# Chapter 6

# Hellinger differentiability

*SECTION 1 relates Hellinger differentiability to the classical regularity conditions for maximum likelihood theory.*

*SECTION 2 discusses connections between Hellinger differentiability and pointwise differentiability of densities, leading to a sufficient condition for Hellinger differentiability.*

*SECTION 3 derives the information inequality, as an illustration of the elegance brought into statistical theory by Hellinger differentiability.*

*SECTION 4 explains why Hellinger differentiability almost implies contiguity for product measures.*

*SECTION 5 explains how one can dispense with the domination assumption when defining Hellinger differentiability, at the cost of a natural extra assumption regarding non-dominated components. The slightly strengthened concept is called Differentiability in Quadratic Mean (DQM) to avoid confusion.*

*SECTION 6 shows that DQM is preserved under measurable maps.*

---

*Final two sections not yet edited.*

---

*SECTION 7 derives some subtle consequences of norm differentiability for unit vectors.*

*SECTION 8 shows that Hellinger differentiability of marginal densities implies existence of a local quadratic approximation to the likelihood ratio for product measures.*

## 1. Heuristics

Modern statistical theory makes clever use of the fact that square roots of probability density functions correspond to unit vectors in spaces of square integrable functions. The Hellinger distance between densities corresponds to the $\mathcal{L}^2$ norm of the difference between the unit vectors. This Chapter explains some of the statistical consequences of differentiability in norm of the square root of the density, a property known as Hellinger differentiability.

Throughout the Chapter, $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ will denote a family of probability measures, on a fixed $(\mathcal{X}, \mathcal{A})$, indexed by a subset $\Theta$ of $\mathbb{R}^k$. In many cases there will exist a dominating sigma-finite measure $\lambda$ with respect to which each $P_\theta$ has a density $f_\theta(x) = dP_\theta/d\lambda$. I will then refer to $\mathcal{P}$ as a ***dominated family*** and write $\| \cdot \|_2$ for the $\mathcal{L}^2(\lambda)$ norm and $\xi_\theta(x)$ for the positive square root of $f_\theta(x)$.

Most results in the Chapter will concern behavior near some arbitrarily chosen point $\theta_0$ of $\Theta$. For simplicity of notation, I will usually assume $\theta_0 = 0$, except in a few basic definitions. Thus an expression such as $(\theta - \theta_0)'\dot{\xi}_{\theta_0}$ will simplify to $\theta'\dot{\xi}_0$, a form that is easier to read and occupies less space on the

page. The simplification involves no loss of theoretical generality, because the same effect could always be achieved by a reparametrization, $\theta := t + \theta_0$.

The traditional regularity conditions for asymptotic statistical theory involve existence of two or three derivatives of density functions, together with domination assumptions to justify differentiation under integral signs. Le Cam (1970) noted that such conditions are unnecessarily stringent. He commented:

> Even if one is not interested in the maximum economy of assumptions one cannot escape practical statistical problems in which apparently "slight" violations of the assumptions occur. For instance the derivatives fail to exist at one point $x$ which may depend on $\theta$, or the distributions may not be mutually absolutely continuous or a variety of other difficulties may occur. The existing literature is rather unclear about what may happen in these circumstances. Note also that since the conditions are imposed upon probability densities they may be satisfied for one choice of such densities but not for certain other choices.

Probably Le Cam had in mind examples such as the double exponential density, $\frac{1}{2}\exp(-|x - \theta|)$, for which differentiability fails at the point $\theta = x$. He showed that the traditional conditions can, for some purposes, be replaced by a simpler assumption of ***Hellinger differentiability***: differentiability in norm of the square root of the density as an element of an $\mathcal{L}^2$ space. The derivation of the information inequality in Section 3 will provide a simple illustration of this point.

<1>     **Definition.**    *Write $\mathcal{L}_+^1(\lambda)$ for the set of nonnegative functions that are integrable with respect to a sigma-finite measure $\lambda$.*

*Say that a set $\mathcal{F} = \{f_\theta : \theta \in \Theta\} \subseteq \mathcal{L}_+^1(\lambda)$, indexed by a subset $\Theta$ of $\mathbb{R}^k$, is **Hellinger differentiable** at a point $\theta_0$ of $\Theta$ if the map $\theta \mapsto \xi_\theta(x) := \sqrt{f_\theta(x)}$ is differentiable in $\mathcal{L}^2(\lambda)$ norm at $\theta_0$, that is, if there exists a vector $\dot{\xi}_{\theta_0}(x)$ of functions in $\mathcal{L}^2(\lambda)$ such that*

<2>      $\xi_\theta(x) = \xi_{\theta_0}(x) + (\theta - \theta_0)'\dot{\xi}_{\theta_0}(x) + r_\theta(x)$        *with $\|r_\theta\|_2 = o(|\theta - \theta_0|)$ near $\theta_0$.*

*Call $\dot{\xi}_{\theta_0}(x)$ the Hellinger derivative at $\theta_0$.*

*In particular, if $f_\theta = dP_\theta/d\lambda$ for a family of probability measures $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, say that $\mathcal{P}$ is Hellinger differentiable at $\theta_0$ if $\mathcal{F}$ is Hellinger differentiable at $\theta_0$.*

Some authors (for example, Bickel, Klaassen, Ritov & Wellner (1993, page 202)) adopt a slightly different definition,

<3>                      $\xi_\theta(x) = \xi_{\theta_0}(x) + \frac{1}{2}(\theta - \theta_0)'\Delta(x)\xi_{\theta_0}(x) + r_\theta(x),$

replacing the Hellinger derivative $\dot{\xi}_{\theta_0}$ by $\frac{1}{2}\Delta(x)\xi_{\theta_0}(x)$. As explained in Section 5, the modification very cleverly adds an extra regularity assumption to the definition. The two definitions are not completely equivalent.

Classical statistical theory, especially when dealing with independent observations from a $P_\theta$, makes heavy use of the function $\ell_\theta(x) := \log f_\theta(x)$, where $f_\theta = dP_\theta/d\lambda$. The vector $\dot{\ell}_\theta(x)$ of partial derivatives with respect to $\theta$ is called the ***score function***. The variance matrix $\mathbb{I}_\theta$ of the score function is called the ***Fisher information matrix*** for the model. The classical regularity conditions justify differentiation under the integral sign to get

<4>                      $$P_\theta \dot{\ell}_\theta(x) = \lambda \dot{f}_\theta(x) = \frac{\partial}{\partial\theta}\lambda f_\theta(x) = 0,$$

whence $\mathbb{I}_\theta := \mathrm{var}_\theta\left(\dot{\ell}_\theta\right) = P_\theta\left(\dot{\ell}_\theta \dot{\ell}_\theta'\right)$.

Under assumptions of Hellinger differentiability, the derivative $\dot{\xi}_\theta$ takes over the role of the score vector. Ignoring problems related to division by zero and distinctions between pointwise and $\mathcal{L}^2(\lambda)$ differentiability, we would have

$$\frac{2\dot{\xi}_\theta(x)}{\xi_\theta(x)} \stackrel{?}{=} \frac{2}{\sqrt{f_\theta(x)}} \frac{\partial}{\partial\theta} \sqrt{f_\theta(x)} = \frac{1}{f_\theta(x)} \frac{\partial f_\theta(x)}{\partial\theta} = \dot{\ell}_\theta(x).$$

Thus the $\Delta$ in the modified definition $\langle3\rangle$ corresponds to the score function.

The equality $\langle4\rangle$ corresponds to the assertion $P_\theta\left(\dot{\xi}_\theta/\xi_\theta\right) = \lambda\left(\xi_\theta\dot{\xi}\right) = 0$, which Section 7 will show to be a consequence of Hellinger differentiability and the fact that $\|\xi_\theta\|_2 =$ for all $\theta$. The Fisher information $\mathbb{I}_\theta$ at $\theta$ corresponds to the matrix

$$P_{\theta_0}\left(\dot{\ell}_\theta\dot{\ell}_\theta'\right) \stackrel{?}{=} 4P_{\theta_0}\left(\dot{\xi}_\theta\dot{\xi}_\theta'/\xi_\theta^2\right) \stackrel{?}{=} 4\lambda\left(\dot{\xi}_\theta\dot{\xi}_\theta'\right).$$

Here I flag both equalities as slightly suspect, not just for the unsupported assumption of equivalence between pointwise and Hellinger differentiabilities, but also because of a possible 0/0 cancellation. For the moment it is better to insert an explicit indicator function, $\{\xi_\theta > 0\}$, to protect against 0/0. To avoid possible ambiguity or confusion, I will write $\mathbb{I}_\theta$ for $4\lambda(\dot{\xi}_\theta\dot{\xi}_\theta')$ and $\mathbb{I}_\theta^\circ$ for $4\lambda(\dot{\xi}_\theta\dot{\xi}_\theta'\{\xi_\theta > 0\})$, to hint at equivalent forms for $\mathbb{I}_\theta$ without yet giving precise conditions under which all three exist and are equal. See Section 4 for an explanation of when the distinction is necessary.

The classical assumptions also justify further interchanges of integrals and derivatives, to derive an alternative representation $\mathbb{I}_\theta = -\mathbb{P}_\theta\ddot{\ell}_\theta$ for the information matrix. It might seem obvious that there can be no analog of this representation for Hellinger differentiability. Indeed, how could an assumption of one-times differentiability, in norm, imply anything about a second derivative? Surprisingly, there is a way, if we think of second derivatives as coefficients of quadratic terms in local approximations. As will be shown in Section 8, the fact that $\|\xi_\theta\|_2 = 1$ for all $\theta$ leads to a quadratic approximation for a log-likelihood ratio—a sort of Taylor expansion to quadratic terms without the usual assumption of twice continuous differentiability. Remarkable.

## 2. A sufficient condition for Hellinger differentiability

How does Hellinger differentiability relate to the classical assumption of pointwise differentiability?

Roughly speaking, the difference between the two concepts is like the difference between convergence in $\mathcal{L}^2$ and convergence almost surely. In fact, it is easy (Problem [1]) to adapt a standard counterexample to show that Hellinger differentiability does not imply pointwise differentiability.

Consider the case where $\Theta$ is one-dimensional, and $f_\theta$ is both Hellinger differentiable and differentiable a.e. [$\lambda$] at $\theta = 0$. Choose a sequence $\{\theta_n\}$ tending to zero so fast that $\sum_n \|r_{\theta_n}\|_2/|\theta_n| < \infty$, which implies $r_{\theta_n}(x) = o(|\theta_n|)$ a.e. [$\lambda$]. For almost all $x$,

$$\xi_{\theta_n}(x) = \xi_0(x) + \theta_n\dot{\xi}_0(x) + o(|\theta_n|)$$
$$\xi_{\theta_n}(x)^2 = \xi_0(x)^2 + \theta_n f_0'(x) + o(|\theta_n|).$$

If $f_0(x) \neq 0$, the second equation can be rewritten as

$$\xi_{\theta_n}(x) = \xi_0(x)\left(1 + \theta_n\frac{f_0'(x)}{\xi_0(x)^2} + o(|\theta_n|)\right)^{1/2} = \xi_0(x) + \tfrac{1}{2}\theta_n\frac{f_0'(x)}{\xi_0(x)} + o(|\theta_n|).$$

It follows (cf. differentiation of $\sqrt{f_\theta(x)}$ by first principles) that $f_0'(x) = 2\xi_0(x)\dot{\xi}_0(x)$. At an $x$ where $f_0(x) = 0$, this argument fails. Instead we would have

$$\xi_{\theta_n}(x)^2 = \theta_n^2 \dot{\xi}_0(x)^2 + o(|\theta_n|^2)$$
$$\xi_{\theta_n}(x)^2 = \theta_n f_0'(x) + o(|\theta_n|).$$

We then deduce that $f_0'(x) = 0$ but apparently we no longer have any control over $\dot{\xi}_0(x)$. However, if 0 is an interior point of the parameter space $\Theta$ we could repeat the argument with $\{\theta_n\}$ replaced by $\{-\theta_n\}$, obtaining for almost all $x$ for which $\xi_0(x) = 0$ that

$$\xi_{\pm\theta_n}(x) = \pm\theta_n\dot{\xi}_0(x) + o(|\theta_n|).$$

Nonnegativity of $\xi_\theta$ would then force $\dot{\xi}_0(x) = 0$.

