

APTS Statistical Inference, Lecture notes

Jonathan Rougier, University of Bristol

December 2012

The first half of this document is a crash-bang-wallop summary of the statistical material that I would present in a comprehensive second undergraduate course on statistical inference at a leading UK university. It covers both Frequentist and Bayesian approaches to estimation and hypothesis testing in some generality, with a linking section on different statistical interpretations of probability.

More advanced material is presented in the second half, covering asymptotic convergence, Decision Theory, Bayes Factors, inferential fallacies, and principles.

Throughout the documents the key results are proved to an acceptable level of rigour; longer proofs can be found in the textbooks listed at the end.

This is the first draft of these notes, and I expect there are mistakes, or passages that are unclear. So comments on this document are particularly welcome; please address them by email to me at j.c.rougier@bristol.ac.uk. This version of the document was created on November 27, 2012, and formatted using the `tufte-handout` class, available at <https://code.google.com/p/tufte-latex/>.

© University of Bristol, 2012.

Contents

1	<i>Probability and statistics background</i>	3
1.1	<i>Probability basics</i>	3
1.2	<i>Statistical models</i>	5
1.3	<i>The exponential family of distributions</i>	8
1.4	<i>Continuous random quantities</i>	9
1.5	<i>The Probability Integral Transform (PIT)</i>	9
1.6	<i>Convergence</i>	10
2	<i>Point and set estimators</i>	11
2.1	<i>Point estimators</i>	11
2.2	<i>Judging point estimators</i>	12
2.3	<i>Set estimators</i>	14
3	<i>Goodness of fit tests</i>	16
3.1	<i>Simple hypothesis</i>	16
3.2	<i>Composite hypothesis</i>	17
	<i>Bibliography</i>	19

1 Probability and statistics background

At the heart of statistical inference is the notion of an *random quantity*, typically denoted X . I strongly advocate the positivist view that X is a set of operations which, if followed, will result in a value. A statement such as $\Pr(X = x)$ is to be read as “the probability that the operations described by X yields the value x ”, where I am adopting the standard convention that small letters denote specific possible values for X . It is a very useful discipline to insist that *random quantities have operational definitions*; it ensures that uncertainty about X is not confounded with ambiguity about the meaning of X .

random quantity

The set of possible outcomes for X is denoted \mathcal{X} , the *sample space* for X . The operational definition for X implies that \mathcal{X} is finite, reflecting the finite precision of our instruments; that is to say, X is *always a discrete random quantity*. It may seem idiosyncratic to rule out continuous random quantities at the start of a course on statistical inference, but in fact many statisticians have advocated something similar. For example, Debabrata Basu:

sample space

“The author holds firmly to the view that this contingent and cognitive universe of ours is in reality only finite and, therefore, discrete. In this essay we steer clear of the logical quick sands of ‘infinity’ and the ‘infinitesimal’. Infinite and continuous models will be used in the sequel, but they are to be looked upon as mere approximations to the finite realities.” (Basu, 1975, footnote, p. 4)

In fact, continuous random quantities will reappear for reasons explained in section 1.4.

A collection of random quantities is denoted $\mathbf{X} := (X_1, \dots, X_n) \in \mathcal{X}$, where very often $\mathcal{X} = \mathcal{X}^n$; likewise, $\mathbf{x} := (x_1, \dots, x_n)$.¹ In lectures and handwritten notes I will use an underscore, e.g. $\underline{\mathbf{X}} \equiv \underline{X}$ and $\underline{\mathbf{x}} \equiv \underline{x}$. Where it is necessary to specify the observed values of \mathbf{X} , these will be denoted \mathbf{x}^{obs} .

¹ Examples where $\mathcal{X} \neq \mathcal{X}^n$ are given in ??.

1.1 Probability basics

At this stage we assume the existence of a *probability mass function (PMF)* f_X with domain \mathcal{X} , such that

probability mass function (PMF)

$$\Pr(\mathbf{X} \in A) := \sum_{\mathbf{x} \in A} f_X(\mathbf{x}) \quad \text{for any } A \subset \mathcal{X}.$$

The three axioms of the probability calculus are automatically satisfied provided that $f_X(\mathbf{x}) \geq 0$ for all \mathbf{x} and $\Pr(\mathbf{X} \in \mathcal{X}) = 1$. For now I skirt over the precise interpretation of f_X , which is discussed in more detail in ??.

If $\mathcal{X} = \mathbf{Y} \times \mathbf{Z}$, then the *marginal PMF* for \mathbf{Z} is defined as

marginal PMF

$$f_Z(\mathbf{z}) := \sum_{\mathbf{y}} f_{\mathbf{Y},\mathbf{Z}}(\mathbf{y}, \mathbf{z}).$$

The *conditional PMF* of \mathbf{Y} given \mathbf{Z} is defined as the function $f_{\mathbf{Y}|\mathbf{Z}}$ satisfying

conditional PMF

$$f_{\mathbf{Y},\mathbf{Z}}(\mathbf{y}, \mathbf{z}) = f_{\mathbf{Y}|\mathbf{Z}}(\mathbf{y} | \mathbf{z}) f_Z(\mathbf{z})$$

for all \mathbf{y} and \mathbf{z} ; clearly such a function always exists and is unique excepting the case where $f_Z(\mathbf{z}) = 0$. It is conventional to state that $f_{Y|Z}$ is undefined in this case, but also overly restrictive, since $f_Z(\mathbf{z}) = 0$ implies that $f_{Y,Z}(\mathbf{y}, \mathbf{z}) = 0$, and so any value for the conditional PMF would do.

The *Law of Total Probability (LTP)* is a simple extension of the definitions,

$$f_Y(\mathbf{y}) = \sum_{\mathbf{z}} f_{Y,Z}(\mathbf{y}, \mathbf{z}) = \sum_{\mathbf{z}} f_{Y|Z}(\mathbf{y} | \mathbf{z}) f_Z(\mathbf{z}). \quad (1)$$

The LTP shows that one has the option of specifying f_Y indirectly, in terms of $f_{Y|Z}$ and f_Z . Breaking down complex probability assessments in this way is often helpful in practice. It also shows that $\min_{\mathbf{z}} f_{Y|Z}(\mathbf{y} | \mathbf{z})$ and $\max_{\mathbf{z}} f_{Y|Z}(\mathbf{y} | \mathbf{z})$ are lower and upper bounds on $f_Y(\mathbf{y})$, because the righthand side of (1) is a convex combination.²

Expectation and variance. Let X be a component of \mathbf{X} . The *expectation* of X is defined as

$$E(X) := \sum_x x f_X(x)$$

and the *variance* of X as

$$\text{Var}(X) := E(\{X - E(X)\}^2).$$

These always exist if X has a bounded sample space, but do not always exist, in more general treatments of random quantities. Note that if $E(X)$ exists, it does not necessarily lie in \mathcal{X} . However, it is a convex combination of \mathcal{X} . Hence if all the elements of \mathcal{X} are non-negative, then $E(X) \geq 0$. This gives rise to the very useful *monotonicity property of expectations*. If X and Y are two components of \mathbf{X} and $f_{X,Y}(x, y) = 0$ unless $x \leq y$, then $E(X) \leq E(Y)$.

