

Influence functions and functional derivatives

Anders Munch

May 11, 2021

Outline

Setting

Motivation

Functional derivatives

Canonical gradient / efficient influence function

Summary of main results

Next step – constructing estimator

References

Disclaimer

- ▶ The note is work in progress, and we have not used it before – you are very welcome to comment on weird/unclear passages.

Disclaimer

- ▶ The note is work in progress, and we have not used it before – you are very welcome to comment on weird/unclear passages.
- ▶ You should see it as a service – some exact mathematical statements are collected there if you care about it, but the important part is the intuition which we talk about today.

Disclaimer

- ▶ The note is work in progress, and we have not used it before – you are very welcome to comment on weird/unclear passages.
- ▶ You should see it as a service – some exact mathematical statements are collected there if you care about it, but the important part is the intuition which we talk about today.
- ▶ Do NOT write like this in your report!

A statistical problem

We call a collection of probability measures \mathcal{P} together with a functional $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ a *statistical problem*.

A statistical problem

We call a collection of probability measures \mathcal{P} together with a functional $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ a *statistical problem*.

Example (Average treatment effect)

We are given n iid. sample of $O \sim P$, with $P \in \mathcal{P}$ and where $O = (X, A, Y)$, with $X \in \mathbb{R}^d$, $A \in \{0, 1\}$, and $Y \in \{0, 1\}$. We want to estimate the average treatment effect

$$\mathbb{E}_P [f(1, X) - f(0, X)],$$

with $f(a, x) := \mathbb{E}_P [Y \mid A = a, X = x]$. The target parameter is

$$\Psi(P) = \mathbb{E}_P [f_P(1, X) - f_P(0, X)].$$

A statistical problem

We call a collection of probability measures \mathcal{P} together with a functional $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ a *statistical problem*.

Example (Average treatment effect)

We are given n iid. sample of $O \sim P$, with $P \in \mathcal{P}$ and where $O = (X, A, Y)$, with $X \in \mathbb{R}^d$, $A \in \{0, 1\}$, and $Y \in \{0, 1\}$. We want to estimate the average treatment effect

$$\mathbb{E}_P [f(1, X) - f(0, X)],$$

with $f(a, x) := \mathbb{E}_P [Y \mid A = a, X = x]$. The target parameter is

$$\Psi(P) = \mathbb{E}_P [f_P(1, X) - f_P(0, X)].$$

Target and nuisance parameters

Target parameter Low-dimensional, scientifically meaningful.

Target and nuisance parameters

Target parameter Low-dimensional, scientifically meaningful.

Nuisance parameters Needed to express the target parameter.

Target and nuisance parameters

Target parameter Low-dimensional, scientifically meaningful.

Nuisance parameters Needed to express the target parameter.

Example (ATE)

The ATE can be written as $\Psi(P) = P[\varphi_1] = P[\varphi_2] = P[\varphi_3]$, for

$$\varphi_1(o; f) := f(1, x) - f(0, x),$$

$$\varphi_2(o; \pi) := \frac{a y}{\pi(x)} - \frac{(1-a) y}{1 - \pi(x)},$$

$$\varphi_3(o; f, \pi) := \varphi_1(o; f) + \varphi_2(o; \pi) - \frac{a f(1, x)}{\pi(x)} + \frac{(1-a) f(0, x)}{1 - \pi(x)}$$

$P[\varphi]$ means

$$P[\varphi] = \mathbb{E}_P [\varphi(O)] = \int \varphi(o) dP(o).$$

High-/infinite-dimensional nuisance parameters

A parametric setting means that \mathcal{P} is finite-dimensional. We are interested in *nonparametric* or *semiparametric* settings which mean that \mathcal{P} is infinite-dimensional.

High-/infinite-dimensional nuisance parameters

A parametric setting means that \mathcal{P} is finite-dimensional. We are interested in *nonparametric* or *semiparametric* settings which mean that \mathcal{P} is infinite-dimensional.

Having our data set and scientific question in mind, why would it be of interest to use infinite-dimensional nuisance parameters?

High-/infinite-dimensional nuisance parameters

A parametric setting means that \mathcal{P} is finite-dimensional. We are interested in *nonparametric* or *semiparametric* settings which mean that \mathcal{P} is infinite-dimensional.

Having our data set and scientific question in mind, why would it be of interest to use infinite-dimensional nuisance parameters?

Trying to control for confounding \implies nice to have:

- ▶ flexible model
- ▶ many covariates

Toy example: Integrated kernel density

\mathcal{P} consist all probability measures with continuous Lebesgue-density (this is an infinite-dimensional space). We want to estimate $F(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$.

