

Lecture 15

In this and the next couple of lectures we will go back to our initial story of estimating "smooth" functionals in infinite dimensional models.

The general set up: X_1, \dots, X_n are a random sample on some sample space (Ω, \mathcal{A}) having distribution P . We assume $P \in \mathcal{P} \leftarrow$ a model and want to estimate $\Psi(P): \mathcal{P} \rightarrow \mathbb{R}$. We saw that in many examples $\mathcal{P} = \{P_\eta: \eta \in \mathbb{H}\}$ where \mathbb{H} is a subset of a normed linear space, and $\Psi(P) = \chi(\eta)$. ~~We saw that~~ We have considered the analysis of quadratic functional of a density as our working example so far. The analysis was initiated by a von-Mises type expansion of the functionals, which requires certain kind of differentiability of the functionals and the distributions under study. We shall now try to make it more concrete. We refer to Aad van der Vaart! Asymptotic Statistics (Chapter 25) for more details.

We first begin by defining tangent space and influence functions when we have X_1, \dots, X_n iid $P \in \mathcal{P}$ and we want to estimate $\Psi(P)$.

Definition (Differentiable Path) A differentiable path is a map $t \rightarrow P_t$ from a neighbourhood of 0, $(0, \varepsilon) \subset (0, +\infty)$ to \mathcal{P} such that for some measurable $g: \Omega \rightarrow \mathbb{R}$,

$$\int \left(\frac{dP_t^{1/2} - dP^{1/2}}{t} - \frac{1}{2} g dP^{1/2} \right)^2 \rightarrow 0 \quad \text{as } t \downarrow 0$$

①

The $\int_0^t g$ is called the score \int_0^t of the submodel $\{P_t: 0 \leq t < \varepsilon\}$ at $t=0$ with $P_0 = P$ the truth.

Indeed the preceding definition is not very enlightening. First of all we have not made clear what the objects $dP_t^{1/2}$, $dP^{1/2}$ are. They can be formalized by introducing Hilbert space of square root of probability measures. (de laan 1960's). A more simple form is

$$\int \left(\frac{P_{tt}^{1/2} - P_t^{1/2}}{t} - \frac{1}{2} g P_t^{1/2} \right)^2 d\mu_t \rightarrow 0 \text{ as } t \downarrow 0$$

where for each t , μ_t is an arbitrary measure relative to which both P and P_t have densities P_t and P_{tt} respectively. ~~It is~~ the choice of μ_t is not important. We shall not pay too much attention to these subtleties and instead try to understand the implications of the definitions. To this end note that

$$\begin{aligned} & \int \left(\frac{dP_t^{1/2} - dP^{1/2}}{t} - \frac{1}{2} g dP^{1/2} \right)^2 \\ &= \int \left(\frac{dP_t^{1/2} - dP^{1/2}}{t dP^{1/2}} - \frac{1}{2} g \right)^2 dP \end{aligned}$$

$$\therefore \text{the definition} \Rightarrow \lim_{t \rightarrow 0} \mathbb{E}_P \left\{ \left(\frac{dP_t^{1/2} - dP^{1/2}}{t dP^{1/2}} - \frac{g}{2} \right)^2 \right\} = 0$$

\Rightarrow We have some kind of differentiability in quadratic means of square root of our densities. This differentiability requires convergence in $L_2(P)$ sense (or $L_2(\mu_t)$ sense), and not pointwise.

$\{P_t: 0 \leq t < \epsilon\}$ that is diff. at quadratic mean at $t=0$ is called Smooth Parametric Submodel of \mathcal{P} .

Definition (Tangent Set) letting $t \rightarrow P_t$ range over a collection of smooth and "regular" submodels we obtain a collection of score functions, which we call a Tangent Set of the model \mathcal{P} at P , denoted by \mathcal{T}_P .

By regular we mean that the Fisher information at $t=0$ of $\{P_t: 0 \leq t < \epsilon\}$ defined as $(\mathbb{E}(g^2))$ is non-zero.

When we consider all possible differentiable paths $t \rightarrow P_t$, we obtain the maximal collection of score functions \rightarrow Maximal Tangent Set.

The following lemmas gives fundamental properties of the score functions.

Lemma 1: Every score function satisfies $\mathbb{E}_P(g) = 0$ and $\mathbb{E}_P g^2 < \infty$.

Proof: For given but arbitrary $(t_n)_{n \geq 1}$ such that $t_n \rightarrow 0$ as $n \rightarrow \infty$, let p_n and p be densities of P_{t_n} and P w.r.t a dominating measure μ . (for example a convex combination of the countably many measures $P_{t_n} + P$).

By definition $(\sqrt{p_n} - \sqrt{p})/t_n$ converges in $L_2(\mu)$ to $\frac{1}{2} g \sqrt{p}$.
 $\Rightarrow \left\| \frac{\sqrt{p_n} - \sqrt{p}}{t_n} - \frac{1}{2} g \sqrt{p} \right\|_{L_2(\mu)}^2 \rightarrow 0$.

Now $(\sqrt{p_n} - \sqrt{p})/t_n \in L_2(\mu) \quad \forall n \geq 1 \Rightarrow \frac{g \sqrt{p}}{2}$ is a limit of $L_2(\mu)$ functions in $L_2(\mu) \Rightarrow g \sqrt{p}/2 \in L_2(\mu) \Rightarrow g \in L_2(\mathcal{P})$.

(3)

$$\text{Now, } \lim_{t_n \rightarrow 0} \|\sqrt{p_n} - \sqrt{p} - \frac{1}{2} g \sqrt{p} t_n\|_{L_2(\mu)} = \lim_{t_n \rightarrow 0} \left\| \frac{\sqrt{p_n} - \sqrt{p}}{t_n} - \frac{1}{2} g \sqrt{p} \right\|_{L_2(\mu)} t_n = 0$$

$$\Rightarrow \text{But } \|\sqrt{p_n} - \sqrt{p} - \frac{1}{2} g \sqrt{p} t_n\|_{L_2(\mu)}$$

$$\geq \|\sqrt{p_n} - \sqrt{p}\|_{L_2(\mu)} - \frac{1}{2} t_n \|g \sqrt{p}\|_{L_2(\mu)}$$

$$\Rightarrow 0 \geq \lim_{n \rightarrow \infty} \left| \|\sqrt{p_n} - \sqrt{p}\|_{L_2(\mu)} - \frac{1}{2} t_n \|g \sqrt{p}\|_{L_2(\mu)} \right|$$

$$\Rightarrow \|\sqrt{p_n} - \sqrt{p}\|_{L_2(\mu)} \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\Rightarrow \sqrt{p_n} \rightarrow \sqrt{p} \text{ in } L_2(\mu)$$

$$\Rightarrow \mathbb{E}_p(g) = \int g dp = \int g p d\mu = \int \frac{1}{2} g \sqrt{p} \times 2\sqrt{p} d\mu$$

$$= \left\langle \frac{1}{2} g \sqrt{p}, 2\sqrt{p} \right\rangle_{L_2(\mu)}$$

$$= \lim_{n \rightarrow \infty} \int \frac{(\sqrt{p_n} - \sqrt{p})}{t_n} \times (\sqrt{p_n} + \sqrt{p}) d\mu$$

$$= \lim_{n \rightarrow \infty} \int \frac{(p_n - p)}{t_n} d\mu = 0 \quad \square \quad \text{(continuity of inner product)}$$

