

Didn't start material until 8 men.

\* Majority of time, Mack just presented foundations

9/5/2011

at the end of the day, you'll have to decide on a target  
param + estimator  
to our  $\Psi$  ( $\hat{\Psi}$ )

So, let  $P \sim P_0 \in \mathcal{M}$

↳ we observe  $O_1, \dots, O_n$  iid from  $P_0$   
our  $\Psi: \mathcal{M} \rightarrow \mathbb{R}^d$

estimator ↗ stat dist mapping  
 $\hat{\Psi}_0 = \Psi(P_0) = \text{TMLE param value}$

Estimator = mapping of data to our param space

↳  $\hat{\Psi}$  aka an algorithm

↳ We have  $P_n =$  the empirical prob dist of  $O_1, \dots, O_n$ , w/  $P_n(O=o_i) = \frac{1}{n}, i=1, \dots, n$   
so for  $P_n$ , we want  $\Psi(P_n)$   
\* if we have randomness in our estimator, recommended to avg across runs  
for  $\hat{\Psi}: \mathcal{M}_{wp} \rightarrow \mathbb{R}^d$   
↑ non-parametric

Example:  $\hat{\Psi}(P) = E_P [E_P[Y|A=1, w] - E_P[Y|A=0, w]]$

↑ think he meant  $E_w$

Note:  $\hat{\Psi}(P_n)$  is a random variable

\* A property we'd like is  $\hat{\Psi}(P_0) = \Psi_0$

↳ Consistency:  $\forall c > 0, P(|\hat{\Psi}(P_n) - \Psi_0| > c) \xrightarrow{n \rightarrow \infty} 0$

↳ Doesn't tell us how fast we're learning though (i.e. how fast approaching)

\* We need Asymp. Linearity:  $\hat{\Psi}(P_n) - \Psi_0 = \underbrace{\frac{1}{n} \sum_i I(O_i)}_{\text{behaves as}} \underbrace{(\hat{\Psi}(P_0))(O_i)}_{\text{rv, mult w/ } \frac{1}{n}}$

empirical mean

still goes to 0 in prob  
↳ cont.

unknown  
 $\downarrow$

cont.  $IC(P_0)$  is a function of  $\theta$ , w/ mean  $\theta$  and finite variance

$\hookrightarrow$  It is called the influence curve of the  $\hat{\psi}$  at  $P_0$

Note: There are params where you won't get  $\hat{\psi}(P_n)$  that are ~~steplike~~ linear

\* So for  $n$  large enough,  $\hat{\psi}(P_n) - \psi_0$  will behave like an empirical mean  $\bar{O}$  and variance of the IC

$$\Rightarrow \sqrt{n}(\hat{\psi}(P_n) - \psi_0) \xrightarrow{D} N(0, \sigma^2 = \text{var}(IC(P_0)(\bar{O}))$$

$$\downarrow \Rightarrow Z_n \xrightarrow{D} Z \equiv P(Z_n \leq x) \xrightarrow{n \rightarrow \infty} P(Z \leq x) \quad \forall x$$

$$\text{if dealing w/ vector, have } \sqrt{n}(\hat{\psi}(P_n) - \psi_0) \xrightarrow{D} N(0, \Sigma^2 = E_0(IC(\bar{O})IC(\bar{O})^T))$$

This allows us to determine  $P(-1.96 < \underbrace{\frac{\sqrt{n}(\hat{\psi}(P_n) - \psi_0)}{\Sigma}}_{\sim Z} < 1.96) \xrightarrow{n \rightarrow \infty} 0.95$  (i.e. 95% CI)

$$P(\hat{\psi}_0 - 1.96 \frac{\Sigma}{\sqrt{n}} < \psi_0 < \hat{\psi}_0 + 1.96 \frac{\Sigma}{\sqrt{n}}) \xrightarrow{n \rightarrow \infty} 0.95$$

\* To understand an estimator, will need to calc. Influence Curve.