Toy example: Integrated kernel density

\mathcal{P} consist all probability measures with continuous Lebesgue-density (this is an infinite-dimensional space). We want to estimate $F(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$. Our target parameter is then $\theta = \Psi(P) = F_P(x)$ which we can express as

$$\Psi(P) = \Psi_0(f) := \int_{-\infty}^x f(z) \, dz, \quad \text{for } P = f \cdot \lambda,$$

because of our assumption about \mathcal{P} .

Toy example: Integrated kernel density

\mathcal{P} consist all probability measures with continuous Lebesgue-density (this is an infinite-dimensional space). We want to estimate $F(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$. Our target parameter is then $\theta = \Psi(P) = F_P(x)$ which we can express as

$$\Psi(P) = \Psi_0(f) := \int_{-\infty}^x f(z) \, dz, \quad \text{for } P = f \cdot \lambda,$$

because of our assumption about \mathcal{P} . We want to use **machine learning** (!) for this problem, so use a kernel estimator, i.e.,

$$\hat{f}_n(x) = \hat{\mathbb{P}}_n[k_h(X, x)] = \frac{1}{n} \sum_{i=1}^n k_h(X_i, x),$$

where k_h is, e.g. $k_h(x, y) = g\left(\frac{x-y}{h}\right)$, with g the density for the standard Gaussian distribution, and the bandwidth h is chosen using cross-validation.

Toy example: Integrated kernel density

\mathcal{P} consist all probability measures with continuous Lebesgue-density (this is an infinite-dimensional space). We want to estimate $F(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$. Our target parameter is then $\theta = \Psi(P) = F_P(x)$ which we can express as

$$\Psi(P) = \Psi_0(f) := \int_{-\infty}^x f(z) \, dz, \quad \text{for } P = f \cdot \lambda,$$

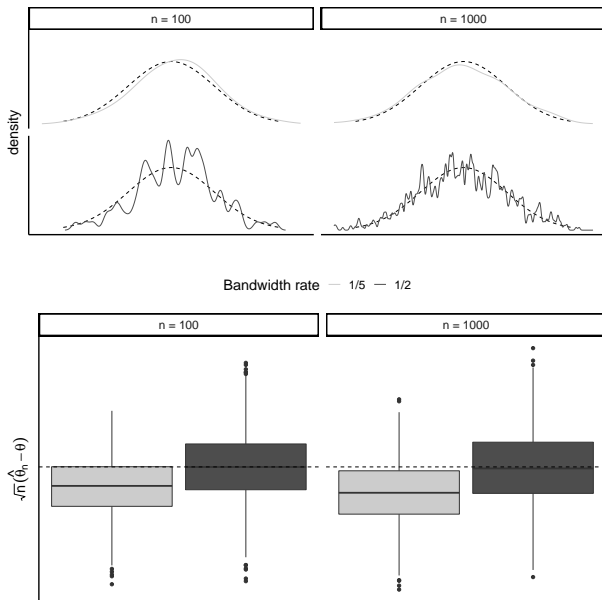
because of our assumption about \mathcal{P} . We want to use **machine learning** (!) for this problem, so use a kernel estimator, i.e.,

$$\hat{f}_n(x) = \hat{\mathbb{P}}_n[k_h(X, x)] = \frac{1}{n} \sum_{i=1}^n k_h(X_i, x),$$

where k_h is, e.g. $k_h(x, y) = g\left(\frac{x-y}{h}\right)$, with g the density for the standard Gaussian distribution, and the bandwidth h is chosen using cross-validation. We then obtain the target estimator $\hat{\theta}_n = \Psi_0(\hat{f}_n)$.

How does this work in practice?

How does this work in practice?



What happened?

What happened?

Consider a general problem (\mathcal{P}, Ψ) for which we can write $\Psi(P) = \Psi_0(P, \nu) = P[\varphi(O, \nu)]$.

What happened?

Consider a general problem (\mathcal{P}, Ψ) for which we can write $\Psi(P) = \Psi_0(P, \nu) = P[\varphi(O, \nu)]$. We have

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \sqrt{n}(\Psi_0(\hat{\mathbb{P}}_n, \hat{\nu}_n) - \Psi_0(P, \nu)) \\ &= \sqrt{n}(\hat{\mathbb{P}}_n[\varphi(O, \hat{\nu}_n)] - P[\varphi(O, \nu)]) \\ &= \sqrt{n}(\hat{\mathbb{P}}_n[\varphi(O, \hat{\nu}_n)] \pm P[\varphi(O, \hat{\nu}_n)] - P[\varphi(O, \nu)]) \\ &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n}\{\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu)\},\end{aligned}$$

with $\mathbb{G}_n := \sqrt{n}(\hat{\mathbb{P}}_n - P)$ the empirical process.