It follows that a tangent set can be identified with a subset of $L_2(p)$ (upto equivalence classes). Tangent sets are of the form linear spaces, and referred to tangent spaces in such cases. Geometrically, we may visualize the model \mathcal{P}_p or rather the corresponding "square root of measures" $dP^{1/2}$, as a subset of the unit ball of a Hilbert space & the set \mathcal{P}_p (or rather the objects $\frac{1}{2} g p^{1/2}$) as its tangent set.

Usually we construct submodels $t \rightarrow P_t$ such that for every x , $g(x) = \frac{\partial}{\partial t} \log dP_t(x) \Big|_{t=0}$. Such pointwise differentiability is not required by our definition. We therefore often make use of the next lemma.

Lemma 2: If p_t is a density w.r.t. fixed measure μ and $t \rightarrow p_t(x)$ is continuously diff. in a neighbourhood of 0 and $t \rightarrow \int \dot{p}_t^2 / p_t d\mu$ is finite and continuous in this neighborhood, then $t \rightarrow P_t$ is a diff. path.

It also turns out that our definition of differentiability is ~~the~~ a reasonable one because it ensures a desired form of local asymptotic normality.

Lemma 3: If $t \rightarrow P_t$ is a differentiable path, then

$$\log \prod_{i=1}^n \frac{dP_{t_i/\sqrt{n}}}{dP}(x_i) = \frac{1}{\sqrt{n}} \sum_{i=1}^n g(x_i) - \frac{1}{2} \mathbb{E}_P(g^2) + o_p(1)$$

Now come to estimation of suitably differentiable functionals. ($P \rightarrow \Psi(P)$)

Differentiable Functionals:

For defining "information" for estimating $\Psi(P)$, only those submodels $t \rightarrow P_t$ along which $t \rightarrow \Psi(P_t)$ is ~~of~~ differentiable is of interest. A minimal requirement is differentiability at $t=0$, but we need more. We follow definition introduced by Levit (1978). Let $\Psi: \mathcal{P} \rightarrow \mathbb{B}$ where \mathbb{B} is a Banach space.

at P

Definition: A map $\Psi: \mathcal{P} \rightarrow \mathbb{B}$ is differentiable relative to a given tangent space $\dot{\mathcal{P}}_P$ if there exists a continuous linear map $\dot{\Psi}_P: L_2(P) \rightarrow \mathbb{B}$ such that $\forall g \in \dot{\mathcal{P}}_P$ and corresponding submodel $t \rightarrow P_t$,

$$\frac{\Psi(P_t) - \Psi(P)}{t} \xrightarrow{\mathbb{B}} \dot{\Psi}_P(g) \text{ i.e. } \left\| \frac{\Psi(P_t) - \Psi(P)}{t} - \dot{\Psi}_P(g) \right\|_{\mathbb{B}} \rightarrow 0$$

It requires that $t \rightarrow \Psi(P_t)$ is differentiable in usual sense with an extra requirement for it to have a special representation. The map $\dot{\Psi}_P$ is much like a Hadamard derivative of Ψ viewed as a map on the space of square root of measures.

Implications: In the case $\mathbb{B} = \mathbb{R}^k$ for some $k \in \mathbb{N}$, we have by Riesz Representation theorem, that $\exists \tilde{\Psi}_P: \Omega \rightarrow \mathbb{R}^k$ s.t. $\Psi_P(g) = \langle \tilde{\Psi}_P, g \rangle_{L_2(P)} = \int \tilde{\Psi}_P g \, dP$ ($\tilde{\Psi}_P \in L_2(P)$). The function $\tilde{\Psi}_P$ is defined uniquely up to inner products with elements of the tangent set (which is not necessarily all of $L_2(P)$). However it is possible to find $\tilde{\Psi}_P$ whose coordinates $\in \text{lin } \dot{\mathcal{P}}_P \rightarrow$ this function is unique and is called the efficient influence function. Any other $\tilde{\Psi}_P$ is simply called an influence function of $\Psi(P)$ at P w.r.t $\dot{\mathcal{P}}_P$. In particular one takes as an influence function any mble $\Psi_P: \Omega \rightarrow \mathbb{R}^k$ whose projection onto $\text{lin } \dot{\mathcal{P}}_P$ is the efficient influence function.

Remark: It can be shown that the maximal tangent set is a cone, which we will assume from now on.

Examples:

① Parametric Model: Consider a parametric model with parameter $\theta \in \Theta \subseteq \mathbb{R}^d$ and Θ is an open subset of \mathbb{R}^d . Assume $\{P_\theta\}_{\theta \in \Theta}$ to be the class of probability distributions and p_θ be density of P_θ w.r.t. some common dominating measure μ .

Suppose \exists vector valued measurable map \dot{l}_θ such that

$$\int \left(p_{\theta+h}^{1/2} - p_\theta^{1/2} - \frac{1}{2} h^T \dot{l}_\theta p_\theta^{1/2} \right)^2 d\mu = o(\|h\|^2) \text{ as } h \rightarrow 0 \text{ in } \mathbb{R}^d.$$

Then a tangent set at P_θ is given by linear space $\{h^T \dot{l}_\theta : h \in \mathbb{R}^m\} \rightarrow$ span of score functions for the coordinates of θ .

If the Fisher Information matrix $I_\theta = P_\theta(\dot{l}_\theta \dot{l}_\theta^T)$ is invertible then every map $\chi: \Theta \rightarrow \mathbb{R}^k$ that is differentiable in an ordinary sense is differentiable as a map $\Psi(P_\theta) = \chi(\theta)$ on the model relative to the tangent space. To see this, note that ~~that~~ the submodel $t \rightarrow \theta + th$ has score $h^T \dot{l}_\theta$ at P_θ and

$$\begin{aligned} \frac{\partial}{\partial t} \chi(\theta + th) \Big|_{t=0} &= \dot{\chi}_\theta h = P_\theta (\dot{\chi}_\theta I_\theta^{-1} \dot{l}_\theta \dot{l}_\theta^T h) \\ &= P_\theta (\underbrace{\dot{\chi}_\theta I_\theta^{-1} \dot{l}_\theta}_{\Psi_{P_\theta}(\dot{g})} \underbrace{h^T \dot{l}_\theta}_{\dot{g}}) \end{aligned}$$

② Nonparametric Model: Suppose \mathcal{P} consists of all probability laws on a given sample space. Then the tangent set at $P = \{g: \int g dP = 0, \int g^2 dP < \infty\} = L_2^0(P)$. This is indeed the maximal tangent set. To see that this is indeed the tangent set we exhibit suitable one dimensional submodels of \mathcal{P} at P s.t. we can get any $g \in L_2^0(P)$ as a score.