The expectation and variance are related by various inequalities, including *Chebyshev's inequality*: if μ is the expectation of X and σ^2 the variance, then

$$\Pr(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

If g is a convex scalar function of x , then *Jensen's inequality* states that

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

(The proofs of these easy results are left as exercises.)

If X and Y are two components of \mathbf{X} , then the *covariance* between X and Y is defined as

$$\text{Cov}(X, Y) := E(\{X - E(X)\}\{Y - E(Y)\}),$$

from which we might also have defined the variance of X as $\text{Cov}(X, X)$. It is a general property that

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

which might be more familiar as the statement that the correlation between X and Y always lies in the interval $[-1, 1]$; see Grimmett and Stirzaker (2001, section 3.6) for a simple proof.

Law of Total Probability (LTP)

² A convex combination of the k elements of set \mathcal{A} has the form $\sum_{j=1}^k p_j a_j$ where each $p_j \geq 0$ and $\sum_j p_j = 1$.

expectation

variance

monotonicity property of expectations

Chebyshev's inequality

Jensen's inequality

covariance

Propositions. A proposition is a statement which is either true or false. Thus propositions

$$a = \text{'I am an astronaut'}$$

is a proposition which is false, and $a(x) = \text{'x is an astronaut'}$ is a proposition whose value is contingent on x . The *indicator function* of a proposition a is indicator function

$$\mathbb{1}(a) := \begin{cases} 0 & \text{if } a \text{ is false} \\ 1 & \text{if } a \text{ is true.} \end{cases}$$

A proposition need not be certain; if $A := a(X)$, then A is a random quantity, and the probability that A is true is equal to the expectation of its indicator function,

$$\Pr(A) = \sum_x f_X(x) \mathbb{1}(a(x)) = E\{\mathbb{1}(A)\}.$$

The following two results are simple but useful inequalities for propositions, which follow from the link between probabilities and expectations of indicator functions. First, if a and b are two propositions and a implies b , then $\Pr(A) \leq \Pr(B)$.

Proof. So a implies b if and only if $\mathbb{1}(a(x)) \leq \mathbb{1}(b(x))$ for all x . Then, using the monotonicity property of expectations,

$$\Pr(A) = E\{\mathbb{1}(A)\} \leq E\{\mathbb{1}(B)\} = \Pr(B).$$

□

Second, if a_1, \dots, a_k is a set of propositions, then

$$\Pr(A_1 \vee \dots \vee A_k) \leq \sum_{i=1}^k \Pr(A_i), \tag{2}$$

where ' \vee ' denotes 'or' (inclusive disjunction). This is the *Bonferroni inequality*.

Bonferroni inequality

Proof. Here is the case with $k = 2$, and $a_1 = a$ and $a_2 = b$.

$$\mathbb{1}(a(x) \vee b(x)) = \max\{\mathbb{1}(a(x)), \mathbb{1}(b(x))\} \leq \mathbb{1}(a(x)) + \mathbb{1}(b(x)).$$

Then, using the monotonicity property of expectations,

$$\Pr(A \vee B) = E\{\mathbb{1}(A \vee B)\} \leq E\{\mathbb{1}(A) + \mathbb{1}(B)\} = \Pr(A) + \Pr(B).$$

The general case follows by induction. □

1.2 Statistical models

The PMF for X is often indexed by a parameter vector $\theta \in \Omega$, written

$$\Pr(X \in A; \theta) = \sum_{x \in A} f_X(x; \theta);$$

other operators dependent on θ , such as the expectation of functions of X , are also indexed by θ . The introduction of a fixed family

f_X and a parameter θ is a common strategy to limit the set of possible probability assignments over \mathcal{X} in order to be consistent with one's judgements about X . Often the nature of f_X is fairly apparent, or conforms to a conventional situation, and only the value of θ is unknown. Interest then focuses on making inferences about θ based on x^{obs} . The combination of f_X and Ω is termed the *statistical model*. To avoid unnecessary complications, the probability assignment in the statistical model is assumed to be an injective function of θ , i.e. two different θ 's cannot give the same probability assignment; in this case the statistical model is said to be *identifiable*.³

I will take Ω to be a convex subset of \mathbb{R}^d . Many useful functions in statistical inference have Ω in their domain, such as the score function, (4), and the likelihood function, (5). For these functions I use the argument $t \in \Omega$. Differentiation with respect to t is denoted using ' ∇ ' (nabla):

$$\nabla := \begin{pmatrix} \partial/\partial t_1 \\ \vdots \\ \partial/\partial t_d \end{pmatrix}, \text{ and } \nabla^2 := \begin{pmatrix} \partial^2/(\partial t_1)^2 & \dots & \partial^2/\partial t_1 \partial t_d \\ \vdots & \ddots & \vdots \\ \partial^2/\partial t_d \partial t_1 & \dots & \partial^2/(\partial t_d)^2 \end{pmatrix}. \quad (3)$$

This is a notationally confusing area. Some authors do not distinguish between θ , the 'true but unknown parameter value' and t , a point in Ω ; this can be confusing for non-experts. The convention on capital letters would suggest that the true parameter be Θ , and the point in Ω could then be θ ; Schervish (1995) follows this consistently, but the size of ' Θ ' makes it quite intrusive. Another possibility is to use θ_0 for the true value and θ for the point, but this gets confusing when the null hypothesis H_0 arrives.

Those x values for which $f_X(x; \theta) > 0$ are termed the *support of X*. Many interesting statistical results concern the special case where the support of X does not depend on the value of θ . These are *regular models*.⁴ One very useful function for parametric statistical models is the *score function*

$$u(x, t) := \nabla \log f_X(x; t). \quad (4)$$

It is straightforward to prove that

$$E\{u(\mathbf{X}; \theta); \theta\} = \mathbf{0}$$

for regular models.

Proof. Start from the identity $\sum_x f_X(x; \theta) = 1$. Then differentiate both sides with respect to θ_j . The righthand side is zero. For a regular model, the lefthand side is

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \sum_x f_X(x; \theta) &= \sum_x \frac{\partial}{\partial \theta_j} f_X(x; \theta) \\ &= \sum_x \frac{\partial}{\partial \theta_j} (\log f_X(x; \theta)) f_X(x; \theta) \\ &= E\{u_j(\mathbf{X}; \theta); \theta\}, \end{aligned}$$

which must therefore equal zero. \square

statistical model

identifiable

³ A stronger condition is that $f_X(x; \theta) \neq f_X(x; \theta')$ unless $\theta = \theta'$, for every $x \in \mathcal{X}$. This type of strong condition is only required in precise proofs of consistency.

support of X

regular models

⁴ Technical note. In fact, regular models are defined to be those for which the order of differentiation with respect to θ and summation/integration with respect to x can be reversed. See Casella and Berger (2002, section 2.4) for details. The invariance of the support of X with respect to θ is the main necessary condition.

score function

The Fisher information matrix is defined as

$$i(\theta) := \text{Var}\{u(\mathbf{X}, \theta); \theta\}$$

(it is a $d \times d$ matrix when θ is a vector); it appears at interesting times in the theory of estimators. It is straightforward to show that for regular models

$$i(\theta) = -E\{\nabla u(\mathbf{X}, \theta)^T; \theta\}.$$

When evaluating $i(\theta)$ it is often easier to compute the expectation of the derivatives of $u(\mathbf{X}, \theta)$ than to compute the variance of $u(\mathbf{X}, \theta)$.

