

# Statistics 2 - Reviewed Notes

Dom Hutchinson

November 28, 2019

## Contents

<b>1</b>	<b>General</b>	<b>2</b>
1.1	Definitions . . . . .	2
1.2	Theorems . . . . .	2
<b>2</b>	<b>Estimation</b>	<b>3</b>
2.1	Likelihood . . . . .	3
2.2	Estimators . . . . .	6
2.3	Confidence Sets . . . . .	6
2.4	Convergence . . . . .	8
2.5	Performance of Estimators . . . . .	10
2.6	Asymptotic Distribution of Estimators . . . . .	13
2.6.1	Confidence Intervals . . . . .	15
2.7	Efficiency of Estimators . . . . .	16
<b>0</b>	<b>Appendix</b>	<b>17</b>
0.1	Notation . . . . .	17
0.2	Definitions . . . . .	17
0.3	Theorems . . . . .	19
0.4	Probability Distributions . . . . .	20
0.5	Identities . . . . .	23
0.5.1	Likelihood . . . . .	23

# 1 General

## 1.1 Definitions

**Definition 1.1 - Probability Space,  $(\Omega, \mathcal{F}, \mathbb{P})$**

A *Probability Space* is a mathematical construct for modelling the real world. A *Probability Space* has three elements

- i)  $\Omega$ , Sample space;
- ii)  $\mathcal{F}$ , Set of events; and,
- iii)  $\mathbb{P}$ , Probability Measure

and must fulfil the following criteria

- i)  $\Omega \in \mathcal{F}$ ;
- ii)  $\forall A \in \mathcal{D} \implies A^c \in \mathcal{F}$ ;
- iii)  $\forall A_0, \dots, A_n \in \mathcal{F} \implies \left( \bigcup_{i=0}^n A_i \right) \in \mathcal{F}$ ;
- iv)  $\mathbb{P}(\Omega) = 1$ ; and,
- v)  $\mathbb{P}\left(\bigcup_{i=0}^n A_i\right) = \sum_{i=0}^n \mathbb{P}(A_i)$  for any  $n$  disjoint  $A_0, \dots, A_n$ .

**Definition 1.2 - Random Variable**

A *Random Variable* is a function which maps an event in the sample space to a value.  $X$  is a random variable if it satisfies the signature

$$X : \Omega \rightarrow \mathbb{R}$$

**Definition 1.3 - Parametric Models**

*Parametric Models* are the class of statistical distributions whose probability mass/density function take parameters. These parameters represent values of interest in the population, such as mean or variance. We generally do not know these values so we estimate them from samples.

**Definition 1.4 - Quantity of Interest**

When analysing distributions it often helps to define *Quantities of Interest* about the distributions (e.g. Mean). These are defined as functions in terms of the parameters  $\tau(\theta)$ . We estimate *Quantities of Interest* by passing estimated values of the parameters  $\hat{\tau} = \tau(\hat{\theta})$ .

## 1.2 Theorems

**Theorem 1.1 - Samples from a Normal Distribution are  $\chi^2$  Distributed**

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$ . Then

$$\begin{aligned} \sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^2} &\sim \chi_n^2 \\ \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2} &\sim \chi_{n-1}^2 \end{aligned}$$

**Theorem 1.2** - Distance between Sample Mean & Population Mean is  $t_r$  Distributed

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$ . Then

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

N.B. We don't need to know  $\sigma^2$  to estimate the distance between  $\bar{X}$  and  $\mu$ .

**Theorem 1.3** - Multidimension Transform of a Random Variable

Consider an  $n$ -dimensional *continuous* random variable  $\mathbf{X} \sim f_{\mathbf{X}}(\cdot)$  which we wish to transform.

Define a continuously differentiable bijective function  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $\mathbf{h} := \mathbf{g}^{-1}$ .

Then if  $\mathbf{Y} := \mathbf{g}(\mathbf{X}) \sim f_{\mathbf{Y}}(\cdot)$  we have

$$f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(\mathbf{h}(\mathbf{y}))J_{\mathbf{h}}(\mathbf{y})$$

where  $J_{\mathbf{h}}(\mathbf{y}) := \left| \det \left( \frac{\partial \mathbf{h}}{\partial \mathbf{y}} \right) \right| = \left| \det \begin{pmatrix} \frac{\partial h_1}{\partial y_1} & \cdots & \frac{\partial h_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_n}{\partial y_1} & \cdots & \frac{\partial h_n}{\partial y_n} \end{pmatrix} \right|$ .

**Theorem 1.4** - Weak Law of Large Numbers

Let  $\{X_n\}_{n \in \mathbb{N}}$  be a sequence of independent & identically distributed random variables.

If  $\mathbb{E}(X_i) = \mu < \infty$  then

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow_{\mathbb{P}} \mu$$

**Theorem 1.5** - Central Limit Theorem

Let  $\{X_n\}_{n \in \mathbb{N}}$  be a sequence of independent & indetically distributed with  $\mathbb{E}(X_i) = \mu < \infty$  and  $\text{Var}(X_i) = \sigma^2 < \infty$ . Then

$$\sqrt{\frac{n}{\sigma^2}}(Z_n - \mu) \rightarrow_{\mathcal{D}} Z \sim \text{Normal}(0, 1)$$

## 2 Estimation

### 2.1 Likelihood

**Definition 2.1** - Likelihood Function

The *Likelihood Function* is a family of functions which measure the likely of a certain realisation of a random variable is given the parameters of a model have a certain value.

$$L(\boldsymbol{\theta}; \mathbf{x}) := C f_{\mathbf{X}}(\mathbf{x}; \boldsymbol{\theta}) \text{ for } C > 0$$

where  $\mathbf{X} \sim f_n(\cdot; \boldsymbol{\theta}^*)$  with  $\boldsymbol{\theta}^*$  unknown and  $\mathbf{x}$  is a realisation of  $\mathbf{X}$ .

N.B. *Likelihood Functions* have signature  $L(\cdot; \mathbf{x}) : \boldsymbol{\theta} \rightarrow [0, \infty)$ .

N.B. This is also known as the *Observed Likelihood Function*.

**Definition 2.2** - Log-Likelihood Function

The *Log-Likelihood Function* is the family of functions which are equivalent to the natural log of the *Likelihood Function*.

$$\ell(\boldsymbol{\theta}; \mathbf{x}) := \ln f_n(\mathbf{x}; \boldsymbol{\theta}) + C \text{ for } \underbrace{C}_{\equiv \ln C} \in \mathbb{R}$$

N.B. This is increasing with  $L(\cdot; \mathbf{x})$ .

**Remark 2.1** - *Likelihood for Independent & Identically Distributed Random Variables*

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta)$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ . Then

$$L_n(\theta; \mathbf{x}) := \prod_{i=1}^n L(\theta; x_i)$$

$$\ell_n(\theta; \mathbf{x}) := \sum_{i=1}^n \ell(\theta; x_i)$$

**Theorem 2.1** - *The Likelihood Function is Invariant under Bijective Transformations which are independent of Model Parameters*

Consider  $\mathbf{X} \sim f_{\mathbf{X}}(\cdot; \theta)$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a bijective function which is independent of  $\theta$ .

Define  $\mathbf{Y} := \mathbf{g}(\mathbf{X}) \sim f_{\mathbf{Y}}(\cdot; \theta)$ . Then

$$f_{\mathbf{Y}}(\mathbf{y}; \theta) \propto f_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{y}); \theta)$$

Hence

$$L_{\mathbf{Y}}(\theta; \mathbf{g}(\mathbf{x})) \propto L_{\mathbf{X}}(\theta; \mathbf{x})$$

**Proof 2.1** - *Theorem 2.1*

Consider  $\mathbf{X} \sim f_{\mathbf{X}}(\cdot; \theta)$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a bijective function which is independent of  $\theta$ .

Define  $\mathbf{h} := \mathbf{g}^{-1}$  and  $\mathbf{Y} := \mathbf{g}(\mathbf{X})$ .

We consider the cases where  $\mathbf{X}$  is discrete & continuous independently

*Discrete Case* Let  $\mathbf{X}$  be a discrete random variable. Then

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}; \theta) &= \mathbb{P}(\mathbf{Y} = \mathbf{y}; \theta) \\ &= \mathbb{P}(\mathbf{g}^{-1}(\mathbf{Y}) = \mathbf{g}^{-1}(\mathbf{y}); \theta) \\ &= \mathbb{P}(\mathbf{h}(\mathbf{Y}) = \mathbf{h}(\mathbf{y}); \theta) \\ &= \mathbb{P} * \mathbf{X} = \mathbf{h}(\mathbf{y}); \theta) \\ &= f_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{y}); \theta) \end{aligned}$$

*Continuous Case* Let  $\mathbf{X}$  be a continuous random variable.

By **Theorem 1.3**

$$f_{\mathbf{Y}}(\mathbf{y}; \theta) = f_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{y}); \theta) J_{\mathbf{g}^{-1}}(\mathbf{y})$$

Since  $J_{\mathbf{g}^{-1}}$  is independent of  $\theta$  this case is solved.

In both cases  $L_{\mathbf{Y}}(\theta; \mathbf{y}) = f_{\mathbf{Y}}(\mathbf{y}; \theta) \propto f_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{y}); \theta) = L_{\mathbf{X}}(\theta; \mathbf{x})$ . □

**Definition 2.3** - *Maximum Likelihood Estimate*

Let  $\mathbf{X} \sim f_n(\cdot; \theta)$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

The *Maximum Likelihood Estimate* of  $\mathbf{X}$  is the value  $\hat{\theta} \in \Theta$  which produces the greatest value of the *Likelihood Function* of  $\mathbf{X}$ .

$$\hat{\theta}_{\text{MLE}}(\mathbf{x}) := \operatorname{argmax}_{\theta} L(\theta; \mathbf{x}) = \operatorname{argmax}_{\theta} \ell(\theta; \mathbf{x})$$

*N.B.* The *Maximum Likelihood Estimate* is not necessarily unique.