What happened?

Consider a general problem (\mathcal{P}, Ψ) for which we can write $\Psi(P) = \Psi_0(P, \nu) = P[\varphi(O, \nu)]$. We have

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \sqrt{n}(\Psi_0(\hat{\mathbb{P}}_n, \hat{\nu}_n) - \Psi_0(P, \nu)) \\ &= \sqrt{n}(\hat{\mathbb{P}}_n[\varphi(O, \hat{\nu}_n)] - P[\varphi(O, \nu)]) \\ &= \sqrt{n}(\hat{\mathbb{P}}_n[\varphi(O, \hat{\nu}_n)] \pm P[\varphi(O, \hat{\nu}_n)] - P[\varphi(O, \nu)]) \\ &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n}\{\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu)\},\end{aligned}$$

with $\mathbb{G}_n := \sqrt{n}(\hat{\mathbb{P}}_n - P)$ the empirical process.

$\mathbb{G}_n[\varphi(O, \hat{\nu}_n)]$ determines the (main) variance

What happened?

Consider a general problem (\mathcal{P}, Ψ) for which we can write $\Psi(P) = \Psi_0(P, \nu) = P[\varphi(O, \nu)]$. We have

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \sqrt{n}(\Psi_0(\hat{\mathbb{P}}_n, \hat{\nu}_n) - \Psi_0(P, \nu)) \\ &= \sqrt{n}(\hat{\mathbb{P}}_n[\varphi(O, \hat{\nu}_n)] - P[\varphi(O, \nu)]) \\ &= \sqrt{n}(\hat{\mathbb{P}}_n[\varphi(O, \hat{\nu}_n)] \pm P[\varphi(O, \hat{\nu}_n)] - P[\varphi(O, \nu)]) \\ &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n}\{\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu)\},\end{aligned}$$

with $\mathbb{G}_n := \sqrt{n}(\hat{\mathbb{P}}_n - P)$ the empirical process.

$\mathbb{G}_n[\varphi(O, \hat{\nu}_n)]$ determines the (main) variance
 $\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu)$ is bias!

What to do? – Taylor expansion

What to do? – Taylor expansion

Assume we could make a Taylor expansion of $\nu \mapsto \Psi_0(P, \nu)$, so that

$$\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu) = D_\nu \Psi_0[\hat{\nu}_n - \nu] + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|_\nu^2).$$

What to do? – Taylor expansion

Assume we could make a Taylor expansion of $\nu \mapsto \Psi_0(P, \nu)$, so that

$$\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu) = D_\nu \Psi_0[\hat{\nu}_n - \nu] + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|_{\mathcal{V}}^2).$$

The decomposition then becomes

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] \tag{1}$$

$$+ D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] \tag{2}$$

$$+ \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_{\mathcal{V}}^2). \tag{3}$$

What to do? – Taylor expansion

Assume we could make a Taylor expansion of $\nu \mapsto \Psi_0(P, \nu)$, so that

$$\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu) = D_\nu \Psi_0[\hat{\nu}_n - \nu] + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|_{\mathcal{V}}^2).$$

The decomposition then becomes

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] \tag{1}$$

$$+ D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] \tag{2}$$

$$+ \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_{\mathcal{V}}^2). \tag{3}$$

(1) can be handled by empirical process theory or sample splitting

What to do? – Taylor expansion

Assume we could make a Taylor expansion of $\nu \mapsto \Psi_0(P, \nu)$, so that

$$\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu) = D_\nu \Psi_0[\hat{\nu}_n - \nu] + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|_{\mathcal{Y}}^2).$$

The decomposition then becomes

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] \tag{1}$$

$$+ D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] \tag{2}$$

$$+ \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_{\mathcal{Y}}^2). \tag{3}$$

(1) can be handled by empirical process theory or sample splitting

(2) is our focus! \rightarrow make sense of this

What to do? – Taylor expansion

Assume we could make a Taylor expansion of $\nu \mapsto \Psi_0(P, \nu)$, so that

$$\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu) = D_\nu \Psi_0[\hat{\nu}_n - \nu] + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|_{\mathcal{V}}^2).$$

The decomposition then becomes

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] \tag{1}$$

$$+ D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] \tag{2}$$

$$+ \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_{\mathcal{V}}^2). \tag{3}$$

(1) can be handled by empirical process theory or sample splitting

(2) is our focus! \rightarrow make sense of this

(3) is specific to the nuisance estimator (and the functional Ψ).