A common situation is where the X 's are *independent and identically distributed (IID)*, in which case

$$f_X(\mathbf{x}; \theta) = \prod_{i=1}^n f_X(x_i; \theta), \quad \text{written } \mathbf{X} \stackrel{\text{iid}}{\sim} f_X(x; \theta),$$

for some marginal PMF f_X . If the X 's are IID and f_X is regular then the score function is a simple sum:

$$u(\mathbf{x}, t) = \nabla \log \prod_{i=1}^n f_X(x_i; t) = \sum_{i=1}^n \nabla \log f_X(x_i; t) = \sum_{i=1}^n u(x_i, t).$$

In this case $i(\theta) = ni_1(\theta)$, where $i_1(\theta) := \text{Var}\{u(X_1, \theta); \theta\}$.

Nuisance parameters. The statistical model $\mathbf{X} \sim f_X(\mathbf{x}; \theta)$ for $\theta \in \Omega$ is meant to be very general, accommodating simple 'student' models such as $(X_1, \dots, X_n) \stackrel{\text{iid}}{\sim} N(x; \mu, \sigma^2)$, but also much more complex models where the IID structure may be buried deep inside the model, and where both \mathbf{X} and θ represent collections of possibly quite heterogeneous quantities.⁵

In all but the simplest statistical models some of the parameters are interesting to the scientist, and others are only there for the purposes of structuring the statistical model. The ones that are *not* interesting are referred to as *nuisance parameters*, and the name is well-earned, because they are indeed a nuisance. In order to discuss nuisance parameters at the appropriate places below, I will partition the parameter space as $\Omega = \Omega' \times \Omega''$, where $\theta = (\theta', \theta'')$ and θ'' are the nuisance parameters.

This division of θ into two parts is made without loss of generality, as long as the techniques being used are *transformation invariant*. Which is to say, if $g : \theta \mapsto \psi$ is a bijective function, then all inferences made about θ from the model $f_X(\mathbf{x}; \theta)$ are equivalent to inferences made about ψ from the model $f_X^\psi(\mathbf{x}; \psi) = f_X(\mathbf{x}; g^{-1}(\psi))$, and *vice versa*. In this case consistent inferences can be made about any function of θ , say $\psi' = g_1(\theta)$, as long as this function can be embedded within a bijective function $\psi = (\psi', \psi'') = (g_1(\theta), g_2(\theta)) = g(\theta)$.

Transformation invariance is also an important general concern, because the precise representation of the parameters in the statistical model is purely a matter of convention. And in fact in many cases the convention has not settled down, and the same

Fisher information matrix

independent and identically distributed (IID)

⁵ A more general treatment would distinguish between the random quantities \mathbf{X} and the observations \mathbf{Y} , where the latter are some function of \mathbf{X} .

nuisance parameters

transformation invariant

distribution can be parameterised in several different ways. In applied work it is always a good idea to state $f_X(x; \theta)$ explicitly as a formula, to remove any ambiguity; a statement such as ‘ X is IID Gamma’ does not perfectly identify the two components of θ .

1.3 The exponential family of distributions

One very common class of statistical models have the general form

$$g_X(x; \psi) = f(x) \exp \{ \psi s(x) - \kappa(\psi) \} \quad \text{for } \psi \in \Psi \subset \mathbb{R},$$

where ψ is a strictly monotone function of θ and s is a scalar function. These are the one parameter *exponential family of distributions*. Common IID distributions based on the Binomial, Poisson, Geometric, Exponential, and Normal (known variance) all belong to this family. One derives a member of the one-parameter exponential family by ‘exponentially tilting’ a base distribution $f(x)$:

$$g_X(x; \psi) \propto f(x) \exp \{ \psi s(x) \},$$

where, in order to sum to one,

$$\kappa(\psi) := \log \sum_x f(x) \exp \{ \psi s(x) \}$$

and the parameter space is defined as $\Psi := \{ \psi : \kappa(\psi) < \infty \}$. It is straightforward to show that Ψ is a convex subset of \mathbb{R} , and that κ is a convex function of ψ (Davison, 2003, sec. 5.2). And also that f_X and g_X have the same support, which does not depend on ψ .

The exponential family generalises to a vector ψ , with the form

$$g_X(x; \psi) \propto f(x) \exp \left\{ \sum_j \psi_j s_j(x) \right\}.$$

This contains many more familiar distributions, including the Normal, Beta, Gamma, Multinomial, and Dirichlet.

The exponential family has special properties, some of which will be mentioned below. These properties are so special, though, that the exponential family does not provide a general basis for an exploration of statistical inference. When computation was a challenge, the exponential family was favoured because of its simplicity; now, however, there is less need for this type of restriction.

In case you are wondering, in these notes there will be no mention at all of *sufficiency*. I regard sufficiency as a very useful computational device in the case where the model admits a fixed length sufficient statistic. But, as the *Pitman-Koopmans-Darmois theorem* states, a support that does not depend on the parameter plus a fixed length sufficient statistic are more-or-less the characterisation of the exponential family; see Schervish (1995, section 2.2.3). So if I judge the exponential family too restrictive for general consideration, then I can ignore sufficiency.⁶

exponential family of distributions

sufficiency

Pitman-Koopmans-Darmois theorem

⁶ A slight regret: there is some beautiful mathematics in sufficiency, for example the Fisher-Neyman factorisation criterion, the Dynkin-Lehmann-Scheffé theorem on minimal sufficient statistics, and the Rao-Blackwell theorem on estimators under convex loss. Never mind.

1.4 Continuous random quantities

As already stated, \mathcal{X} is taken to be finite, i.e. X is a discrete random quantity. However, *continuous random quantities* with uncountable state spaces are still useful, as an approximation; and because they arise naturally in the Bayesian approach, where the parameter itself (which is generally continuous) acquires a distribution. In the case that \mathbf{X} is continuous $f_{\mathbf{X}}$ is taken to be a *probability density function (PDF)* with the property that

continuous random quantities

probability density function (PDF)

$$\Pr(\mathbf{X} \in A; \theta) = \int_A f_{\mathbf{X}}(x; \theta) dx$$

where $A \subset \Omega$. Discrete and continuous random quantities can be unified within *measure theory*, but we are not going to worry about measure theory in these notes; see Grimmett and Stirzaker (2001) and then Rosenthal (2006) or Kingman and Taylor (1966) for an accessible introduction.

measure theory

One useful result for continuous \mathbf{X} is the *transformation of the PDF*. If $\mathbf{y} = g(\mathbf{x})$ where g is a differentiable bijective function, then

transformation of the PDF

$$f_Y(\mathbf{y}) = f_X(\mathbf{x}) |\nabla g(\mathbf{x})^T|^{-1} \quad \text{where } \mathbf{x} := g^{-1}(\mathbf{y})$$

for $\mathbf{y} \in g(\mathcal{X})$, and zero otherwise; here $|\nabla g^T|$ is known as the *Jacobian* of the transformation, and is the determinant of the square matrix of first derivatives, remembering that g is a vector-valued function in general, with one component for each x_i .