**Theorem 2.2** - *Maximum Likelihood Estimate of Reparameterisation*

Define random variable  $\tau = g(X)$  where  $g : \mathbb{R} \rightarrow \mathbb{R}$ . Then

$$\hat{\tau}_{\text{MLE}} = \tau(\hat{\theta}_{\text{MLE}})$$

**Proof 2.2** - *Theorem 2.2*

*This is a proof by contradiction.*

Suppose  $\exists \tau^* \in G$  st  $\tilde{f}(x; \tau^*) > \tilde{f}(x; \tau)$ .

We know that  $\forall \theta \in \Theta$ ,  $f(x; \theta) = \tilde{f}(x; g(\theta))$  and  $\forall \tau \in G$ ,  $f(x; g^{-1}(\tau)) = \tilde{f}(x; \tau)$ .

We deduce that

$$\begin{aligned} f(x; g^{-1}(\tau^*)) &= \tilde{f}(x; \tau^*) \\ &> \tilde{f}(x; \hat{\tau}) \text{ by assumption} \\ &= f(x; g^{-1}(\hat{\tau})) \\ &= f(x; \hat{\theta}) \end{aligned}$$

This contradicts the assumption that  $\hat{\theta}$  is an maximum likelihood estimate of  $\theta$ .  $\square$

**Remark 2.2 - Finding Maximum Likelihood Estimates - Multivariate**

Let  $X \sim f_X(\cdot; \theta)$  be continuous random variable where  $f_X(\cdot)$  is differentiable and  $\theta$  is an  $n$ -dimensional parameter.

Let  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

To find a *Maximum Likelihood Estimate* for  $\theta$

- i) Find the gradient of  $\ell(\theta; \mathbf{x})$  wrt  $\theta$ .

$$\nabla \ell(\theta; \mathbf{x}) := \left( \frac{\partial}{\partial \theta_1} \ell(\theta; \mathbf{x}) \quad \dots \quad \frac{\partial}{\partial \theta_n} \ell(\theta; \mathbf{x}) \right)$$

- ii) Equate  $\nabla \ell(\theta; \mathbf{x})$  to the zero-vector and solve for each  $\theta$  to find extrema of  $\ell$ .

$$\nabla \ell(\theta; \mathbf{x}) = \mathbf{0}$$

- iii) Calculate the *Hessian* of  $\ell(\theta; \mathbf{x})$

$$\nabla^2 \ell(\theta; \mathbf{x}) = \begin{pmatrix} \frac{\partial^2}{\partial \theta_1^2} \ell(\theta; \mathbf{x}) & \dots & \frac{\partial^2}{\partial \theta_1 \partial \theta_n} \ell(\theta; \mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial \theta_n \partial \theta_1} \ell(\theta; \mathbf{x}) & \dots & \frac{\partial^2}{\partial \theta_n^2} \ell(\theta; \mathbf{x}) \end{pmatrix}$$

- iv) Test each extremum  $\hat{\theta}$  to see if it is a maximum

If  $\det(H(\hat{\theta})) > 0$  and  $\frac{\partial}{\partial \theta^2} \ell(\hat{\theta}; \mathbf{x}) < 0$  then  $\hat{\theta}$  is a local maximum.

*i.e.* Check  $H(\hat{\theta})$  is *negative definite*.

**Definition 2.4 - Likelihood Ratio**

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta^*)$  for  $\theta^* \in \Theta$  and  $\{\hat{\theta}_i\}_{i=1}^n$  be a sequence of consistent *Maximum Likelihood Estimators* of  $\theta^* \in \Theta$ .

We define the *Likelihood Ratio* as

$$\Lambda_n(\mathbf{x}) := \frac{L(\theta^*; \mathbf{x})}{L(\hat{\theta}_n; \mathbf{x})} \in [0, 1] \text{ for } \mathbf{x} \in \mathcal{X}^n$$

**Theorem 2.3 - Asymptotic Distribution of Likelihood Ratio**

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta^*)$  for  $\theta^* \in \Theta$  and  $\{\hat{\theta}_i\}_{i=1}^n$  be a sequence of consistent *Maximum Likelihood Estimators* of  $\theta^* \in \Theta$ .

Suppose the conditions of **Theorem 2.13** hold (*i.e.*  $X_n$  is asymptotically normal). Then

$$-2 \ln \Lambda_n(\mathbf{X}_n) \rightarrow_{\mathcal{D}(\cdot; \theta^*)} \chi_1^2$$

## 2.2 Estimators

### Definition 2.5 - Estimation

Let  $\mathbf{X} \sim f_n(\cdot; \theta^*)$  with  $\theta^* \in \Theta$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

As *Estimation* of model parameter  $\theta^*$  is a statistic,  $\hat{\theta}(\mathbf{x}) = T(\mathbf{x})$ , which is intended to approximate the true value of  $\theta^*$ .

*N.B.* Interchangeable with *Estimate*.

### Definition 2.6 - Estimator

Let  $\mathbf{X} \sim f_n(\cdot; \theta^*)$  with  $\theta^* \in \Theta$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

An *Estimator* of model parameter  $\theta^*$  is the random variable  $\hat{\theta} := \hat{\theta}(\mathbf{X})$  where  $\hat{\theta}(\mathbf{x})$  is an *estimation* of  $\theta^*$ .

### Definition 2.7 - Bias

The *Bias* of an *Estimator*,  $\hat{\theta}$ , is its expected error.

*i.e.* By how much an estimator consistently deviates from the true value of the parameter).

Let  $\theta^*$  be the true value of parameter  $\theta$ . Then

$$\begin{aligned} \text{Bias}(\hat{\theta}; \theta^*) &:= \mathbb{E}(\hat{\theta} - \theta^*; \theta^*) \\ &= \mathbb{E}(\hat{\theta}; \theta^*) - \theta^* \end{aligned}$$

*N.B.* An *Estimator* is *Unbiased* if  $\forall \theta \in \Theta \text{ Bias}(\hat{\theta}; \theta) = 0 \iff \mathbb{E}(\hat{\theta}; \theta) = \theta$ .

### Definition 2.8 - Mean Square Error

The *Mean Square Error* of an *Estimator*,  $\hat{\theta}$ , measures the average of its square error.

Let  $\theta^*$  be the true value of parameter  $\theta$ . Then

$$\begin{aligned} \text{MSE}(\hat{\theta}; \theta^*) &:= \mathbb{E}[(\hat{\theta}(\mathbf{X}) - \theta^*)^2; \theta^*] \\ &= \text{Var}(\hat{\theta}; \theta^*) + \text{Bias}(\hat{\theta}; \theta^*)^2 \end{aligned}$$

### Definition 2.9 - Distribution of an Estimator

Let  $\mathbf{X} \sim f_n(\cdot; \theta^*)$  with  $\theta^* \in \Theta \subseteq \mathbb{R}$ .

Let  $\hat{\theta}(\mathbf{X})$  be a real-valued *Estimator* of  $\theta^*$ . Then

$$\begin{aligned} F_{\hat{\theta}(\mathbf{X})}(t; \theta^*) &:= \mathbb{P}(\hat{\theta}(\mathbf{X}) \leq t; \theta^*) \\ &= \int_{\mathcal{X}^n} \mathbb{1}\{\hat{\theta}(\mathbf{x}) \leq t\} f_n(\mathbf{x}; \theta^*) d\mathbf{x} \end{aligned}$$

*N.B.* The distribution of an *Estimator* depends on the true value of the parameter it is estimating.

*N.B.* As sample size increases the distribution of an estimator should converge to a more standard distribution.

## 2.3 Confidence Sets

### Definition 2.10 - Random Interval

TODO

### Definition 2.11 - Observed Confidence Interval

Let  $\mathbf{X}$  be a random variable,  $\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})]$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

$\mathcal{I}(\mathbf{x}) = [L(\mathbf{x}), U(\mathbf{x})]$  is an *Observed Confidence Interval*.

### Definition 2.12 - Coverage of an Interval

Let  $\mathbf{X} \sim f_n(\cdot; \theta)$  for  $\theta \in \Theta = \mathbb{R}$ .

Define  $L : \mathcal{X}^n \rightarrow \Theta$  &  $U : \mathcal{X}^n \rightarrow \Theta$  st  $\forall \mathbf{x} \in \mathcal{X}^n, L(\mathbf{x}) < U(\mathbf{x})$ .

The *Coverage* of the *Random Interval*  $\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})]$  at  $\theta$  is defined to be

$$C_{\mathcal{I}}(\theta) := \mathbb{P}(\theta \in [L(\mathbf{X}), U(\mathbf{X})]; \theta)$$

*N.B.* *Coverage* is the probability that a realisation of a random variable lies in a given random interval for a given parameter value.

**Definition 2.13 - Confidence Interval**

Let  $\alpha \in [0, 1]$  and  $\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})]$  be a random interval.

We say that  $\mathcal{I}(\mathbf{X})$  is a  $1 - \alpha$  *Confidence Interval* if

$$\forall \theta \in \Theta, C_{\mathcal{I}}(\theta) \geq 1 - \alpha$$

*N.B.*  $\mathcal{I}(\mathbf{X})$  is an Exact *Confidence Interval* if  $\forall \theta \in \Theta, C_{\mathcal{I}}(\theta) = 1 - \alpha$ .

**Proposition 2.1 - Transformed Confidence Interval**

Let  $\mathbf{X} \sim f(\cdot; \theta^*)$  for  $\theta^* \in \Theta$  and  $\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})]$  be a confidence interval for  $\theta^*$ .

Let  $\tau := g(\theta)$  be a bijective, continuously differential function. If

- $g(\cdot)$  is **increasing** then  $[L(\mathbf{x}), U(\mathbf{x})] = [g(L(\mathbf{x})), g(U(\mathbf{x}))]$ .
- $g(\cdot)$  is **decreasing** then  $[L(\mathbf{x}), U(\mathbf{x})] = [g(U(\mathbf{x})), g(L(\mathbf{x}))]$ .

**Proposition 2.2 - Confidence Interval for Reparameterisations**

Let  $\mathbf{X}_n \sim f(\cdot; \theta^*)$  for  $\theta^* \in \Theta \subseteq \mathbb{R}$  and  $\tau_n := g(\theta)$  be a bijective & continuously differentiable function.

