

I. ESTIMATION OF PARAMETERS

1. PARAMETERS; LIKELIHOOD

To do Statistics – to handle the mathematics and data analysis of situations involving randomness – we need to *model* the situation. Here we confine ourselves to models that can be specified by a *parameter*, θ , which will be *finite-dimensional*. Often, θ will be one-dimensional. Usually, the dimensionality will be quite low (at most 5 or 6, say), unless the parameters are vectors or scalars (which will be the case with Multivariate Analysis, Ch. VII. When infinitely many dimensions are needed, one speaks instead of a *non-parametric* model; see Ch. IV. Sometimes, one has a compound model, with a parametric part and a non-parametric part; one speaks then of a *semi-parametric model*.

Things should be kept as simple as possible (but not simpler!) So we should always work with as few parameters as possible – or, in the lowest possible number of dimensions. If we are unsure about what this is, we need to formulate a question on this, and test it on the data. This is the context of *Hypothesis testing*, Ch. II.

We deal with a probability distribution, F , describable by a parameter θ . Our data consists of a random variable X , or random variables X_1, \dots, X_n , drawn from this distribution. A *statistic* is just a function of the data – something we can calculate when we have done our sampling and obtained our data; an *estimator* of θ is a statistic used to estimate a parameter θ . Often our data X_1, \dots, X_n will be *independent and identically distributed (iid)*; we call them independent *copies* drawn from F , or independent *draws* from F . We shall use the same letter F for the probability distribution or law (a measure), and the corresponding probability distribution function (a function); F will be a Lebesgue-Stieltjes measure (function) in the language of Measure Theory (Stochastic Processes, Ch. I – SP I). By the Lebesgue decomposition theorem,

$$F = F_{ac} + F_d + F_{cs} = F_{ac} + F_s,$$

where F_{ac} is the absolutely continuous component (w.r.t. Lebesgue measure; write f for its Radon-Nikodym derivative, called the (probability) *density*

(function) of F, X), F_d is the discrete component (probability mass $m_n > 0$ at a finite or countable set of points x_n), and F_{cs} is the continuous singular component. We often combine the last two, into the *singular* component, F_s . In this course, without further comment, we shall always be dealing with the absolutely continuous case, with density f , *or* with the discrete case, in which case (partly to simplify notation, partly to emphasise that here the base or reference measure is counting measure rather than Lebesgue measure) we write $f(x_n)$ for the probability mass m_n at the point x_n .

The most basic questions to ask about a random variable are ‘how big is it’ (on average), and this is measured by the *mean*,

$$\mu \quad \text{or} \quad \mu_X := E[X],$$

and ‘how variable (or how random) is it’, which is measured by the *variance*

$$\sigma^2, \quad \text{or} \quad \sigma_X^2 := E[(X - E[X])^2] = E[X^2] - [EX]^2.$$

We write

$$\mu_2 := E[X^2], \quad \mu_n := E[X^n] \quad (n = 1, 2, \dots).$$

Our first task is usually to estimate the mean, and we like to be ‘right on average’. We call an estimator S for θ *unbiased* if

$$ES = \theta;$$

otherwise it has *bias* $ES - \theta$. For the mean, we have an obvious estimator, the *sample mean* \bar{X} . This is unbiased, and by the Strong Law of Large Numbers (SLLN – Stochastic Processes), $\bar{X} \rightarrow \mu$ ($n \rightarrow \infty$) a.s.; we say that \bar{X} is *consistent* for μ (we ‘get the right answer in the limit’). For the variance, matters are somewhat more complicated. The sample variance

$$S^2 := \frac{1}{n} \sum_1^n (X_k - \bar{X})^2 = \overline{X^2} - [\bar{X}]^2$$

is consistent, as by SLLN

$$S^2 \rightarrow E[X^2] - [EX]^2 = \sigma^2 \quad (n \rightarrow \infty).$$

However, it is biased, and to obtain the unbiased version we have to divide by $n - 1$ instead of n (as the authors of many textbooks do for this reason –

always check!) For,

$$\begin{aligned}
nS^2 &= \sum_{k=1}^n (X_k - \bar{X})^2 \quad (\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i) \\
&= \sum_k X_k^2 - \frac{2}{n} \sum_{ik} X_k X_i + n \cdot \frac{1}{n^2} \sum_{ij} X_i X_j \\
&= \sum_k X_k^2 - \frac{1}{n} \sum_{ik} X_k X_i.
\end{aligned}$$