In summary: If the $f_\theta(x)$ are pointwise differentiable at $\theta = 0$ for almost all $x$ and if 0 is an interior point of $\Theta$ then the only possible candidate (up to an almost sure equivalence) for the Hellinger derivative at 0 is

$$\dot{\xi}_0(x) = \tfrac{1}{2}\frac{f_0'(x)}{\xi_0(x)}$$

What more do we need in order to show that this $\dot{\xi}_0$ is, in fact, an $\mathcal{L}^2(\lambda)$ derivative of $\theta_\theta$ at $\theta = 0$? The answer requires careful attention to the problem of when functions of a real variable can be recovered as integrals of their derivatives.

<5>   **Definition.**   *A real valued function $H$ defined on an interval $[a, b]$ of the real line is said to be **absolutely continuous** if to each $\epsilon > 0$ there exists a $\delta > 0$ such that $\sum_i |H(b_i) - H(a_i)| < \epsilon$ for all finite collections of nonoverlapping subintervals $[a_i, b_i]$ of $[a, b]$ for which $\sum_i (b_i - a_i) < \delta$.*

*Absolute continuity of a function defined on the whole real line is taken to mean absolute continuity on each finite subinterval.*

The following connection between absolute continuity and integration of derivatives is one of the most celebrated results of classical analysis (UGMTP §3.4).

<6>   **Theorem.**   *A real valued function $H$ defined on an interval $[a, b]$ is absolutely continuous if and only if the following three conditions hold.*

*(i) The derivative $H'(t)$ exists at Lebesgue almost all points of $[a, b]$.*

*(ii) The derivative $H'$ is Lebesgue integrable*

*(iii) $H(t) - H(a) = \int_a^t H'(s)\,ds$ for each $t$ in $[a, b]$*

Put another way, a function $H$ is absolutely continuous on an interval $[a, b]$ if and only if there exists an integrable function $h$ for which

<7>
$$H(t) = \int_a^t h(s)\,ds \qquad \text{for all } t \text{ in } [a, b]$$

The function $H$ must then have derivative $h(t)$ at almost all $t$. As a systematic convention we could take $h$ equal to the measurable function

$$\dot{H}(t) = \begin{cases} H'(t) & \text{at points } t \text{ where the derivative exists,} \\ 0 & \text{elsewhere.} \end{cases}$$

I will refer to $\dot{H}$ as the **density** of $H$. Of course it is actually immaterial how $\dot{H}$ is defined on the Lebesgue negligible set of points at which the derivative does not exist, but the convention helps to avoid ambiguity.

Now consider a *nonnegative* function $H$ that is differentiable at a point $t$. If $H(t) > 0$ then the chain rule of elementary calculus implies that the function

$2\sqrt{H}$ is also differentiable at $t$, with derivative $H'(t)/\sqrt{H(t)}$. At points where $H(t) = 0$, the question of differentiability becomes more delicate, because the map $y \mapsto \sqrt{y}$ is not differentiable at the origin. If $t$ is an internal point of the interval and $H(t) = 0$ then we must have $H'(t) = 0$. Thus $H(y) = o(|y - t|)$ near $t$. If $\sqrt{H}$ had a derivative at $t$ then $\sqrt{H(y)} = o(|y - t|)$ near $t$, and hence $H(y) = o(|y - t|^2)$. Clearly we need to take some care with the question of differentiability at points where $H$ equals zero.

Even more delicate is the fact that absolute continuity of a nonnegative function $H$ need not imply absolute continuity of the function $\sqrt{H}$, without further assumptions—even if $H$ is everywhere differentiable (Problem [2]).

<8>   **Lemma.**   *Suppose a nonnegative function $H$ is absolutely continuous on an interval $[a, b]$, with density $\dot{H}$. Let $\Delta(t) := \frac{1}{2}\dot{H}(t)\{H(t) > 0\}/\sqrt{H(t)}$. If $\int_a^b |\Delta(t)|\, dx < \infty$ then $\sqrt{H}$ is absolutely continuous, with density $\Delta$, that is,*

$$\sqrt{H(t)} - \sqrt{H(a)} = \int_a^t \Delta(s)\, ds \qquad \text{for all } t \text{ in } [a, b]$$

*Proof.*   Fix an $\eta > 0$. The function $H_\eta := \eta + H$ is bounded away from zero, and hence $\sqrt{H_\eta}$ has derivative $H_\eta' = H'/(2\sqrt{H + \eta})$ at each point where the derivative $H'$ exists. Moreover, absolute continuity of $\sqrt{H_\eta}$ follows directly from the Definition <5>, because

$$|\sqrt{H_\eta(b_i)} - \sqrt{H_\eta(a_i)}| = \frac{|H_\eta(b_i) - H_\eta(a_i)|}{\sqrt{H_\eta(b_i)} + \sqrt{H_\eta(a_i)}} \leq \frac{|H(b_i) - H(a_i)|}{2\sqrt{\eta}}$$

for each interval $[a_i, b_i]$. From Theorem <6>, for each $t$ in $[a, b]$,

$$\sqrt{H(t) + \eta} - \sqrt{H(a) + \eta} = \int_a^t \frac{\dot{H}(s)}{2\sqrt{H(s) + \eta}}\, ds.$$

As $\eta$ decreases to zero, the left-hand side converges to $\sqrt{H(t)} - \sqrt{H(a)}$. The integrand on the right-hand side converges to $\Delta(s)$ at points where $H(s) > 0$. For almost all $s$ in $\{H = 0\}$ the derivative $H'(s)$ exists and equals zero; the integrand converges to $0 = \Delta(s)$ at those points. By Dominated Convergence, the right-hand side converges to $\int_a^t \Delta(s)\, ds$.   □

The integral representation for the square root of an absolutely continuous function is often the key to proofs of Hellinger differentiability. For simplicity of notation, the following sufficient condition is stated only for a one-dimensional $\Theta$ with $0$ as an interior point.

<9>   **Theorem.**   *Suppose $\mathcal{F} = \{f_\theta(x) : |\theta| < \delta\} \subseteq \mathcal{L}_+^1(\lambda)$ for some $\delta > 0$. Suppose also that*

(i) *the map $(x, \theta) \mapsto f_\theta(x)$ is product measurable;*

(ii) *for $\lambda$ almost all $x$, the function $\theta \mapsto f_\theta(x)$ is absolutely continuous on $[-\delta, \delta]$, with almost sure derivative $\dot{f}_\theta(x)$;*

(iii) *for $\lambda$ almost all $x$, the function $\theta \mapsto f_\theta(x)$ is differentiable at $\theta = 0$;*

(iv) *for each $\theta$ the function $\dot{\xi}_\theta(x) := \frac{1}{2}\dot{f}_\theta(x)\{f_\theta(x) > 0\}/\sqrt{f_\theta(x)}$ belongs to $\mathcal{L}^2(\lambda)$ and $\lambda\dot{\xi}_\theta^2 \to \lambda\dot{\xi}_0^2$ as $\theta \to 0$.*

*Then $\mathcal{F}$ has Hellinger derivative $\dot{\xi}_0(x)$ at $\theta = 0$.*

> REMARK.   Assumption (iii) might appear redundant, because (ii) implies differentiability of $\theta \mapsto f_\theta(x)$ at Lebesgue almost all $\theta$, for $\lambda$-almost all $x$. A mathematical optimist (or Bayesian) might be prepared to gamble that $0$ does not belong to the bad negligible set; a mathematical pessimist might prefer Assumption (iii).

*Proof.*    As before, write $\xi_\theta(x)$ for $\sqrt{f_\theta(x)}$ and define $r_\theta(x) := \xi_\theta(x) - \xi_0(x) - \theta\dot{\xi}_0(x)$. We need to prove that $\lambda r_\theta^2 = o(|\theta|^2)$ as $\theta \to 0$.

Assumption (i) and the convention about densities imply joint measurability of $(x, \theta) \mapsto \dot{f}_\theta(x)$.

For simplicity of notation, consider only positive $\theta$. The arguments for negative $\theta$ are analogous. Write $\mathfrak{m}$ for Lebesgue measure on $[-\delta, \delta]$.

With no loss of generality (or by a suitable decrease in $\delta$) we may assume that $\lambda\dot{\xi}_\theta^2$ is bounded, so that, by Tonelli, $\infty > \mathfrak{m}^\theta\lambda^x\dot{\xi}_\theta(x)^2 = \lambda^x\mathfrak{m}^\theta\dot{\xi}_\theta(x)^2$, implying $\mathfrak{m}^\theta\dot{\xi}_\theta(x)^2 < \infty$ a.e. $[\lambda]$. From Lemma <8> it then follows that

$$\frac{\xi_\theta(x) - \xi_0(x)}{\theta} = \frac{1}{\theta}\int_0^\theta \dot{\xi}_s(x)\,ds \qquad \text{a.e. } [\lambda].$$

By Jensen's inequality for the uniform distribution on $[0, \theta]$, and (iv),

<10>
$$\lambda\left|\frac{\xi_\theta(x) - \xi_0(x)}{\theta}\right|^2 \leq \frac{1}{\theta}\int_0^\theta \lambda\dot{\xi}_s(x)^2\,ds \to \lambda\dot{\xi}_0^2 \qquad \text{as } \theta \to 0.$$

Define nonnegative, measurable functions

$$g_\theta(x) := 2\,|\xi_\theta(x) - \xi_0(x)|^2\,/\theta^2 + 2\dot{\xi}_0(x)^2 - |r_\theta(x)/\theta|^2\,.$$

By (iii), $r_\theta(x)/\theta \to 0$ at almost all $x$ where $\xi_0(x) > 0$, and hence $g_\theta(x) \to 4\dot{\xi}_0(x)^2$. At almost all points where $\xi_0(x) = 0$ we have $\dot{\xi}_0(x) = 0$, so that $\xi_\theta(x) = r_\theta(x)$ and $g_\theta(x) \geq 0$. Thus $\liminf g_\theta(x) \geq 4\dot{\xi}_0(x)^2$ a.e. $[\lambda]$. By Fatou's Lemma (applied along subsequences), followed by an appeal to <10>,

$$4\lambda\dot{\xi}_0^2 \leq \liminf_{\theta \to 0} \lambda g_\theta \leq 4\lambda\dot{\xi}_0^2 - \limsup_{\theta \to 0} \lambda\,|r_\theta(x)/\theta|^2\,.$$

☐    That is, $\lambda r_\theta^2 = o(\theta^2)$, as required for Hellinger differentiability.

