

9/8/2011 * we want estimators w/ bias $\leq \frac{1}{n}$

For a rv $O \sim P_0 \in \mathcal{M}$, we wish to learn about $\psi: \mathcal{M} \rightarrow \mathbb{R}$

$\psi_0 =$ true value of param of interest

param is pathwise differentiable at $P \Rightarrow$



Action

A.1

Cramer-Rao

we want it differentiable at every P
 ↳ needs to be within model

For class of 1-dim submodel (generates rich enough space)

$$\left\{ \begin{array}{l} P(\epsilon) : \epsilon \in (-\delta, \delta) \\ P(\epsilon=0) = P \end{array} \right\} \subset \mathcal{M}$$

* Want to choose as rich a class of submodels as possible

Pathwise derivative : $\frac{\partial}{\partial \epsilon} \psi(P(\epsilon))|_{\epsilon=0} = \lim_{\epsilon \rightarrow 0} \frac{\psi(P(\epsilon)) - \psi(P)}{\epsilon}$ (we want this to exist for every path)

$$S(O) = \text{Score of param model} = \frac{\partial \psi(P(\epsilon))}{\partial \epsilon}|_{\epsilon=0}$$

↓ if exists
 $= E_P [D(P)(O) S(O)]$ for some $D(P)$

Hilbert spaces $L^2(P)$, $\langle h_1, h_2 \rangle = E_P [h_1(O) h_2(O)]$
 linear space $= \text{Cov}_P(h_1, h_2)$

we can define norm $\|h\| = \sqrt{\langle h, h \rangle}$

$\frac{\partial}{\partial \epsilon} \psi(P(\epsilon))$ can be thought of as mapping $L^2(P)$ to real space \mathbb{R}

$$A: L^2(P) \rightarrow \mathbb{R}$$

Hilbert space

A is linear bounded ($A(h_1, h_2) = A(h_1) + A(h_2)$)

$$\Rightarrow A(h) = \langle O, h \rangle = E_P [D(O), h(O)]$$

aka Riesz-Figyel theorem

Ex. Any $\mathcal{M} = NP$ model

$O =$ real valued

$$\psi(P) = \int O P(O) dM(O)$$

$$P(\epsilon) = (1+\epsilon h)P, \quad \left\{ \begin{array}{l} E_P h(O) = 0 \\ \|h\|_\infty < \infty \end{array} \right.$$

↑
 is submodel

↗ cont.

$$\text{cont. } \frac{\partial}{\partial \epsilon} \psi(P(\epsilon)) = \frac{\partial}{\partial \epsilon} \int O(1 + \epsilon h(O)) P(O) d\mu(O)$$

$$= \int O \cdot h(O) P(O) \mu(O)$$

$$= \langle f, h \rangle$$

\hookrightarrow needs to have mean 0

$$= \langle f - E_P(f), h \rangle$$

$\Delta(P)$

\Rightarrow inner product of $\Delta(P)$ and h (the score)

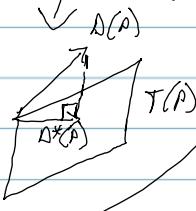
Note: $\Delta(P) \in L^2(P)$

Hilbert space

Let $T(P)$ be the tangent space of this class of submodels through P
 = closure of the linear span of all scores

$$\frac{\partial}{\partial \epsilon} \psi(P(\epsilon))|_{\epsilon=0} = E_P(\Delta^*(P)(O) S(O)) \text{ for a unique gradient } \Delta^* \in T(P)$$

Note: if gradient Δ , $\pi[\Delta|T(P)] = \Delta^*$
 i.e. projection of Δ onto a tangent space.



Def'n Projection: $\pi[\Delta|T(P)]$ is defined to be
 ① $\pi[\Delta|T(P)] \in T(P)$
 ② $\Delta - \pi[\Delta|T(P)] \perp T(P)$

$$\text{i.e. } \langle \Delta - \pi[\Delta|T(P)], h \rangle = 0 \quad \forall h \in T(P)$$

Cx on $\Delta = (w, A, y)$ projection

$$E[h(w, a, y)|w] = \pi(h|T_w)$$

$$\{h(w) \cdot h \in L^2(P)\}$$

$$\langle h(w, a, y) - E(h|w), h(w) \rangle$$

$$= E[h(w, a, y) - E(h|w), h(w)] = 0$$

$$= E[h(w, a, y)h(w)] - E[E(h|w)h(w)] = 0$$

* the efficient influence curve is not influenced in any way by placing restrictions on global space.

Generalized Cramer-Rao lower bound

Assume $O \sim P(c_0) \in \{P(\epsilon) : \epsilon\}$ (Note: $c_0 = 0$, don't know this yet)

Target param: $\psi(P(c_0))$

$$\text{C.Rao: } \left(\frac{\partial}{\partial c_0} \psi(P(c_0)) \right)^2 = \frac{[E_P(\Delta^*(P)(O) S(O))]^2}{E_P[S^2(O)]}$$

$$I(\epsilon_0) = E[S^2(O)]$$

Var Score

$$\text{Generalized Cr-R lower bound} \cdot \sup_{S \in T(P)} \frac{(E_P D^*(P) S(O))^2}{E_P S^2(O)} \leq \underbrace{\sup_{S \in T(P)} \frac{E(D^{**}) \cdot E(S^2)}{E(S^{**})}}_{= E(D^*(P))^2(\alpha)} = \boxed{\text{Var}(D^*)}$$

\hookrightarrow least favorable submodel

Cauchy-Schwarz Inequality: $E_P [D^* S] \leq \sqrt{E(D^{**})} \sqrt{E(S^2)}$
 i.e. $\langle D^*, S \rangle \leq \|D\| \cdot \|S\|$
 \approx inner product.

Note: can choose score
 $S(O) = D^*(P)$
 i.e. equal to canonical gradient

Example:

$$O = (w, A, Y) \sim P_0 \in \mathcal{M} = \text{NP model}$$

$$\psi(P) = E_P E_P [Y | A=1, w] = \sum_w \sum_y y P(Y | A=1, w) P(w)$$

$$\text{Recall: } P(w, a, y) = P(w) P(a|w) P(y|a, w)$$

For pathwise deriv.: need: $P_c(w) = (1 + c S_w(w)) P(w)$, $S_w(w)$ has mean 0

$$\sum_w (S_w(w)) P(w) = 0$$

$$P_c(a|w) = (1 + c S_w(a|w)) P(a|w), E(S_A(A|w)) = 0$$

$$P_c(y|a, w) = (1 + c S_y(y|a, w)) P(y|a, w), E(S_Y(Y|A, w)) = 0$$

$$\downarrow \frac{\partial}{\partial c} \psi /$$

fluctuations