Now if $i = k$ $E[X_i X_k] = E[X_k^2] = \mu_2$, and if $i \neq k$ $E[X_i X_k] = E[X_i] \cdot E[X_k] = \mu^2$, by independence. So

$$nE[S^2] = n\mu_2 - \frac{1}{n}[n\mu_2 + n(n-1)\mu^2] = (n-1)[\mu_2 - \mu^2] = (n-1)\sigma^2.$$

So

$$E[S^2] = \frac{n-1}{n}\sigma^2 : \quad E\left[\frac{n}{n-1}S^2\right] = \sigma^2,$$

or

$$E[S_u^2] = \sigma^2, \quad S_u^2 := \frac{1}{n-1} \sum_1^n (X_k - \bar{X})^2.$$

Here S_u^2 is called the *unbiased* (version of the) sample variance.

We recapitulate from Introductory Statistics Ch. II (IS II).

Likelihood.

We write θ for a parameter (scalar or vector), and write such examples as $f(x|\theta)$, which we will call the *density* (w.r.t. Lebesgue measure in the first three examples, counting measure in the fourth – see SP I). Here x is the *argument* of a function, the density function.

If we have n independent copies sampled from this density, the joint density is the product of the marginal densities:

$$f(x_1, \dots, x_n|\theta) = f(x_1|\theta) \dots f(x_n|\theta), \quad (*)$$

which we may abbreviate to

$$f(., \dots, .|\theta) = f(.\theta) \dots f(.\theta),$$

DATA.

Now suppose that the numerical values of the random variables in our

data set are x_1, \dots, x_n . Fisher's great idea of 1912 was to put the data x_i where the arguments x_i were in (*). He called this (later, 1921 on) the *likelihood*, L – a function of the parameter θ :

$$L(\theta) := f(x_1, \dots, x_n | \theta) = f(x_1 | \theta) \dots f(x_n | \theta). \quad (L)$$

The data point will tend to be concentrated where the probability is concentrated. Fisher advocated choosing as our estimate of the (unknown, but non-random) parameter θ , the value(s) $\hat{\theta}$ (or $\hat{\theta}_n$) for which the likelihood $L(\lambda)$ is maximised. This gives the *maximum likelihood estimator* (MLE); the method is the *Method of Maximum Likelihood*. It is intuitive, simple to use and very powerful – ‘everyone’s favourite method of estimating parameters’.

It is often more convenient to use the *log-likelihood*,

$$\ell := \log L,$$

and maximise that instead (as log is increasing, maximising L and ℓ are the same).

Examples.

1. *Normal*, $N(\mu, \sigma)$ (or $N(\mu, \sigma^2)$).

As in IS II, the MLEs are

$$\hat{\mu} = \bar{X}, \quad \hat{\sigma}^2 = S^2 (= \frac{1}{n} \sum_1^n (X_k - \bar{X})^2).$$

But by above, this is biased: to obtain an unbiased estimator for σ^2 , we have to use S_u^2 and divide by $n - 1$ instead of n . So desirable properties of estimators (e.g. being MLE and unbiased, as here) may be incompatible.

Note. If we use X here (in X_1, \dots, X_n , \bar{X} etc.), we are thinking of the X s as random variables (“before sampling”). If we use the corresponding lower-case letters, we are thinking of them as data – the numerical values obtained (“after sampling”). We shall feel free to use either, depending on convenience – but the second is customary in Statistics, so we shall use it as our default option here.

We quote (see e.g. [BF] Th. 2.4) that for $N(\mu, \sigma^2)$

- (i) \bar{X} and S^2 are independent;
- (ii) $\bar{X} \sim N(\mu, \sigma^2/n)$;
- (iii) $nS^2/\sigma^2 \sim \chi^2(n - 1)$.