When  $\mathbf{X}_n$  is a regular statistical model we have

$$\sqrt{n\tilde{I}(\tau^*)}(\hat{\tau}_n - \tau^*) \rightarrow_{\mathcal{D}(\cdot; \tau^*)} Z \sim \text{Normal}(0, 1)$$

which leads to the *Confidence Interval*

$$\tilde{\mathcal{I}}(\mathbf{X}) := [\tilde{L}(\mathbf{X}), \tilde{U}(\mathbf{X})] \text{ where } \tilde{L}(\mathbf{X}) = \hat{\tau}_n - z_{\alpha/2} \sqrt{\frac{g'(\theta^*)^2}{nI(\theta^*)}} \text{ and } \tilde{U}(\mathbf{X}) = \hat{\tau}_n + z_{\alpha/2} \sqrt{\frac{g'(\theta^*)^2}{nI(\theta^*)}}$$

*N.B.* This confidence interval is **not** necessarily the same as transforming  $[L(\mathbf{x}), U(\mathbf{x})]$  directly.

**Proposition 2.3 - Confidence Intervals with unknown variance,  $\sigma^2$**

When variance,  $\sigma^2$ , is unknown we can define a consistent sequence of estimators  $\{\hat{\sigma}_n^2\}_{n \in \mathbb{N}}$

$$\hat{\sigma}_n^2 := \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \hat{\mu}_n)^2$$

**Definition 2.14 - Wald's Approach**

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta^*)$  for  $\theta^* \in \Theta \subset \mathbb{R}$ .

Using *Wald's Approach* we can define a confidence interval for  $\theta^*$  using the asymptotic distribution of the *Maximum Likelihood Estimator* for  $\theta^*$ .

$$\mathcal{I}(\tau^*) := [L(\mathbf{X}), U(\mathbf{X})] \text{ where } L(\mathbf{x}) := \hat{\theta}_n - z_{\alpha/2} \sqrt{nI(\theta^*)} \text{ and } U(\mathbf{x}) = \hat{\theta}_n + z_{\alpha/2} \sqrt{nI(\theta^*)}$$

*N.B.* This definition ensures that as  $\mathbb{P}(\theta \in \mathcal{I}(\mathbf{X})) \xrightarrow[n \rightarrow \infty]{} 1 - \alpha$ .

**Remark 2.3 - Limitations of Wald's Approach**

Let  $\mathcal{I}(\theta^*)$  be a *Confidence Interval* defined using *Wald's Approach*.

There are certain limitations of *Wald's Approach*

- i) It is possible  $\exists \theta \notin \mathcal{I}(\theta^*)$  st  $\exists \theta' \in \mathcal{I}(\theta^*)$  where  $L(\theta; \mathbf{x}) > L(\theta'; \mathbf{x})$ .
- ii) It is possible  $\exists \theta \in \mathcal{I}(\theta^*)$  where  $L(\theta; \mathbf{x}) = 0$ .
- iii) *Wald Confidence Intervals* are not invariant under reparameterisation.

**Definition 2.15 - Confidence Set**

Let  $\mathbf{X}_n \stackrel{\text{iid}}{\sim} f_n(\cdot; \theta^*)$  for  $\theta^* \in \Theta$  and  $\hat{\theta}_n$  be an estimator of  $\theta$ .

*Confidence Sets* for  $\theta^*$  are the possible values of  $\theta$  whoses likelihood is close to that of the *Maximum Likelihood Estimate* of  $\theta$ .

*Confidence Sets* are not necessarily contiguous.

$$C(\mathbf{X}_n) := \left\{ \theta \in \Theta : \ell(\hat{\theta}_n; \mathbf{X}_n) - \ell(\theta; \mathbf{X}_n) \leq \frac{1}{2} \chi_{1,\alpha}^2 \right\} \subseteq \Theta$$

*Confidence Interval* sets are asymptotically  $1 - \alpha$  for  $\theta^*$  since

$$\mathbb{P}(\theta^* \in C(\mathbf{X}_n); \theta^*) \xrightarrow{n \rightarrow \infty} 1 - \alpha$$

*N.B.* This definition and result are applications of **Definition 2.4** & **Theorem 2.3**.

*N.B.* *Confidence Sets* are hard to define explicitly without a computer.

**Theorem 2.4 - Confidence Set of Reparameterisation**

Let  $\mathbf{X} \sim f(\cdot; \theta^*)$  for  $\theta^* \in \Theta$  and  $\tau := g(\theta)$  where  $g : \Theta \rightarrow G$  is a bijection.

Let  $C(\mathbf{x})$  be a confidence set for  $\theta^*$  and  $\tilde{C}(\mathbf{x})$  be a confidence set for  $\tau^*$ . Then

$$\forall \mathbf{x} \in \mathcal{X}^n, \theta^* \in \Theta \text{ we have } \theta \in C(\mathbf{x}) \iff g(\theta) \in \tilde{C}(\mathbf{x})$$

$$\text{N.B. } \tilde{C}(\mathbf{x}) := \left\{ \theta \in \Theta : \tilde{\ell}_n(\hat{\theta}_n; \mathbf{x}) - \tilde{\ell}(\theta; \mathbf{x}) \leq \frac{1}{2} \chi_{1,\alpha}^2 \right\}.$$

**Proof 2.3 - Theorem 2.4**

Let  $\mathbf{x} \in \mathcal{X}^n$  be arbitrary.

Everything rests on the observation that

$$\forall \theta \in \Theta, \ell(\theta; \mathbf{x}) = \ln f(\mathbf{x}; \theta) = \ln f(\mathbf{x}; g(\theta)) = \tilde{\ell}(g(\theta); \mathbf{x})$$

and similiary

$$\forall \tau \in G, \tilde{\ell}(\tau; \mathbf{x}) = \ln \tilde{f}(\mathbf{x}; \tau) = \ln f(\mathbf{x}; g^{-1}(\tau)) = \ell(g^{-1}(\tau); \mathbf{x})$$

Note that  $g(\hat{\theta}_n)$  is the *Maximum Likelihood Estimate* of  $\tau$ .

Assume  $\theta \in C(\mathbf{x})$ . Then

$$-2 \left[ \ell(\theta; \mathbf{x}) - \ell(\hat{\theta}_n; \mathbf{x}) \right] \leq \chi_{1,\alpha}^2$$

Thus

$$-2 \left[ \tilde{\ell}(g(\theta); \mathbf{x}) - \tilde{\ell}(g(\hat{\theta}_n); \mathbf{x}) \right] \leq \chi_{1,\alpha}^2$$

Thus  $g(\theta) \in \tilde{C}(\mathbf{x})$ .

So  $\theta \in C(\mathbf{x}) \implies g(\theta) \in \tilde{C}(\mathbf{x})$ .

Similarly, assume that  $g(\theta) \in \tilde{C}(\mathbf{x})$ . Thus

$$-2 \left[ \ell(\theta; \mathbf{x}) - \ell(\hat{\theta}_n; \mathbf{x}) \right] \leq \chi_{1,\alpha}^2$$

Thus  $\theta \in C(\mathbf{x})$ .

So  $\theta \in C(\mathbf{x}) \iff g(\theta) \in \tilde{C}(\mathbf{x})$ .



For the last part, this correspondence implies that

$$\{\mathbf{x} \in \chi^n; \theta^* \in C(\mathbf{x})\} = \{\mathbf{x} \in \chi^2 : g(\theta^*) \in \tilde{C}(\mathbf{x})\}$$

Thus, we can conclude from the equivalence of the events

$$\{\theta^* \in C(\mathbf{X}) = \{g(\theta^*) \in \tilde{C}(\mathbf{X})\}$$

**Remark 2.4 - Confidence Set Rule of Thumb**

Under the conditions of **Theorem 2.3** there is a rule of thumb that

$$\mathbb{P}(\theta^* \in C(\mathbf{x})) \approx 0.95 \text{ where } C \approx \left\{ \theta \in \Theta : \ell(\hat{\theta}_n; \mathbf{x}) - \ell(\theta; \mathbf{x}) \leq 2 \right\}$$

**Definition 2.16 - Wilk's Approach**

TODO - Is this just confidence sets?

## 2.4 Convergence

**Definition 2.17 - Convergence**

Let  $\{z_n\}_{n \in \mathbb{N}}$  be a deterministic sequence of real values and  $z \in \mathbb{R}$ .

We say  $\{z_n\}$  converges to limit  $z$  if

$$\forall \varepsilon > 0 \exists n_0 \in \mathbb{N} \text{ st } \forall n \geq n_0 \quad |z_n - z| \leq \varepsilon$$

*N.B.* This is the same for vectors.

**Definition 2.18 - Convergence in Probability**

Let  $\{Z_n\}_{n \in \mathbb{N}}$  be a sequence of random variables and  $Z$  be a random variable in the same probability space.

We say that  $\{Z_n\}_{n \in \mathbb{N}}$  Converges in Probability to  $Z$  if

$$\forall \varepsilon > 0 \quad \lim_{n \rightarrow \infty} \mathbb{P}(|Z_n - Z| > \varepsilon) = 0$$

*N.B.* This is denoted as  $Z_n \rightarrow_{\mathbb{P}} Z$ .

**Definition 2.19 - Convergence in Distribution**

Let  $\{Z_n\}_{n \in \mathbb{N}}$  be a sequence of random variables and  $Z$  be a random variable, not necessarily in the same probability space.

We say  $\{Z_n\}_{n \in \mathbb{N}}$  Converges in Distribution to  $Z$  if

$$\forall z \in Z \text{ where } F_Z(z) \text{ is continuous } \lim_{n \rightarrow \infty} F_{Z_n}(z) = F_Z(z)$$

*i.e.*  $F_{X_n}$  converges in value to  $F_X$  as  $n \rightarrow \infty$ .

*N.B.* This is denoted as  $Z_n \rightarrow_{\mathcal{D}} Z$ .

**Definition 2.20 - Convergence in Quadratic Mean**

Let  $\{Z_n\}_{n \in \mathbb{N}}$  be a sequence of random variables and  $Z$  be a random variable, not necessarily in the same probability space.