For a bounded $g \in L_2^0(P)$, consider

$$p_t(x) = c(t) \exp(tg(x)) p(x) \quad \begin{array}{l} p \text{ being} \\ \text{density of } P \\ \text{w.r.t some} \\ \text{fixed measure} \\ \mu. \end{array}$$

$$\Rightarrow g(x) = \left. \frac{\partial}{\partial t} \log p_t(x) \right|_{t=0}$$

One can also show that g is a score in an L_2 -sense (by direct calculation using boundedness of g)

For unbounded functions g these submodels might not be well defined. However one can always use

$$p_t(x) = c(t) f(tg(x)) p(x) \quad \text{with}$$

$$f(0) = f'(0) = 1$$

e.g: $f(x) = 2(1 + e^{-2x})^{-1}$ works.

For differentiability of $\Psi(P)$ at P we need for any differentiable submodel $\{P_t\}$

$$\frac{\Psi(P_t) - \Psi(P)}{t} \rightarrow \int \dot{\Psi}_P g dP \quad \text{for some } \dot{\Psi}_P.$$

e.g: When P is on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ and $\Psi(P) = \int x dP$, $\dot{\Psi}_P(x) = x$

⑧

(check by direct calculation)

~~Score and information operators~~

Score and information operators

Consider first the situation that $\mathcal{P} = \{P_\eta : \eta \in H\}$ and η is itself a probability measure on some measurable space. We are interested in estimation of $\psi(P_\eta) = \chi(\eta)$ for $\chi: H \rightarrow \mathbb{R}^d$ on the model H . We will work $k=1$ here.

The model H gives rise to a tangent set \dot{H}_η at η . If the map $\eta \rightarrow P_\eta$ is differentiable in an appropriate sense then the corresponding derivative will map every score $b \in \dot{H}_\eta$ into a score g for the model \mathcal{P} at P_η . Let's make this precise.

Assume that a smooth parametric submodel $t \rightarrow \eta_t$ at η induces a smooth parametric submodel $t \rightarrow P_{\eta_t}$ at P_η where $\eta = \eta_0$, and that the score functions b of the submodel $t \rightarrow \eta_t$ and g of $t \rightarrow P_{\eta_t}$ are related by

$$g = A_\eta b$$

for some operator $A_\eta : \dot{H}_\eta \subset L_2(\eta) \rightarrow L_2(P_\eta)$.

We assume that

$A_\eta: \text{lin } \dot{H}_\eta \subset L_2(\eta) \rightarrow L_2(P_\eta)$ is continuous & linear.

Next assume that the function $\eta \rightarrow \chi(\eta)$ is differentiable with influence function $\tilde{\chi}_\eta$ w.r.t. \dot{H}_η .

Then by definition the function $\psi(P_\gamma) = x(\gamma)$ is pathwise differentiable w.r.t the tangent set $\dot{x}_{P_\gamma} = A_\gamma \dot{H}_\gamma$ iff \exists a function $\tilde{\psi}_{P_\gamma}$ such that

$$\langle \tilde{\psi}_{P_\gamma}, A_\gamma b \rangle_{P_\gamma} = \frac{\partial}{\partial t} \psi(P_{\gamma_t}) \Big|_{t=0} = \frac{\partial}{\partial t} x(\gamma_t) \Big|_{t=0} = \langle \tilde{x}_\gamma, b \rangle_\gamma \quad \forall b \in \dot{H}_\gamma$$

This equation can be written in terms of the adjoint score operator $A_\gamma^* : L_2(P_\gamma) \rightarrow \overline{\text{lin } \dot{H}_\gamma}$. By definition of the adjoint

$$\langle h, A_\gamma b \rangle_{P_\gamma} = \langle A_\gamma^* h, b \rangle_\gamma \quad \text{for } h \in L_2(P_\gamma), b \in \dot{H}_\gamma$$

Note that we define A_γ^* to have range $\overline{\text{lin } \dot{H}_\gamma}$ so that it is the adjoint of $A_\gamma : \dot{H}_\gamma \rightarrow L_2(P_\gamma)$. This is the adjoint of an extension $A_\gamma : \underline{L_2(\dot{H}_\gamma)} \rightarrow L_2(P_\gamma)$ followed by an orthogonal projection onto $\overline{\text{lin } \dot{H}_\gamma}$.

Some Results from Functional Analysis: Every continuous linear map $A : H_1 \rightarrow H_2$ between Hilbert spaces H_1 and H_2 has an adjoint map $A^* : H_2 \rightarrow H_1$ which is a continuous linear map & uniquely determined by

$$\langle A^* h_2, h_1 \rangle_1 = \langle h_2, A h_1 \rangle_2 \quad \text{for all } h_1 \in H_1, h_2 \in H_2.$$

If A is considered to be the restriction to $H_1 \subset \tilde{H}_1$ of a continuous linear map $\tilde{A} : \tilde{H}_1 \rightarrow H_2$ where \tilde{H}_1 an Hilbert that contains H_1 isometrically, then $A^* = \Pi \tilde{A}^*$ where $\Pi : \tilde{H}_1 \rightarrow H_1$ the orthogonal projection on H_1 .

Going back to our problem, we have

$$\langle h, A_\gamma b \rangle_{P_\gamma} = \langle A_\gamma^* h, b \rangle_\gamma \quad \forall h \in L_2(P_\gamma), b \in \dot{H}_\gamma$$

But pathwise differentiability requires

$$\langle \tilde{\Psi}_{P_\gamma}, A_\gamma b \rangle_{P_\gamma} = \langle \tilde{x}_\gamma, b \rangle_\gamma \quad \forall b \in \dot{H}_\gamma$$

But

$$\langle \tilde{\Psi}_{P_\gamma}, A_\gamma b \rangle_{P_\gamma} = \langle A_\gamma^* \tilde{\Psi}_{P_\gamma}, b \rangle_\gamma \quad \forall b \in \dot{H}_\gamma$$

$$\Rightarrow A_\gamma^* \tilde{\Psi}_{P_\gamma} = \tilde{x}_\gamma \dots (*)$$

$\Rightarrow \psi(P_\gamma) = x(\gamma)$ is differentiable w.r.t. to the tangent set $\dot{x}_{P_\gamma} = A_\gamma \dot{H}_\gamma$ iff (*) can be solved for $\tilde{\Psi}_{P_\gamma}$.