Proof. Consider, just for simplicity, the case of scalar X and increasing g . Denote the probability distribution function of X as F_X , and similarly for Y . Then

$$F_Y(y) := \Pr\{Y \leq y\} = \Pr\{g(X) \leq y\} = \Pr\{X \leq g^{-1}(y)\} = F_X(x)$$

with $x := g^{-1}(y)$. Differentiate with respect to y to find the PDF, giving

$$f_Y(y) = f_X(x) \frac{dx}{dy} = f_X(x) \left(\frac{dy}{dx} \right)^{-1} = f_X(x) \left| \frac{dy}{dx} \right|^{-1},$$

the last equality following because g is increasing. This is easily generalised to the case where g is decreasing, for which $F_Y(y) = 1 - F_X(x)$. The full result, for vector \mathbf{X} and \mathbf{Y} , is effectively the change of variables formula for multivariate integration. \square

It is common to use f to denote either a probability mass function or a PDF, because many results are the same in either case. Below, though, I will use f for the probability mass function and π for the PDF, notably in ?? and beyond, where I write π_θ for the PDF of θ .

1.5 The Probability Integral Transform (PIT)

Let $X \in \mathcal{X} \subset \mathbb{R}$ be a scalar random quantity with distribution function $F_X : \mathcal{X} \rightarrow [0, 1]$. Define a new random quantity $Y := F_X(X)$;

i.e. Y is the random quantity one gets by putting X into its own distribution function. It is a very useful fact that Y has a *sub-uniform distribution*, i.e.

$$F_Y(u) \leq u \quad \text{for all } u \in [0, 1],$$

and that $F_Y(u) = u$ if there exists an $x \in \mathcal{X}$ such that $u = F_X(x)$.

Proof. First, consider the case where $u = F_X(x)$ for some $x \in \mathcal{X}$:

$$F_Y(u) = \Pr\{F_X(X) \leq F(x)\} = \Pr\{X \leq x\} = F_X(x) = u.$$

The ‘cancellation’ of F at the second equality occurs because of the bijective relationship between x and $F(x)$ for $x \in \mathcal{X}$.⁷ This proves the second part of the claim: in the case where X is a continuous random quantity, the points u in $(0, 1)$ are in a bijective relationship with the points x in \mathcal{X} , and Y is uniformly distributed.

Otherwise, let x and x' be two consecutive values in \mathcal{X} , with $u = F_X(x)$ and $u' = F_X(x')$, and let $u + \delta$ be some value in the open interval (u, u') . Then

$$Y \leq u + \delta \implies X \leq x$$

and so $F_Y(u + \delta) \leq F_X(x) = u$. But we must also have $F_Y(u + \delta) \geq F_Y(u) = u$. Therefore we conclude that $F_Y(u + \delta) = u$, and hence $F_Y(u + \delta) < u + \delta$. \square

So the distribution function of Y looks like a staircase where each step starts from the 45° line drawn from $(0, 0)$ to $(1, 1)$; see Figure 1. For a continuous random quantity the steps are infinitesimally small, and the distribution function and the 45° line coincide.

1.6 Convergence

There are several different types of convergence for sequences of random quantities. X_n converges in probability to Y , written $X_n \xrightarrow{P} Y$, if

$$\Pr(|X_n - Y| \geq \varepsilon) \longrightarrow 0 \quad \text{as } n \rightarrow \infty$$

for all $\varepsilon > 0$. If X_1, X_2, \dots is an IID sequence where each X_i has finite expectation μ , then

$$n^{-1}(X_1 + \dots + X_n) \xrightarrow{P} \mu.$$

This is the *Weak Law of Large Numbers (WLLN)*. It is easily proved using Chebyshev’s inequality in the case where $\text{Var}(X_i)$ is finite. In the WLLN convergence is to a constant, rather than to a random quantity.

Convergence in distribution is a weaker form of convergence (i.e. it is implied by convergence in probability). X_n converges in distribution to Y , written $X_n \xrightarrow{D} Y$, if

$$F_{X_n}(y) \longrightarrow F_Y(y) \quad \text{as } n \rightarrow \infty$$

sub-uniform distribution

⁷ Technical note: here we can ignore points in \mathcal{X} that have zero probability.

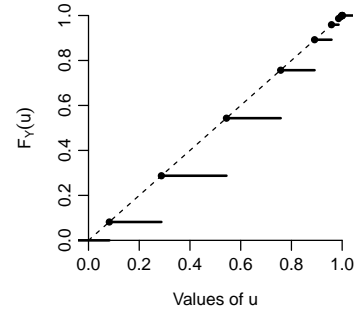


Figure 1: Distribution function of $Y = F_X(X)$, where $X \sim \text{Poisson}(x; \lambda = 2.5)$.

converges in probability

Weak Law of Large Numbers (WLLN)

converges in distribution

for all y for which F_Y is continuous. The most famous law in statistics, the *Central Limit Theorem (CLT)*, states that if X_1, X_2, \dots is an IID sequence with finite variance, then

Central Limit Theorem (CLT)

$$\frac{n^{-1}(X_1 + \dots + X_n) - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} N(0,1),$$

where μ and σ^2 are the mean and variance of each X_i . For convergence in distribution it is common to write a distribution on the right-hand side, as shown above, rather than a Y with a specified distribution.

2 Point and set estimators

Any function of \mathbf{x} is termed a *statistic*. A statistic designed to estimate θ is termed an *estimator*; typically these can be divided into *point estimators*, which map from \mathcal{X} to a point in Ω , and *set estimators*, which map from \mathcal{X} to a set in Ω . If s is a statistic, then $S := s(\mathbf{X})$ is a random quantity, whose distribution depends on θ . So, just to clarify, s is a function with domain \mathcal{X} , $s(\mathbf{x})$ is a value in the range of s , and S is a random quantity with sample space equal to the range of s . Where it is necessary to identify a particular point in the co-domain of s I will use s' ; e.g. in expressions of the form $\Pr(S \leq s'; \theta)$.⁸

statistic

point estimators

set estimators

If s is an estimator for θ , then one important principle in statistics states that s should be judged by comparing S and θ across the range of possible values of θ . This is the basis of the Frequentist approach to statistics, which is discussed in this section and the two that follow. The Bayesian approach is introduced in ??.

⁸ That s is a function and s' is a value is not ideal, but I think it is the least-bad compromise consistent with standard notation. Another option would be to write s everywhere as $s(\cdot)$, but to me this seems redundant and ugly.