Importantly, the rate $\sqrt{n}\|\hat{\nu}_n - \nu\|_{\mathcal{V}} = \mathcal{O}_P(n^{-1/4})$ is sufficient.

Defining a functional derivative

What is a derivative?

Defining a functional derivative

What is a derivative?

A linear approximation $\dot{\Psi}_x$ to the map Ψ at $x \in \mathcal{M}$, i.e.,

$$\left\| \Psi(x + \varepsilon_n h_n) - \Psi(x) - \dot{\Psi}_x(\varepsilon_n h_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

Defining a functional derivative

What is a derivative?

A linear approximation $\dot{\Psi}_x$ to the map Ψ at $x \in \mathcal{M}$, i.e.,

$$\left\| \Psi(x + \varepsilon_n h_n) - \Psi(x) - \dot{\Psi}_x(\varepsilon_n h_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

This expression also makes sense for functionals (or operators) Ψ .

Defining a functional derivative

What is a derivative?

A linear approximation $\dot{\Psi}_x$ to the map Ψ at $x \in \mathcal{M}$, i.e.,

$$\left\| \Psi(x + \varepsilon_n h_n) - \Psi(x) - \dot{\Psi}_x(\varepsilon_n h_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

This expression also makes sense for functionals (or operators) Ψ .

- For which h_n should this hold? Along “lines”, “paths”, or “uniformly” (h_n fixed, converging, or bounded)?

Defining a functional derivative

What is a derivative?

A linear approximation $\dot{\Psi}_x$ to the map Ψ at $x \in \mathcal{M}$, i.e.,

$$\left\| \Psi(x + \varepsilon_n h_n) - \Psi(x) - \dot{\Psi}_x(\varepsilon_n h_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

This expression also makes sense for functionals (or operators) Ψ .

- ▶ For which h_n should this hold? Along “lines”, “paths”, or “uniformly” (h_n fixed, converging, or bounded)?
- ▶ Which norm on \mathcal{M} should we use?

Defining a functional derivative

What is a derivative?

A linear approximation $\dot{\Psi}_x$ to the map Ψ at $x \in \mathcal{M}$, i.e.,

$$\left\| \Psi(x + \varepsilon_n h_n) - \Psi(x) - \dot{\Psi}_x(\varepsilon_n h_n) \right\| = \mathcal{O}(\varepsilon_n),$$

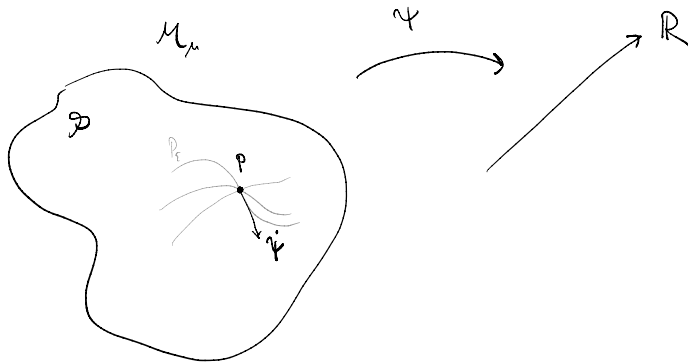
when $\varepsilon_n \rightarrow 0$.

This expression also makes sense for functionals (or operators) Ψ .

- ▶ For which h_n should this hold? Along “lines”, “paths”, or “uniformly” (h_n fixed, converging, or bounded)?
- ▶ Which norm on \mathcal{M} should we use?
- ▶ In which space should we represent \mathcal{P} ?

Pathwise Hadamard differentiability

Think of the gradient of a map defined on a manifold (surface).



Canonical gradient

Definition (Canonical gradient)

Let (\mathcal{P}, Ψ) be a statistical problem, with $\mathcal{P} \subset \mathcal{M}_\mu$, and $\dot{\mathcal{P}}_P$ the tangent space of \mathcal{P} at $P \in \mathcal{P}$. If $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is Hadamard differentiable at P tangential to $\dot{\mathcal{P}}_P$, we refer to the Hadamard derivative $\dot{\Psi}_P$ as the *canonical gradient of the statistical problem*.

Canonical gradient

Definition (Canonical gradient)

Let (\mathcal{P}, Ψ) be a statistical problem, with $\mathcal{P} \subset \mathcal{M}_\mu$, and $\dot{\mathcal{P}}_P$ the tangent space of \mathcal{P} at $P \in \mathcal{P}$. If $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is Hadamard differentiable at P tangential to $\dot{\mathcal{P}}_P$, we refer to the Hadamard derivative $\dot{\Psi}_P$ as the *canonical gradient of the statistical problem*.