We say  $\{Z_n\}_{n \in \mathbb{N}}$  Converges in Quadratic Mean to  $Z$  if

$$\lim_{n \rightarrow \infty} \mathbb{E}[(Z_n - Z)^2] = 0$$

*N.B.* This is denoted as  $Z_n \rightarrow_{\text{qm}} Z$ .

**Theorem 2.5** -  $Z_n \rightarrow_{\mathbb{P}} Z \implies Z_n \rightarrow_{\mathcal{D}} Z$

**Theorem 2.6** -  $Z_n \rightarrow_{qm} Z \implies Z_n \rightarrow_{\mathbb{P}} Z$

**Theorem 2.7** -  $Z_n \rightarrow_{\mathbb{P}} a \iff Z_n \rightarrow_{\mathcal{D}} a$  for  $a \in \mathbb{R}$

**Theorem 2.8** - *Continuous Mapping Theorem*

Let  $\{Z_n\}_{n \in \mathbb{N}}$  be a sequence of random variables and  $Z$  be a random variable.

Let  $g : Z \rightarrow G$  be a function which maps from the space of random variable  $Z$  to a space  $G$ .  
Then

i) If  $Z_n \rightarrow_{\mathbb{P}} Z$  then  $g(Z_n) \rightarrow_{\mathbb{P}} g(Z)$ .

ii) If  $Z_n \rightarrow_{\mathcal{D}} Z$  then  $g(Z_n) \rightarrow_{\mathcal{D}} g(Z)$ .

**Theorem 2.9** - *Slutsky's Theorem*

Let  $\{Y_n\}_{n \in \mathbb{N}}$  &  $\{Z_n\}_{n \in \mathbb{N}}$  be sequences of random variables,  $Y$  be a random variable &  $c \in \mathbb{R} \setminus \{0\}$ .  
If  $Y_n \rightarrow_{\mathcal{D}} Y$  and  $Z_n \rightarrow_{\mathcal{D}} c$ . Then

i)  $Y_n + Z_n \rightarrow_{\mathcal{D}} Y + c$ .

ii)  $Y_n Z_n \rightarrow_{\mathcal{D}} Y c$ .

iii)  $\frac{Y_n}{Z_n} \rightarrow_{\mathcal{D}} \frac{Y}{c}$ .

**Definition 2.21** - *Consistent Sequence of Estimators*

Let  $\mathbf{X}_n \sim f_n(\cdot; \theta)$  be a random vector and  $\{\hat{\theta}_n(\cdot) : \mathcal{X}^n \rightarrow \Theta\}_{n \in \mathbb{N}}$  be a sequence of estimators for  $\theta$ .

We say  $\{\hat{\theta}_n\}$  is *Consistent* if

$$\forall \theta \in \Theta \quad \hat{\theta}_n(\mathbf{X}_n) \rightarrow_{\mathbb{P}(\cdot; \theta)} \theta$$

**Theorem 2.10** -  $\hat{\theta}_n \rightarrow_{qm} \theta \implies \{\hat{\theta}_n\}$  is consistent

## 2.5 Performance of Estimators

**Remark 2.5** - *Measuring Performance of an Estimator*

We measure the performance of an estimator  $\hat{\theta}$  in terms of variance since its mean should be  $\theta^*$  and is thus a bad measure.

Lower variance indicates better performance.

**Definition 2.22** - *Fisher Information Regularity Conditions*

Define  $\Theta \subset \mathbb{R}$  and  $f(x; \theta)$  be a probability mass/density function.

If a model fulfils the following criteria then it is sufficiently *regular* for *Fisher Information* to be drawn from it

i)  $\forall x \in \mathcal{X}$  both  $L'(\theta; x) = \frac{d}{d\theta} f(x; \theta)$  and  $L''(\theta; x) = \frac{d^2}{d\theta^2} f(x; \theta)$  exist.

ii)  $\forall \theta \in \Theta$  the set  $S := \{x \in \mathcal{X} : f(x; \theta) > 0\}$  is independent of  $\theta \in \Theta$ .

iii) The identity below exists

$$\int_S \frac{d}{d\theta} f(x; \theta) dx = \frac{d}{d\theta} \int_S f(x; \theta) dx = 0$$

*N.B.* Statistical Models which fulfil all these criteria are described as *Regular*.

**Definition 2.23** - *Score Function - Single Random Variable*

Let  $X \sim f(\cdot; \theta)$  for some  $\theta \in \Theta$  and  $x$  be a realisation of  $X$ .

The *Score Function* measures the sensitivity of the likelihood function wrt the parameter it is estimating.

$$\ell'(\theta; x) := \frac{d}{d\theta} \ell(\theta; x) = \frac{\frac{d}{d\theta} f(x; \theta)}{f(x; \theta)}$$

**Definition 2.24** - *Score Function - Independent & Identically Distributed Random Variables*

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta)$  with  $\theta \in \Theta$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

$$\ell'_n(\theta; \mathbf{x}) := \sum_{i=1}^n \frac{d}{d\theta} \ell(\theta; x_i)$$

**Proof 2.4** - *By Regularity Conditions*  $\mathbb{E}(\ell'(\theta; X); \theta) = 0 \forall \theta \in \Theta$

$$\begin{aligned} \mathbb{E}(\ell'(\theta; X); \theta) &= \int_S \frac{\frac{d}{d\theta} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx \\ &= \int_S \frac{d}{d\theta} f(x; \theta) dx \\ &= \frac{d}{d\theta} \int_S f(x; \theta) dx \\ &= \frac{d}{d\theta} (1) \\ &= 0 \forall \theta \in \Theta \end{aligned}$$

**Definition 2.25** - *Fisher Information - Single Random Variable*

Let  $X \sim f(\cdot; \theta)$  be an sufficiently regular (see **Definition 2.14**) observable random variable with  $\theta$  unknown.

*Fisher Information* measures the amount of information  $X$  carries about  $\theta$ .

$$\begin{aligned} I(\theta) &:= \mathbb{E}(\ell'(\theta; X)^2; \theta) \\ &= \text{Var}(\ell'(\theta; X); \theta) \text{ by } \mathbf{Proof 2.3} \end{aligned}$$

*N.B.* This is the expectation of the score, squared  $\equiv$  The second moment of the score.

**Definition 2.26** - *Fisher Information - Independent & Identically Distributed Random Variables*

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta)$  with  $\theta \in \Theta$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

$$\begin{aligned} I_n(\theta) &:= \mathbb{E}(\ell'_n(\theta; \mathbf{X})^2; \theta) \\ &= \text{Var}(\ell'_n(\theta; \mathbf{X}); \theta) \\ &= nI(\theta) \end{aligned}$$

**Definition 2.27** - *Observed Fisher Information*

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} f(\cdot; \theta^*)$  be a random  $n$ -dimensional vector.

The *Observed Fisher Information* at  $\theta$  is

$$nJ_n(\theta) = -\ell''(\theta; \mathbf{X}) = -\sum_{i=1}^n \ell''(\theta; X_i)$$

*N.B.*  $\mathbb{E}(J_n(\theta^*; \theta^*)) = I(\theta^*)$ . This is a deterministic value, not an expectation like *Fisher Information*.

**Theorem 2.11 - Fisher Information of Reparameterisation**

Let  $X \sim f(\cdot; \theta)$  for  $\theta \in \Theta \subseteq \mathbb{R}$  and  $\tau := g(\theta)$  be a bijective & continuously differentiable function. Consider the reparameterisation  $\tilde{f}(x; \tau) := f(x; g(\theta)) = f(x; g^{-1}(\tau))$ . The Fisher Information for this reparameterisation,  $\tilde{f}$  is given by

$$\tilde{I}(\tau) = \frac{I(\theta)}{g'(\theta)^2}$$

**Proof 2.5 - Theorem 2.9**

Since  $\tilde{f}(x; \tau) = f(x; g^{-1}(\tau))$  the log-likelihood for  $\tau$  is

$$\tilde{\ell}(\tau; x) = \ln \tilde{f}(x; \tau) = \ln f(x; g^{-1}(\tau))$$

The score is therefore

$$\begin{aligned} \tilde{\ell}'(\tau; x) &= \frac{d}{d\tau} \ln f(x; g^{-1}(\tau)) \\ &= \frac{d}{d\theta} \ln f(x; g^{-1}(\tau)) \times \frac{d}{d\tau} g^{-1}(\tau) \\ &= \ell'(g^{-1}(\tau); x) \times \frac{1}{g'(g^{-1}(\tau))} \\ &= \frac{\ell'(\theta; x)}{g'(\theta)} \end{aligned}$$

No we use the definition of Fisher Information

$$\begin{aligned} \tilde{I}(\tau) &= \mathbb{E}(\tilde{\ell}'(\tau; X)^2; \tau) \\ &= \mathbb{E}\left(\frac{\ell'(\theta; X)^2}{g'(\theta)^2}; \theta\right) \\ &= \frac{1}{g'(\theta)^2} \mathbb{E}(\ell'(\theta; X)^2; \theta) \\ &= \frac{I(\theta)}{g'(\theta)^2} \end{aligned}$$

**Theorem 2.12 - Alternative Expression of Fisher Information**

Let  $X \sim f(\cdot; \theta)$  be a sufficiently regular random variable. Then

$$\text{if } \forall \theta \in \Theta \int_{\mathcal{X}} \frac{d^2}{d\theta^2} f(x; \theta) dx = \frac{d}{d\theta} \int_{\mathcal{X}} \frac{d}{d\theta} f(x; \theta) dx \text{ then } I(\theta) = -\mathbb{E}\left(\frac{d^2}{d\theta^2} \ell(\theta; X); \theta\right)$$

**Proof 2.6 - Theorem 2.9**

By the Quotient Rule

$$\begin{aligned} \frac{d^2}{d\theta^2} \ell(\theta; x) &= \frac{d}{d\theta} \frac{\frac{d}{d\theta} f(x; \theta)}{f(x; \theta)} \\ &= \frac{\frac{d^2}{d\theta^2} f(x; \theta)}{f(x; \theta)} - \left(\frac{\frac{d}{d\theta} f(x; \theta)}{f(x; \theta)}\right)^2 \end{aligned}$$