So (by definition of the Student t -distribution)

$$\sqrt{n-1} \cdot \frac{\sqrt{n}(\bar{X} - \mu)/\sigma}{\sqrt{n}S/\sigma} = \sqrt{n-1}(\bar{X} - \mu)/S \sim t(n-1).$$

Note that σ (a nuisance parameter if we are interested in the mean) cancels.

As in IS II:

2. *Poisson* $P(\lambda)$: $\bar{\lambda} = \bar{x}$.
3. *Exponential* $E(\lambda)$: $\bar{\lambda} = 1/\bar{x}$.

The first example is a two-parameter problem, the next two are one-parameter problems. But the first example contains two one-parameter sub-problems:

- 1a. *Normal* $N(\mu, \sigma^2)$, σ *known*. The calculation above gives $\hat{\mu} = \bar{x}$ again. Note that $\hat{\mu} \sim N(\mu, \sigma^2/n)$ (whether or not σ is known).
- 1b. *Normal* $N(\mu, \sigma^2)$, μ *known*. The calculation above gives

$$\hat{\sigma}^2 = \frac{1}{n} \sum_1^n (x_i - \mu)^2.$$

This is now a statistics, as μ is known – call it S_μ^2 . Then (recall that $\chi^2(r)$ is the distribution of the sum of the squares of r copies of standard normals)

$$nS_\mu^2/\sigma^2 \sim \chi^2(n).$$

By contrast, in Ex. 1,

$$nS^2/\sigma^2 \sim \chi^2(n-1)$$

(see e.g. [BF], Th. 2.4). We shall see other differences in Ch. II on Hypothesis Testing: the tests used vary depending on what is known.

2. THE CRAMÉR-RAO INEQUALITY

As above: we like parameter estimates to be unbiased ("get it right on average"). We also like estimates to be precise ("have values close together" – as little randomness as possible). We can think of precision as the reciprocal of the variance, so we like maximum precision, or minimum variance. Thus an ideal estimator is minimum-variance unbiased (MVU), and we shall study such estimators below.

But before we do this, it is important to consider the trade-off between precision and bias. Consider, by analogy, setting the sights for a rifle. Bias concerns whether the weapon fires, say, too high or to the right, etc. Precision concerns the grouping of a number of shots. One would prefer a precision weapon firing a bit high to a blunderbuss, with its shots all over the place but 'right on average'. One can formalise this, using the language of Decision Theory, but we shall not do this.

We now focus on MVU estimators. The remarkable thing is that there

are theoretical limits to the accuracy they can attain.

As above, we have a joint density $f = f(x_1, \dots, x_n; \theta)$, which we write as $f = f(x; \theta)$. This integrates to 1: $\int f(x; \theta) dx = 1$ (where dx is n -dimensional Lebesgue measure), which we abbreviate to

$$\int f = 1.$$

We assume throughout that $f(x; \theta)$ is smooth enough for use to differentiate under the integral sign (w.r.t. dx , understood) w.r.t. θ , twice. Then

$$\int \frac{\partial f}{\partial \theta} = \frac{\partial}{\partial \theta} \int f = \frac{\partial}{\partial \theta} 1 = 0 : \quad \int \left(\frac{1}{f} \frac{\partial f}{\partial \theta} \right) \cdot f = 0 : \quad \int \left(\frac{\partial}{\partial \theta} \log f \right) \cdot f = 0.$$

Now $E[g(X)] = \int g(x) f(x; \theta) dx = \int g f$, so in probabilistic language this says

$$E\left[\frac{\partial \log L}{\partial \theta}\right] = 0 : \quad E\left[\frac{\partial \ell}{\partial \theta}\right] = 0 : \quad E[\ell'(\theta)] = 0.$$