<11>    **Example.**    Let $q$ be a probability density with respect to Lebesgue measure $\mathfrak{m}$ on the real line. Suppose $q$ is absolutely continuous, with density $\dot{q}$ for which $\mathbb{I}_q := \mathfrak{m}\big(\{q > 0\}\dot{q}^2/q\big) < \infty$. Define $Q_\theta$ to have density $f_\theta(x) := q(x - \theta)$ with respect to $\lambda$, for each $\theta$ in $\mathbb{R}$. The conditions of Theorem <9> are satisfied, with

$$\dot{\xi}_\theta(x) = -\tfrac{1}{2}\frac{\dot{q}(x - \theta)}{\sqrt{q(x - \theta)}}\{q(x - \theta) > 0\} \qquad \text{and} \qquad 4\mathfrak{m}\dot{\xi}_\theta^2 \equiv \mathbb{I}_q.$$

The family $\mathcal{Q} := \{Q_\theta : \theta \in \mathbb{R}\}$ is Hellinger differentiable at $\theta = 0$. In fact, the same argument works at every $\theta$; the family is everywhere Hellinger differentiable, with Hellinger derivative $\dot{\xi}_\theta$ at $\theta$.

It is traditional to call $\mathbb{I}_q$ the Fisher information for $q$, even though it would be more more precise to call it the Fisher information for the shift family
☐    generated by $q$.

## 3.   Information inequality

The information inequality for the model $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ bounds the variance of a statistic $T(x)$ from below by an expression involving the expected value of the statistic and the Fisher information: under suitable regularity conditions,

$$\text{var}_\theta(T) \geq \dot{\gamma}_\theta'\mathbb{I}_\theta^{-1}\dot{\gamma}_\theta \qquad \text{where } \gamma_\theta := P_\theta T(x) \text{ and } \dot{\gamma}_\theta := \frac{d}{d\theta}\gamma_\theta.$$

The classical proof of the inequality imposes assumptions that derivatives can be passed inside integral signs, typically justified by more primitive assumptions involving pointwise differentiability of densities and domination assumptions about their derivatives.

By contrast, the proof of the information inequality based on an assumption of Hellinger differentiability replaces the classical requirements by simple properties of $\mathcal{L}^2(\lambda)$ norms and inner products. The gain in elegance and economy of assumptions illustrates the typical benefits of working with Hellinger differentiability. The main technical ideas are captured by the following Lemma. Once again, with no loss of generality I consider only behavior at $\theta = 0$.

> REMARK.    The measure $P_\theta$ might itself be a product measure, representing the joint distribution of a sample of independent observations from some distribution $\mu_\theta$. As shown by Problem [5], Hellinger differentiability of $\theta \mapsto \mu_\theta$ at $\theta = 0$ would then imply Hellinger differentiability of $\theta \mapsto P_\theta$ at $\theta = 0$. We could substitute an explicit product measure for $P_\theta$ in the next Lemma, but there would be no advantage to doing so.

<12>  **Lemma.**  *Suppose a dominated family $\mathcal{P}$ has Hellinger derivative $\dot{\xi}_0$ at 0 and that $\sup_{\theta \in U} P_\theta T(x)^2 < \infty$, for some neighborhood $U$ of 0. Then the function $\theta \mapsto \gamma_\theta := P_\theta^x T(x)$ has derivative $\dot{\gamma}_0 = 2\lambda(\xi_0 \dot{\xi}_0 T)$ at 0.*

> REMARK.    Notice that $P_\theta T$ is well defined throughout $U$, because of the bound on the second moment. Also $\left(\lambda|\xi_0 \dot{\xi}_0 T|\right)^2 \le \left(\lambda \xi_0^2 T^2\right)\left(\lambda|\dot{\xi}_0|^2\right) < \infty$.

*Proof.*    Write $C^2$ for $\sup_{\theta \in U} P_\theta T(x)^2$, so that $\|\xi_\theta T\|_2 \le C$ for each $\theta$ in $U$. For simplicity, I consider only the one-dimensional case. The proof for $\mathbb{R}^k$ differs only notationally.

The proof is easy if $T$ is bounded by a finite constant $K$. By expanding $\xi_\theta^2$ using <2> then invoking the Cauchy-Schwarz inequality, we get

$$
\begin{aligned}
&|\gamma_\theta - \gamma_0 - 2\theta\lambda(\xi_0 \dot{\xi}_0 T)| \\
&= |\lambda(R_\theta T)| \qquad \text{where } R_\theta(x) := \xi_\theta^2 - \xi_0^2 - 2\theta\xi_0\dot{\xi}_0 \\
&= \lambda \left|\theta^2 \dot{\xi}_0^2 + r_\theta^2 + 2\xi_0 r_\theta + 2\theta\dot{\xi}_0 r_\theta\right| |T| \\
&\le K\left(\theta^2\|\dot{\xi}_0\|_2^2 + \|r_\theta\|_2^2 + 2|\theta|\|\dot{\xi}_0\|_2 \|r_\theta\|_2\right) + 2\|r_\theta\|_2 \|\xi_0 T\|_2 \\
&= o(|\theta|).
\end{aligned}
$$

<13>

The result for bounded $T$ suggests that we break the general $T$ into two pieces, $T\{|T| \le K\} + T\{|T| > K\}$. We could then argue as above for the contribution from $\{|T| \le K\}$. For the other contribution we would need to invoke the integrability of $\dot{\xi}_0^2$ to show that

<14>
$$
\|\dot{\xi}_0\{|T| > K\}\|_2^2 = \lambda\left(\dot{\xi}_0^2\{|T| > K\}\right) \to \qquad \text{as } K \to \infty.
$$

The argument would require some delicacy to ensure that whenever we invoked Cauchy-Schwarz we split into a product of two terms terms that are controlled by <14>, the Hellinger differentiability, or the boundedness of the second moment.

It is slightly more elegant to dispose of some contributions from $r_\theta$ before splitting $T$. Note that

$$
\begin{aligned}
R_\theta(x) &= (\xi_\theta - \xi_0)(\xi_\theta + \xi_0) - 2\theta\xi_0\dot{\xi}_0 \\
&= \theta\dot{\xi}_0(\xi_\theta - \xi_0) + r_\theta(\xi_\theta + \xi_0).
\end{aligned}
$$

For the $r_\theta$ contribution to $|\lambda(R_\theta T)|$ we have

$$
\left|\lambda\left(r_\theta(\xi_\theta + \xi_0)T\right)\right| \le \|r_\theta\|_2 \left(\|\xi_\theta T\|_2 + \|\xi_0 T\|_2\right) = o(|\theta|).
$$

For the other contribution we already have a $|\theta|$. We only need another $o(1)$ factor. Split according to whether $|T| \le K$ or not then bound by

$$
\begin{aligned}
&\lambda|\dot{\xi}_0\{|T| > K\}T(\xi_\theta - \xi_0)| + \lambda|\dot{\xi}_0\{|T| \le K\}T(\xi_\theta - \xi_0)| \\
&\le \|\dot{\xi}_0\{|T| > K\}\|_2 \left(\|\xi_\theta T\|_2 + \|\xi_0 T\|_2\right) + K\|\dot{\xi}_0\|_2\|\xi_\theta - \xi_0\|_2
\end{aligned}
$$

☐   Choose $K$ to make the first term suitably small then let $\theta$ tend to 0.

<15>   **Corollary.**   *In addition to the conditions of the Lemma, suppose $\mathbb{I}_0 := 4\lambda(\dot{\xi}_0\dot{\xi}_0')$ is nonsingular. Then $\mathrm{var}_0 T \geq \dot{\gamma}_0 \mathbb{I}_0^{-1} \dot{\gamma}_0$.*

*Proof.*   The special case of the Lemma where $T \equiv 1$ gives $\lambda(\xi_0\dot{\xi}_0) = 0$. Let $\alpha$ be a fixed vector in $\mathbb{R}^k$. Deduce that

$$
\begin{aligned}
(\alpha'\dot{\gamma}_0)^2 &= 4\left(\lambda\left(\alpha'\dot{\xi}_0\right)(T - \gamma_0)\xi_0\right)^2 \\
&\leq 4\alpha'\lambda(\dot{\xi}_0\dot{\xi}_0')\alpha\,\lambda\left(\xi_0^2(T - \gamma_0)^2\right) \qquad \text{by Cauchy-Schwarz} \\
&= \left(\alpha'\mathbb{I}_0\alpha\right)\mathbb{P}_0\left(T - \gamma_0\right)^2
\end{aligned}
$$

☐   Choose $\alpha := \mathbb{I}_0^{-1}\dot{\gamma}_0$ to complete the proof.

Variations on the information inequality lead to other useful lower bounds for variances and mean squared errors of statistics.

<16>   **Example.**   [The Van Trees inequality] Let $\Theta$ be an open subset of $\mathbb{R}$. Suppose $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ is dominated by $\lambda$, with Hellinger derivative $\dot{\xi}_\theta$ existing at each $\theta$ and information function $\mathbb{I}(\theta) = 4\lambda(\dot{\xi}_\theta^2)$. Suppose also that $q$ is an absolutely continuous probability density (with respect to Lebesgue measure $\mathfrak{m}$ on $\mathcal{B}(\mathbb{R})$), which satisfies the assumptions of Example <11> and vanishes outside an interval $[a, b] \subseteq \Theta$.

Let $T(x)$ be an estimator for the unknown parameter $\theta$, with expected value $\gamma(\theta) = P_\theta T = \lambda f_\theta(x)T(x)$. Suppose $\sup_{\theta \in K} P_\theta T(x)^2 < \infty$ for each compact subset $K$ of $\Theta$.

Let $\mathcal{Q}$ be the shift family generated by $q$, with densities $\eta_\alpha(\theta)^2 := q(\theta - \alpha)$ with respect to $\mathfrak{m}$. Remember that $\mathcal{Q}$ has Hellinger derivative

$$
\dot{\eta}_\alpha(\theta) = -\tfrac{1}{2}\frac{\dot{q}(\theta - \alpha)}{\sqrt{q(\theta - \alpha)}}\{q(\theta - \alpha) > 0\} \qquad \text{and} \qquad 4\mathfrak{m}\dot{\eta}_\alpha^2 \equiv \mathbb{I}_q.
$$

Consider the one-parameter family of probabilities $\mathcal{G} := \{\mathbb{G}_\alpha : |\alpha| < \delta\}$, for some small, positive $\delta$, defined by densities

$$
g_\alpha(x, \theta) := q(\theta - \alpha)f_{\theta-\alpha}(x) \qquad \text{with respect to } \lambda \otimes \mathfrak{m}.
$$

Under slightly awkward assumptions, the family $\mathcal{G}$ has Hellinger derivative

$$
\Delta(x, \theta) = \dot{\eta}_0(\theta)\xi_\theta(x) - \eta_0(\theta)\dot{\xi}_\theta(x) \qquad \text{at } \alpha = 0.
$$

Lemma <12> shows that the function $g(\alpha) := \mathbb{G}_\alpha\left(T(x) - \theta\right)$ has derivative

<17>   $$2\lambda \otimes \mathfrak{m}\left(\eta_0(\theta)\xi_\theta(x)\Delta(x, \theta)\left(T(x) - \theta\right)\right)$$

with corresponding information inequality

$$
\mathfrak{m}^\theta q(\theta)P_\theta\left(T(x) - \theta\right)^2 \geq \frac{1}{\mathbb{I}_q + \mathfrak{m}^\theta q(\theta)\mathbb{I}(\theta)},
$$

This result is also known as the ***van Trees inequality***. It has many statistical applications. See Gill & Levit (1995) for details.