Consequently

$$\begin{aligned} \mathbb{E}\left(\frac{d^2}{d\theta^2} \ell(\theta; X); \theta\right) &= \int_S \frac{\frac{d^2}{d\theta^2} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx - \int_S \left(\frac{\frac{d}{d\theta} f(x; \theta)}{f(x; \theta)}\right)^2 f(x; \theta) dx \\ &= \int_S \frac{d^2}{d\theta^2} f(x; \theta) dx - \int_S \ell'(\theta; x)^2 f(x; \theta) dx \\ &= 0 - \mathbb{E}(\ell'(\theta; X)^2; \theta) \\ &= -I(\theta) \\ \implies I(\theta) &= -\mathbb{E}\left(\frac{d^2}{d\theta^2} \ell(\theta; X); \theta\right) \end{aligned}$$

□

**Theorem 2.13** - *Distribution of Maximum Likelihood Estimators for Regular Models*

Let  $\mathbf{X}_n \stackrel{\text{iid}}{\sim} f_n(\cdot; \theta^*)$  be a sufficiently regular statistically model and  $\{\hat{\theta}_n\}_{n \in \mathbb{N}}$  be a consistent sequence of *Maximum Likelihood Estimators* for  $\theta^*$ . Then

$$\sqrt{nI(\theta^*)}(\hat{\theta}_n - \theta^*) \rightarrow_{\mathcal{D}(\cdot; \theta^*)} Z \sim \text{Normal}(0, 1)$$

Here  $I(\theta^*)$  is unknown so we replace it with

i)  $I(\hat{\theta}_n)$  when

(a)  $I(\theta)$  is continuous in a neighbourhood of  $\theta^*$ ;

(b) And, the interval  $[L(\mathbf{X}), U(\mathbf{X})]$  with  $L(\mathbf{x}) := \hat{\theta}_n - z_{\alpha/2} \sqrt{nI(\hat{\theta}_n)}$  and  $U(\mathbf{x}) := \hat{\theta}_n + z_{\alpha/2} \sqrt{nI(\hat{\theta}_n)}$  is an asymptotically exact  $1 - \alpha$  confidence interval for  $\theta^*$ .

ii)  $J_n(\hat{\theta}_n) := -\frac{1}{n} \sum_{i=1}^n \ell''(\hat{\theta}_n; X_i)$  when

(a)  $\hat{\theta}_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} \theta^*$ ;

(b)  $I(\theta) = -\mathbb{E}(\ell''(\theta; X); \theta) \forall \theta \in \Theta$ ;

(c)  $\exists C : \mathcal{X} \rightarrow [0, \infty)$  st  $\mathbb{E}(C(X_1); \theta^*) < \infty$ ,  $\Xi \subset \Theta$  is an open set containing  $\theta^*$  and  $\Delta(\cdot) : \Xi \rightarrow [0, \infty)$  is continuous at 0 st  $\Delta(0) = 0$ , and st  $\forall \theta, \theta^*, x \in \Xi^2 \times \mathcal{X}$

$$|\ell''(\theta; x) - \ell''(\theta'; x)| \leq C(x)\Delta(\theta - \theta')$$

(d) And, the interval  $[L(\mathbf{X}), U(\mathbf{X})]$  with  $L(\mathbf{x}) := \hat{\theta}_n - z_{\alpha/2} \sqrt{nJ_n(\hat{\theta}_n)}$  and  $U(\mathbf{x}) := \hat{\theta}_n + z_{\alpha/2} \sqrt{nJ_n(\hat{\theta}_n)}$  is an asymptotically exact  $1 - \alpha$  confidence interval for  $\theta^*$

**Theorem 2.14** - *Cramer-Rao Inequality*

Let *Cramer-Rao Inequality* provides us with a *lower bound* for the performance of all estimators.

Let  $\mathbf{X}_n \stackrel{\text{iid}}{\sim} f(\cdot; \theta)$  be a sufficiently regular random vector and  $\hat{\theta}_n(\cdot)$  be an estimator of  $\theta$  with expectation  $m_1(\theta) := \mathbb{E}(\hat{\theta}_n(\mathbf{X}_n); \theta)$ .

$$\text{if } \forall \theta \in \Theta, \underbrace{\frac{d}{d\theta} \int \hat{\theta}_n(\mathbf{x}) f_n(\mathbf{x}; \theta) d\mathbf{x}}_{\mathbb{E}(\hat{\theta}_n)} = \int \hat{\theta}_n(\mathbf{x}) \frac{d}{d\theta} f_n(\mathbf{x}; \theta) d\mathbf{x}$$

Then

$$\forall \theta \in \Theta, \text{Var}(\hat{\theta}_n(\mathbf{X}_n); \theta) \geq \frac{m'_1(\theta)^2}{nI(\theta)}$$

**Proof 2.7** - *Cramer-Rao Inequality*

We notice that

$$\begin{aligned} m'(\theta) &= \frac{d}{d\theta} \mathbb{E}(\hat{\theta}_n(\mathbf{X}_n); \theta) \\ &= \frac{d}{d\theta} \int_{S^n} \hat{\theta}_n(\mathbf{x}_n) f_n(\mathbf{x}_n; \theta) d\mathbf{x}_n \end{aligned}$$

The clever part of this proof is to observe that

$$\begin{aligned} \text{Var}(\hat{\theta}_n(\mathbf{X}_n); \theta) nI(\theta) &= \text{Var}(\hat{\theta}_n(\mathbf{X}_n); \theta) \text{Var}(\ell_n(\theta; \mathbf{X}_n); \theta) \\ &\geq \text{Cov}(\hat{\theta}_n(\mathbf{X}_n), \ell'_n(\theta; \mathbf{X}_n); \theta)^2 \text{ by Covariance Inequality} \end{aligned}$$

Thus

$$\begin{aligned}
\text{Cov}(\hat{\theta}_n(X_n), \ell'_n(\theta; \mathbf{X}_n); \theta)^2 &= \mathbb{E}(\hat{\theta}_n(X_n) \ell'_n(\theta; \mathbf{X}_n); \theta) - \mathbb{E}(\hat{\theta}_n(\mathbf{X}_n); \theta) \mathbb{E}(\ell'_n(\theta; \mathbf{X}_n); \theta) \\
&= \mathbb{E}(\hat{\theta}_n(X_n) \ell'_n(\theta; \mathbf{X}_n); \theta) - \mathbb{E}(\hat{\theta}_n(\mathbf{X}_n); \theta) \times 0 \\
&= \mathbb{E}(\hat{\theta}_n(X_n) \ell'_n(\theta; \mathbf{X}_n); \theta) \\
&= \int_{S^n} \hat{\theta}_n(\mathbf{x}_n) \ell'_n(\theta; \mathbf{x}_n) f_n(\mathbf{x}_n; \theta) d\mathbf{x}_n \\
&= \int_{S^n} \hat{\theta}_n(\mathbf{x}_n) \frac{\frac{d}{d\theta} f_n(\mathbf{x}_n; \theta)}{f_n(\mathbf{x}_n; \theta)} f_n(\mathbf{x}_n; \theta) d\mathbf{x}_n \\
&= \int_{S^n} \hat{\theta}_n(\mathbf{x}_n) \frac{d}{d\theta} f_n(\mathbf{x}_n; \theta) \\
&= \frac{d}{d\theta} \int_{S^n} \hat{\theta}_n(\mathbf{x}_n) f_n(\mathbf{x}_n; \theta) d\mathbf{x}_n \text{ by regularity assumption} \\
&= m'(\theta) \\
\implies \text{Var}(\hat{\theta}_n(X_n); \theta) n I(\theta) &\geq m'(\theta)^2
\end{aligned}$$

**Remark 2.6 - Cramer-Rao Inequality with an Unbiased Estimator**

Let  $\hat{\theta}_n$  be an unbiased estimator of  $\theta$  (i.e.  $m_1(\theta) = \theta$ ). Then

$$\text{Var}(\hat{\theta}_n(\mathbf{X}_n); \theta) = \text{MSE}(\hat{\theta}_n(\mathbf{X}_n); \theta) \geq \frac{1}{nI(\theta)}$$

## 2.6 Asymptotic Distribution of Estimators

**Theorem 2.15 - Asymptotic Distribution of Maximum Likelihood Estimators**

Suppose that  $\mathbf{X}_n \stackrel{\text{iid}}{\sim} f(\cdot; \theta^*)$  for some  $\theta^* \in \Theta$  and assume that

- i) The sequence of maximum likelihood estimators  $\{\hat{\theta}_n(\mathbf{X}_n)\}$  is consistent;
- ii) The *Fisher Information Regularity Conditions* (**Definition 6.2**) hold and  $I(\theta^*) = -\mathbb{E}[\ell''(\theta; X); \theta] > 0$ .
- iii)  $\exists C : \mathcal{X} \rightarrow [0, \infty)$  such that  $\mathbb{E}[C(X_1); \theta^*] < \infty$  and  $\Delta : \Xi \rightarrow [0, \infty)$ , where  $\Xi \subset \Theta$  st  $\theta^* \in \Xi$ , that is continuous at 0 st  $\Delta(0) = 0$ , such that

$$\forall (\theta, \theta', x') \in \chi^2 \times \mathcal{X}, \quad |\ell''(\theta; x) - \ell''(\theta'; x)| \leq C(x) \Delta(\theta - \theta')$$

Then  $\forall \theta^* \in \Theta$

$$\sqrt{nI(\theta^*)}(\hat{\theta}_n(\mathbf{X}_n) - \theta^*) \rightarrow_{\mathcal{D}(\cdot; \theta^*)} Z \sim \text{Normal}(0, 1)$$

**Proof 2.8 - Theorem 2.11**

By **Theorem 2.11**  $\ell'_n(\hat{\theta}_n; \mathbf{X}) = \ell'_n(\theta^*; \mathbf{X}) + (\hat{\theta}_n - \theta^*)[\ell''_n(\theta^*; \mathbf{X}) + R_n]$  where  $\frac{1}{n}R_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} 0$ .

Since  $\hat{\theta}_n$  is the maximum likelihood estimator & the *Fisher Information Regularity Conditions* hold, the score at  $\ell'(\hat{\theta}_n; X) = 0$ .