We now introduce the (Fisher) *score function*

$$s(\theta) := \ell'(\theta) :$$

$$E[s(\theta)] = 0. \tag{a}$$

Differentiate under the integral sign wrt θ again:

$$\frac{\partial}{\partial \theta} \int \left(\frac{1}{f} \frac{\partial f}{\partial \theta} \right) \cdot f = 0, \quad \int \frac{\partial}{\partial \theta} \left[\left(\frac{1}{f} \frac{\partial f}{\partial \theta} \right) \cdot f \right] = 0 :$$

$$\int \left[\left(\frac{1}{f} \frac{\partial f}{\partial \theta} \right) \frac{\partial f}{\partial \theta} + f \frac{\partial}{\partial \theta} \left(\frac{1}{f} \frac{\partial f}{\partial \theta} \right) \right] = 0.$$

As the bracket in the second term is $\partial \log f / \partial \theta$, this says

$$\int \left[\left(\frac{1}{f} \frac{\partial f}{\partial \theta} \right)^2 + \frac{\partial}{\partial \theta} \left(\frac{\partial \log f}{\partial \theta} \right) \right] f = 0, \quad \int \left[\left(\frac{\partial \log f}{\partial \theta} \right)^2 + \frac{\partial^2}{\partial \theta^2} (\log f) \right] f = 0,$$

or as above

$$E\left[\left(\frac{\partial}{\partial \theta} \log L\right)^2 + \frac{\partial^2}{\partial \theta^2} \log L\right] = 0 : \quad E[\{\ell'(\theta)\}^2 + \ell''(\theta)] = 0 :$$

$$E[s(\theta)^2 + s'(\theta)] = 0. \tag{b}$$

We write

$$I(\theta) := E[\{\ell'(\theta)\}^2] = -E[\ell''(\theta)] : \quad I(\theta) = E[s^2(\theta)] = -E[s'(\theta)], \quad (c)$$

and call $I(\theta)$ the (Fisher) *information* on θ (in the sample (x_1, \dots, x_n)). By (a) and (c):

Proposition. The score function $s(\theta) := \ell'(\theta)$ has mean 0 and variance $I(\theta)$.

When x_1, \dots, x_n are independent, the joint density is the product of the marginal densities; so the log-likelihoods ℓ add; so the informations $-E[\ell''] = -E[s']$ (from (c)) add: *the information in a sample of size n is n times the information per reading*. Also from (c), $s^2 \geq 0$, so $E[s^2] \geq 0$: *information is non-negative*. These two properties suggest that the term information is indeed well chosen.

Theorem (Cramér-Rao Inequality, or Information Inequality, H. Cramér (1946), C. R. Rao (1945)). Let $Y = u(\mathbf{X})$ be any unbiased estimator of θ . Then the minimum variance bound for $\text{var } Y$ is

$$\text{var } Y \geq 1/I(\theta, \mathbf{X}) = 1/(nI(\theta)),$$

where $I(\theta)$ is the information per reading.

Proof. As $Y = u(\mathbf{X})$ is unbiased,

$$\theta = E[u(\mathbf{X})] = \int u(\mathbf{x})f(\mathbf{x}; \theta)d\mathbf{x} = \int uf.$$

$\partial/\partial\theta$:

$$1 = \frac{\partial}{\partial\theta} \int uf = \int u \left(\frac{1}{f} \frac{\partial f}{\partial\theta} \right) f = \int u (\partial \log f / \partial\theta) f :$$

$$1 = E[u \partial \log L / \partial\theta] = E[u\ell'] = E[us].$$

By (a), (b) and (c),

$$\text{var } s = \text{var } \ell' = E[(\ell')^2] = I(\theta; \mathbf{X}), = I(\theta),$$

say. The correlation coefficient is

$$\rho := \rho(u, s) = \frac{\text{cov}(u, s)}{\sqrt{\text{var } u} \sqrt{\text{var } s}} = \frac{E[us] - E[u]E[s]}{\sqrt{\text{var } u} \sqrt{I}} = \frac{1}{\sqrt{\text{var } u} \sqrt{I}},$$

as $E[s] = 0$, $E[us] = 1$. But

$$\rho^2 \leq 1$$

(correlation bound: Cauchy-Schwarz Inequality). So

$$\text{var} u \geq 1/I. \quad //$$

Defn. We call an estimator *efficient* if it is unbiased and its variance achieves the CR lower bound, *asymptotically efficient* if its bias tends to 0 and its variance achieves the CR bound asymptotically.