In fact, because of the separation of variables in the function $T(x) - \theta$, it is not essential that $\mathcal{G}$ be Hellinger differentiable We can establish the value of the expression in <17> by two separate calculations. By Fubini,

$$
\begin{aligned}
&2\lambda \otimes \mathfrak{m}\left(\eta_0(\theta)\xi_\theta(x)\Delta(x, \theta)T(x)\right) \\
&= \mathfrak{m}^\theta\lambda^x\left(2\eta_0(\theta)\dot{\eta}_0(\theta)f_\theta(x)T(x) - q(\theta)2\xi_\theta(x)\dot{\xi}_\theta(x)T(x)\right) \\
&= -\mathfrak{m}^\theta\left(\dot{q}_0(\theta)\{q(\theta) > 0\}\gamma(\theta)\right) - \mathfrak{m}^\theta\left(q(\theta)\dot{\gamma}(\theta)\right)
\end{aligned}
$$

Remember that $\dot{q} = 0$ almost everywhere on $\{q = 0\}$ because $q \geq 0$. we can discard the indicator function from the first expression, leaving

$$
-\mathfrak{m}\left(\dot{q}(\theta)\gamma(\theta) + q(\theta)\dot{\gamma}(\theta)\right)
$$

The expression to be integrated is the almost sure derivative of the absolutely continuous function $q(\theta)\gamma(\theta)$, which vanishes outside the interval $[a, b]$. The contribution from $T(x)$ to $<17>$ is zero.

Similarly

$$2\lambda \otimes \mathfrak{m}\left(\eta_0(\theta)\xi_\theta(x)\Delta(x, \theta)\theta\right)$$
$$= \mathfrak{m}^\theta\left(\dot{q}(\theta)\theta\right) - 0 \qquad \text{because } \lambda\xi_\theta(x)\dot{\xi}_\theta(x) = 0$$
$$= -\mathfrak{m}^\theta(q(\theta)1) \qquad \text{by absolute continuity of } q(\theta)\theta \text{ on } [a, b]$$
$$= -1$$

Thus

$$2\lambda \otimes \mathfrak{m}\left(\eta_0(\theta)\xi_\theta(x)\Delta(x, \theta)\left(T(x) - \theta\right)\right) = 1.$$

By the Cauchy-Schwarz inequality,

$$4\lambda \otimes \mathfrak{m}\left(q(\theta)f_\theta(x)\left(T(x) - \theta\right)^2\right)\lambda \otimes \mathfrak{m}\left(\Delta(x, \theta)^2\right) \geq 1.$$

Finally, note that

$$4\lambda \otimes \mathfrak{m}\left(\Delta(x, \theta)^2\right)$$
$$= 4\lambda \otimes \mathfrak{m}\left(\dot{\eta}_0(\theta)^2 f_\theta(x) - 2\dot{\eta}_0(\theta)\eta_0(\theta)\dot{\xi}_\theta(x)\xi_\theta(x) + q(\theta)\dot{\xi}_\theta(x)^2\right)$$
$$= \mathbb{I}_q - 0 + \mathfrak{m}^\theta q(\theta)\mathbb{I}(\theta),$$

$\square$  which completes the direct proof of the van Trees inequality.

## 4.  Possible trouble at the boundary

At the end of Section 1, in order to postpone a postone difficulty with division by zero, I introduced temporary notation to distinguish between between two candidates for the title of information function under a Hellinger differentiability assumption. For simplicity of notation, consider the case of $\theta$ equal to zero. The two candidates are then $\mathbb{I}_0 := 4\lambda(\dot{\xi}_0\dot{\xi}_0')$ and $\mathbb{I}_0^\circ := 4\lambda(\dot{\xi}_0\dot{\xi}_0'\{\xi_0 > 0\})$. Their difference is a nonnegative definite matrix,

$$B_0 := \mathbb{I}_0 - \mathbb{I}_0^\circ = 4\lambda\left(\dot{\xi}_0\dot{\xi}_0'\{\xi_0 = 0\}\right).$$

If $0$ is an interior point of $\Theta$, nonnegativity of $\xi_\theta$ in a neighborhood of $0$ forces $\dot{\xi}_0(x) = 0$ for $\lambda$-almost all $x$ in $\{\xi_0 = 0\}$. (The argument from the start of Section 2 generalizes easily to higher dimensions.) Only if $0$ is a boundary point of $\Theta$ might $B_0$ be nonzero.

For $x$ in the set $\{\xi_0 = 0\}$ we have

$$\xi_\theta(x) = 0 + \theta'\dot{\xi}_0(x) + r_\theta(x),$$

which implies

$$P_\theta\{\xi_0 = 0\} = \lambda\left(\theta'\dot{\xi}_0 + r_\theta\right)^2\{\xi_0 = 0\} = \theta'B_0\theta + o(|\theta|^2)$$

The quantity on the left-hand side equals $P_\theta^\perp(\mathfrak{X})$, the total mass of the part of $P_\theta$ that is singular with respect to $P_0$.

<18>    **Example.**    Define $\mathcal{P} := \{P_\theta : 0 \le \theta \le 1\}$ via the densities

$$f_\theta(x) = \xi_\theta(x)^2 := (1 - \theta^2)(1 - |x|)^+ + \theta^2(1 - |x - 2|)^+$$

with respect to Lebesgue measure $\lambda$ on $[-1, 3]$. The densities $f_0$ and $f_1$ have disjoint support, and $\xi_\theta = (1 - \theta^2)^{1/2}\xi_0 + \theta\xi_1$. By direct calculation

$$\lambda\,|\xi_\theta(x) - \xi_0(x) - \theta\xi_1(x)|^2 = \left(\sqrt{1 - \theta^2} - 1\right)^2 = O(\theta^4).$$

Thus $\mathcal{P}$ is Hellinger differentiable at $\theta = 0$ with $\mathcal{L}^2(\lambda)$ derivative $\dot\xi_0 := \sqrt{f_1}$, but $P_\theta\{f_0 = 0\} = \theta^2$. The random variable $Z_n$ is equal to zero a.e. $[P_0^n]$, and
□    $\mathbb{I}_0 = 1$, and $\mathbb{I}_0 = 0$.

In general, if $B_0$ is nonzero there might be sequences $\{\theta_n\}$ in $\Theta$ approaching 0 at a $1/\sqrt{n}$ rate for which $\theta_n/|\theta_n| \to \delta$ with $c := \delta'B_0\delta \neq 0$. For such $\theta_n$ we would have $nP_{\theta_n}^\perp(\mathcal{X}) = nP_{\theta_n}\{\xi_0 = 0\} \to c$. As you will see in Section 8, this possibility will cause awkward behavior for the the likelihood ratio $dP_{\theta_n}^n/dP_0^n$, an awkwardness essentially due to a failure of contiguity—see Problem [11].

## 5.    An intrinsic characterization of Hellinger differentiability

For the definition of Hellinger differentiability, the choice of dominating measure $\lambda$ for the family of probability measures $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ is somewhat arbitrary. In fact, there is really no need for a single dominating $\lambda$, provided we guard against bad behavior by $P_\theta^\perp$, the part of $P_\theta$ that is singular with respect to $P_0$. As you saw in Section 4, the assumption $P_\theta^\perp\mathcal{X} = o(|\theta|^2)$ is needed to ensure some asymptotic difficulties related to failure of contiguity. We lose little by building the assumption into the definition. Following Le Cam (1986, Section 17.3) and Le Cam & Yang (2000, Section 7.2), I will call the slightly stronger property **_differentiability in quadratic mean (DQM)_**, to stress that the definition requires a little more than Hellinger differentiability.

The definition makes no assumption that the family of probability measures $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ is dominated. Instead it is expressed directly in terms of the Lebesgue decomposition of $P_\theta$ with respect to $P_{\theta_0}$, for a fixed $\theta_0$ in $\Theta$. As before, I will assume $\theta_0 = 0$ to simplify notation. Remember that $P_\theta = \tilde{P}_\theta + P_\theta^\perp$, where the absolutely continuous part $\tilde{P}_\theta$ has a density $p_\theta$ with respect to $P_0$ and the singular part $P_\theta^\perp$ concentrates on a $P_0$-negligible set $\mathcal{N}_\theta$,

$$P_\theta g = \tilde{P}_\theta\left(g(x)p_\theta(x)\{x \in \mathcal{N}_\theta^c\}\right) + P_\theta^\perp\left(g(x)\{x \in \mathcal{N}_\theta\}\right),$$

at least for nonnegative measurable functions $g$ on $\mathcal{X}$.

<19>    **Definition.**    *Say that $\mathcal{P}$ is differentiable in quadratic mean (DQM) with score function $\Delta$ at 0 if*

(i)  $P_\theta^\perp(\mathcal{X}) = o(|\theta|^2)$ *as* $|\theta| \to 0$,
(ii)  $\Delta$ *is a vector of $\mathcal{L}^2(P_0)$ functions for which*

$$\sqrt{p_\theta(x)} = 1 + \tfrac{1}{2}\theta'\Delta(x) + r_\theta(x) \qquad \text{with } P_0\left(r_\theta^2\right) = o(|\theta|^2) \text{ near 0.}$$

REMARK.    Some authors (for example, Bickel et al. 1993, page 457) use the term DQM as a synonym for differentiability in $L^2$ norm. The factor of $1/2$ simplifies some calculations, by making the vector $\Delta$ correspond to the score function at 0.

When $\mathcal{P}$ is dominated by a sigma-finite measure, the definition agrees with the definition of Hellinger differentiability under the assumption (i), which is needed for contiguity of product measures.

<20>  **Theorem.** *Suppose $\mathcal{P}$ is dominated by a sigma-finite measure $\lambda$, with corresponding densities $f_\theta(x)$.*

(i) *Suppose $P_\theta\{f_0 = 0\} = o(|\theta|^2)$ and, for some vector $\dot\xi$ of functions in $\mathcal{L}^2(\lambda)$,*

$$\sqrt{f_\theta(x)} = \sqrt{f_0} + \theta'\dot\xi(x) + R_\theta(x) \qquad \text{where } \lambda\left(R_\theta^2\right) = o(|\theta|^2) \text{ with } \theta = 0.$$

*Then $\mathcal{P}$ satisfies the DQM condition at 0, with $\Delta := 2\{f_0 > 0\}\dot\xi/\sqrt{f_0}$ and $r_\theta := \{f_0 > 0\}R_\theta/\sqrt{f_0}$.*

(ii) *If $\mathcal{P}$ satisfies the DQM condition at 0 then it is also Hellinger differentiable at 0, with $\mathcal{L}^2(\lambda)$ derivative $\dot\xi := \frac{1}{2}\Delta\sqrt{f_0}$.*

*Proof.* For the Lebesgue decomposition we can take $p_\theta := \{f_0 > 0\}f_\theta/f_0$ and $\mathcal{N}_\theta := \{f_0 = 0\}$. Thus $P_\theta^\perp\mathcal{X} = \lambda f_\theta\{f_0 = 0\}$.

If $\mathcal{P}$ is Hellinger differentiability, as in (i), then

$$P_0\left|\sqrt{p_\theta} - 1 - \tfrac{1}{2}\theta'\Delta\right|^2 = \lambda f_0\left|\{f_0 > 0\}\sqrt{f_\theta/f_0} - 1 - \tfrac{1}{2}\theta'\dot\xi\{f_0 > 0\}/\sqrt{f_0}\right|^2$$

$$= \lambda\left(\{f_0 > 0\}\left|\sqrt{f_\theta} - \sqrt{f_0} - \theta'\dot\xi\right|^2\right) = o(|\theta|^2).$$

Conversely, if $\mathcal{P}$ satisfies DQM then

$$\lambda\left|\sqrt{f_\theta} - \sqrt{f_0} - \theta'\dot\xi\right|^2 = \lambda\{f_0 = 0\}\left(\sqrt{f_\theta} - 0\right)^2$$

$$+ \lambda\{f_0 > 0\}\left|\sqrt{f_0 p_\theta} - \sqrt{f_0} - \tfrac{1}{2}\theta'\Delta\sqrt{f_0}\right|^2$$

$$= o(|\theta|^2) + P_0\left|\sqrt{p_\theta} - 1 - \tfrac{1}{2}\theta'\Delta\right|^2 = o(|\theta|^2).$$

$\square$

The definition of DQM has some advantages over the definition of Hellinger differentiability, even beyond the elimination of the dominating measure $\lambda$. For $\theta$ near zero, $p_\theta \approx 1$, a simplification that has subtle consequences, as illustrated by the next Section.

