approx}}{\sim} \mathcal{N}(\theta, \tau_2^2/n). \end{aligned}$$

The **Asymptotic Relative Efficiency** (ARE) of $\hat{\theta}_{n_2}$ with respect to $\hat{\theta}_{n_1}$ is given by

$$\text{ARE}(\hat{\theta}_{n_2}) = \tau_1^2/\tau_2^2$$

[Background Notes, page 46].

Example 58 (Asymptotic Distribution of Bernoulli MLE). Example taken from [Kro13, page 177]. For $X_1, \dots, X_n \overset{\text{iid}}{\sim} \text{Ber}(p)$, the MLE for p is

$$\hat{p}_n = \bar{x} = \frac{1}{n} \sum_i x_i.$$

To compute the information number (see Definition 24) for p , note that regularity conditions hold so that

$$\mathring{\mathcal{J}}(p) = \text{Var}_p(S(p))$$

where

$$\begin{aligned} S(p) &= \frac{d}{dp} \ln(p^x(1-p)^{1-x}) \\ &= \frac{d}{dp} [x \cdot \ln(p) + (1-x) \ln(1-p)] \\ &= \frac{x}{p} - \frac{(1-x)}{1-p} = \frac{x - \theta}{p(1-p)}. \end{aligned}$$

This means that

$$\mathring{\mathcal{J}}(p) = \text{Var}_p(S(p))$$

$$\begin{aligned}
&= \text{Var}_p \left(\frac{X - \theta}{p(1-p)} \right) \\
&= \text{Var}_p \left(\frac{X}{p(1-p)} \right) \\
&= \frac{p(1-p)}{p^2(1-p)^2} = \frac{1}{p(1-p)}.
\end{aligned}$$

Theorem 54 states that

$$\hat{p}_n \overset{\text{approx}}{\sim} \mathbf{N} \left(p, \frac{p(1-p)}{n} \right).$$

Expectation Maximization Algorithm. The expectation-maximization (EM) algorithm is a broadly applicable approach to the iterative computation of maximum likelihood (ML) estimates, useful in a variety of incomplete data problems, where algorithms such as the Newton-Raphson method may turn out to be more complicated. On each iteration of the EM algorithm, there are two steps called the Expectation step or the E -step and the Maximization step or the M -step. Because of this, the algorithm is called the EM algorithm.

Formulation of the EM Algorithm. We let \mathbf{Y} be the random vector corresponding to the observed data \mathbf{y} having p.d.f. postulated as $g(\mathbf{y}; \Psi)$, where $\Psi = (\Psi_1, \Psi_2, \dots, \Psi_d)^\top$ is a vector of unknown parameters with parameter space Ω . We let $g_c(\mathbf{x}; \Psi)$ denote the p.d.f. of the random vector \mathbf{X} corresponding to the complete data vector \mathbf{x} . Then the complete-data log likelihood function that could be formed for Ψ if \mathbf{x} were fully observable is given by

$$\ln L_c(\Psi) = \ln g_c(\mathbf{x}; \Psi).$$

Formally, we have two sample space X and Y and a many-to-one mapping X to Y . Instead of observing the complete-data vector $\mathbf{x} \in X$, we observe the incomplete-data vector $\mathbf{y} = \mathbf{y}(\mathbf{x}) \in Y$. It follows that

$$g(\mathbf{y}; \Psi) = \int_{X(\mathbf{y})} g_c(\mathbf{x}; \Psi) d\mathbf{x}$$

where $X(\mathbf{y})$ is the subset of X determined by the equation $\mathbf{y} = \mathbf{y}(\mathbf{x})$. The EM algorithm approaches the problem of solving the incomplete-data likelihood equation

$$\nabla_{\Psi} \ln L(\Psi) = 0$$

indirectly by proceeding iteratively in terms of the complete-data log likelihood function, $\ln L_c(\Psi)$. As it is unobservable, it is replaced by its conditional expectation given \mathbf{y} , using the current fit for \mathbf{y} . More specifically, let $\Psi^{(0)}$ be some initial value for Ψ . Then on the first iteration, the E -step requires the calculation of

$$Q(\Psi; \Psi^{(0)}) = \mathbb{E}_{\Psi^{(0)}} [\ln L_c(\Psi) | \mathbf{y}].$$

The M -step requires the maximization of $Q(\Psi; \Psi^{(0)})$ with respect to Ψ over the parameter space Ω . That is, we choose $\Psi^{(1)}$ such that

$$Q(\Psi^{(1)}; \Psi^{(0)}) \geq Q(\Psi; \Psi^{(0)})$$

for all $\Psi \in \Omega$. The E - and M -steps are then carried out again, but this time with $\Psi^{(0)}$ replaced by the current fit $\Psi^{(1)}$. On the $(k+1)^{th}$ iteration, the E - and M -steps are defined as follows:

Definition 59 (E -step). Calculate $Q(\Psi; \Psi^{(k)})$ as

$$Q(\Psi; \Psi^{(k)}) = \mathbb{E}_{\Psi^{(k)}} [\ln L_c(\Psi) | \mathbf{y}].$$

Definition 60 (M -step). Choose $\Psi^{(k+1)}$ to be any value of $\Psi \in \Omega$ that maximises $Q(\Psi; \Psi^{(k)})$, that is,

$$Q(\Psi^{(k+1)}; \Psi^{(k)}) \geq Q(\Psi; \Psi^{(k)})$$

for all $\Psi \in \Omega$.

The E – and M – steps are alternated repeatedly until the difference

$$L(\Psi^{(k+1)}) - L(\Psi^{(k)})$$

changes by an arbitrarily small amount in the case of convergence of the sequence of likelihood value $L(\Psi^{(k)})$. Another way of expressing Definition 59 is to say that $\Psi^{(k+1)}$ belongs to

$$\mathcal{M}(\Psi^{(k)}) = \underset{\Psi}{\operatorname{argmax}} Q(\Psi; \Psi^{(k)}),$$

which is the set of points that maximise $Q(\Psi; \Psi^{(k)})$.

Example 61 (A Multinomial Example). Example taken from [McL08, page] [Kro13, page 185]. We consider first the multinomial example that DLR used to introduce the EM algorithm and that has been subsequently used many times in the literature to illustrate various modifications and extensions of this algorithm. The data relates to a problem of estimation of linkage in genetics where an observed data vector of frequencies

$$\mathbf{y} = (y_1, y_2, y_3, y_4)^\top$$

is postulated to arise from a multinomial distribution with four cells with cell probabilities

$$\frac{1}{2} + \frac{1}{4}\Psi, \frac{1}{4}(1 - \Psi), \frac{1}{4}(1 - \Psi), \text{ and } \frac{1}{4}\Psi$$

with $0 \leq \Psi \leq 1$. The parameter Ψ is to be estimated on the basis of the observed information \mathbf{y} . The probability of the observed data \mathbf{y} is given by

$$g(\mathbf{y}; \Psi) = \frac{n!}{y_1!y_2!y_3!y_4!} \left(\frac{1}{2} + \frac{1}{4}\Psi\right)^{y_1} \left(\frac{1}{4} - \frac{1}{4}\Psi\right)^{y_2} \left(\frac{1}{4} - \frac{1}{4}\Psi\right)^{y_3} \left(\frac{1}{4}\Psi\right)^{y_4}.$$

Suppose now that the first of the original four multinomial cells, which has an associated probability of $\frac{1}{2} + \frac{1}{4}\Psi$, could be split into two subcells have probabilities $\frac{1}{2}$ and $\frac{1}{4}$, respectively, and let y_{11} and y_{12} be the corresponding split of y_1 , where

$$y_1 = y_{11} + y_{12}.$$

Thus, the observed vector of frequencies \mathbf{y} is viewed as being incomplete and the complete-data vector is taken to be

$$\mathbf{x} = (y_{11}, y_{12}, y_2, y_3, y_4)^\top.$$

The cell frequencies in \mathbf{x} are assumed to arise from a multinomial distribution having five cells with probabilities

$$\frac{1}{2}, \frac{1}{4}\Psi, \frac{1}{4}(1 - \Psi), \frac{1}{4}(1 - \Psi), \text{ and } \frac{1}{4}\Psi.$$

In this framework, y_{11} and y_{12} are regared as the unobservable or missing data since we only get their sum y_1 . The complete-data log likelihood is then

$$g_c(\mathbf{y}; \Psi) = C(\mathbf{x}) \left(\frac{1}{2}\right)^{y_{11}} \left(\frac{1}{4}\Psi\right)^{y_{12}} \left(\frac{1}{4} - \frac{1}{4}\Psi\right)^{y_2} \left(\frac{1}{4} - \frac{1}{4}\Psi\right)^{y_3} \left(\frac{1}{4}\Psi\right)^{y_4}$$