An efficient (= minimum-variance unbiased, MVU) estimator is also called a *best* estimator.

When dealing with regression (Ch. V), we shall often meet *linear* estimators; the above then become *BLUES* (best linear unbiased estimators).

3. LARGE-SAMPLE PROPERTIES OF MAXIMUM-LIKELIHOOD ESTIMATORS

We assume the following regularity conditions:

- (i) differentiability under the integral sign twice (as before);
- (ii) finite positive Fisher information per reading $I(\theta)$;
- (iii) In some neighbourhood N of the true parameter value θ_0 ,

$$\left| \frac{\partial^3}{\partial \theta^3} \log f(\mathbf{x}; \theta) \right| \leq H(\mathbf{x}), \quad \text{where} \quad \sup_{\theta \in N} E_{\theta} H(\mathbf{X}) \leq M < \infty.$$

Theorem (Cramér, 1946). Under the above regularity conditions, the MLE $\hat{\theta}$ of the true parameter value θ_0 is

- (i) *strongly consistent*: $\hat{\theta} \rightarrow \theta_0$ as $n \rightarrow \infty$, a.s.,
- (ii) *asymptotically efficient*: $\text{var } \hat{\theta} \sim 1/(nI(\theta))$, the Cramér-Rao lower bound;
- (iii) *asymptotically normal*: $\sqrt{nI(\theta_0)}(\hat{\theta} - \theta_0) \rightarrow \Phi = N(0, 1)$ ($n \rightarrow \infty$).

Proof. (i) We use Taylor's theorem to expand the score function $s = \ell'$ about $\theta = \theta_0$:

$$s(\theta) = s(\theta_0) + (\theta - \theta_0)s'(\theta_0) + \frac{1}{2}(\theta - \theta_0)^2 s''(\theta^*),$$

for some θ^* between θ_0 and θ . Since $\ell(x_1, \dots, x_n; \theta) = \sum_1^n \ell(x_i; \theta)$, and similarly for $s = \ell'$, this says (on dividing by n)

$$s(\mathbf{x}; \theta)/n = \frac{1}{n} \sum_1^n s(x_i; \theta) = \frac{1}{n} \sum_1^n s(x_i; \theta_0) + (\theta - \theta_0) \cdot \frac{1}{n} \sum_1^n s'(x_i; \theta_0)$$

$$+\frac{1}{2}(\theta - \theta_0)^2 \cdot \frac{1}{n} \sum_1^n s''(x_i; \theta^*). \quad (*)$$

The first term on the RHS is an average of iid rvs with mean $Es(\theta_0) = 0$, by (a) of I.2. So by SLLN, this $\rightarrow 0$ a.s. as $n \rightarrow \infty$. Similarly, by SLLN the second term on RHS converges a.s. to

$$(\theta - \theta_0)s'(\theta_0) = -I(\theta_0)(\theta - \theta_0).$$

The third term on RHS of (*) is bounded by $\frac{1}{2}M(\theta - \theta_0)^2$, by our regularity assumption (iii) (as $s'' = \ell'''$). For θ close enough to θ_0 , this is negligible wrt the second term. So $\text{RHS} \sim \text{second term}$:

$$\text{RHS} \sim -I(\theta_0) \cdot (\theta - \theta_0).$$

Since $I(\theta_0) \in (0, \infty)$, by (ii), the *sign* of the RHS is thus *opposite* to that of $\theta - \theta_0$, for large enough n and θ close enough to θ_0 . For such n and θ , RHS *changes sign* in every neighbourhood of θ_0 (just take θ through θ_0). But LHS = RHS, so the LHS too change sign in every neighbourhood of θ_0 , for large enough n . This says that there is a root $\hat{\theta} = \hat{\theta}_n$ of the *likelihood equation*

$$s(\mathbf{x}; \theta) = 0 \quad (LE)$$

in every neighbourhood of θ_0 .