## 6.  Preservation of DQM under measurable maps

Suppose $T$ is a measurable map from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{T}, \mathcal{B})$. The distribution of $T$ under $P_\theta$ is a probability measure on $\mathcal{B}$, the image measure $Q_\theta := T P_\theta$ on $\mathcal{B}$, defined by $Q_\theta g := P_\theta^x g(Tx)$ for each $g$ in $\mathcal{M}^+(\mathcal{T}, \mathcal{B})$.

For each $h$ in $\mathcal{L}^1(P_0)$ we can define a (signed) measure $\nu_h$ on $\mathcal{B}$ by

$$\nu_h(g) := P_0^x\left(h(x)g(Tx)\right) \qquad \text{for each } g \text{ in } \mathcal{M}^+(\mathcal{T}, \mathcal{B}).$$

The measure $\nu_h$ is absolutely continuous with respect to $Q_0$. The Kolmogorov conditional expectation $P_0(h \mid T = t)$ is just the density $d\nu_h/dQ_0$. (See UGMTP §5.6.) I will also denote it by $\pi_t(h)$. Thus, at least for each bounded, $\mathcal{B}$-measurable real function $g$ on $\mathcal{T}$,

<21>  $$P_0^x\left(h(x)g(Tx)\right) = Q_0^t\left(g(t)\pi_t(h)\right)$$

Now suppose $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ is DQM at $\theta = 0$ with score function $\Delta$. For the moment, suppose each $P_\theta$ is dominated by $P_0$ with density $p_\theta(x)$. Then each member of $\mathcal{Q} := \{Q_\theta : \theta \in \Theta\}$ is dominated by $Q_0$, with density $dQ_\theta/dQ_0 = P_0(p_\theta \mid T = t) = \pi_t(p_\theta)$. Under DQM,

$$\sqrt{\pi_t(p_\theta)} = \left(\pi_t\left(1 + \tfrac{1}{2}\theta'\Delta + r_\theta\right)^2\right)^{1/2}$$

$$= \left(1 + \theta'\pi_t\Delta + \ldots\right)^{1/2} = 1 + \tfrac{1}{2}\theta'\pi_t\Delta + \ldots$$

If all the omitted terms can be ignored, in an $\mathcal{L}^2(Q_0)$ sense, then $\mathcal{Q}$ would be Hellinger differentiable at 0, with $\mathcal{L}^2(Q_0)$-derivative $\pi_t(\Delta)$. The next Theorem makes this heuristic argument rigorous, even without a domination assumption on $\mathcal{P}$.

<22>   **Theorem.**   *Suppose $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ is DQM with score function $\Delta$ at 0. Suppose $T$ is a measurable map from $(\mathcal{X}, \mathcal{A})$ into $(\mathcal{T}, \mathcal{B})$. Then $\{T P_\theta : \theta \in \Theta\}$ is DQM at 0, with score function $P_0(\Delta \mid T = t)$.*

*Proof.*   To simplify notation, I will assume $\Theta$ is one-dimensional. No extra conceptual difficulties arise in higher dimensions.

Define $Q_\theta := T P_\theta$ and $\tilde{Q}_\theta := T \tilde{P}_\theta$. In fact $\tilde{Q}_\theta$ might not be the part of $Q_\theta$ that is singular with respect to $Q_0$, because $T P_\theta^\perp$ might contain a component dominated by $Q_0$. As you will see at the end of the proof, there is little harm in ignoring the possible contribution from $T P_\theta^\perp$ for the moment.

Write $\xi_\theta$ for $\sqrt{p_\theta}$, where $p_\theta = d\tilde{P}_\theta / dP_0$. By definition of DQM,

$$\xi_\theta(x) = 1 + \tfrac{1}{2}\theta\Delta(x) + r_\theta(x) \qquad \text{with } P_0 r_\theta^2 = o(\theta^2).$$

Use a bar to denote "averaging" with respect to $\pi_t$:

$$\bar{\Delta}(t) := \pi_t(\Delta), \qquad \bar{r}_\theta(t) := \pi_t(r_\theta), \qquad \bar{\xi}_\theta(t) := \pi_t(\xi_\theta) = 1 + \tfrac{1}{2}\theta\bar{\Delta}(t) + \bar{r}_\theta(t).$$

Notice that $\bar{\xi}_\theta$ converges in $Q_0$ probability to 1 as $\theta \to 0$ because

$$Q_0\bar{\Delta}^2 \le Q_0\pi_t\Delta^2 = P_0\Delta^2 < \infty$$
$$Q_0\bar{r}_\theta^2 \le Q_0\pi_t r_\theta^2 = P_0 r_\theta^2 = o(\theta^2).$$

Also, the function $\xi_\theta$ has a small conditional variance:

$$\sigma_\theta^2(t) := \pi_t\left(\xi_\theta - \bar{\xi}_\theta\right)^2$$
$$= \pi_t\left(\tfrac{1}{2}\theta(\Delta - \bar{\Delta}) + r_\theta - \bar{r}_\theta\right)^2$$
$$\le 2(\tfrac{1}{2}\theta)^2\pi_t(\Delta - \bar{\Delta})^2 + 2\pi_t\left(r_\theta - \bar{r}_\theta\right)^2$$

<23>
$$\le \theta^2 J(t) + \epsilon_\theta(t) \qquad \text{where } J(t) := \tfrac{1}{2}\pi_t\Delta^2 \text{ and } \epsilon_\theta(t) = 2\pi_t r_\theta(t)^2$$

so that

$$Q_0\sigma_\theta^2(t) \le \theta^2 Q_0 J + Q_0\epsilon_\theta = \tfrac{1}{2}\theta^2 P_0\Delta^2 + 2P_0 r_\theta^2 = O(\theta^2) + o(\theta^2).$$

REMARK.      The cancellation of the leading constants when $\bar{\xi}_\theta$ is subtracted from $\xi_\theta$ seems to be vital to the proof. For general Hellinger differentiability, the cancellation does not occur.

The density of $\tilde{Q}_\theta$ with respect to $Q_0$ equals

$$\eta_\theta^2(t) := \pi_t(\xi_\theta^2) = \sigma_\theta^2(t) + \bar{\xi}_\theta(t)^2.$$

Consequently, $\eta_\theta(t) \ge \bar{\xi}_\theta(t)$ or, equivalently,

$$\eta_\theta(t) = \bar{\xi}_\theta(t) + \delta_\theta(t) \qquad \text{for some } \delta_\theta \ge 0$$
$$= 1 + \tfrac{1}{2}\theta\bar{\Delta}(t) + \bar{r}_\theta(t) + \delta_\theta(t)$$

To establish DQM for $\mathcal{Q}$ we need to show that

$$Q_0\left(\eta_\theta - 1 - \tfrac{1}{2}\theta\bar{\Delta}\right)^2 = o(\theta^2).$$

It is enough, therefore, to show that both $Q_0\bar{r}_\theta^2$ and $Q_0\delta_\theta^2$ are of order $o(\theta^2)$.

The $\bar{r}_\theta$ is easily handled via a conditional Jensen inequality:

$$Q_0\bar{r}_\theta^2 = Q_0\left(\pi_t r_\theta\right)^2 \le Q_0\pi_t r_\theta^2 = P_0 r_\theta^2 = o(\theta^2).$$

The argument for $\delta_\theta$ is a little more delicate. From the bound <23> and the equality $\left(\bar{\xi}_\theta + \delta_\theta\right)^2 = \eta_\theta^2 = \sigma_\theta^2 + \bar{\xi}_\theta^2$ we have

<24>
$$2\bar{\xi}_\theta(t)\delta_\theta(t) + \delta_\theta^2(t) = \sigma_\theta^2 \le \theta^2 J(t) + \epsilon_\theta(t).$$

For a fixed small $\beta > 0$, define

$$A_\theta = \{t : \bar{\bar{\xi}}_\theta(t) \geq 1/2, \ \sigma_\theta^2(t) \leq \beta\}.$$

Notice that $Q_0 A_\theta^c \to 0$ as $\theta \to 0$. For $t$ in $A_\theta$ we have $\delta_\theta(t)^2 \leq \sigma_\theta^4(t) \leq \beta\sigma_\theta^2(t)$; and for $t$ in $A_\theta^c$ we have $\delta_\theta^2(t) \leq \sigma_\theta^2$. Thus

$$\begin{aligned}
Q_0\delta_\theta^2 &\leq \beta Q_0\sigma_\theta^2 A_\theta + Q_0\sigma_\theta^2 A_\theta^c \\
&\leq \beta\left(\theta^2 Q_0 J + Q_0\epsilon_\theta\right) + \theta^2 Q_0(J A_\theta^c) + Q_0\epsilon_\theta \\
&\leq \beta\theta^2 Q_0 J + o(\theta^2) + \theta^2 o(1) + o(\theta^2)
\end{aligned}$$

As $\beta$ could be chosen arbitrarily small, it follows that $Q_0\delta_\theta^2 = o(\theta^2)$.

Finally, let me take care of any part of $T P_\theta^\perp$ that might have been dominated by $Q_0$. The density for the part of $Q_\theta$ with respect to $Q_0$ might actually equal $\eta_\theta^2 + s_\theta$, where $s_\theta \geq 0$ and $Q_0 s_\theta \leq (T P_\theta^\perp)(\mathcal{T}) = P_\theta^\perp(\mathcal{X}) = o(\theta^2)$. That is, we need to replace $\sigma_\theta^2$ by $\sigma_\theta^2 + s_\theta$, which adds only another $o(\theta^2)$ terms to the bounds in the previous paragraph. The modification has an asymptotically negligible effect on the argument. The family $\{Q_\theta : \theta \in \Theta\}$ inherits DQM from $\square$  the family $\{\tilde{Q}_\theta : \theta \in \Theta\}$.

---

> Discuss consequences for loss of information under measurable maps.
> Sufficiency. cf. Kagan & Shepp (2005)

---

## 7.   Differentiability of unit vectors

Suppose $\tau$ is a map from $\mathbb{R}^k$ into some inner product space $\mathcal{H}$ (such as $\mathcal{L}^2(\lambda)$). Suppose also that $\tau$ is differentiable (in norm) at $\theta_0$,

$$\tau_\theta = \tau_{\theta_0} + (\theta - \theta_0)'\dot{\tau}_{\theta_0} + r_\theta \qquad \text{with } \|r_\theta\| = o(|\theta - \theta_0|) \text{ near } \theta_0.$$

For simplicity of notation, suppose $\theta_0 = 0$.

The Cauchy-Schwarz inequality gives $|\langle\tau_0, r_\theta\rangle| \leq \|\tau_0\|\|r_\theta\| = o(|\theta|)$. It would usually be a blunder to assume naively that the bound must therefore be of order $O(|\theta|^2)$; typically, higher-order differentiability assumptions are needed to derive approximations with smaller errors. However, if $\|\tau_\theta\|$ is constant—that is, if $\tau_\theta$ is constrained to take values lying on the surface of a sphere—then the naive assumption turns out to be no blunder. Indeed, in that case, it is easy to show that in general $\langle\tau_0, r_\theta\rangle$ equals a quadratic in $\theta$ plus an error of order $o(|\theta|^2)$. The sequential form of the assertion will be more convenient for the calculations in Section 8.