Thus, the complete-data log likelihood is, therefore

$$\ln L_c(\Psi) = (y_{12} + y_4) \ln \Psi + (y_2 + y_3) \ln(1 - \Psi) + c.$$

for some constant c not involving Ψ . Let $\Psi^{(0)}$ be the value specified initially for Ψ . Then on the first iteration of the EM algorithm, the E -step requires the computation of the conditional expectation of $L_c(\Psi)$ given \mathbf{y} , using $\Psi^{(0)}$, which can be written as

$$Q(\Psi; \Psi^{(0)}) = \mathbb{E}_{\Psi^{(0)}} [\ln L_c(\Psi) \mid \mathbf{y}].$$

As $\ln L_c(\Psi)$ is a linear function of the unobservable data y_{11} and y_{12} for this problem, the E -step is effected simply by replacing y_{11} and y_{12} by their current conditional expectations given the observed data \mathbf{y} . Considering the random variable Y_{11} corresponding to y_{11} , it is easy to verify that conditional on \mathbf{y} , effectively y_1, Y_{11} has a binomial distribution with sample size y_1 and probability parameter

$$\frac{1}{2} / \left(\frac{1}{2} + \frac{1}{4} \Psi^{(0)} \right),$$

where $\Psi^{(0)}$ is used in place of the unknown parameter Ψ . Thus the initial conditional expectation of Y_{11} given y_1 is

$$\mathbb{E}_{\Psi^{(0)}} [Y_{11} \mid y_1] = y_{11}^{(0)},$$

where

$$\begin{aligned} y_{11}^{(0)} &= \frac{1}{2} y_1 - y_{11}^{(0)} \\ &= \frac{1}{4} y_1 \Psi^{(0)} / \left(\frac{1}{2} + \frac{1}{4} \Psi^{(0)} \right) \end{aligned}$$

The M-step is undertaken on the first iteration by choosing $\Psi^{(1)}$ to be the value of Q that maximizes $Q(\Psi; \Psi^{(0)})$ with respect to Ψ . Since this Q -function is given simply by replacing the unobservable frequencies y_{11} and y_{12} with their current conditional expectations $y_{11}^{(0)}$ and $y_{12}^{(0)}$ in the complete-data log likelihood, $\Psi^{(1)}$ is obtained taking the derivative of $\ln L_c(\Psi)$ to find its maximising value, that is,

$$\begin{aligned} \frac{\partial}{\partial \Psi} \ln L_c(\Psi) = 0 &= (y_{12}^{(0)} + y_4) \frac{1}{\Psi} - (y_2 + y_3) \frac{1}{1 - \Psi} \\ \Psi (y_2 + y_3) &= (1 - \Psi) (y_{12}^{(0)} + y_4) \\ \Psi (y_{12}^{(0)} + y_2 + y_3 + y_4) &= (y_{12}^{(0)} + y_4) \\ \Psi &= \frac{y_{12}^{(0)} + y_4}{y_{12}^{(0)} + y_2 + y_3 + y_4} \\ &= \frac{y_{12}^{(0)} + y_4}{n - y_{11}^{(0)}}. \end{aligned}$$

It follows on so alternating the E - and M -steps on the $(k+1)^{th}$ iteration of the EM algorithm that

$$\Psi^{(k+1)} = \frac{y_{12}^{(k)} + y_4}{n - y_{11}^{(k)}}$$

where

$$\begin{aligned} y_{11}^{(k)} &= \frac{1}{2} y_1 / \left(\frac{1}{2} + \frac{1}{4} \Psi^{(k)} \right) \\ y_{12}^{(k)} &= y_1 - y_{11}^{(k)} \end{aligned}$$

HYPOTHESIS TESTING

Methods of Finding Tests.

Definition 62 (Hypothesis). A **hypothesis** is a statement about a population parameter [Cas01, page 373].

Definition 63 (Null and Alternative Hypothesis). The complementary hypotheses in a hypothesis testing problem are called the **null hypothesis** and the **alternative hypothesis**. They are denoted by H_0 and H_1 respectively [Cas01, page 373].