Equivalently, this happens iff $\tilde{x}_\gamma \in \text{Range}(A_\gamma^*)$. Since A_γ^* is not necessarily onto $\overline{\text{lin } \dot{H}_\gamma}$, this is indeed a condition!

Now two solutions of (*) can only differ by an element of the null space $\mathcal{N}(A_\gamma^*) = \mathcal{R}(A_\gamma)^\perp$ where $\mathcal{R}(A_\gamma)$ is the range $A_\gamma: \text{lin } \dot{H}_\gamma \rightarrow L_2(P_\gamma)$. Thus there is at most one solution $\tilde{\Psi}_{P_\gamma}$ that is contained in $\overline{\mathcal{R}(A_\gamma)} = \overline{\text{lin } A_\gamma \dot{H}_\gamma}$.

If \tilde{x}_γ is contained in the smaller range of $A_\gamma^* A_\gamma$ then equation (*) can be solved &

$$\tilde{\Psi}_{P_\gamma} = A_\gamma (A_\gamma^* A_\gamma)^- \tilde{x}_\gamma$$

$A_\gamma^* A_\gamma$ is called an information operator & $(A_\gamma^* A_\gamma)^-$ is a "generalized inverse". The following lemma shows that this attractive solution is available for any functional x if the range of the score operator is closed, a situation which fails often (unfortunately).

Results from Functional Analysis

Lemma 1: Let $A: H_1 \rightarrow H_2$ be a continuous linear map between Hilbert spaces H_1 and H_2 . Then the following are equivalent.

- (i) $\mathcal{R}(A)$ is closed.
- (ii) $\mathcal{R}(A^*)$ is closed.
- (iii) $\mathcal{R}(A^*A)$ is closed.
- (iv) $\mathcal{R}(A^*A) = \mathcal{R}(A^*)$.

Lemma 2: Let $A: H_1 \rightarrow H_2$ be a continuous linear map between two Hilbert spaces H_1 and H_2 . Then

- (i) $\mathcal{N}(A) = \mathcal{R}(A^*)^\perp$
- (ii) $\mathcal{N}(A^*) = \mathcal{R}(A)^\perp$

Also, $A^*A: H_1 \rightarrow H_1$ is one-to-one and onto and has a continuous inverse iff A is ~~not~~ one-to-one and onto and $\mathcal{R}(A)$ is closed iff A^*A is one-to-one and onto.

More generally: So far we have assumed that η is a probability distribution. However a close inspection of the arguments reveal that this is not necessary.

More generally let $\mathcal{P} = \{P_\eta : \eta \in H\}$ where H is an arbitrary set. Let H_η be a subset of a Hilbert space that indexes "directions" b in which η can be approximated within H .

Suppose \exists continuous linear map/operators $A_\eta: \text{lin } H_\eta \rightarrow L_2(P_\eta)$ and $\dot{x}_\eta: \text{lin } H_\eta \rightarrow \mathbb{R}^2 \forall b \in H_\eta \exists$ a path $t \rightarrow \eta_t$ s.t. as $t \downarrow 0$

$$\int \left[\frac{dP_{\eta_t}^{1/2} - dP_\eta^{1/2}}{t} - \frac{1}{2} A_\eta b dP_\eta^{1/2} \right]^2 \rightarrow 0$$

$$\mathbb{E} \frac{x(\eta_t) - x(\eta)}{t} \rightarrow \dot{x}_\eta b$$

Shun by Riesz representation theorem, $\exists \tilde{x}_\eta \in \overline{\text{lin } H_\eta}$ s.t.
 $\tilde{x}_\eta b = \langle \tilde{x}_\eta, b \rangle$. Hereafter proceeding as before, we can
obtain the following theorem.

Theorem: The map $\psi: \mathcal{P} \rightarrow \mathbb{R}$ given by $\psi(P_\eta) = x(\eta)$ is
differentiable at P_η relative to the tangent set $\overline{A_\eta H_\eta}$
iff $\tilde{x}_\eta \in \mathcal{R}(A_\eta^*)$ for $A_\eta^*: \text{L}_2(P_\eta) \rightarrow \overline{\text{lin } H_\eta}$.

Proof: By assumption $A_\eta H_\eta$ is tangent set. The rest are
similar arguments as before.

There are two ~~very~~ reasons why one might fail to
have a solution to (*)

$$A_\eta^* \tilde{\psi}_{P_\eta} = \tilde{x}_\eta, \quad \tilde{x}_\eta \in \overline{\text{lin } H_\eta}$$

$$A_\eta: H_\eta \rightarrow \text{L}_2(P_\eta)$$

$$A_\eta^*: \text{L}_2(P_\eta) \rightarrow \overline{\text{lin } H_\eta}$$

Firstly $\mathcal{R}(A^*) \not\subseteq \overline{\text{lin } H_\eta}$. Since $b \perp \mathcal{R}(A^*)$ iff $b \in \mathcal{N}(A_\eta)$,
this can only happen when A_η is not $\perp\text{-}\perp$. This
means two directions b may lead to same score for $A_\eta b$
 \Rightarrow information matrix ^{for} corresponding to ~~the~~ two dimensional
model is singular. An intuitive explanation is that
the parameter is not locally identifiable & hence we
have a problem.

Secondly, $\mathcal{R}(A^*)$ may be ~~closed~~ dense but not closed.
There for any $\tilde{x}_\eta \in \mathcal{R}(A_\eta^*)$ that are
arbitrarily close to \tilde{x}_η but (*) cannot be solved.
Such situations can give rise to problematic situations.

Theorem: If $\tilde{x}_\eta \notin \mathcal{R}(A_\eta^*)$ then

(i) \nexists estimators for $x(\eta)$ that are "regular" at P_η .

(ii) $\sup_{b \in H_\eta} \frac{\langle \tilde{x}_\eta, b \rangle_\eta^2}{\|A_\eta b\|_{P_\eta}^2} = +\infty$ (Has implications in \sqrt{n} -consistent
estimation).

Example: (Information loss Model)

Suppose that a typical observation is distributed as a measurable transformation $X = m(Y)$ where Y is not observed. Suppose now that $Y \sim \eta$, $\eta \in H$. This yields a natural parametrization of the distribution of $X \sim P_\eta$. Assume the form of m is known.