<25>   **Lemma.**   *Let $\{\alpha_n\}$ be a sequence of constants tending to zero. Let $\tau_0, \tau_1, \ldots$ be elements of norm one for which $\tau_n = \tau_0 + \alpha_n W + \rho_n$, with $W$ a fixed element of $\mathcal{H}$ and $\|\rho_n\| = o(\alpha_n)$. Then $\langle\tau_0, W\rangle = 0$ and $2\langle\tau_0, \rho_n\rangle = -\alpha_n^2\|W\|^2 + o(\alpha_n^2)$.*

*Proof.*   Because both $\tau_n$ and $\tau_0$ have unit length,

$$\begin{aligned}
0 = \|\tau_n\|^2 - \|\tau_0\|^2 = \ &2\alpha_n\langle\tau_0, W\rangle && \text{order } O(\alpha_n) \\
&+ 2\langle\tau_0, \rho_n\rangle && \text{order } o(\alpha_n) \\
&+ \alpha_n^2\|W\|^2 && \text{order } O(\alpha_n^2) \\
&+ 2\alpha_n\langle W, \rho_n\rangle + \|\rho_n\|^2 && \text{order } o(\alpha_n^2).
\end{aligned}$$

The $o(\alpha_n)$ and $o(\alpha_n^2)$ rates of convergence in the second and fourth lines come from the Cauchy-Schwarz inequality. The exact zero on the left-hand side of the equality exposes the leading $2\alpha_n\langle\tau_0, W\rangle$ as the only $O(\alpha_n)$ term on the

right-hand side. It must be of smaller order, $o(\alpha_n)$ like the other terms, which can happen only if $\langle \tau_0, W \rangle = 0$, leaving

$$0 = 2\langle \tau_0, \rho_n \rangle + \alpha_n^2 \|W\|^2 + o(\alpha_n^2),$$

□    as asserted.

> REMARK.    Without the fixed length property, the difference $\|\tau_n\|^2 - \|\tau_0\|^2$ might contain terms of order $\alpha_n$. The inner product $\langle \tau_0, \rho_n \rangle$, which inherits $o(\alpha_n)$ behaviour from $\|\rho_n\|$, might then not decrease at the $O(\alpha_n^2)$ rate.

<26>    **Corollary.**   *If $\mathcal{P}$ has a Hellinger derivative $\dot{\xi}_{\theta_0}$ at 0, and if 0 is an interior point of $\Theta$, then $\lambda\left(\xi_0 \dot{\xi}_0\right) = 0$ and $8\lambda\left(\xi_0 r_\theta\right) = -\theta' \mathbb{I}_0 \theta + o(|\theta|^2)$ near 0.*

*Proof.*   Start with the second assertion, in its equivalent form for sequences $\theta_n \to 0$. Write $\theta_n$ as $|\theta_n| u_n$, with $u_n$ a unit vector in $\mathbb{R}^k$. By a subsequencing argument, we may assume that $u_n \to u$, in which case,

$$\xi_{\theta_n} = \xi_0 + |\theta_n| u_n' \dot{\xi}_0 + r_{\theta_n} = \xi_0 + |\theta_n| u' \dot{\xi}_0 + \left(r_{\theta_n} + |\theta_n|(u_n - u)' \dot{\xi}_0\right).$$

Invoke the Lemma (with $W = u' \dot{\xi}_0$) to deduce that $u' \lambda\left(\xi_0 \dot{\xi}_0\right) = 0$ and

$$-4|\theta_n|^2 \lambda\left(u' \dot{\xi}_0\right)^2 + o(|\theta_n|^2) = 8\lambda\left(\xi_0\left(r_{\theta_n} + |\theta_n|(u_n - u)' \dot{\xi}_0\right)\right)$$
$$= 8\lambda\left(\xi_0 r_{\theta_n}\right) + 8|\theta_n|(u_n - u)' \lambda\left(\xi_0 \dot{\xi}_0\right).$$

Because 0 is an interior point, for every unit vector $u$ there are sequences $\theta_n \to 0$ through $\Theta$ for which $u = \theta_n / |\theta_n|$. Thus $u' \lambda\left(\xi_0 \dot{\xi}_0\right) = 0$ for every unit vector $u$, implying that $\lambda\left(\xi_0 \dot{\xi}_0\right) = 0$. The last displayed equation reduces the □    sequential analog of the asserted approximation.

> REMARK.    If 0 were not an interior point of the parameter space, there might not be enough directions $u$ along which $\theta_n \to 0$ through $\Theta$, and it might not follow that $\lambda(\xi_0 \dot{\xi}_0) = 0$. Roughly speaking, the set of such directions is called the ***contingent*** of $\Theta$ at $\theta_0$. If the contingent is 'rich enough', we do not need to assume that 0 is an interior point. See Le Cam & Yang (2000, Section 7.2) and Le Cam (1986, page 575) for further details.

## 8.    Quadratic approximation for log likelihood ratios

Suppose observations $\{x_i\}$ are drawn independently from the distribution $P_0$. Under the classical regularity conditions, the log of the likelihood ratio $dP_\theta^n / dP_0^n = \prod_{i \le n} f_\theta(x_i) / f_0(x_i)$ has a local quadratic approximation in $1/\sqrt{n}$ neighborhoods of 0, under $P_0^n$. (Remember that, in general, $d\mathbb{Q}/d\mathbb{P}$ denotes the density with respect to $\mathbb{P}$ of the part of $\mathbb{Q}$ that is absolutely continuous with respect to $\mathbb{P}$.) For example, the following result (for one dimension) was proved in Chapter 7.

<27>    **Theorem.**   *Let $\mathbb{P}_n := P_0^n$ and $\mathbb{Q}_n := P_{\theta_n}^n$, for $\theta_n := \delta_n / \sqrt{n}$ with $\{\delta_n\}$ bounded. Suppose the map $\theta \mapsto f_\theta$ is twice differentiable in a neighborhood $U$ of 0 with:*
> *(i) $\theta \mapsto \ddot{f}_\theta(x)$ is continuous at 0;*
> *(ii) there exists a $\lambda$-integrable function $M(x)$ with $\sup_{\theta \in U} |\ddot{f}_\theta(x)| \le M(x)$ a.e. $[P_0]$;*
> *(iii) $P_0^x\left(\dot{f}_\theta(x)/f_0(x)\right)^2 \to P_0^x\left(\dot{f}_0(x)/f_0(x)\right)^2 =: \mathbb{I}_0 < \infty$ as $\theta \to 0$;*
> *(iv) $P_\theta\{f_0 = 0\} = o(\theta^2)$ as $\theta \to 0$.*
>
> *Then $P_0 \dot{\ell}_0(x) = 0 = P_0\left(\ddot{f}(x)/f_0(x)\right)$ and, under $\{\mathbb{P}_n\}$,*

$$\frac{d\mathbb{Q}_n}{d\mathbb{P}_n} = \left(1 + o_p(1)\right) \exp\left(\delta_n Z_n - \tfrac{1}{2}\delta_n^2 \mathbb{I}_0\right),$$

*where* $Z_n := \sum_{i \leq n} \dot{\ell}_0(x_i)/\sqrt{n} \rightsquigarrow N(0, \mathbb{I}_0)$. *Consequently*, $\mathbb{Q}_n \lhd \mathbb{P}_n$.

The method of proof consisted of writing the likelihood ratio as

$$\prod_{i \leq n} \left(1 + \epsilon_n(x_i)\right) \qquad \text{where } \epsilon_n(x) := \{f_0(x) > 0\}\left(f_{\theta_n}(x) - f_0(x)\right)/f_0(x),$$

then showing that, under $\mathbb{P}_n$,

   (a) $\max_{i \leq n} |\epsilon_n(x_i)| = o_p(1)$,
   (b) $\sum_{i \leq n} \epsilon_n(x_i) = \delta_n Z_n + o_p(1)$,
   (c) $\sum_{i \leq n} \epsilon_n(x_i)^2 = \delta_n^2 \mathbb{I}_0 + o_p(1)$.

Result (a) plus the fact that $\sum_{i \leq n} \epsilon_n(x_i)^2 = O_p(1)$ implied that

<28>
$$\prod_{i \leq n} \left(1 + \epsilon_n(x_i)\right) = \left(1 + o_p(1)\right) \exp\left(\sum_{i \leq n} \epsilon_n(x_i) - \tfrac{1}{2}\epsilon_n(x_i)^2\right),$$

from which the final assertion followed.

Le Cam (1970) established a similar quadratic approximation under an assumption of Hellinger differentiability. The method of proof is very similar to the method just outlined, but with a few very subtle differences. Remember that $\mathbb{I}_0 := 4\lambda\left(\dot{\xi}_\theta \dot{\xi}_\theta'\right)$ and $\mathbb{I}_0^\circ := 4\lambda\left(\dot{\xi}_0 \dot{\xi}_0'\{\xi_\theta > 0\}\right)$.

<29>   **Theorem.**   *Suppose $\mathcal{P}$ is Hellinger differentiable at 0, with $\mathcal{L}^2(\lambda)$ deriva-tive $\dot{\xi}_0$. Let $\mathbb{P}_n := P_0^n$ and $\mathbb{Q}_n := P_{\theta_n}^n$, with $\theta_n := \delta_n/\sqrt{n}$ for a bounded sequence $\{\delta_n\}$. Then, under $\{\mathbb{P}_n\}$,*

$$\frac{d\mathbb{Q}_n}{d\mathbb{P}_n} = \left(1 + o_p(1)\right)\exp\left(\delta_n' Z_n - \tfrac{1}{4}\delta_n'(\mathbb{I}_0 + \mathbb{I}_0^\circ)\delta_n\right),$$

*where*

$$Z_n := 2n^{-1/2}\sum_{i \leq n}\{\xi_0(x_i) > 0\}\dot{\xi}_0(x_i)/\xi_0(x_i) \rightsquigarrow N(0, \mathbb{I}_0^\circ).$$

> REMARK.     It is traditional to absorb the $1+o_p(1)$ factor for the likelihood ratio into the exponent. One then has some awkwardness with the right-hand side of the approximation at samples for which the left-hand side is zero. The awkwardness occurs with positive $\mathbb{P}_n$ probability if $P_{\theta_0}\{f_{\theta_n} = 0\} > 0$.

*Proof.*   I will give the proof only for the one-dimensional case. The proof for the multi-dimensional case is analogous.

Write $\tau_n$ for $\xi_{\theta_n}$, and $\rho_n$ for $r_{\theta_n}$, and $L_n$ for $d\mathbb{Q}_n/d\mathbb{P}_n$. By Hellinger differentiability,

$$\tau_n(x) = \xi_0(x) + n^{-1/2}\delta_n\dot{\xi}_0(x) + \rho_n(x) \qquad \text{with } \lambda\rho_n^2 = o(\theta_n^2).$$

Define

<30>
$$\eta_n(x) := \{\xi_0(x) > 0\}\frac{\tau_n(x) - \xi_0(x)}{\xi_0(x)} = \frac{\delta_n}{\sqrt{n}}D(x) + R_n(x),$$

where

$$D(x) := \{\xi_0(x) > 0\}\dot{\xi}_0(x)/\xi_0(x) \qquad \text{and} \qquad R_n(x) := \{\xi_0(x) > 0\}\rho_n(x)/\xi_0(x).$$

The indicator functions have no effect within the set $A_n := \cap_{i \leq n}\{\xi_0(x_i) > 0\}$, which has $\mathbb{P}_n$-probability one, but they will protect against $0/0 \stackrel{?}{=} 1$ when converting from $P_0$- to $\lambda$-integrals. On the set $A_n$,

$$\sqrt{L_n} = \prod_{i \leq n} \tau_n(x_i)/\xi_0(x_i) = \prod_{i \leq n}\left(1 + \eta_n(x_i)\right).$$

For almost the same reason as in the proof of Theorem <27>, we need to show that

   (i) $\max_{i \leq n}|\eta_n(x_i)| = o_p(1)$,

(ii) $\sum_{i \le n} \eta_n(x_i) = \frac{1}{2}\delta_n Z_n - \frac{1}{8}\delta_n^2 \mathbb{I}_0 + o_p(1)$,

(iii) $\sum_{i \le n} \eta_n(x_i)^2 = \frac{1}{4}\delta_n^2 \mathbb{I}_0^\circ + o_p(1)$.