2.1 Point estimators

There are many ways of constructing estimators for θ , but the *maximum likelihood estimator (MLE)* is often preferred for both theoretical and practical reasons, particularly when we are confident about the adequacy of the statistical model.

maximum likelihood estimator (MLE)

The MLE $\hat{\theta}$ satisfies

$$f_{\mathbf{X}}(\mathbf{x}; \hat{\theta}(\mathbf{x})) \geq f_{\mathbf{X}}(\mathbf{x}; t) \quad \text{for all } t \in \Omega.$$

If $\theta = (\theta_1, \theta_2)$, and the MLE of θ is $\hat{\theta}(\mathbf{x}) = (\hat{\theta}_1(\mathbf{x}), \hat{\theta}_2(\mathbf{x}))$, then the MLE of θ_1 is defined to be $\hat{\theta}_1(\mathbf{x})$. Now we can show that MLEs are transformation-invariant.⁹

Proof. Let $g : \theta \mapsto \psi$ be a bijective function, so that $f_{\mathbf{X}}^{\psi}(\mathbf{x}; \psi) = f_{\mathbf{X}}(\mathbf{x}; g^{-1}(\psi))$. Then

⁹ There is a more general definition of the likelihood function of θ_1 which preserves the property of transformation-invariance of the MLE, but also works in the case of non-bijective functions of θ ; see Casella and Berger (2002, section 7.2.2).

$$\begin{aligned} f_{\mathbf{X}}^{\psi}(\mathbf{x}; \hat{\psi}) &\geq f_{\mathbf{X}}^{\psi}(\mathbf{x}; v) \quad \text{for all } v \in g(\Omega) \\ \iff f_{\mathbf{X}}(\mathbf{x}; g^{-1}(\hat{\psi})) &\geq f_{\mathbf{X}}(\mathbf{x}; t) \quad \text{for all } t \in \Omega \\ \iff g^{-1}(\hat{\psi}) &= \hat{\theta} \\ \iff \hat{\psi} &= g(\hat{\theta}). \end{aligned}$$

This argument also extends to parts of ψ , as long as they can be embedded in a bijective function. Thus, let $\psi = (\psi_1, \psi_2) = (g_1(\theta), g_2(\theta))$ where g is bijective as before. Then $\hat{\psi} = (g_1(\hat{\theta}), g_2(\hat{\theta}))$ as already shown, and hence $\hat{\psi}_1 = g_1(\hat{\theta})$. \square

In general, $\hat{\theta}(x)$ is an implicit function of x , and the value of the MLE is found by maximising the *likelihood function*,

$$L(t) := f_X(\mathbf{x}^{\text{obs}}; t). \quad (5)$$

The value $\hat{\theta}(\mathbf{x}^{\text{obs}}) = \operatorname{argmax}_{t \in \Omega} L(t)$ is termed the *Maximum Likelihood (ML) estimate*. For regular models, the MLE should satisfy the first order conditions

$$u(x, \hat{\theta}(x)) = \mathbf{0},$$

excepting complications that arise at the boundary of Ω —which would usually be evidence that the model was inadequate. The situation is more complicated for non-regular models. For example, if $\mathbf{X} \stackrel{\text{iid}}{\sim} U(x; 0, \theta)$, the MLE is well defined but the likelihood function is not differentiable at its maximum.

Predictions. Suppose that we would like to predict the value of an unobserved X , based upon the \mathbf{x}^{obs} that we have. All such predictions are ultimately functions of θ . For example,

$$\Pr(X_{n+1} = x \mid \mathbf{X} = \mathbf{x}^{\text{obs}}; \theta) = \frac{f_{X, X_{n+1}}(\mathbf{x}^{\text{obs}}, x; \theta)}{f_X(\mathbf{x}^{\text{obs}}; \theta)}. \quad (6)$$

As long as the function on the right-hand side is or can be extended to a bijective function of θ , then the MLE of the probability of $X_{n+1} = x$ must be $f_{X, X_{n+1}}(x, x; \hat{\theta}(x)) / f_X(x; \hat{\theta}(x))$. Thus we have the general rule, that ML predictions about unobserved X are made by plugging in the value of the MLE. In the special case where $(\mathbf{X}, X_{n+1}) \stackrel{\text{iid}}{\sim} f_X(x; \theta)$, eq. (6) simplifies to

$$\Pr(X_{n+1} = x \mid \mathbf{X} = \mathbf{x}^{\text{obs}}; \theta) = f_X(x; \theta),$$

and hence the MLE is just $f_X(x; \hat{\theta}(x))$.

2.2 Judging point estimators

Suppose some statistic $s : \mathcal{X} \rightarrow \Omega$ is claiming to be a point estimator for θ —how do we judge whether it is a good estimator or a poor one? One property has already been mentioned: that the estimator should be transformation-invariant, a property that the MLE possesses. There have been many suggestions for other attractive properties of estimators. I just mention a couple here.

I will use ‘estimator’ for both the function s and also for the random quantity $S := s(\mathbf{X})$. The estimator S is *unbiased* if

$$\text{bias}(S; \theta) := E(S; \theta) - \theta$$

is zero for all $\theta \in \Omega$. This is a superficially attractive criterion but leads to daft results even in simple cases, where it exists, and

likelihood function

Maximum Likelihood (ML) estimate

unbiased estimator

often does not exist. Thus the unique unbiased estimator of θ in $X \sim \text{Geometric}(x; \theta)$ is daft (Cox and Hinkley, 1974, section 8.2),¹⁰ and the unbiased estimator of θ in $X \sim \text{Exponential}(x; \theta)$ does not exist (Schervish, 1995, section 5.1).¹¹ Furthermore, unbiasedness is not consistent with transformation-invariance. For example, if S is an unbiased estimator of scalar positive θ , then $1/S$ is not an unbiased estimator of $1/\theta$, because $E\{1/S; \theta\} \neq 1/E\{S; \theta\} = 1/\theta$.

A better criterion is that S has a small *mean squared error (MSE)*,

$$\text{MSE}(S; \theta) := E\{(S - \theta)^2; \theta\} = \text{bias}(S; \theta)^2 + \text{Var}(S; \theta), \quad (7)$$

where the second expression follows from a simple rearrangement. The squared error is one example of a *loss function* for point estimation, and the MSE is then the *risk*; see ??.

There is an interesting relationship between the bias and the variance, the *Cramér-Rao lower bound (CRLB)*. Taking θ scalar, for simplicity,

$$\text{Var}\{S; \theta\} \geq \frac{\{1 + \nabla \text{bias}(S; \theta)\}^2}{i(\theta)} \quad (8)$$

where $\nabla \text{bias} := d \text{bias} / d\theta$; the Fisher Information $i(\theta)$ was defined in section 1.2. The CRLB holds for regular models.