Definition 64 (Hypothesis Test). A **hypothesis testing procedure** or **hypothesis test** is a rule that specifies

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

[Cas01, page 374].

Definition 65 (Acceptance and Rejection Regions). The subset of the sample space for which H_0 will be rejected is called the **rejection region** or **critical region**. The complement of the rejection region is called the **acceptance region** [Cas01, page 374].

Definition 66 (Acceptance and Rejection Regions). The subset of the sample space for which H_0 will be rejected is called the **rejection region** or **critical region**. The complement of the rejection region is called the **acceptance region** [Cas01, page 374].

Definition 67 (Likelihood Ratio Test (LRT)). 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 $0 \leq c \leq 1$ [Cas01, page 375].

Example 68 (Normal LRT). Example taken from [Cas01, page 375]. Let X_1, X_2, \dots, X_n be a random sample from a $N(\theta, 1)$ population. Consider testing $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$. Here θ_0 is a number fixed by the experimenter prior the experiment. Since there is only one value of θ specified by H_0 , the numerator of $\lambda(\mathbf{x})$ is $L(\theta_0 | \mathbf{x})$. In Example 32, the (unrestricted MLE) we found to be \bar{X} , the sample mean. Thus the denominator of $\lambda(\mathbf{x})$ is $L(\bar{x} | \mathbf{x})$. So the LRT statistic is

$$\begin{aligned} \lambda(\mathbf{x}) &= \frac{(2\pi)^{-n/2} \exp \left[-\sum_{i=1}^n (x_i - \theta_0)^2 / 2 \right]}{(2\pi)^{-n/2} \exp \left[-\sum_{i=1}^n (x_i - \bar{x})^2 / 2 \right]} \\ &= \exp \left[\left(-\sum_{i=1}^n (x_i - \theta_0)^2 + \sum_{i=1}^n (x_i - \bar{x})^2 \right) / 2 \right]. \end{aligned}$$

The expression for $\lambda(\mathbf{x})$ can be simplified by noting that

$$\sum_{i=1}^n (x_i - \theta_0)^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \theta_0)^2.$$

Thus the LRT statistic is $\lambda(\mathbf{x}) = \exp \left[-n(\bar{x} - \theta_0)^2 / 2 \right]$.

Theorem 69. If $T(\mathbf{X})$ is a sufficient statistic for θ and $\lambda^*(t)$ and $\lambda(x)$ are the LRT statistics based on T and \mathbf{X} , respectively, then $\lambda^*(T(\mathbf{X})) = \lambda(x)$ for every x in the sample space [Cas01, page 376].

Definition 70 (Type I and Type II Errors). Suppose R denotes the rejection of H_0 for a test. then for $\theta \in \Theta_0$, the test will make a mistake if $\mathbf{x} \in R$, so the probability of a **Type I Error** is $\mathbb{P}_\theta(\mathbf{X} \in R)$. For $\theta \in \Theta_0^c$, the probability of a **Type II Error** is $\mathbb{P}(\mathbf{X} \in R^c)$ [Cas01, page 383].

Definition 71 (Power Function). The **power function** of a hypothesis test with rejection region R is the function of θ defined by $\beta(\theta) = \mathbb{P}_\theta(\mathbf{X} \in R)$ [Cas01, page 383].

Example 72 (Normal Power Function). Example taken from [Cas01, page 375]. Let X_1, X_2, \dots, X_n be a random sample from a $N(\theta, \sigma^2)$ population. AN LRT of $H_0 : \theta \leq \theta_0$ versus $H_1 : \theta > \theta_0$ is a test that rejects H_0 if $(\bar{X} - \theta_0) / (\sigma/\sqrt{n}) > c$. That constant c can be any positive number. The power function of this test is

$$\begin{aligned}\beta(\theta) &= \mathbb{P}\left(\frac{\bar{X} - \theta_0}{\sigma/\sqrt{n}} > c\right) \\ &= \mathbb{P}\left(\frac{\bar{X} - \theta}{\sigma/\sqrt{n}} > c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right) \\ &= \mathbb{P}\left(Z > c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right),\end{aligned}$$

where Z is a standard normal random variable, since $(\bar{X} - \theta) / (\sigma/\sqrt{n}) \sim N(0, 1)$. As θ increases from $-\infty$ to ∞ , it is easy to see that this normal probability increases from 0 to 1. Therefore, it follows that $\beta(\theta)$ is an increasing function of θ , with

$$\lim_{\theta \rightarrow -\infty} \beta(\theta) = 0, \quad \lim_{\theta \rightarrow \infty} \beta(\theta) = 1, \quad \text{and } \beta(\theta_0) = \alpha \text{ if } \mathbb{P}(Z > c) = \alpha.$$