The analog of $<28>$, with $\eta_n$ replacing $\epsilon_n$, will then give

$$\sqrt{L_n} = \left(1 + o_p(1)\right) \exp\left(\sum_{i \le n} \eta_n(x_i) - \frac{1}{2}\sum_{i \le n} \eta_n(x_i)^2\right),$$

from which the assertion of the Theorem follows by squaring both sides.

REMARK.     Notice that (ii) differs significantly from its analog (b) for the proof of Theorem $<27>$, through the addition of a constant term. However, the difference is compensated by a halving of the corresponding constant in (iii), as compared with (c). The differences occur because, on the set $\{f_0(x) > 0\}$,

$$\epsilon_n(x) = \frac{\tau_n(x)^2 - \xi_0(x)^2}{\xi_0(x)^2} = \frac{\tau_n(x) - \xi_0(x)}{\xi_0(x)} \frac{2\xi_0(x) + \tau_n(x) - \xi_0(x)}{\xi_0(x)} = 2\eta_n(x) + \eta_n(x)^2.$$

Thus

$$\sum_{i \le n} \epsilon_n(x_i) = 2\sum_{i \le n} \eta_n(x_i) + \sum_{i \le n} \eta_n(x_i)^2 = \delta_n Z_n - \frac{1}{4}\delta_n^2 \mathbb{I}_0 + \frac{1}{4}\mathbb{I}_0^\circ + o_p(1).$$

As you will see in Section 4, the conditions of Theorem $<27>$ actually imply $\mathbb{I}_0 = \mathbb{I}_0^\circ$, a condition equivalent to the contiguity $\mathbb{Q}_n \lhd \mathbb{P}_n$.

Assertions (i), (ii), and (iii) will follow from $<30>$, via simple probability facts, including: if $Y_1, Y_2, \ldots$ are independent, identically distributed random variables with $\mathbb{P}|Y_1|^r < \infty$ for some constant $r \ge 1$ then $\max_{i \le n}|Y_i| = o_p(n^{1/r})$. (The proof appeared as a Problem to Chapter 7.)

First note that

$$P_0 D(x) = \lambda\left(\xi_0(x)^2 \frac{\dot\xi_0(x)}{\xi_0(x)}\{\xi_0(x) > 0\}\right) = \lambda\left(\xi_0 \dot\xi_0\right) = 0 \qquad \text{by Corollary } <26>,$$

$$P_0 D(x)^2 = \lambda\left(\xi_0(x)^2 \frac{\dot\xi_0(x)^2}{\xi_0(x)^2}\{\xi_0(x) > 0\}\right) = \lambda\left(\dot\xi_0^2\{\xi_0(x) > 0\}\right) = \frac{1}{4}\mathbb{I}_0^\circ,$$

$$8 P_0 R(x) = 8\lambda\left(\xi_0(x)\rho_n(x)\right) = -\delta_n^2 \mathbb{I}_0/n + o(1/n),$$

$$P_0 R(x)^2 \le \lambda\rho_n(x)^2 = o(1/n).$$

From the expressions involving $D$ we get

$$Z_n = 2\sum_{i \le n} D(x_i)/\sqrt{n} \rightsquigarrow N(0, \tfrac{1}{4}\mathbb{I}^\circ),$$

$$n^{-1}\sum_{i \le n} D(x_i)^2 = \tfrac{1}{4}\mathbb{I}^\circ + o_p(1),$$

$$\max_{i \le n}|D(x_i)| = o_p(n^{1/2}).$$

From the expressions involving $R_n$ we get

$$\mathbb{P}_n\left(\sum_{i \le n} R_n(x_i)\right) = -\delta_n^2 \mathbb{I}_0 + o(1),$$

$$\text{var}\left(\sum_{i \le n} R_n(x_i)\right) \le \sum_{i \le n} \mathbb{P}_n R_n(x_i)^2 \to 0,$$

which together imply that

$$\sum_{i \le n} R(x_i) = -\tfrac{1}{8}\delta_n^2 \mathbb{I}_0 + o_p(1),$$

$$\left(\max_{i \le n}|R_n(x_i)|\right)^2 \le \sum_{i \le n} R_n(x_i)^2 = o_p(1).$$

Assertions (i), (ii), and (iii) now follow easily.
For (i):

$$\max_{i \leq n} |\eta_n(x_i)| \leq |\delta_n| \max_{i \leq n} \frac{|D(x_i)|}{\sqrt{n}} + \max_{i \leq n} |R_n(x_i)| = o_p(1).$$

For (ii):

$$\sum_{i \leq n} \eta_n(x_i) = \tfrac{1}{2}\delta_n \sum_{i \leq n} \frac{D(x_i)}{\sqrt{n}} + \sum_{i \leq n} R_n(x_i) = \tfrac{1}{2}\delta_n Z_n - \tfrac{1}{8}\delta_n^2 \mathbb{I}_0 + o_p(1).$$

For (iii):

$$\left| \left( \sum_{i \leq n} \eta_n(x_i)^2 \right)^{1/2} - \left( \delta_n^2 \sum_{i \leq n} \frac{D(x_i)^2}{n} \right)^{1/2} \right| \leq \left( \sum_{i \leq n} R_n(x_i)^2 \right)^{1/2} = o_p(1),$$

implying that

$$\sum_{i \leq n} \eta_n(x_i)^2 = \delta_n^2 \sum_{i \leq n} \frac{D(x_i)^2}{n} + o_p(1) = \tfrac{1}{4}\delta_n^2 \mathbb{I}_0^\circ + o_p(1).$$

☐   The asserted quadratic approximation follows.

## 9.   Problems

> Problems not yet checked.

[1]   Let $\lambda$ denote Lebesgue measure on $\mathcal{B}[0,1]$. Let $\{B_k : k \in \mathbb{N}\}$ be a sequence of sets with $\lambda B_{2k} = 1/k = \lambda B_{2k+1}$ and

$$\lambda\{x : x \in B_{2k} \cup B_{2k+1} \text{ for infinitely many } k \} = 1.$$

Define $f_0(x) := 1$ for all $x$. Define $a_k := k^{-1/3}$. For $a_{k+1} \leq \theta < a_k$ define
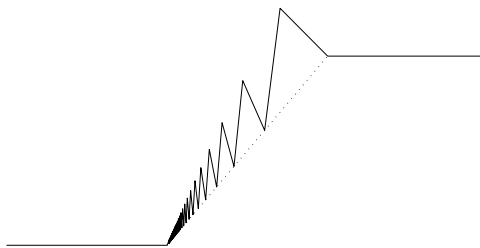
$$f(x,\theta) := 1 + \tfrac{1}{2}\{x \in B_{2k}\} - \tfrac{1}{2}\{x \in B_{2k+1}\}$$

  (i)   Show that $\lambda \left( \sqrt{f_\theta} - \sqrt{f_0} \right)^2 = O(\theta^3)$.
  (ii)   Deduce that $\{f_\theta : 0 \leq \theta < 1\}$ has Hellinger derivative 0 at $\theta = 0$.
  (iii)   Show that the function $\theta \mapsto f(x,\theta)$ is discontinuous at $\theta = 0$, for every $x$. Deduce that none of the functions is differentiable at $\theta = 0$.
  (iv)   Modify the construction to give a family with a nonzero Hellinger derivative at $\theta = 0$ for which none of the $f(x, \cdot)$ are differentiable at $\theta = 0$.

[2]   (Construction of an absolutely continuous density whose square root is not absolutely continuous.) For $i \geq 3$ define

$$\alpha_i = \frac{1}{i(\log i)^2} \qquad \text{and} \qquad \beta_i = \frac{1}{i(\log i)^5},$$

Define $B_i = 2\sum_{j \geq i} \beta_j$. Define functions

$$H_i(t) = \alpha_i \left( 1 - |t - B_i - \beta_i|/\beta_i \right)^+ \qquad \text{and} \qquad H(t) = (1 \wedge t)^+ + \sum_{i \geq 3} H_i(t).$$

(i) Show that $B_i$ decreases like $(\log i)^{-4}$.

(ii) Use the fact that $\sum_i \alpha_i < \infty$ to prove that $H$ is absolutely continuous.

(iii) Show that $\alpha_i / B_i \to 0$, then deduce that $H$ has derivative 1 at 0.

(iv) Show that

$$\sqrt{H(B_{i-1} - \beta_i)} - \sqrt{H(B_{i-1})} = \frac{\alpha_i - \beta_i}{\sqrt{H(B_{i-1} + \beta_i)} + \sqrt{H(B_{i-1})}},$$

which decreases like $1/i$, then deduce that

$$\sum_{i=k}^{k+m} |\sqrt{H(B_{i-1} - \beta_i)} - \sqrt{H(B_{i-1})}|$$

can be made arbitrarily large while keeping $\sum_{i=k}^{k+m} |\beta_i|$ arbitrarily small. Deduce that $\sqrt{H}$ is not absolutely continuous.

(v) Show, by an appropriate "rounding off of the corners" at each point where $H$ has different left and right derivatives followed by some smooth truncation and rescaling, that there exists an absolutely continuous, everywhere differentiable probability density function $f$ for which $\sqrt{f}$ is not absolutely continuous.

[3] Let $f_\theta(x) = \frac{1}{2} \exp(-|x - \theta|)$, for $\theta \in \mathbb{R}$ (the double-exponential location family of densities with respect to Lebesgue measure).

(i) Show that $\int \sqrt{f_\theta(x) f_{\theta+\delta}(x)} \, dx = (1 + \delta/2) \exp(-\delta/2)$.

(ii) Deduce that the density $f_\theta$ is Hellinger differentiable at every $\theta$.

(iii) Show that $\theta \mapsto f_\theta(x)$ is not differentiable, for each fixed $x$, at $\theta = x$.

(iv) Prove Hellinger differentiability by a direct Dominated Convergence argument, without the explicit calculation from (i).

(v) Prove Hellinger differentiability by an appeal to Example <11>, without the explicit calculation from (i).

[4] Suppose $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ is a familiy of densities indexed by a subset $\Theta$ of $\mathbb{R}^k$. Suppose 0 is an interior point of $\Theta$ and that $\mathcal{F}$ is Hellinger differentiable at $\theta = 0$, with derivative $\Delta$. Show that $\Delta(x) = 0$ almost everywhere on $\{f_0 = 0\}$. Hint: Approach 0 from each direction in $\mathbb{R}^k$. Deduce that both $\mathbb{P}_0 \Delta \{f_0 = 0\}$ and $\mathbb{P}_0 \Delta^2 \{f_0 = 0\}$ equal zero.

[5] Suppose $\mathcal{F} = \{f_t(x) : t \in T\}$ is a family of probability densities with respect to a measure $\lambda$, $\mathcal{G} = \{g_s(x) : s \in S\}$ is a family of probability densities with respect to a measure $\mu$. Suppose $\mathcal{F}$ is Hellinger differentiable at $t = 0$ and $\mathcal{G}$ is Hellinger differentiable at $s = 0$. Show that the family of densities $\{f_s(x) g_t(y) : (s, t) \in S \otimes T\}$ with respect to $\lambda \otimes \mu$ is Hellinger differentiable at $(s, t) = (0, 0)$. Hint: Use Cauchy-Schwarz to bound contributions from most of the cross-product terms in the expansion of $\sqrt{f_t(x) g_s(y)}$.