There is cool property that this type of censoring mechanisms entail: if $t \rightarrow \eta_t$ is a differentiable submodel of H , then the induced submodel $t \rightarrow P_{\eta_t}$ is a differentiable of $\{P_\eta : \eta \in H\}$. Moreover, the score $g = A_\eta b$ for induced model $t \rightarrow P_{\eta_t}$ can be obtained from score b of the model $t \rightarrow \eta_t$ by taking conditional expectations.

$$A_\eta b(x) = \mathbb{E}_\eta(b(Y) | X=x).$$

lemma: Suppose that $\{\eta_t : 0 < t < 1\}$ is a collection of probability measures on a measurable space (Y, \mathcal{B}) , s.t. for some measurable $b: Y \rightarrow \mathbb{R}$,

$$\int \left(\frac{d\eta_t^{1/2} - d\eta^{1/2}}{t} - \frac{1}{2} b d\eta^{1/2} \right)^2 \rightarrow 0$$

For a measurable map $m: Y \rightarrow \mathcal{X}$ let P_η be the distribution of $m(Y)$ if $Y \sim \eta$ & let $A_\eta b(x) = \mathbb{E}(b(Y) | m(Y)=x)$

Then

$$\int \left(\frac{dP_{\eta_t}^{1/2} - dP_\eta^{1/2}}{t} - \frac{1}{2} A_\eta b dP_\eta^{1/2} \right)^2 \rightarrow 0$$

Proof: Assume for the sake of notational simplicity that η_t and η have densities h_t and h w.r.t. some dominating measure ν . (If this is not the case choose $\nu = \nu_t = \frac{1}{2}(\eta + \eta_t)$ and add an extra t -subscript throughout the following arguments).

Further more assume b is uniformly bounded by M (o.w. truncate b at $M_t \rightarrow \infty$ & add t^2 's in the following chain of arguments accordingly).

then $\int u_t^2 \nu \rightarrow 0$ where

$$u_t = \left(\frac{h_t^{\frac{1}{2}} - h^{\frac{1}{2}}}{t} - \frac{1}{2} b h^{\frac{1}{2}} \right)$$

Define μ to be the law of $X = m(Y)$ where Y has ν as its distribution. Then

$$p_t(x) = \mathbb{E}_\nu (h_t(Y) | X=x) \text{ \& } p(x) = \mathbb{E}_\nu (h(Y) | X=x)$$

are densities of P_{η_t} and P_η w.r.t. μ . For P_η this follows because

$$\begin{aligned} P_\eta(A) &= \int \mathbb{1}_A(m(y)) d\eta(y) = \mathbb{E}_\nu (\mathbb{1}_A(m(Y)) h(Y)) \\ &= \mathbb{E}_\nu (\mathbb{1}_A(X) \mathbb{E}_\nu (h(Y) | X)) \\ &= \int \mathbb{1}_A(x) p(x) d\mu(x) \end{aligned}$$

By a similar argument we have a.s. under P_η

$$A_\eta b(x) p(x) = \mathbb{E}_\nu (b(Y) h(Y) | X)$$

Now, from definition of u_t ,

$$h_t = h + t b h + t^2 u_t^2 + t \left(t u_t b h^{\frac{1}{2}} + 2 u_t h^{\frac{1}{2}} + \frac{1}{4} t b^2 h \right)$$

Evaluating these at Y and taking conditional expectations with respect to X ,

$$p_t = p + t A_\eta b p + c_t + d_t$$

$$c_t(x) = t^2 \mathbb{E}_\nu (u_t^2(Y) | X)$$

$$d_t(x) = t \mathbb{E}_\nu \left(\left(t u_t b h^{\frac{1}{2}} + 2 u_t h^{\frac{1}{2}} + \frac{1}{4} t b^2 h \right) (Y) | X \right)$$

$$|d_t(x)|^2 = t^2 \mathbb{E}_\nu \left(\left(t u_t b h^{\frac{1}{2}} + 2 u_t h^{\frac{1}{2}} + \frac{1}{4} t b^2 h \right) (\gamma) | x \right)^2$$

$$\lesssim t^2 \mathbb{E}_\nu \left(\left(u_t h^{\frac{1}{2}} (tM+1) + t M^2 h \right) (\gamma) | x \right)$$

$$\lesssim t^2 \left(\mathbb{E}_\nu (u_t^2(\gamma) | x) (tM+1)^2 + t^2 M^4 p(x) \right) p(x) \quad (\text{by Cauchy Schwarz inequality})$$

Now using ~~the~~ lemma 3 below we have

$$\left[\frac{p_t^{\frac{1}{2}} - p^{\frac{1}{2}}}{t} - \frac{1}{2} (A_\eta b) p^{\frac{1}{2}} \right]^2 \lesssim \mathbb{E}_\nu (u_t^2(\gamma) | x) (tM+1)^2 + t^2 M^4 p(x) + \mathbb{E}_\nu (u_t^2(\gamma) | x) + \left| \frac{1}{\sqrt{1-Mt}} - 1 \right|^2 M^2 p(x)$$

on the set $A = \{p > 0\}$.

The integral over the set A of the above ~~is~~ relative to μ converges to 0 as $t \rightarrow 0$.

Finally $\eta(m^{-1}(A^c)) = P_\eta(A^c) = 0 \Rightarrow P_{\eta_t}(A^c) = \eta_t(m^{-1}(A^c)) = 0 (t^2)$ because $\nu(t u_t \perp B) = \eta_t(B)$ if $\eta(B) = 0$. Thus the integral over A^c converges to 0 as well.

lemma 3: For $a, b, c, d \in \mathbb{R}$ with $a > 0$, $b/a \leq \varepsilon < 1$, $c \geq 0$ and

$$a + b + c + d \geq 0,$$

$$\left| \sqrt{a+b+c+d} - \sqrt{a} - \frac{1}{2} \frac{b}{\sqrt{a}} \right| \leq \frac{3d^2}{a(1-\varepsilon)} + 3c + \left| \frac{1}{\sqrt{1-\varepsilon}} - 1 \right|^2 \frac{b^2}{a}.$$

Now view $A_\eta: L_2(\eta) \rightarrow L_2(P_\eta)$. To get its adjoint note

$$\langle g, A_\eta b \rangle_{L_2(P_\eta)} = \mathbb{E} (g(x) A_\eta b(x))$$

$$= \mathbb{E} (g(x) \mathbb{E}_\eta (b(\gamma) | x))$$

$$= \mathbb{E} (g(x) b(\gamma))$$

$$= \mathbb{E} (\mathbb{E}_\eta (g(x) | \gamma) b(\gamma))$$

$$= \langle A_\eta^* g, b \rangle_{L_2(\eta)}$$

$$A_\eta^* g = \mathbb{E}_\eta (g(x) | \gamma = \cdot) : L_2(P_\eta) \rightarrow L_2(\eta)$$

Although the score operator is defined for all $L_2(\gamma)$ we need its restriction to $\text{lin } H_\gamma$. Hence forth we need to project the conditional expectation $\mathbb{E}_\gamma(g(X)|Y=y)$ onto $\text{lin } H_\gamma$.