Suppose the experimenter wishes to have a maximum Type I Error probability of 0.1. Suppose, in addition, the experimenter wishes to have a maximum Type II Error probability of 0.2 if $\theta \geq \theta_0 + \sigma$. We now show how to choose c and n to achieve these goals, using a test that rejects $H_0 : \theta \leq \theta_0$ if $(\bar{X} - \theta_0) / (\sigma/\sqrt{n}) > c$. As noted above, the power function of such a test is

$$\beta(\theta) = \mathbb{P}\left(Z > c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right).$$

Because $\beta(\theta)$ is increasing in θ , the requirements will be met if

$$\beta(\theta_0) = 0.1 \quad \text{and} \quad \beta(\theta_0 + \sigma) = 0.8.$$

By choosing $c = 1.28$, we achieve $\beta(\theta_0) = \mathbb{P}(Z > 1.28) = 0.1$, regardless of n . Now we wish to choose n so that $\beta(\theta_0 + \sigma) = \mathbb{P}(Z > 1.28 - \sqrt{n}) = 0.8$. But, $\mathbb{P}(Z > -0.84) = 0.8$. So setting $1.28 - \sqrt{n} = -0.84$ and solving for n yields $n = 4.49$. Of course n must be an integer. SO choosing $c = 1.28$ and $n = 5$ yields a test with error probabilities controlled as specified by the experimenter.

BAYESIAN INFERENCE

Monte Carlo Sampling.

Definition 73 (Empirical distribution). Let x_1, \dots, x_n be an iid real-valued sample from a cdf F . The function

$$F_n(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{x_i \leq x\}}(x), \quad x \in \mathbb{R}$$

is called the **empirical cdf** of the data [Kro13, page 196].

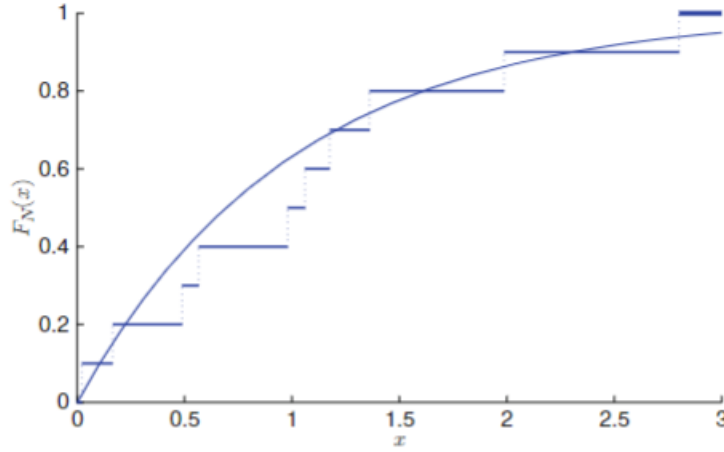


FIGURE 1. The empirical cdf 73 for a sample size 10 from a $\text{Exp}(0.2)$ distribution as well as the true cdf. Image from [Kro13, page 197].

Note that for the ordered sample $x_{(1)} < x_{(2)} < \dots < x_{(N)}$,

$$F_N(x_{(i)}) = \frac{i}{N}.$$

assuming for the sake of simplicity that all $\{x_i\}$ take on different values. If instead of deterministic $\{x_i\}$ we take random X_i , then $F_N(x)$ becomes random as well. To distinguish between the deterministic and the random case, let us denote the random empirical cdf by $\hat{F}_N(x)$. We now have

$$\mathbb{P} \left[\hat{F}_N(x) = \frac{i}{N} \right] = \mathbb{P} [X_{(i)} \leq x, X_{(i+1)} > x] = \binom{N}{i} (F(x))^i (1 - F(x))^{N-i}$$