[6] Suppose $\mathcal{F} = \{f_\theta : \theta \in \mathbb{R}^k\}$ has Hellinger derivative $\Delta$ at $\theta_0$. Show that $\mathcal{F}$ is also differentiable in $\mathcal{L}^1$ norm with derivative $\Delta_1 = 2\sqrt{f_{\theta_0}} \Delta$, that is, show

$$\lambda |f_\theta - f_{\theta_0} - (\theta - \theta_0)' \Delta_1| = o(|\theta - \theta_0|) \qquad \text{near } \theta_0$$

[7] If $\mathcal{F}$ is $\mathcal{L}^1$ differentiable and $\lambda \dot{f}^2/f_0 < \infty$ is $\mathcal{F}$ also Hellinger differentiable? [Expand.]

[8] Let $\mathbb{P}_\theta$ be the probability measure defined by the density $f_\theta(\cdot)$. A simple application of the Cauchy-Schwarz inequality shows that

$$H(\mathbb{P}_\theta, \mathbb{P}_{\theta_0})^2 = (\theta - \theta_0)'\lambda \left(\dot{\xi}(x)\dot{\xi}(x)'\right)(\theta - \theta_0) + o(|\theta - \theta_0|^2).$$

Provided the matrix $\Gamma = \lambda\left(\Delta(x)\Delta(x)'\right)$ is nonsingular, it then follows that there exist nonzero constants $C_1$ and $C_2$ for which

$$C_1|\theta - \theta_0| \le H(\mathbb{P}_\theta, \mathbb{P}_{\theta_0}) \le C_2|\theta - \theta_0| \qquad \text{near } \theta_0.$$

If such a pair of inequalities holds, with fixed strictly positive constants $C_1$ and $C_2$, throughout some subset of $\Theta$, then Hellinger distance plays the same role as ordinary Euclidean distance on that set.

[9] Suppose $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ is a family of probability densities with respect to a measure $\lambda$, with index set $\Theta$ a subset of the real line. As in Theorem <9>, suppose

$$\sqrt{f_{\theta+\beta}(x)} - \sqrt{f_\theta(x)} = \int_\theta^{\theta+\beta} \Delta_t(x)\, dt \qquad \text{mod}[\lambda], \quad \text{for } |\beta| \le \delta,\, a \le \theta \le b,$$

with $\sup_t \lambda \Delta_t^2 = C < \infty$, where $[a-\delta, b+\delta] \subseteq \Theta$. Let $\mathcal{Q} = \{q_\alpha : -\delta < \alpha < \delta\}$ be a family of probability densities with respect to Lebesgue measure $\mu$ on $[a, b]$, each bounded by a fixed constant $K$, and with Hellinger derivative $\dot{\eta}$ at $\alpha = 0$. Create a new family $\mathcal{P} = \{p_{\alpha,\beta}(x, \theta) : \max(|\alpha|, |\beta|) < \delta\}$ of probability densities $p_{\alpha,\beta}(x, \theta) = q_\alpha(\theta)f_{\theta+\beta}(x)$ with respect to $\lambda \otimes \mu$.

  (i) Show that $\mathcal{P}$ is Hellinger differentiable at $\alpha = 0$, $\beta = 0$ with derivative having components $\dot{\eta}\sqrt{f_\theta}$ and $\sqrt{f_0}\Delta_\theta$.

  (ii) Try to relax the assumptions on $\mathcal{Q}$.

[10] Suppose $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, with $0 \in \Theta \subseteq \mathbb{R}^k$, is a dominated family of probability measures on a space $\mathcal{X}$, having densities $f_\theta(x)$ with respect to a sigma-finite measure $\lambda$. Define $\mathcal{U}$ as the set of unit vectors

$$\mathcal{U} = \{u : \text{there exists a sequence } \{\theta_i\} \text{ in } \Theta \text{ such that } \theta_i/|\theta_i| \to u \text{ as } i \to \infty\}$$

    Write $\tilde{P}_\theta$ for the part of $P_\theta$ that is absolutely continuous with respect to $P_0$ and $P_\theta^\perp = P_\theta - \tilde{P}_\theta$ for the part that is singular with respect to $P_0$. Write $\tilde{p}_\theta$ for the density $d\tilde{P}_\theta/dP_0$.

    The following are equivalent.

  (i) For some vector $\dot{\xi}$ of functions in $\mathcal{L}^2(\lambda)$,

$$\sqrt{f_\theta} = \sqrt{f_0} + \theta'\dot{\xi} + r_\theta \qquad \text{where } \lambda\left(r_\theta^2\right) = o(|\theta|^2) \text{ near } \theta = 0,$$

    and $u'\dot{\xi} = 0$ a.e. $[\lambda]$ on $\{f_0 = 0\}$, for each $u$ in $\mathcal{U}$.

  (ii) For some vector $\Delta$ of functions in $\mathcal{L}^2(P_0)$,

$$\sqrt{f_\theta} = \sqrt{f_0} + 2\theta'\Delta\sqrt{f_0} + R_\theta \qquad \text{where } \lambda\left(R_\theta^2\right) = o(|\theta|^2) \text{ near } \theta = 0,$$

    and $P_\theta^\perp \mathcal{X} = o(|\theta|^2)$.

  (iii) For some vector $\tilde{\Delta}$ of functions in $\mathcal{L}^2(P_0)$,

$$\sqrt{\tilde{p}_\theta} = 1 + 2\theta'\tilde{\Delta} + \tilde{r}_\theta \qquad \text{where } P_0\left(\tilde{r}_\theta^2\right) = o(|\theta|^2) \text{ near } \theta = 0,$$

    and $P_\theta^\perp \mathcal{X} = o(|\theta|^2)$.

[11] Suppose $\mathcal{P}$ is Hellinger differentiable at 0 and that $\theta_n = \delta_n/\sqrt{n}$ is a sequence in $\Theta$ with $\{\delta_n\}$ bounded. Let $P_{\theta_n}^\perp$ be the part of $P_{\theta_n}$ that is singular with respect to $P_0$.

  (i) Show that $P_{\theta_n}^n \lhd P_0^n$ if and only if $nP_{\theta_n}^\perp(\mathcal{X}) \to 0$.

(ii) If $P_{\theta_n}^n \lhd P_0^n$, show that

$$dP_{\theta_n}^n / dP_0^n = \left(1 + o_p(1; \mathbb{P}_n)\right) \exp\left(\delta_n' Z_n - \tfrac{1}{2}\delta_n' \mathbb{I}_0 \delta_n\right).$$

[12]    Generalize the van Trees inequality to the case where $q$ has support that is a closed subset of $\Theta$. Approximate $q$ by elements of $\mathcal{L}_+^1(\mathfrak{m})$ with compact support then pass to the llimit.

## 10.    Notes

<div style="border:1px solid">Incomplete</div>

     I borrowed the exposition Section 7 from Pollard (1997). The essential argument is fairly standard, but the interpretation of some of the details is novel. Compare with the treatments of Le Cam (1970, and 1986 Section 17.3), Ibragimov & Has'minskii (1981, page 114), Millar (1983, page 105), Le Cam & Yang (1990, page 101), or Strasser (1985, Chapter 12).

     Hájek (1962) used Hellinger differentiability to establish limit behaviour of rank tests for shift families of densities. Most of results in Section 2 are adapted from the Appendix to Hájek (1972), which in turn drew on Hájek & Šidák (1967, page 211) and earlier work of Hájek. For a proof of the multivariate version of Theorem <9> see Bickel et al. (1993, page 13). A reader who is puzzled about all the fuss over negligible sets, and behaviour at points where the densities vanish, might consult Le Cam (1986, pages 585–590) for a deeper discussion of the subtleties.

     The proof of the information inequality (Lemma <12>) is adapted from Ibragimov & Has'minskii (1981, Section 1.7), who apparently gave credit to Blyth & Roberts (1972), but I could find no mention of Hellinger differentiability in that paper.

     Gilles Stoltz explained to me how the van Trees inequality could be derived without the full Hellinger differentiability of the family $\mathcal{G}$.

*Reference?*

     Cite van der Vaart (1988, Appendix A3) and Bickel et al. (1993, page 461) for Theorem <22>. Ibragimov & Has'minskii (1981, page 70) asserted that the result follows by "direct calculations". Indeed my proof uses the same truncation trick as in the proof of Lemma <12>, which is based on the argument of Ibragimov & Has'minskii (1981, page 65). Le Cam & Yang (1988, Section 7) deduced an analogous result (preservation of DQM under restriction to sub-sigma-fields) by an indirect argument using equivalence of DQM with the existence of a quadratic approximation to likelihood ratios of product measures (an LAN condition).

<div align="center">REFERENCES</div>

Bickel, P. J., Klaassen, C. A. J., Ritov, Y. & Wellner, J. A. (1993), *Efficient and Adaptive Estimation for Semiparametric Models*, Johns Hopkins University Press, Baltimore.

Blyth, C. & Roberts, D. (1972), On inequalities of Cramér-Rao type and admissibility proofs, *in* L. Le Cam, J. Neyman & E. L. Scott, eds, 'Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability', Vol. I, University of California Press, Berkeley, pp. 17–30.

Gill, R. & Levit, B. (1995), 'Applications of the van Trees inequality: a Bayesian Cramér-Rao bound', *Bernoulli* **1**, 59–79.

Hájek, J. (1962), 'Asymptotically most powerful rank-order tests', *Annals of Mathematical Statistics* **33**, 1124–1147.

Hájek, J. (1972), Local asymptotic minimax and admissibility in estimation, *in* L. Le Cam, J. Neyman & E. L. Scott, eds, 'Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability', Vol. I, University of California Press, Berkeley, pp. 175–194.

Hájek, J. & Šidák, Z. (1967), *Theory of Rank Tests*, Academic Press. Also published by Academia, the Publishing House of the Czechoslavak Academy of Sciences, Prague.

Ibragimov, I. A. & Has'minskii, R. Z. (1981), *Statistical Estimation: Asymptotic Theory*, Springer, New York.

Kagan, A. & Shepp, L. A. (2005), 'A sufficiency paradox: an insufficient statistic preserving the Fisher information', *The American Statistician* **59**, 54–56.

Le Cam, L. (1970), 'On the assumptions used to prove asymptotic normality of maximum likelihood estimators', *Annals of Mathematical Statistics* **41**, 802–828.

Le Cam, L. (1986), *Asymptotic Methods in Statistical Decision Theory*, Springer-Verlag, New York.

Le Cam, L. & Yang, G. L. (1988), 'On the preservation of local asymptotic normality under information loss', *Annals of Statistics* **16**, 483–520.

Le Cam, L. & Yang, G. L. (1990), *Asymptotics in Statistics: Some Basic Concepts*, Springer-Verlag.

Le Cam, L. & Yang, G. L. (2000), *Asymptotics in Statistics: Some Basic Concepts*, 2nd edn, Springer-Verlag.

Millar, P. W. (1983), 'The minimax principle in asymptotic statistical theory', *Springer Lecture Notes in Mathematics* **976**, 75–265.

Pollard, D. (1997), Another look at differentiability in quadratic mean, *in* D. Pollard, E. Torgersen & G. L. Yang, eds, 'A Festschrift for Lucien Le Cam', Springer-Verlag, New York, pp. 305–314.

Strasser, H. (1985), *Mathematical Theory of Statistics: Statistical Experiments and Asymptotic Decision Theory*, De Gruyter, Berlin.

van der Vaart, A. (1988), *Statistical estimation in large parameter spaces*, Center for Mathematics and Computer Science. CWI Tract 44.