Hence,  $0 = \ell''(\hat{\theta}_n; X) = \ell'_n(\theta; X) + (\hat{\theta}_n - \theta^*)\{\ell''(\theta; X) + R_n\}$ .

Rearranging & rescaling by  $\sqrt{n}$  gives

$$\sqrt{n}(\hat{\theta}_n - \theta^*) = \frac{\frac{1}{\sqrt{n}}\ell'(\theta^*; X)}{-\frac{1}{\sqrt{n}}\{\ell''(\theta^*; X) + R_n\}} =: \frac{U_n}{V_n - \frac{R_n}{n}}$$

Recall that  $\ell'_n(\theta^*; X) = \sum_{i=1}^n \ell'(\theta; X_i)$  and  $\ell''_n(\theta^*; X) = \sum_{i=1}^n \ell''(\theta^*; X_i)$ .

Since  $\mathbb{E}(\ell'(\theta^*; X_i); \theta^*) = 0$  and  $\text{Var}(\ell'(\theta^*; X_i); \theta^*) = I(\theta^*)$

$\implies U_n \rightarrow_{\mathcal{D}(\cdot; \theta^*)} U \sim \text{Normal}(0, I(\theta^*))$  by the *Central Limit Theorem*.

We observed that  $V_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} I(\theta^*)$  by the *Weak Law of Large Numbers* since  $\mathbb{E}(-\ell''(\theta^*; X_i); \theta^*) = I(\theta^*)$ .

It follows that  $V_n - \frac{1}{n}R_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} I(\theta^*)$  by *Slutsky's Theorem*.

Using *Slutsky's Theorem* again

$$\sqrt{n}(\hat{\theta}_n - \theta^*) = \frac{U_n}{V_n - \frac{1}{n}R_n} \rightarrow_{\mathcal{D}(\cdot; \theta^*)} \frac{\sqrt{I(\theta^*)}}{I(\theta^*)} Z \text{ where } Z \sim \text{Normal}(0, 1)$$

We can rewrite this as

$$\sqrt{nI(\theta^*)}(\hat{\theta}_n - \theta^*) \rightarrow_{\mathcal{D}(\cdot; \theta^*)} Z \sim \text{Normal}(0, 1)$$

**Theorem 2.16** - *Convergence of Score of Maximum Likelihood Estimators*

Under the conditions in **Theorem 2.11**, with  $\hat{\theta}_n$  a Maximum Likelihood Estimator

$$\ell'_n(\hat{\theta}_n; \mathbf{X}) = \ell'_n(\theta^*; \mathbf{X}) + (\hat{\theta}_n - \theta^*)[\ell''_n(\theta^*; \mathbf{X}) + R_n]$$

where  $\frac{1}{n}R_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} 0$ .

**Proof 2.9** - *Theorem 2.12*

*This is an non-examinable, sketch proof of Theorem 8.2.*

By the regularity conditions and the mean value theorem

$$\frac{\ell'_n(\theta; \mathbf{x}) - \ell'_n(\theta^*; \mathbf{x})}{\theta - \theta^*} = \ell''_n(\tilde{\theta}; \mathbf{x})$$

for some  $\tilde{\theta} \in (\theta, \theta^*)$ . Hence, we deduce that

$$\begin{aligned} \ell'_n(\theta; \mathbf{x}) - \ell'_n(\theta^*; \mathbf{x}) &= (\theta - \theta^*)\ell''_n(\tilde{\theta}; \mathbf{x}) \\ &= (\theta - \theta^*)\{\ell''_n(\theta^*; \mathbf{x}) + [\ell''_n(\tilde{\theta}; \mathbf{x}) - \ell''_n(\theta^*; \mathbf{x})]\} \\ &= (\theta - \theta^*)\{\ell''_n(\theta; \mathbf{x}) + R_n(\theta, \theta^*, \mathbf{x})\} \end{aligned}$$

Now we replace  $\theta$  with the maximum likelihood estimator  $\hat{\theta}_n := \hat{\theta}_n(\mathbf{X})$ . We find

$$\ell'(\hat{\theta}_n; \mathbf{X}) = \ell'_n(\theta^*; \mathbf{X}) + (\hat{\theta}_n - \theta^*)\{\ell''_n(\theta^*; \mathbf{X}) + R_n(\hat{\theta}_n, \theta^*, \mathbf{x})\}$$

and we need to analyse  $R_n$ .

Since  $\hat{\theta}_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} \theta^*$  we can take  $n$  large enough that  $\mathbb{P}(\hat{\theta}_n \in \Xi; \theta^*)$  with arbitrarily high probability.

On the event  $\{\hat{\theta} \in \Xi\}$  and we have  $\{\tilde{\theta}_n \in \Xi\}$  since  $\tilde{\theta}_n \in (\hat{\theta}_n, \theta^*)$  and

$$\begin{aligned} |\frac{1}{n}R_n| &= \frac{1}{n}|\ell''_n(\tilde{\theta}_n; \mathbf{X}) - \ell''_n(\theta^*; \mathbf{X})| \\ &= \frac{1}{n} \left| \sum_{i=1}^n \ell''(\tilde{\theta}_n; X_i) - \ell''(\theta^*; X_i) \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n |\ell''(\tilde{\theta}_n; X_i) - \ell''(\theta^*; X_i)| \\ &\leq \Delta(\tilde{\theta}_n - \theta^*) \left\{ \frac{1}{n} \sum_{i=1}^n C(X_i) \right\} \end{aligned}$$

from the smoothness condition on  $\ell''$ .

From the *Weak Law of Large Numbers*

$$\frac{1}{n} \sum_{i=1}^n C(X_i) \rightarrow_{\mathbb{P}(\cdot; \theta^*)} \mathbb{E}(C(X_1); \theta^*) < \infty$$

and from the consistency of  $\{\hat{\theta}_n\}$  and  $\{\tilde{\theta}_n\}$  and continuity of  $\Delta(\cdot)$  we have by the *Continuous Mapping Theorem*

$$\Delta(\tilde{\theta}_n - \theta^*) \rightarrow_{\mathbb{P}(\cdot; \theta^*)} 0$$

Hence,  $\frac{1}{n}R_n \rightarrow_{\mathbb{P}(\cdot; \theta^*)} 0$  □

### 2.6.1 Confidence Intervals

**Theorem 2.17** - *Convergence in Distribution of Confidence Intervals*

Let  $\mathbf{X} \sim f(\cdot; \theta^*)$  with  $\theta \in \Theta$  and define  $\{\hat{\theta}_n\}_{n \in \mathbb{N}}$  be a consistent sequence of estimators of  $\theta^*$ . Suppose that  $\{\hat{\theta}_n\}$  is asymptotically normal in the sense that

$$\exists \sigma^2 > 0 \text{ st } \frac{\hat{\theta}_n(\mathbf{X}) - \theta^*}{\sqrt{\sigma^2/n}} \rightarrow_{\mathcal{D}(\cdot; \theta^*)} Z \sim \text{Normal}(0, 1)$$

Then

$\forall \alpha \in (0, 1)$ ,  $\mathcal{I}_n(\mathbf{X}) = [L_n(\mathbf{X}), U_n(\mathbf{X})]$  is an asymptotically exact  $1 - \alpha$  confidence interval

where  $L_n(\mathbf{x}) := \hat{\theta}_n(\mathbf{x}) - z_{\alpha/2} \sqrt{\frac{\sigma^2}{n}}$  and  $U_n(\mathbf{x}) := \hat{\theta}_n(\mathbf{x}) + z_{\alpha/2} \sqrt{\frac{\sigma^2}{n}}$ .

**Proof 2.10** - *Theorem 2.13*

Let  $\{W_n\}_{n \in \mathbb{N}}$  be defined by  $W_n := \frac{\hat{\theta}_n(X) - \theta^*}{\sqrt{\sigma^2/n}}$ .

Since  $W_n \rightarrow_{\mathcal{D}(\cdot; \theta^*)} Z \sim \text{Normal}(0, 1)$  we have

$$\begin{aligned} \mathbb{P}(-z_{\alpha/2} \leq W_n \leq z_{\alpha/2}) &= F_{W_n}(z_{\alpha/2}) - F_{W_n}(-z_{\alpha/2}) \\ &\xrightarrow{n \rightarrow \infty} \Phi(z_{\alpha/2}) - \Phi(-z_{\alpha/2}) \\ &= 1 - \alpha \end{aligned}$$

Similary to before we have the equivalence of events

$$\{-z_{\alpha/2} \leq W_n \leq z_{\alpha/2}\} = \left\{ \hat{\theta}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \theta^* \leq \hat{\theta}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right\}$$

So  $\lim_{n \rightarrow \infty} \mathbb{P} \left( \hat{\theta}_n(X) - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \theta^* \leq \hat{\theta}_n(X) + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}; \theta^* \right) = 1 - \alpha$ . □

### 2.7 Efficiency of Estimators

**Definition 2.28** - *Efficient Estimator*

Let  $\hat{\theta}$  be an estimator of parameter  $\theta$ .

$\hat{\theta}$  is said to be an *Efficient Estimator* if its variance is equal to the *Cramer-Rao Lower Bound*  $\forall \theta^*$ .

$$\forall \theta^*, \text{Var}(\hat{\theta}; \theta^*) = \frac{m'(\theta^*)^2}{nI(\theta)}$$

**Definition 2.29** - *Asymptotically Efficient Sequence of Estimators*

Let  $\mathbf{X} \sim f(\cdot; \theta)$  for  $\theta \in \Theta$  and  $\{\hat{\theta}_n(\mathbf{X})\}_{n \in \mathbb{N}}$  be a sequence of estimators.