(This ends our discussion on existence of first order influence functions)

Morally Correct Calculations of Influence Functions

We work with an example, that of estimation of a treatment effect functional.

Suppose your typical observation consists of (Y, A, X) where $A \in \{0, 1\}$ corresponds to a treatment, $Y \in \mathbb{R}$ is an outcome and $X \in [0, 1]$ is a covariate. We want to study the treatment effect of A on Y . One common way of quantifying this through the "variance weighted treatment effect"

$$\tau = \frac{\mathbb{E}(\text{Var}(A|X) c(X))}{\mathbb{E}(\text{Var}(A|X))}$$

$$\text{where } c(X) = \mathbb{E}(Y|A=1, X) - \mathbb{E}(Y|A=0, X)$$

is treatment effect in the stratum of X under no ~~measure~~ unmeasured confounding. While considering fixed treatment effect $c(x) = \varphi^*$ for all x , $\tau = \varphi^*$, justifying its value as a marker for treatment effect. τ also φ^* under a commonly studied model in this area: $\mathbb{E}(Y|A, X) = \varphi^* A + b(X)$.

How to estimate τ ?

$$\text{Note that } \tau = \frac{\mathbb{E}(\text{Var}(A|X) c(X))}{\mathbb{E}(\text{Var}(A|X))} = \frac{\mathbb{E} \text{Cov}(Y, A|X)}{\mathbb{E} \text{Var}(A|X)}$$

To see this note that

$$\begin{aligned} \text{Cov}(Y, A|X) &= \mathbb{E}(AY|X) - \mathbb{E}(A|X) \mathbb{E}(Y|X) \\ &= (\mathbb{E}(Y|A=1, X) - \mathbb{E}(Y|X)) \mathbb{E}(A|X) \end{aligned}$$

Similarly,

$$\begin{aligned} \text{Cov}(Y, 1-A|X) &= (\mathbb{E}(Y|A=0, X) - \mathbb{E}(Y|X)) \mathbb{E}(1-A|X) \\ &= (\mathbb{E}(Y|A=1, X) - c(X) - \mathbb{E}(Y|X)) (1 - \mathbb{E}(A|X)) \end{aligned}$$

$$= \frac{\text{cov}(Y, A|X) (1 - E(A|X)) - c(X) (1 - E(A|X))}{E(A|X)}$$

$$= \frac{\text{cov}(Y, A|X) (1 - E(A|X)) - c(X) \text{var}(A|X)}{E(A|X)}$$

$$\Rightarrow \text{cov}(Y, 1-A|X) E(A|X) = \text{cov}(Y, A|X) (1 - E(A|X)) - c(X) \text{var}(A|X)$$

$$\Rightarrow -\text{cov}(Y, A|X) E(A|X) = \text{cov}(Y, A|X) (1 - E(A|X)) - c(X) \text{var}(A|X)$$

$$\Rightarrow \text{cov}(Y, A|X) = c(X) \text{var}(A|X)$$

therefore its enough to understand estimation of $E \text{cov}(Y, A|X)$. (the estimation of $E \text{var}(A|X)$ is similar with "Y=A")

$$\text{But } E \text{cov}(Y, A|X) = E(E(AY|X) - E(A|X) E(Y|X)) = E(AY) - E(E(A|X) E(Y|X))$$

now $E(AY)$ can be efficiently estimated by $\frac{1}{n} \sum_{i=1}^n A_i Y_i$.

We therefore focus on the estimation of

$$E(E(A|X) E(Y|X)) = \int a(x) b(x) g(x) dx$$

$$\text{where } a(x) = E(A|X=x)$$

$$b(x) = E(Y|X=x)$$

$g(x) \leftarrow$ prob. density of X on $[0,1]$ w.r.t. λ .

An influence function for $\Psi(P_\eta) = \chi(\eta)$ where $P_\eta = P(a, b, c, g)$ and $\chi(\eta) = \int a b g$ will

operationalize our von-Mises calculus. In particular the likelihood of the model indexed by (a, b, c, g) is given by (at a single obs. (A, Y, Z)) if $Y \in \{0,1\}$

$$a(x)^A (1-a(x))^{1-A} (c(1-a)+b)(x)^{AY} (1-c(1-a)-b)(x)^{(1-Y)A} \times (-ca+b)(x)^{Y(1-A)} (1+ca-b)(x)^{(1-Y)(1-A)} = f_\eta(Y, A, X) \text{ (say)}$$

Then an influence function of $\Psi(P_\eta) = x(\eta)$ at P_η is given by $(Y - b(x))(A - a(x)) - x(\eta)$, w.r.t. to the tangent space given by perturbations in given directions. Namely, we consider submodels (one-dimensional) in given directions a', b', c' , and g' as follows:

$$a_t = a + t a'$$

$$b_t = b + t b'$$

$$c_t = c + t c'$$

$$g_t = g(1 + t g'), \quad \int g' g \neq 0$$

Then one can calculate the corresponding likelihood $p_{\eta_t} = p(a_t, b_t, c_t, g_t)$ and the corresponding score in the given direction (a', b', c', g') is calculated by $\frac{d}{dt} \log p_{\eta_t} \Big|_{t=0} := \delta_\eta(a', b', c', g')$ (say). For any such score one can verify by tedious but straight forward calculations that

$$\begin{aligned} \frac{d}{dt} x(\eta_t) \Big|_{t=0} &= \frac{d}{dt} \int a_t b_t g_t \Big|_{t=0} \\ &= \mathbb{E}_\eta \left(\left\{ (Y - b(x))(A - a(x)) - x(\eta) \right\} \delta_\eta(a', b', c', g') \right) \end{aligned}$$

This then implies that $(Y - b(x))(A - a(x)) - x(\eta)$ is an influence function at P_η w.r.t. the tangent space obtained above.

Remark: Indeed to justify that δ_η is a score and $(Y - b(x))(A - a(x)) - x(\eta)$ is a valid influence function one needs to be more rigorous. However, we will see that this heuristic and morally correct calculation suffices for us since all we care about is obtaining an estimator of $x(\eta)$ and correct its bias suitably. However for systematic handling of the procedure one definitely needs to pin down the gory details.