Proof of (8), with scalar θ . Correlation coefficients lie in the interval $[-1, 1]$, or, in our case,

$$\text{Cov}\{S, u(\mathbf{X}, \theta); \theta\}^2 \leq \text{Var}\{S; \theta\} \text{Var}\{u(\mathbf{X}, \theta); \theta\};$$

see section 1.1. As $E\{u(\mathbf{X}, \theta); \theta\} = 0$ for regular models,

$$\begin{aligned} \text{Cov}\{S, u(\mathbf{X}, \theta); \theta\} &= E\{Su(\mathbf{X}, \theta); \theta\} \\ &= \sum_{\mathbf{x}} s(\mathbf{x}) \frac{d}{d\theta} \left\{ \log f_{\mathbf{X}}(\mathbf{x}; \theta) \right\} f_{\mathbf{X}}(\mathbf{x}; \theta) \\ &= \frac{d}{d\theta} \sum_{\mathbf{x}} s(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}; \theta) \\ &= \frac{d}{d\theta} E\{S; \theta\} \\ &= \frac{d}{d\theta} \{\theta + \text{bias}(S; \theta)\} = 1 + \nabla \text{bias}(S; \theta). \end{aligned}$$

The result follows from the inequality at the top of the proof.

The generalisation to vector θ is not completely straightforward, excepting the obvious case where $s(\mathbf{x})$ is an estimator of one component of θ . See Cox and Hinkley (1974, section 8.3). \square

This proof also demonstrates that an estimator attains the CRLB if and only if

$$u(\mathbf{x}, \theta) = a(\theta)(s(\mathbf{x}) - \theta), \quad (9)$$

where a is any function of θ (and other values treated as known, such as n), which is necessary and sufficient to ensure that the correlation between S and $u(\mathbf{X}, \theta)$ is ± 1 for any given θ , neglecting the trivial case where $s(\mathbf{x})$ is a constant. Estimators that attain the CRLB are termed *efficient estimators*. Integrating (9) and rearrang-

¹⁰ It is $\hat{\theta}(x) = \mathbb{1}(x = 1)$.

¹¹ A much more technical point is that from the point of view of Decision Theory, some unbiased estimators, where they exist, are *inadmissible*. This is demonstrated by the *Stein effect*, discussed below in ??.

mean squared error (MSE)

Cramér-Rao lower bound (CRLB)

efficient estimators

ing shows that an estimator can attain the CRLB if and only if the statistical model has the general form

$$\log f_{\mathbf{X}}(\mathbf{x}; \theta) = A(\theta)s(\mathbf{x}) + B(\theta) + C(\mathbf{x}),$$

i.e. if it is a subset of the exponential family of distributions; see section 1.3. It follows from the special properties of this family that where an estimator is efficient it is also the MLE.

The MSE is an intuitive criterion with which to judge point estimators, but it is important to be aware that unbiased and efficient estimators do not necessarily minimise the MSE for regular models. This is because some biased estimators increase the first term in (7) but decrease the second term by more; the CRLB shows that these estimators must have $\nabla \text{bias}(S; \theta) < 0$. The Stein effect provides a good illustration; see ??.

2.3 Set estimators

Set estimators of θ are also useful: crudely, their volume represents the variability in the estimator that arises from having only a finite sample. Let \mathcal{C} be a function from \mathcal{X} to a subset of Ω , i.e. $\mathcal{C}(\mathbf{x}) \subset \Omega$. The *coverage* of \mathcal{C} is the probability that $\mathcal{C}(X)$ will contain θ , which is a function of θ . \mathcal{C} is a *level* $(1 - \alpha)$ *confidence set* for θ if its coverage is at least $1 - \alpha$ for all θ . That is,

$$\Pr\{\theta \in \mathcal{C}(X); \theta\} \geq 1 - \alpha \quad \text{for all } \theta \in \Omega.$$

Here $(1 - \alpha)$ is termed the *nominal level* of the confidence set. When the coverage equals the nominal level for all θ , \mathcal{C} is said to be an *exact confidence set*. Typically, though, this only happens in one or two special cases, or in the limit as $n \rightarrow \infty$ for IID X 's from regular models; see below, around (10). In the case where θ is a scalar parameter and $\mathcal{C}(\mathbf{x})$ is convex for all \mathbf{x} , \mathcal{C} is a *confidence interval* for θ .

In the case of nuisance parameters, i.e. where $\theta = (\theta', \theta'')$ but only θ' is interesting, a level $(1 - \alpha)$ confidence set for θ' is any function \mathcal{C}' from \mathcal{X} to subsets of Ω' , with the property that

$$\Pr\{\theta' \in \mathcal{C}'(X); \theta\} \geq 1 - \alpha \quad \text{for all } \theta \in \Omega.$$

It is reassuring that such *marginal confidence sets* can always be constructed from full confidence sets. Define the projection \mathbb{P} as $\mathbb{P}(\theta) = \theta'$. If \mathcal{C} is a $(1 - \alpha)$ confidence set for θ , then $\mathcal{C}' := \mathbb{P}\mathcal{C}$ is a $(1 - \alpha)$ confidence set for θ' .

Proof. $\theta \in \mathcal{C}(\mathbf{x})$ implies $\theta' \in \mathcal{C}'(\mathbf{x})$, and hence, for all θ ,

$$1 - \alpha \leq \Pr\{\theta \in \mathcal{C}(X); \theta\} \leq \Pr\{\theta' \in \mathcal{C}'(X); \theta\},$$

using the result on propositions given in ??.

Of course, just because such confidence sets exist, does not mean that they have attractive properties (like having small volumes). This is discussed further in ??.

coverage of \mathcal{C}
level $(1 - \alpha)$ confidence set

nominal level
exact confidence set
confidence interval

marginal confidence sets

□

This definition of confidence sets is all very well in theory, but set estimators with a known confidence level are extremely hard to construct in practice, except in very simple statistical models, and for uninformative confidence sets.¹² On the other hand, for any reasonable confidence set, it is extremely hard to work out the infimum of the coverage over Ω , and it is even harder to reverse the problem, and identify a confidence set for which the infimum of the coverage is above a target level. It is harder still to do this in the presence of nuisance parameters.

In general, the best that can be done is construct confidence sets that have approximately the right coverage in special cases, notably where the X 's are IID and the statistical model is regular. The simplest approach is to use asymptotic theory (i.e. $n \rightarrow \infty$). Here it is important to insert a *caveat*. Resorting to asymptotic constructions is an *act of desperation*, because in reality datasets are stubbornly finite, and often quite small. It is a regrettable feature of Frequentist inference that, except in a few special cases, asymptotic arguments are all that are available. The Bayesian approach makes no particular use of asymptotic arguments, but makes additional demands on the scientist that some statisticians are unwilling to accept, as discussed in ??.