Since this neighbourhood can be arbitrarily small, we must have

$$\hat{\theta} = \hat{\theta}_n \rightarrow \theta_0 \quad (n \rightarrow \infty) \quad a.s.,$$

proving the strong consistency of the MLE $\hat{\theta}$ and (i).

Now put $\theta = \hat{\theta}$ in (*). The LHS is 0, by definition of MLE (recall $s = \ell' = (\log L)'$). The third term on RHS is negligible wrt the second term, because of the extra factor $\theta - \theta_0 \rightarrow 0$, by (i). So we can neglect this term, leaving

$$0 \sim \frac{1}{n} \sum_1^n s(x_i; \theta_0) + (\hat{\theta} - \theta_0) \cdot \frac{1}{n} \sum_1^n s'(x_i; \theta_0).$$

Rearranging,

$$\sqrt{nI(\theta_0)}(\hat{\theta} - \theta_0) \sim \frac{[\sum_1^n s(x_i; \theta_0)]/\sqrt{nI(\theta_0)}}{[\sum_1^n s'(x_i; \theta_0)]/(nI(\theta_0))} \quad (n \rightarrow \infty). \quad (**)$$

By CLT, the numerator on RHS $\rightarrow \Phi = N(0, 1)$, as $s(x_i; \theta_0)$ are iid with mean 0 and variance $I(\theta_0)$. By LLN, the denominator on RHS $\rightarrow 1$ a.s., as $-s'(x_i; \theta_0)$ are iid with mean $I(\theta_0)$ (by I.2, Prop.). Combining, RHS $\rightarrow \Phi$. So LHS $\rightarrow \Phi$:

$$\sqrt{nI(\theta_0)}(\hat{\theta} - \theta_0) \rightarrow \Phi = N(0, 1) : \quad \text{var } (\hat{\theta} - \theta_0) = \text{var } \hat{\theta} \sim 1/(nI(\theta_0)),$$

the Cramér-Rao bound, so $\hat{\theta}$ is asymptotically efficient, proving (ii), and $\hat{\theta}$ is asymptotically normal, proving (iii). //

The regularity conditions above can be weakened, at the cost of a harder proof, but *some* regularity conditions are needed. The phrase "under suitable regularity conditions" recurs with remorseless regularity in textbook treatments of large-sample properties of MLEs.

Example. Uniform $U(a, b)$, or $U(\theta - \frac{1}{2}, \theta + \frac{1}{2})$. See IS II, where we showed that here the MLEs converge at a different rate (n , not \sqrt{n}) and to a different limit (exponential, or symmetric exponential, not normal).

Vector Parameters.

If $\theta = (\theta_1, \dots, \theta_r)$ is an r -dimensional (vector) parameter, one can proceed as above. We now obtain the (Fisher) *information matrix*

$$I(\theta) = (I_{ij}(\theta))_{i,j=1}^r,$$

where we use suffix notation for partial differentiation ($g_i := \partial g / \partial \theta_i$, etc.) and

$$I_{ij}(\theta) := E[\ell_i(\theta)\ell_j(\theta)] = E[-\ell_{ij}(\theta)].$$

Under regularity conditions as above (we assume the information matrix is positive definite, so we can invert it), we again obtain consistency, asymptotic efficiency and asymptotic normality:

$$\hat{\theta} \sim N_r(\theta, n^{-1}I^{-1}(\theta)).$$

Stochastic process versions.

The true context for results such as the above is not random variables as above, but stochastic processes. Infinite-dimensional versions are possible (and needed), in which conclusions are drawn, not for one time-point at a time, but for infinitely many together – say, all $t \in [0, 1]$, or all $t \geq 0$. We shall develop these ideas later (Day 3). Meanwhile, we mention our main

source for such things,

[vdVW] Aad van der VAART and Jon A. WELLNER, *Weak convergence and empirical processes, with applications to statistics*, Springer, 1996.