The sequence  $\{\hat{\theta}_n\}$  is *Asymptotically Efficient* if either

i) its *Mean-Squared Error* converges in value to the *Cramer-Rao Lower Bound*

$$\forall \theta \in \Theta, n\text{MSE}(\hat{\theta}_n(\mathbf{X}_n); \theta) \xrightarrow{n \rightarrow \infty} \frac{1}{I(\theta)}$$

ii) Or,  $\hat{\theta}_n$  *Converges in Distribution* to a standard Normal

$$\forall \theta \in \Theta, \sqrt{nI(\theta)}(\hat{\theta} - \theta) \rightarrow_{\mathcal{D}(\cdot; \theta)} Z \sim \text{Normal}(0, 1)$$

**Remark 2.7** - *Under the conditions of Theorem 2.11 Maximum Likelihood Estimators are Asymptotically Efficient*



## 0 Appendix

### 0.1 Notation

**Notation 0.1 - Convergence**

$\{z_n\}_{n \in \mathbb{N}} \rightarrow z$  denotes that the sequence of deterministic values  $\{z_n\}_{n \in \mathbb{N}}$  converges in value to  $z \in \mathbb{R}$ .

$\{Z_n\}_{n \in \mathbb{N}} \rightarrow_{\mathbb{P}} Z$  denotes that the sequence of random variables  $\{Z_n\}_{n \in \mathbb{N}}$  converges in probability to random variable  $Z$ .

$\{Z_n\}_{n \in \mathbb{N}} \rightarrow_{\mathbb{P}(\cdot; \theta)} Z$  denotes that the sequence of random variables  $\{Z_n\}_{n \in \mathbb{N}}$  converges in probability to random variable  $Z$ , dependent upon parameter  $\theta$ .

$\{Z_n\}_{n \in \mathbb{N}} \rightarrow_{\mathcal{D}} Z$  denotes that the sequence of random variables  $\{Z_n\}_{n \in \mathbb{N}}$  converges in distribution to random variable  $Z$ .

**Notation 0.2 - Gamma Function**

$$\Gamma(x) := \int_0^{\infty} t^{x-1} e^{-t} dt$$

### 0.2 Definitions

**Definition 0.1 - Correlation**

Let  $X$  &  $Y$  be random variables.

*Correlation* is a measure of dependence between two random variables

$$\text{Corr}(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \in [-1, 1]$$

**Definition 0.2 - Covariance**

*Covariance* is a measure of the joint variability of two random variables.

Consider random variable  $X$  &  $Y$

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

If  $X$  &  $Y$  are independent then  $\text{Cov}(X, Y) = 0$ .

By definition of *Covariance*  $\text{Cov}(X, X) = \text{Var}(X)$ .

**Definition 0.3 - Estimation**

Let  $\mathbf{X} \sim f_n(\cdot; \theta^*)$  with  $\theta^* \in \Theta$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

As *Estimation* of model parameter  $\theta^*$  is a statistic,  $\hat{\theta}(\mathbf{x}) = T(\mathbf{x})$ , which is intended to approximate the true value of  $\theta^*$ .

*N.B.* Interchangeable with *Estimate*.

**Definition 0.4 - Estimator**

Let  $\mathbf{X} \sim f_n(\cdot; \theta^*)$  with  $\theta^* \in \Theta$  and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ .

An *Estimator* of model parameter  $\theta^*$  is the random variable  $\hat{\theta} := \hat{\theta}(\mathbf{X})$  where  $\hat{\theta}(\mathbf{x})$  is an *estimation* of  $\theta^*$ .

**Definition 0.5 - Expectation**

*Expectation* is the mean value for a random variable.

Consider *continuous* random variable  $X$  with pdf  $f_X$  and *discrete* random variable  $Y$  with pmf  $f_Y$ . Then

$$\mathbb{E}(X) := \int_{\mathbb{R}} x f_X(x) dx \quad \text{and} \quad \mathbb{E}(Y) := \sum_{y \in \mathcal{Y}} y p_Y(y)$$

For a function  $g : \mathbb{R} \rightarrow \mathbb{R}$  we have

$$\mathbb{E}(g(X)) := \int_{\mathbb{R}} g(x)f_X(x)dx \quad \text{and} \quad \mathbb{E}(g(Y)) := \sum_{y \in \mathcal{Y}} g(y)p_Y(y)$$

For linear transformations of a random variable  $Z$  we find

$$\mathbb{E}(aZ + b) = a\mathbb{E}(Z) + b \quad \text{for } a, b \in \mathbb{R}$$

**Definition 0.6 - Five-Number Summary**

The *Five-Number Summary* of a sample contains the sample's: median; lower hinge; upper hinge; minimum value; & maximum value.

**Definition 0.7 - Hinges**

*Hinges* describe the spread of data in a sample, while trying to ignore extreme data. The *Lower Hinge*,  $H_1$ , is the median of the set containing the median & values with rank less than the sample median. The *Upper Hinge*,  $H_3$ , is the median of the set containing the median & values with rank greater than the sample median.

**Definition 0.8 - Median**

The *Median* is the central value of a data set.

Consider a data set  $x_0, \dots, x_n$

- If  $\exists m \in \mathbb{N}$  st  $n = 2m + 1$  (i.e.  $n$  is odd) then the median is  $x_{(m+1)}$ .
- Else  $\exists m \in \mathbb{N}$  st  $n = 2m$  (i.e.  $n$  is even) then the median is  $x_{(m+1)}$ .

**Definition 0.9 - Moments**

The *Moments* of a random variable  $X$  are the expected values of powers of  $X$ .

$$n^{\text{th}} \text{ moment of } X := \mathbb{E}(X^n)$$

N.B.  $\mathbb{E}(X^n) \neq \mathbb{E}(X)^n$ .

**Definition 0.10 - Order Statistic**

An *Order Statistic* is a data set where the data has been placed in increasing order of value, not time. We use  $x_{(i)}$  to denote the  $i^{\text{th}}$  lowest value in  $(x_0, \dots, x_n)$ .

**Definition 0.11 - Quartiles**

*Quartiles* describe the spread of data in a sample. The *Lower Quartile*,  $Q_1$ , is the median of the set of values with rank less than the sample median. The *Upper Quartile*,  $Q_3$ , is the median of the set of values with rank greater than the sample median.

N.B. These sets do not contain the median.

**Definition 0.12 - Sample Mean**

The *Sample Mean* is the mean value of all data points within a sample. Consider a sample  $\{x_1, \dots, x_n\}$

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$$

**Definition 0.13 - Sample Variance**

*Sample Variance* is a measure of spread of data in a sample around the sample mean. For a sample  $\{x_1, \dots, x_n\}$

$$s^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \left( \sum_{i=1}^n x_i^2 \right) - n\bar{x}^2 \right)$$

**Definition 0.14 - Statistic**

Let  $\mathbf{x}$  be some data.

A *Statistic* is any function of the data,  $T(\mathbf{x})$ .

*N.B.* *Statistics* are independent of unknown model parameters.

**Definition 0.15 - Trimmed Sample Mean**

The *Trimmed Sample Mean* is the average value of a subset of data points within a sample. The subset is defined to ignore the  $\frac{\Delta}{2}\%$  largest & smallest values of the sample. For a  $\Delta\%$  trimmed mean we define

$$\bar{x}_{\Delta} := \frac{1}{n - 2k} \sum_{i=k+1}^{n-k-1} x_i \text{ with } k = \left\lfloor \frac{n\Delta}{100} \right\rfloor$$

**Definition 0.16 - Variance**

*Variance* measures how far a set of random numbers are spread from their average value.

Consider random variable  $X$

$$\text{Var}(X) := \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}(X^2) - \mathbb{E}(X)^2$$

For linear transformation of a random variable  $X$  we find

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

For a linear transformation of two random variables  $X$  &  $Y$  we have

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y) \quad \text{for } a, b \in \mathbb{R}$$

**Definition 0.17 - Skew**

*Skew* describes the spread of values in a sample which are less than the median, relative to the spread of values greater than the median. A sample is *Left-Skewed* if  $|H_3 - H_2| < |H_1 - H_2|$ . A sample is *Right-Skewed* if  $|H_3 - H_2| > |H_1 - H_2|$ .

**0.3 Theorems****Theorem 0.1 - Cauchy-Schwarz Inequality**

Let  $X$  &  $Y$  be real-valued random variables in the same probability space. Then

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

**Theorem 0.2 - Chebyshev's Inequality**

Let  $X$  be a random variable.

Define  $\mu := \mathbb{E}(X)$  and  $\sigma^2 := \text{Var}(X)$ . Then

$$\forall a > 0 \quad \mathbb{P}(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2}$$

**Theorem 0.3 - Covariance Inequality**

Let  $X$  &  $Y$  be real-valued random variables in the same probability space. Then

$$\text{Cov}(X, Y)^2 \leq \text{Var}(X)\text{Var}(Y)$$

**Theorem 0.4 - Joint Probability Density of Simple Random Sample**

Let  $\mathbf{X}_1, \dots, X_n$  be a set of independent random variables with pdfs  $f_{X_1}, \dots, f_{X_n}$ , respectively,

and  $x_1, \dots, x_n$  be a realisation of  $X_1, \dots, X_n$ .

The probability of obtaining  $x_1, \dots, x_n$  is

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i; \theta)$$

**Theorem 0.5 - Markov's Inequality**

Let  $X \sim f_X(\cdot)$  be a non-negative continuous random variable. Then

$$\forall a > 0 \quad \mathbb{P}(X \geq a) \leq \frac{\mathbb{E}(X)}{a}$$

## 0.4 Probability Distributions

**Definition 0.18 - Bernoulli Distribution**

Let  $X \sim \text{Bernoulli}(p)$ .

A *discrete* random variable which takes 1 with probability  $p$  & 0 with probability  $(1 - p)$ . Then

$$\begin{aligned} p_X(k) &= \begin{cases} 1 - p & \text{if } k = 0 \\ p & \text{if } k = 1 \\ 0 & \text{otherwise} \end{cases} \\ P_X(k) &= \begin{cases} 0 & \text{if } k < 0 \\ 1 - p & \text{if } k \in [0, 1) \\ 1 & \text{otherwise} \end{cases} \\ \mathbb{E}(X) &= p \\ \text{Var}(X) &= p(1 - p) \\ \mathcal{M}_X(t) &= (1 - p) + pe^t \end{aligned}$$

N.B. Often we define  $q := 1 - p$  for simplicity.

**Definition 0.19 - Binomial Distribution**

Let  $X \sim \text{Binomial}(n, p)$ .