Here is one important example, based on the score function. If the X 's are IID and the statistical model is regular, then the CLT implies that $u(\mathbf{X}, \theta) \xrightarrow{D} N_d(\mathbf{0}, i(\theta))$, and this in turn implies that¹³

$$w(\mathbf{X}, \theta) := u(\mathbf{X}, \theta)^T i(\theta)^{-1} u(\mathbf{X}, \theta) \xrightarrow{D} \chi_d^2, \quad (10)$$

remembering that $u(\mathbf{x}, \theta)$ is a d -dimensional vector of partial derivatives, and $i(\theta)$ is a $d \times d$ variance matrix. Then

$$\mathcal{C}(\mathbf{x}) := \{\theta : w(\mathbf{x}, \theta) \leq F_{\chi_d^2}^{-1}(1 - \alpha)\} \quad (11)$$

is asymptotically an exact $(1 - \alpha)$ confidence set for θ , where $F_{\chi_d^2}^{-1}$ is the quantile function of the χ_d^2 distribution.

Proof.

$$\begin{aligned} \Pr\{\theta \in \mathcal{C}(\mathbf{X}); \theta\} &= \Pr\{w(\mathbf{X}, \theta) \leq F_{\chi_d^2}^{-1}(1 - \alpha); \theta\} \\ &= \Pr\{F_{\chi_d^2}(w(\mathbf{X}, \theta)) \leq 1 - \alpha; \theta\} \longrightarrow 1 - \alpha \end{aligned}$$

for all θ , applying the PIT (section 1.4) and (10). \square

This confidence set is transformation-invariant.

Proof. Take θ scalar, for simplicity. Apply the chain rule to show that $u^\psi(\mathbf{x}, \psi) = u(\mathbf{x}, \theta)(g'(\theta))^{-1}$, where $g : \theta \mapsto \psi$ is bijective and differentiable, $g' := (d/d\theta)g$ and $\theta = g^{-1}(\psi)$. Then it follows that $w(\mathbf{x}, \theta) = w^\psi(\mathbf{x}, \psi)$ for all \mathbf{x} .

The general proof is the same, but more fiddly. \square

¹² I.e., the equivalent of the uninformative statistical tests mentioned below in ??.

¹³ This is a standard result in Normal distribution theory; see Mardia *et al.* (1979, chapter 3).

One loose end needs to be tidied up. The Fisher information $i(\theta)$ in (10) can be replaced by an estimate, for convenience, and the asymptotic properties still hold as long as the estimator converges appropriately for large n . One suggestion is the *expected Fisher Information*, $i(\hat{\theta}(\mathbf{x}^{\text{obs}}))$. As shown in ??, $\hat{\theta}(\mathbf{X}) \xrightarrow{P} \theta$ for IID X 's from a regular model, and so $i(\hat{\theta}(\mathbf{X})) \xrightarrow{P} i(\theta)$, because i is a continuous function of θ . In practice, however, a different approximation is often preferred, the *observed Fisher Information*, denoted J , which is defined below in ?. Efron and Hinkley (1978) provide a theoretical justification for preferring J , which is that J is closer to $\text{Var}(u(\mathbf{X}, \theta) \mid s(\mathbf{X}) = s^{\text{obs}}; \theta)$, where s is an experimental ancillary statistic (see ?).

expected Fisher Information

Thus we have established, at least in the case where the X 's are IID and the statistical model is regular, that transformation-invariant confidence sets for θ (and consequently for θ') can be approximated, and that this approximation ought to be reasonable when n is large.

What about when n is not large? In that case the nominal coverage of a set such as \mathcal{C} in (11) may differ substantially from the actual coverage, in a way that varies with θ . A computer-intensive resampling approach termed *the bootstrap* provides an elegant correction for this, notably the *prepivotting* approach of Beran (1987). The bootstrap is one of the two most important innovations in computational statistics in the last forty years.

the bootstrap

prepivotting

The other being Markov Chain Monte Carlo (MCMC); see ?.

This is not the place to review the extensive and still-developing literature on the bootstrap; instead, see the review of Davison *et al.* (2003, section 4)—and the other papers in this number of *Statistical Science*, celebrating the silver anniversary of the bootstrap—and Young and Smith (2005, chapter 11); the broader perspective in Efron (1998) is also interesting.

3 Goodness of fit tests

Tests of goodness of fit are also known as (pure) *significance tests*. One starts with a hypothesis about the data \mathbf{X} , and then examines whether the observations \mathbf{x}^{obs} are consistent with that hypothesis. Cox and Hinkley (1974, chapter 3) provide more details.

significance tests

3.1 Simple hypothesis

A *simple hypothesis* completely defines the distribution of \mathbf{X} :

simple hypothesis

$$H_0 : \mathbf{X} \sim f_{\mathbf{X}}(\mathbf{x}).$$

A simple hypothesis is evaluated according to a *P-value*, which is any statistic with a sub-uniform distribution under H_0 . That is, a *P-value* is a statistic p with the property that if $P := p(\mathbf{X})$ then

P-value

$$\Pr(P \leq u; H_0) \leq u \quad \text{for all } u \in [0, 1].$$

Thus, a small $p^{\text{obs}} := p(\mathbf{x}^{\text{obs}})$, say $p^{\text{obs}} = 1\%$, indicates an outcome that would be improbable under H_0 .

The standard way to construct a P -value is using a statistic s with the property that large values of $s(\mathbf{x})$ indicate a worrying—or at least interesting—departure from the hypothesis H_0 . Letting $S := s(\mathbf{X})$, the P -value is defined as

$$p(s') := \Pr(S \geq s'; H_0).$$

Under H_0 the random quantity $P = p(S)$ has a *sub-uniform* distribution in general, and a uniform distribution if S is continuous.

Proof. This proof uses a nifty trick from Casella and Berger (2002, section 8.3.4). Let G be the distribution function of $-S$ under H_0 . Then

$$p(s') = \Pr(S \geq s'; H_0) = \Pr(-S \leq -s'; H_0) = G(-s').$$

Then since $P = p(S) = G(-S)$, the result follows from the PIT, see section 1.5. □

Note that one can get a large P -value by choosing a poor test statistic. For example, $s(\mathbf{x}) = a$ where a is any scalar constant will have a P -value of 1, no matter what the value of \mathbf{x} . So there is no basis to claim, without additional information, that a large P -value supports H_0 .

It is important to be able to choose the test statistic in the light of possible alternatives to H_0 . But where this is challenging, one useful portmanteau test statistic for a simple hypothesis is *Box's test statistic* (Box, 1980). This is $s(\mathbf{x}) = -f_{\mathbf{X}}(\mathbf{x})$. Thus Box's P -value would be

Box's test statistic

$$p^{\text{obs}} = \Pr(f_{\mathbf{X}}(\mathbf{X}) \leq f_{\mathbf{X}}(\mathbf{x}^{\text{obs}}); H_0).$$

A small p^{obs} indicates that \mathbf{x}^{obs} is in the tail of the distribution specified by H_0 .