[Kro13, page 198]. The above equation can be summarized as $N\hat{F}_N(x) \sim \text{Bin}(N, F(x))$. Consequently,

$$\mathbb{E} [\hat{F}_N(x)] = F(x)$$

$$\text{Var} [\hat{F}_N(x)] = F(x) (1 - F(x)).$$

Moreover, by the law of large numbers and the central limit theorem, we have

$$\mathbb{P} \left[\lim_{N \rightarrow \infty} \hat{F}_N(x) = F(x) \right] = 1,$$

and

$$\lim_{N \rightarrow \infty} \mathbb{P} \left[\frac{\hat{F}_N(x) - F(x)}{\sqrt{F(x)(1 - F(x)/N)}} \right] = \Phi(z).$$

Definition 74 (Confidence Intervals). Let X_1, \dots, X_n be random variables with a joint distribution depending on a parameter $\theta \in \Theta$. Let $T_1 < T_2$ be functions of the data but not of θ . The random interval (T_1, T_2) is called a **stochastic confidence interval** for θ with confidence $1 - \alpha$ if

$$\mathbb{P}_\theta [T_1 < \theta < T_2] \geq 1 - \alpha \quad \text{for all } \theta \in \Theta.$$

If t_1 and t_2 are the observed values of T_1 and T_2 , then the interval (t_1, t_2) is called the **numerical confidence interval** for θ with confidence $1 - \alpha$ [Kro13, page 128].

An approximate $1 - \alpha$ confidence interval for $F(x)$ is

$$F_N(x) \pm z_{1-\alpha/2} \sqrt{\frac{F_N(x)(1 - F_N(x))}{N}}$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the standard normal distribution. Moreover for the ordered sample $x_{(1)} < x_{(2)} < \dots < x_{(N)}$, an approximate $1 - \alpha$ confidence interval for $F(x_{(i)})$ is

$$\frac{i}{N} \pm z_{1-\alpha/2} \sqrt{\frac{i(1 - i/n)}{N^2}}$$

[Kro13, page 198].

Simultaneous Confidence Bands.

Definition 75 (Kolmogorov–Smirnov statistic). For a continuous F , the statistic of

$$D_n = \sup_{x \in \mathbb{R}} |\hat{F}_N(x) - F(x)|$$

is called the **Kolmogorov–Smirnov statistic** [Kro13, page 200].

Note that this statistic does not depend on F and can be used to test whether iid samples X_1, \dots, X_n come from a specified distribution as well as constructing a simultaneous confidence band for F . The empirical cdf can be used to estimate any cdf, both discrete and continuous, but it is always ‘jumpy’ which may be undesirable, for example when the underlying distribution is continuous. When we have continuous data, we may want a continuous density estimate instead. This leads to our next topic of density estimation.

Kernel Density Estimation.

Definition 76 (Kernel Density Estimator). Let $X_1, \dots, X_n \stackrel{iid}{\sim} f$ from a continuous distribution with cdf F . A **Kernel density estimator** (KDE), also known as a sliding window estimator, with bandwidth h is given by

$$\hat{f}(x; h) = \frac{1}{N} \sum_{i=1}^N K\left(\frac{x_i - x}{h}\right).$$

The function $K(\cdot)$ is the kernel function (“window shape”) and h is the window (“window size”).

There are many choice of K , but we shall just focus on the Gaussian kernel.

Definition 77 (Gaussian Kernel). *The **Gaussian Kernel** uses a kernel of*

$$K(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$$

so that

$$\hat{f}(x; h) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{x_i - x}{h}\right)^2\right)$$

[Kro13, page 201].

Essentially, the Gaussian kernel is just an equally weighted $\frac{1}{N}$ normal mixture model. It has N Gaussian "humps", centered at each observation and has standard deviation h . As it turns out, the choice of the kernel is not that crucial, but the choice of the bandwidth is very important.

REFERENCES

- [Cas01] George and Berger Casella Roger, *Statistical Inference*, Cengage, Mason, OH, 2001 (eng).
- [Kro13] Dirk P and C.C. Chan Kroese Joshua, *Statistical Modeling and Computation*, Springer New York, New York, NY, 2013 (eng).
- [McL08] Geoffrey John and Krishnan McLachlan T. (Thriyambakam) and McLachlan, *The EM algorithm and extensions* / Geoffrey J. McLachlan, Thriyambakam Krishnan., Wiley series in probability and statistics, Wiley-Interscience, 2008.