In particular, [vdVW] contains detailed accounts of *M-estimators* (3.2) ('M for maximum' – generalising MLEs), and *Z-estimators* (3.3) ('Z for zero': the MLE is a zero of the *likelihood equation* $\ell' = 0$).

Iterative solution of the Likelihood Equation

It may not be possible to solve the Likelihood Equation $\ell' = 0$ (*LE*) in closed form. In such cases, we have to proceed as elsewhere in Mathematics – in particular, in Numerical Analysis – and proceed iteratively.

To assess the problem, begin by drawing a rough graph of ℓ . By looking for sign changes, and using trial values, it is usually possible (without excessive effort) to find a rough approximation to the desired root (there may – will in general – be multiple roots, but usually the root we need will be clear enough from context). Call this trial value t . Then (with $s = \ell'$)

$$0 = s(\hat{\theta}) = s(t) + ((\hat{\theta} - t)s'(\theta^*)),$$

with θ^* between t and $\hat{\theta}$. Solving,

$$\hat{\theta} = t - s(t)/s'(\theta^*). \quad (*)$$

We now have a choice about how to proceed. We know that $\hat{\theta}$ is (strongly) consistent, $\hat{\theta} \rightarrow \theta_0$, so $\hat{\theta} \sim \theta_0$, so with a good enough starting value t , also $t \sim \hat{\theta}(\sim \theta_0)$ and $\theta^* \sim \hat{\theta}(\sim \theta_0)$.

Newton-Raphson iteration.

This is also known as the *tangent approximation*. It relies on replacing a function by its tangent near a point. If x_n is near a root of

$$f(x) = 0,$$

then a better approximation is

$$x_{n+1} := x_n - f(x_n)/f'(x_n).$$

So starting from the approximation t , replacing $s'(\theta^*)$ in $(*)$ by $s'(t)$ gives a better approximation; this is the Newton-Raphson method.

Fisher's method of scoring.

Here we replace $s'(\theta^*)$ by $E[s'(t)] = 1/I(t)$ (we know $\theta^* \sim t$ by above;

formally, $s'(t)$ is non-random (so equal to its expectation), but in practice we may have found it by data-dependent [random] means, so we reduce to the non-random case by taking its expectation). Then our next (better) approximation is

$$\hat{\theta} \sim t - s(t)/E[s'(t)] = t + s(t)I(t).$$

This is Fisher's *method of scoring*.

As always with iterations: to implement this numerically, one needs a "do-loop" (while do, else stop).

Exercise. Implement this in C++ (the "official programming language" for this course), for the Cauchy location family (Problems 3 Q3). First, choose a μ (arbitrarily – or, by sampling from a chosen distribution). Then, sample from the Cauchy distribution with this μ . Then, perform the above iterations (by either, or better still both, of the Newton-Raphson and Fisher methods), to estimate this μ from the data.

Reparametrisation and the Delta Method.

Suppose we are using parameter θ , but wish to change to some alternative parametrisation, $g(\theta)$, where g is continuously differentiable. A CLT for θ such as

$$\sqrt{n}(T_n - \theta) \rightarrow N(0, \sigma(\theta)^2)$$

(as holds above, with T_n the MLE $\hat{\theta}$ based on a sample of size n and $\sigma^2(\theta) = 1/I(\theta)$) transforms into a CLT for $g(\theta)$:

$$\sqrt{n}(g(T_n) - g(\theta)) \rightarrow N(0, [g'(\theta)\sigma(\theta)]^2).$$

For,

$$g(T_n) - g(\theta) = (T_n - \theta)(g'(\theta) + \epsilon_n),$$

with ϵ_n a (random) error term. One can show this to be negligible for large n , so

$$g(T_n) - g(\theta) \sim (T_n - \theta)g'(\theta).$$

Since $\text{var}(cX) = c^2 \text{var } X$, the result follows.

This is called the *delta method*, and is often useful. It can be extended from random variables to stochastic processes (i.e., from one or finitely many to infinitely many dimensions), and we shall meet it again later.