The P -value for a simple hypothesis can easily be computed by *Monte Carlo integration*. One samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(r)} \stackrel{\text{iid}}{\sim} f_{\mathbf{X}}(\mathbf{x})$, and then computes

Monte Carlo integration

$$\hat{P}(\mathbf{x}^{\text{obs}}; \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(r)}) := r^{-1} \sum_{j=1}^r \mathbb{1}(s^{(j)} \geq s^{\text{obs}}),$$

where $s^{(j)} := s(\mathbf{x}^{(j)})$ and $s^{\text{obs}} := s(\mathbf{x}^{\text{obs}})$. According to the WLLN,

$$\hat{P}(\mathbf{x}^{\text{obs}}; \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(r)}) \xrightarrow{P} E\{\mathbb{1}(S \geq s^{\text{obs}}); H_0\} = p^{\text{obs}}$$

as $r \rightarrow \infty$, and confidence intervals for the estimator with finite r can be established using the CLT.

3.2 Composite hypothesis

The hypothesis might also be a *composite hypothesis*,

composite hypothesis

$$H_0 : \exists \theta \in \Omega_0 \text{ such that } \mathbf{X} \sim f_{\mathbf{X}}(\mathbf{x}; \theta),$$

which involves additional unknown parameters θ . With $\theta = (\theta', \theta'')$, where θ'' are nuisance parameters, often only θ' is specified, in which case $\Omega_0 = \{\theta'_0\} \times \Omega''$. For a composite hypothesis, one common definition of the P -value is

$$p_{\Omega_0}(s') := \sup_{t \in \Omega_0} p(s'; t)$$

where $p(s'; t)$ is a P -value for the simple hypothesis $X \sim f_X(x; t)$. This has a sub-uniform distribution under all points in Ω_0 , i.e.

$$\Pr\{p_{\Omega_0}(S) \leq u; \theta\} \leq u \quad \text{for all } \theta \in \Omega_0.$$

Proof. Follows from the fact that $p_{\Omega_0}(s') \leq u$ implies $p(s'; \theta) \leq u$ for all $\theta \in \Omega_0$. Therefore

$$\Pr\{p_{\Omega_0}(S) \leq u; \theta\} \leq \Pr\{p(S; \theta) \leq u; \theta\} \leq u$$

for all $\theta \in \Omega_0$. □

One special case is worth mentioning, when θ is a scalar parameter. If the statistical model has a Monotone Likelihood Ratio (MLR) in the test statistic s , then the result that the power function is non-decreasing in θ implies that

$$\sup_{t \in \Omega_0} p(s'; \theta) = p(s'; \sup \Omega_0).$$

This case shows that this definition of a P -value for a composite H_0 often gives large P -values from which nothing interesting can be inferred, and modifications are required unless Ω_0 is quite constrained.

As an alternative, *Pearson's chi-squared test* (or goodness of fit test) works for both simple and composite null hypotheses, but only in the case where the X 's are IID from a regular model. This involves a particular test statistic, in which the sample space \mathcal{X} is first partitioned into k bins, and then

$$s(\mathbf{x}) = \sum_{j=1}^k \frac{(o_j - e_j)^2}{e_j}$$

where o_j is the observed number in the j th bin, and e_j is the expected number in the j th bin under H_0 . For a composite null hypothesis with a d -dimensional unknown parameter, each e_j is computed using the value of the MLE:

$$e_j = n \sum_{x \text{ in bin } j} f_X(x; \hat{\theta}(\mathbf{x}^{\text{obs}})).$$

Then $S \xrightarrow{D} \chi_{k-d-1}^2$ under H_0 (not proved here), from which a P -value can be computed.¹⁴ A commonly-used heuristic is to take the asymptotic limit as reliable provided that $e_j \geq 5$ for all j .

Pearson's chi-squared test has the disadvantage that the result will depend on the partition of \mathcal{X} . It is also not very powerful, in that there is no option to tune the choice of test statistic to a particular departure from the statistical model. But it is available in standard statistical software, so it crops up a lot in applied statistics.

See ahead to ?? for the appropriate definitions and results.

Pearson's chi-squared test

¹⁴ Or $T \xrightarrow{D} \chi_{k-1}^2$ for a simple hypothesis. The rule is to subtract one from the degrees of freedom of the χ^2 distribution for each parameter that is replaced by the value of its MLE.

A final caveat. A statistical hypothesis, whether simple or composite, can never really be 'true'. A P -value tries to address the question of whether a proposed hypothesis is adequate, but it is hard to do this without an alternative model to contrast H_0 with, as addressed in ?? and ??. If you have a large enough n then your P -value is very likely to be small, unless you allow the model complexity (i.e. the number of parameters) to increase with n .

Just as it is important not to get excited when your P -value is large, it is also important not to get excited when your P -value is small, if your n is large. In this situation, your model may be perfectly adequate for your purposes, even though it has not captured every wrinkle in the observations. For example, it is a mistake to conclude that paranormal effects exist just because your P -value is 0.0003 in a test with $n = 104,490,000$; see Jefferys (1990).

References

- D. Basu, 1975. Statistical information and likelihood. *Sankhyā*, **37**(1), 1–71. With discussion. 3
- R. Beran, 1987. Prepivoting to reduce level error of confidence sets. *Biometrika*, **74**(3), 457–468. 16
- G.E.P. Box, 1980. Sampling and Bayes' inference in scientific modelling and robustness. *Journal of the Royal Statistical Society, Series A*, **143**(4), 383–430. With discussion. 17
- G. Casella and R.L. Berger, 2002. *Statistical Inference*. Pacific Grove, CA: Duxbury, 2nd edition. 6, 11, 17
- D.R. Cox and D.V. Hinkley, 1974. *Theoretical Statistics*. London: Chapman and Hall. 13, 16
- A.C. Davison, 2003. *Statistical Models*. Cambridge, UK: Cambridge University Press. 8
- A.C. Davison, D.V. Hinkley, and G.A. Young, 2003. Recent developments in bootstrap methodology. *Statistical Science*, **18**(2), 141–157. 16
- B. Efron, 1998. R.A. Fisher in the 21st century. *Statistical Science*, **13**(2), 95–114. With discussion, 114–122. 16
- B. Efron and D.V. Hinkley, 1978. Assessing the accuracy of the maximum likelihood estimator: Observed versus expected fisher information. *Biometrika*, **65**(3), 457–482. 16
- G.R. Grimmett and D.R. Stirzaker, 2001. *Probability and Random Processes*. Oxford University Press, 3rd edition. 4, 9
- W.H. Jefferys, 1990. Bayesian analysis of random event generator data. *Journal of Scientific Exploration*, **4**(2), 153–169. Available online, <http://bayesrules.net/papers/reg.pdf>. 19
- J.F.C. Kingman and S.J. Taylor, 1966. *Introduction to Measure and Probability*. Cambridge UK: Cambridge University Press. 9
- K.V. Mardia, J.T. Kent, and J.M. Bibby, 1979. *Multivariate Analysis*. London: Harcourt Brace & Co. 15
- J.S. Rosenthal, 2006. *A first look at rigorous probability theory*. Singapore: World Scientific Publishing Co. Pte. Ltd., 2nd edition. 9
- M.J. Schervish, 1995. *Theory of Statistics*. New York: Springer. Corrected 2nd printing, 1997. 6, 8, 13

G.A. Young and R.L. Smith, 2005. *Essentials of Statistical Inference*. Cambridge UK: Cambridge University Press. 16