Now one can check by tedious but straightforward calculations that

$$\frac{d}{dt} x(\eta_t) \Big|_{t=0} = \frac{d}{dt} x((a_t, b_t, c_t, g_t)) \Big|_{t=0}$$

$$= \frac{d}{dt} \int a_t b_t g_t \Big|_{t=0} = \mathbb{E}_\eta \left(\left\{ \begin{array}{l} (Y-b(x))(A-ax) \\ -x(\eta) \end{array} \right\} B_\eta(a', b', c', g') \right)$$

Therefore $(Y-b(x))(A-ax) - x(\eta)$ is indeed a ~~sub~~ w.r.t. the tangent space generated by these submodels.

Therefore our estimator in first order is

$$\hat{x}_{n,1} = x(\hat{\eta}) + \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{b}(X_i)) (A_i - \hat{a}(X_i)) - x(\hat{\eta})$$

where $x(\hat{\eta}) = \int \hat{a} \hat{b} \hat{g}$ and $\hat{a}, \hat{b}, \hat{g}$ are estimated from X_{n+1}, \dots, X_{2n} .

$$\text{Therefore } \hat{x}_{n,1} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{b}(X_i)) (A_i - \hat{a}(X_i))$$

lets understand the bias of this estimator. By direct calculations,

$$\mathbb{E}_{\eta,1}(\hat{x}_{n,1}) - x(\eta) = \int (\hat{a} - a)(\hat{b} - b)g$$

where $\mathbb{E}_{\eta,1}$ denotes expectation with second half of the sample kept fixed.

Our second order von-Mises

expansion aims to correct for this bias. Indeed if \exists a kernel $K(x,y)$

s.t. $\int h(x) K(x,y) dy = h(y)$ a.e. y and all $h \in L_2[0,1]$ we

can correct this bias by considering

$$- \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \frac{(Y_{i_1} - \hat{b}(X_{i_1}))}{\sqrt{\hat{g}(X_{i_1})}} K(X_{i_1}, X_{i_2}) \frac{(A_{i_2} - \hat{a}(X_{i_2}))}{\sqrt{\hat{g}(X_{i_2})}}$$

because this has expectation $-\int (b - \hat{b})(a - \hat{a})g$.

But \nexists such a kernel with as an element of

$L_2([0,1] \times [0,1])$, so we work with partial representation.

That is we correct our bias using projection kernels onto finite but growing with n dimensional spaces. In particular, recall our definition of $V_j = \{ \phi_{j,k}, \psi_{j,k} | z \in \mathcal{I}_0, k \in \mathcal{Z}_j \}$ to be the linear span of the ~~first few~~ ^{first few} basis functions from an MRA generated by (ϕ, ψ) . Then our bias corrected ~~the~~ second order estimator becomes:

$$\hat{\chi}_{n,2} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{b}(X_i)) (A_i - \hat{a}(X_i)) - \frac{1}{n(n-1)} \sum_{i \neq i_2} \frac{(Y_{i_1} - \hat{b}(X_{i_1}))}{\hat{g}(X_{i_1})^{\frac{1}{2}}} K_{V_j}(X_{i_1}, X_{i_2}) \frac{(A_{i_2} - \hat{a}(X_{i_2}))}{\hat{g}(X_{i_2})^{\frac{1}{2}}}$$

We next show that $\hat{\chi}_{n,2}$ has desirable bias and variance properties.

$$\text{Bias: } \left| \mathbb{E}_{\eta}(\hat{\chi}_{n,2}) - \chi(\eta) \right|$$

$$= \left| \mathbb{E}_{\eta,2} \int (b - \hat{b})(x) (a - \hat{a})(x) g(x) - \mathbb{E}_{\eta} \left\{ \frac{(Y_1 - \hat{b}(X_1))}{\hat{g}(X_1)^{\frac{1}{2}}} K_{V_j}(X_1, X_2) \frac{(A_2 - \hat{a}(X_2))}{\hat{g}(X_2)^{\frac{1}{2}}} \right\} \right|$$

For notational brevity, put $\delta b = b - \hat{b}$,
 $\delta a = a - \hat{a}$,
 $\delta g = g - \hat{g}$.

$$\Rightarrow \mathbb{E}_{\eta} \left(\frac{(Y_1 - \hat{b}(X_1))}{\hat{g}(X_1)^{\frac{1}{2}}} K_{V_j}(X_1, X_2) \frac{(A_2 - \hat{a}(X_2))}{\hat{g}(X_2)^{\frac{1}{2}}} \right)$$

$$= \mathbb{E}_{\eta,2} \left\{ \iint \frac{\delta b(x_1) g(x_1)}{\sqrt{\hat{g}(x_1)}} K_{V_j}(x_1, x_2) \frac{\delta a(x_2) g(x_2)}{\sqrt{\hat{g}(x_2)}} dx_1 dx_2 \right\}$$

$$= \mathbb{E}_{\eta,2} \int \frac{\delta b(x_1) g(x_1)}{\sqrt{\hat{g}(x_1)}} \Pi_{V_j} \left(\frac{\delta a g}{\sqrt{\hat{g}}} \right) (x_1) dx_1$$

$$= \mathbb{E}_{\eta,2} \int \frac{\delta b(x_1) \delta a(x_2) g(x_1) g(x_2)}{\sqrt{\hat{g}(x_1)} \sqrt{\hat{g}(x_2)}} dx_1 - \mathbb{E}_{\eta,2} \int \frac{\delta b(x_1) g(x_1)}{\sqrt{\hat{g}(x_1)}} \Pi_{V_j} \left(\frac{\delta a g}{\sqrt{\hat{g}}} \right) (x_1) dx_1$$

$$= \mathbb{E}_{\eta,2} \int \delta b(x_1) \delta a(x_2) g(x_1) dx_1 + \mathbb{E}_{\eta,2} \int \delta b(x_1) \delta a(x_1) g^2(x_1) \left(\frac{1}{\hat{g}(x_1)} - \frac{1}{g(x_1)} \right) dx_1$$

$$- \mathbb{E}_{\eta,2} \int \frac{\delta b(x_1) g(x_1)}{\sqrt{\hat{g}(x_1)}} \Pi_{V_j} \left(\frac{\delta a g}{\sqrt{\hat{g}}} \right) (x_1) dx_1$$

$$= \mathbb{E}_{\eta, 2} \int \delta b \delta a g + \mathbb{E}_{\eta, 2} \int \delta b \delta a \delta g \frac{g}{\hat{g}}$$

$$- \mathbb{E}_{\eta, 2} \int \Pi_{V_j^\perp} \left(\frac{\delta b g}{\sqrt{\hat{g}}} \right) \Pi_{V_j^\perp} \left(\frac{\delta a g}{\sqrt{\hat{g}}} \right)$$

$$\therefore \text{Bias: } \left| \mathbb{E}_{\eta, 2} \int \delta b \delta a \delta g \frac{g}{\hat{g}} - \mathbb{E}_{\eta, 2} \int \Pi_{V_j^\perp} \left(\frac{\delta b g}{\sqrt{\hat{g}}} \right) \Pi_{V_j^\perp} \left(\frac{\delta a g}{\sqrt{\hat{g}}} \right) \right|$$

$$\leq \mathbb{E}_{\eta, 2} \left(\left\| \frac{g}{\hat{g}} \right\|_\infty \|\delta g\|_2 \|\delta a\|_4 \|\delta b\|_4 \right)$$

$$- \mathbb{E}_{\eta, 2} \left\| \Pi_{V_j^\perp} \left(\frac{\delta b g}{\sqrt{\hat{g}}} \right) \right\|_2 \left\| \Pi_{V_j^\perp} \left(\frac{\delta a g}{\sqrt{\hat{g}}} \right) \right\|_2$$