A *discrete* random variable modelled by a *Binomial Distribution* on  $n$  independent events and rate of success  $p$ .

$$\begin{aligned} p_X(k) &= \binom{n}{k} p^k (1 - p)^{n-k} \\ P_X(k) &= \sum_{i=1}^k \binom{n}{i} p^i (1 - p)^{n-i} \\ \mathbb{E}(X) &= np \\ \text{Var}(X) &= np(1 - p) \\ \mathcal{M}_X(t) &= [(1 - p) + pe^t]^n \end{aligned}$$

N.B. If  $Y := \sum_{i=1}^n X_i$  where  $\mathbf{X} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$  then  $Y \sim \text{Binomial}(n, p)$ .

**Definition 0.20 -  $\chi^2$  Distribution**

Let  $X \sim \chi_r^2$ .

A *continuous* random variable modelled by the  $\chi^2$  Distribution with  $r$  degrees of freedom. Then

$$\begin{aligned} f_X(x) &= \frac{1}{2^{r/2} \Gamma(r/2)} x^{\frac{r}{2}-1} e^{-\frac{x}{2}} \\ F_X(x) &= \frac{1}{\Gamma(k/2)} \gamma\left(\frac{r}{2}, \frac{x}{2}\right) \\ \mathbb{E}(X) &= r \\ \text{Var}(X) &= 2r \\ \mathcal{M}_X(t) &= \mathbb{1}\{t < \frac{1}{2}\} (1 - 2t)^{-\frac{r}{2}} \end{aligned}$$

*N.B.* If  $Y := \sum_{i=1}^k Z_i^2$  with  $\mathbf{Z} \stackrel{\text{iid}}{\sim} \text{Normal}(0, 1)$  then  $Y \sim \chi_k^2$ .

**Definition 0.21 - Exponential Distribution**

Let  $X \sim \text{Exponential}(\lambda)$ .

A *continuous* random variable modelled by a *Exponential Distribution* with rate-parameter  $\lambda$ . Then

$$\begin{aligned} f_X(x) &= \mathbf{1}\{t \geq 0\} \cdot \lambda e^{-\lambda x} \\ F_X(x) &= \mathbf{1}\{t \geq 0\} \cdot (1 - e^{-\lambda x}) \\ \mathbb{E}(X) &= \frac{1}{\lambda} \\ \text{Var}(X) &= \frac{1}{\lambda^2} \\ \mathcal{M}_X(t) &= \mathbf{1}\{t < \lambda\} \frac{\lambda}{\lambda - t} \end{aligned}$$

*N.B.* Exponential Distribution is used to model the wait time between decays of a radioactive source.

**Definition 0.22 - Gamma Distribution**

Let  $X \sim \Gamma(\alpha, \beta)$ .

A *continuous* random variable modelled by a *Gamma Distribution* with shape parameter  $\alpha > 0$  & rate parameter  $\beta$ . Then

$$\begin{aligned} f_X(x) &= \frac{1}{\Gamma(\alpha)} \beta^\alpha x^{\alpha-1} e^{-\beta x} \\ F_X(x) &= \frac{\Gamma(\alpha)}{\Gamma(\alpha)} (\alpha, \beta x) \\ \mathbb{E}(X) &= \frac{\alpha}{\beta} \\ \text{Var}(X) &= \frac{\alpha}{\beta^2} \\ \mathcal{M}_X(t) &= \mathbf{1}\{t < \beta\} \left(1 - \frac{t}{\beta}\right)^{-\alpha} \end{aligned}$$

*N.B.* There is an equivalent definition of a *Gamma Distribution* in terms of a shape & scale parameter. The scale parameter is 1 over the rate parameter in this definition.

**Definition 0.23 - Normal Distribution**

Let  $X$  be a continuous random variable modelled by a *Normal Distribution* with mean  $\mu$  & variance  $\sigma^2$ .

Then

$$\begin{aligned} f_X(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ F_X(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^x e^{-\frac{(y-\mu)^2}{2\sigma^2}} dy \\ \mathbb{E}(X) &= \mu \\ \text{Var}(X) &= \sigma^2 \\ \mathcal{M}_X(\theta) &= e^{\mu\theta + \sigma^2\theta^2(1/2)} \end{aligned}$$

**Definition 0.24 - Pareto Distribution**

Let  $X \sim \text{Pareto}(x_0, \theta)$ .

A *continuous* random variable modelled by a *Pareto Distribution* with minimum value  $x_0$  &

shape parameter  $\alpha > 0$ . Then

$$\begin{aligned} f_X(x) &= \frac{\alpha x_0^\alpha}{x^{\alpha+1}} \\ F_X(x) &= 1 - \left(\frac{x_0}{x}\right)^\alpha \\ \mathbb{E}(X) &= \begin{cases} \infty & \alpha \leq 1 \\ \frac{\alpha x_0}{\alpha - 1} & \alpha > 1 \end{cases} \\ \text{Var}(X) &= \begin{cases} \infty & \alpha \leq 2 \\ \frac{x_0^2 \alpha}{(\alpha - 1)^2 (\alpha - 2)} & \alpha > 2 \end{cases} \\ \mathcal{M}_X(t) &= \mathbf{1}_{\{t < 0\}} \alpha (-x_0 t)^\alpha \Gamma(-\alpha, -x_0 t) \end{aligned}$$

**Definition 0.25 - Poisson Distribution**

Let  $X \sim \text{Poisson}(\lambda)$ .

A *discrete* random variable modelled by a *Poisson Distribution* with rate parameter  $\lambda$ . Then

$$\begin{aligned} p_X(k) &= \frac{e^{-\lambda} \lambda^k}{k!} \quad \text{for } k \in \mathbb{N}_0 \\ P_X(k) &= e^{-\lambda} \sum_{i=1}^k \frac{\lambda^i}{i!} \\ \mathbb{E}(X) &= \lambda \\ \text{Var}(X) &= \lambda \\ \mathcal{M}_X(t) &= e^{\lambda(e^t - 1)} \end{aligned}$$

*N.B.* Poisson Distribution is used to model the number of radioactive decays in a time period.

**Definition 0.26 -  $t$ -Distribution**

Let  $X \sim t_r$ .

A *continuous* random variable with  $r$  degrees of freedom. Then

$$\begin{aligned} f_X(k) &= \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} \\ \mathbb{E}(X) &= \begin{cases} 0 & \text{if } \nu > 1 \\ \text{undefined} & \text{otherwise} \end{cases} \\ \text{Var}(X) &= \begin{cases} \frac{\nu}{\nu-2} & \text{if } \nu > 2 \\ \infty & 1 < \nu \leq 2 \\ \text{undefined} & \text{otherwise} \end{cases} \\ \mathcal{M}_X(t) &= \text{undefined} \end{aligned}$$

*N.B.* Let  $Y \sim \text{Normal}(0, 1)$  &  $Z \sim \chi_r^2$  be independent random variables then  $X := \frac{Y}{\sqrt{Z/r}} \sim t_r$ .

**Definition 0.27 - Uniform Distribution - Uniform**

Let  $X \sim \text{Uniform}(a, b)$ .

A *continuous* random variable with lower bound  $a$  & upper bound  $b$ . Then

$$\begin{aligned} f_X(x) &= \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases} \\ F_X(x) &= \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a, b] \\ 1 & \text{otherwise} \end{cases} \\ \mathbb{E}(X) &= \frac{1}{2}(a+b) \\ \text{Var}(X) &= \frac{1}{12}(b-a)^2 \\ \mathcal{M}_X(t) &= \begin{cases} \frac{e^{tb} - e^{ta}}{t(b-a)} & t \neq 0 \\ 1 & t = 0 \end{cases} \end{aligned}$$

## 0.5 Identities

### 0.5.1 Likelihood

#### Proposition 0.1 - Binomial

Let  $X \sim \text{Binomial}(n, p)$  with  $n$  &  $p$  unknown and  $x$  be a realisation of  $X$ . Then

$$\begin{aligned} L(n, p; x) &\propto \binom{n}{x} p^x (1-p)^{n-x} \\ \ell(n, p; \mathbf{x}) &= \ln \binom{n}{x} + x \ln p + (n-x) \ln(1-p) + C \\ \hat{n}_{\text{MLE}} &= \frac{x}{\hat{p}} \\ \hat{p}_{\text{MLE}} &= \frac{x}{\hat{n}} \end{aligned}$$

#### Proposition 0.2 - Normal

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$  with  $\mu$  &  $\sigma^2$  unknown and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ . Then

$$\begin{aligned} L(\mu, \sigma^2; \mathbf{x}) &\propto (\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) \\ \ell(\mu, \sigma^2; \mathbf{x}) &= -n \ln \sigma^2 - \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + C \\ \hat{\mu}_{\text{MLE}} &= \bar{\mathbf{x}} \\ \hat{\sigma}_{\text{MLE}}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \end{aligned}$$

#### Proposition 0.3 - Poisson

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$  with  $\lambda$  unknown and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ . Then

$$\begin{aligned} L(\lambda; \mathbf{x}) &\propto e^{-\lambda n} \lambda^{n\bar{x}} \\ \ell(\lambda; \mathbf{x}) &= -\lambda n + n\bar{x} \ln \lambda + C \\ \hat{\lambda}_{\text{MLE}} &= \bar{x} \end{aligned}$$

#### Proposition 0.4 - Uniform

Let  $\mathbf{X} \stackrel{\text{iid}}{\sim} \text{Uniform}(a, b)$  with  $a$  &  $b$  unknown and  $\mathbf{x}$  be a realisation of  $\mathbf{X}$ . Then

$$\begin{aligned} L(a, b; \mathbf{x}) &\propto \begin{cases} \frac{1}{(b-a)^n} & a \leq x_i \leq b \forall x_i \in \mathbf{x} \\ 0 & \text{otherwise} \end{cases} \\ \ell(a, b; \mathbf{x}) &= \begin{cases} -\ln(b-a) & a \leq x_i \leq b \forall x_i \in \mathbf{x} \\ 0 & \text{otherwise} \end{cases} \\ \hat{a}_{\text{MLE}} &= \min\{x_i : x_i \in \mathbf{x}\} \\ \hat{b}_{\text{MLE}} &= \max\{x_i : x_i \in \mathbf{x}\} \end{aligned}$$