$\hookrightarrow$  can estimate it w/  $\hat{IC}$

$$\hookrightarrow \text{gives us var-covar matrix est of } \hat{\Sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{IC}(\bar{O}_i) \hat{IC}(\bar{O}_i)^T$$

$\Rightarrow$  To learn about IC's and how to calculate them,  
 read all 3 Appendix A's.

$$P_n f = \int f(\bar{O}) dP_n(\bar{O}) = \frac{1}{n} \sum_i f(O_i)$$

$$\hookrightarrow \hat{\psi}(P_n) = \hat{\psi}(P_n f : f \in \mathcal{F})$$

Ex. Take correlation:

We only need to estimate 5 values from the dist to est this:

$$E(X), E(Y), E(X^2), E(Y^2), \text{cov}(X, Y) \quad \Rightarrow \text{corr}(X, Y) = f(E(X), E(Y), \dots)$$

$\hookrightarrow$  The more adaptive the estimator, the more it depends on empirical data

$\mathcal{F}$  can be a big family of functions of  $\bar{O}$

$$\text{eg. } \mathcal{F} = \{I(O_i = o) : o\}$$

Suppose our target param  $\Psi(P) = E_p [E_w(Y|A=1, w) - E_w(Y|A=0, w)]$   
 $\hookrightarrow$  says we need all empirical probabilities

For what class of func would: (for probability to consider)

$$\sup_{f \in \mathcal{F}} |(P_n - P_0)f| \xrightarrow{n \rightarrow \infty} 0$$

$\hookrightarrow$  cont.

cont.  $G_n = (\sqrt{n}(P_n - P_0) f \cdot f)$

$\forall$  random func b/c for a given  $P_n$ , can calc for all values of  $f$ .  
 $\Rightarrow$  Gaussian Process  $(G) = (G_f : f \in F)$   
 true when certain entropy conditions hold.

Note:  $G_n = \text{Sluens class}$

$F = \text{Donsker class}$

\* listen to Mark prove consistency again: used  $\hat{\Psi}(P_n) = \hat{\Psi}(A_f : f \in F)$   
 $\hookrightarrow$  should look at section A-1

$\forall$   $\hat{\Psi}$  is differentiable, then we have an inheritance

$$\hat{\Psi}(P_n) - \hat{\Psi}(P_0) \approx \hat{\Psi}'(P_0)(P_n - P_0)$$

vectors

$$\text{recall: } f(x_n) - f(x) \approx f'(x)(x_n - x)$$

that's what we do above

only instead of one var, we apply over whole vector

$$= \sum_{f \in F} \frac{d}{dP_0 f} \hat{\Psi}(P_0, f : f \in F) \cdot (P_n - P_0)(f)$$

$$= \frac{1}{n} \sum_{i=1}^n [f(O_i) - P_0 f] \text{ ie the emp. mean}$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{f \in F} \frac{d}{dP_0 f} \hat{\Psi}(P_0) [f(O_i) - P_0 f]$$

$$\text{The influence curve: } IC(P_0)(O_i)$$

$$\hookrightarrow \text{thus, the } IC(P_0)(O_i) = \hat{\Psi}'(P_0)(f(O_i) - P_0 f : f \in F)$$

\* can be written as sum of partial de

$$\text{Note: } \psi'(P_0)(h) = \frac{d}{d\varepsilon} \hat{\Psi}(P_0 + \varepsilon h) |_{\varepsilon=0}$$

$h$  being the  $f(O_i) - P_0 f : f \in F$  vector  
 Gateau derivative (aka directional derivative)

\* If we can apply certain chars to our estimator,  
 then our estimator will be asymptotic linear

Good Properties of IC

1) CI's can be computed

2) Robust stats

3) Allows us to see what behaviours are needed for good estimators

✓

4) Suppose we have two estimators that are both asymptotic linear

$\hookrightarrow$  the one w/ the lower var(IC) is the more efficient one.

$\hookrightarrow$  so IC allows us to compare estimators

Note: We can calc. efficient IC by taking derivative of  $\Psi$