Now its time to assume some regularity properties on a, b, g . In particular assume $a \in H(\alpha, M)$, $b \in H(\beta, M)$, and $g \in H(\gamma, M)$ on $[0, 1]$ where $\gamma > 2 \max\{\alpha, \beta\}$. ~~and $\gamma > 2 \max\{\alpha, \beta\}$~~

Also assume that $\|a\|_\infty \vee \|b\|_\infty \vee \|g\|_\infty \leq B_U$ and $g(x) \geq B_0 > 0$ for all x . Then we have estimators $\hat{a}, \hat{b}, \hat{g}$ of a, b, g s.t. the following hold:

$$(i) \mathbb{E}_\eta \|\delta a\|_\infty \leq c^* \left(\frac{n}{\log n} \right)^{-\frac{\alpha}{2\alpha+1}}, \quad \mathbb{E}_\eta \|\delta b\|_\infty \leq c^* \left(\frac{n}{\log n} \right)^{-\frac{\beta}{2\beta+1}},$$

$$\mathbb{E}_\eta (\|\delta g\|_\alpha) \leq c^* \left(\frac{n}{\log n} \right)^{-\frac{\gamma}{2\gamma+1}} \quad \forall \eta \in \mathcal{P}(\alpha, \beta, \gamma) \text{ for } n \geq n^*$$

$$(ii) \frac{B_U}{2} \leq \hat{g}(x) \leq 2B_U \text{ for all } x. \quad \forall \eta \in \mathcal{P}(\alpha, \beta, \gamma), \text{ for } n \geq n^*$$

$$(iii) \sup_{\eta \in \mathcal{P}(\alpha, \beta, \gamma)} \mathbb{P}_\eta (\hat{a} \in H(\alpha, c^{**})) \geq 1 - \frac{1}{n^3}$$

$$\sup_{\eta \in \mathcal{P}(\alpha, \beta, \gamma)} \mathbb{P}_\eta (\hat{b} \in H(\beta, c^{**})) \geq 1 - \frac{1}{n^3} \quad \forall n \geq n^*$$

$$\sup_{\eta \in \mathcal{P}(\alpha, \beta, \gamma)} \mathbb{P}_\eta (\hat{g} \in H(\gamma, c^{**})) \geq 1 - \frac{1}{n^3}$$

where c^*, c^{**}, n^* depends on B_U and M , and

$$\mathcal{P}(\alpha, \beta, \gamma) = \left\{ (a, b, c, g) : \|a\|_\infty \vee \|b\|_\infty \vee \|g\|_\infty \leq B_U, g(x) \geq B_0 \quad \forall x \in [0, 1], \int g = 1 \right\}$$

This implies that,

$$|Bias| \leq c(B_U, B_L, M) \left(\frac{n}{\log n} \right)^{-\left(\frac{\alpha}{2\alpha+1} + \frac{\beta}{2\beta+1} + \frac{\nu}{2\nu+1}\right)} + c(B_U, B_L, M) \left\{ 2^{-2j\alpha + \frac{1}{n^\beta}} \right\}^{\frac{1}{2}}$$

$$\leq c(B_U, B_L, M) \left\{ n^{-\frac{2\alpha+2\beta}{2\alpha+2\beta+1}} + \left(2^{-2j\alpha + \frac{1}{n^\beta}} \right)^{\frac{1}{2}} \left(2^{-2j\beta + \frac{1}{n^\beta}} \right)^{\frac{1}{2}} \right\} \left\{ 2^{-2j\beta + \frac{1}{n^\beta}} \right\}^{\frac{1}{2}}$$

(since $\nu > 2(\alpha + \beta)$)

□

Now we go the variance calculation.

We additionally assumed $\|a\|_\infty, \|b\|_\infty \leq B_U$.

This implies that we can add an extra property to our list (i) - (iii) on page (21) and obtain estimator \hat{a}, \hat{b} which satisfy

$$\sup_{\eta \in \mathcal{V}(\alpha, \beta, \nu)} \mathbb{P}_\eta \left(\max(\|\hat{a}\|_\infty, \|\hat{b}\|_\infty) \leq 2B_U \right) \geq 1 - \frac{1}{n^3} \quad \forall n \geq n^*$$

(provided we have (φ, ψ) with sufficient regularity & compact support)

If finally we assume $|Y| \leq B_U$ a.s.

we can do our variance calculation by standard ~~trick~~ Hoeffding's decomposition and get

$$\text{Var}(\hat{x}_{n,2}) \leq c(B_U) \left\{ \frac{1}{n} + \frac{2^j}{n^2} \right\} \quad (\text{check!})$$

The right bias variance is then obtained by $2^j \sim n$ if $\frac{\alpha + \beta}{2} \geq \frac{1}{4}$ and $2^j \sim n^{\frac{2}{2\alpha+2\beta+1}}$ if $\frac{\alpha + \beta}{2} < \frac{1}{4}$

This gives us the following result.

Theorem: Assume $|Y| \leq B_U$ a.s. and $|A| \leq B_U$ a.s. Then \exists an estimator \hat{x}_n of $\psi(P_\eta) = \mathbb{E}(w(A, Y | X))$ s.t. if $\nu > 2 \max\{\alpha, \beta\}$ and $n \geq n^*$ (depending on M, B_U, B_L)

$$\sup_{\eta \in \mathcal{V}(\alpha, \beta, \nu)} \mathbb{E}_\eta \left(\hat{x}_n - \psi(P_\eta) \right)^2 \leq c(B_U, B_L, M) n^{-\frac{4\alpha+4\beta}{2\alpha+2\beta+1}}$$

$$= c(B_U, B_L, M) n^{-\frac{8\delta_{\text{var}}}{4\delta_{\text{var}}+1}}$$

where $\delta_{\text{var}} = \frac{\alpha + \beta}{2} \leq \frac{1}{4}$

$$= c(B_U, B_L, M) \frac{1}{n^2} \quad \delta_{\text{var}} \geq \frac{1}{4}$$