Characterizing property

With $\Gamma_P := \overline{\text{span}}\{\dot{\ell}_0\} \subset \mathcal{L}_P^2$, where $\dot{\ell}_0 = \partial_0 \log(P_\varepsilon)$ is the score function of the sub-model P_ε , there exists a unique element $\varphi_P \in \Gamma_P$ such that

$$\partial_0 \Psi(P_\varepsilon) = \langle \varphi_P, \dot{\ell}_0 \rangle_P$$

holds for any differentiable submodel P_ε with score function $\dot{\ell}_0$.

Canonical gradient for the ATE

Example (ATE)

When we make no assumptions about \mathcal{P} , the canonical gradient for the ATE problem

$$\begin{aligned}\varphi_{\mathcal{P}}(o; f, \pi) &:= f(1, x) - f(0, x) \\ &\quad + \frac{a y}{\pi(x)} - \frac{(1-a) y}{1 - \pi(x)} \\ &\quad - \frac{a f(1, x)}{\pi(x)} + \frac{(1-a) f(0, x)}{1 - \pi(x)} \\ &\quad - \Psi(\mathcal{P})\end{aligned}$$

Canonical gradient for the ATE

Example (ATE)

When we make no assumptions about \mathcal{P} , the canonical gradient for the ATE problem

$$\begin{aligned}\varphi_{\mathcal{P}}(o; f, \pi) &:= f(1, x) - f(0, x) \\ &\quad + \frac{a y}{\pi(x)} - \frac{(1-a) y}{1 - \pi(x)} \\ &\quad - \frac{a f(1, x)}{\pi(x)} + \frac{(1-a) f(0, x)}{1 - \pi(x)} \\ &\quad - \Psi(\mathcal{P})\end{aligned}$$

One way to show this is to first show that the tangent space $\Gamma_{\mathcal{P}}$ is the full subset $\mathbb{H}_0 \subset \mathcal{L}_{\mathcal{P}}^2$ of zero-mean functions, and then show that $\partial_0 \Psi(\mathcal{P}_{\varepsilon}) = \langle \varphi_{\mathcal{P}}, \dot{\ell}_0 \rangle_{\mathcal{P}}$ for all $\mathcal{P}_{\varepsilon}$ (see for instance Kennedy [2016]).

Neyman orthogonality

Theorem (Neyman orthogonality)

If $\Psi(P) = \Psi_0(P, \nu) = P[\varphi(O, \nu(P))]$ and $\varphi(\cdot, \nu) - P[\varphi(O, \nu)]$ is the canonical gradient of (\mathcal{P}, Ψ) then $D_\nu \Psi_0 = 0$.

Neyman orthogonality

Theorem (Neyman orthogonality)

If $\Psi(P) = \Psi_0(P, \nu) = P[\varphi(O, \nu(P))]$ and $\varphi(\cdot, \nu) - P[\varphi(O, \nu)]$ is the canonical gradient of (\mathcal{P}, Ψ) then $D_\nu \Psi_0 = 0$.

Debiasing

The *first order* bias, coming from $\Psi_0(P, \hat{\nu}_n) - \Psi_0(P, \nu)$, is removed.

Efficiency

Definition (RAL estimators)

An estimator $\hat{\theta}_n$ of the parameter $\theta = \Psi(P)$ under the model \mathcal{P} , is called *asymptotically linear* with *influence function* $IF(\cdot, P) \in \mathcal{L}_P^2$, if $P[IF(O, P)] = 0$ for all $P \in \mathcal{P}$, and

$$\hat{\theta}_n - \theta = \hat{\mathbb{P}}_n[IF(O, P)] + o_P(n^{-1/2}).$$

Efficiency

Definition (RAL estimators)

An estimator $\hat{\theta}_n$ of the parameter $\theta = \Psi(P)$ under the model \mathcal{P} , is called *asymptotically linear* with *influence function* $IF(\cdot, P) \in \mathcal{L}_P^2$, if $P[IF(O, P)] = 0$ for all $P \in \mathcal{P}$, and

$$\hat{\theta}_n - \theta = \hat{\mathbb{P}}_n[IF(O, P)] + o_P(n^{-1/2}).$$

Theorem (Efficient influence function)

The RAL estimator with lowest possible asymptotic variance has the canonical gradient as its influence function.

Constructing estimators: Solve the efficient score equation

Find a parametrization $\Psi(P) = P[\varphi(O, \nu)]$ such that φ is the (canonical) gradient. Then by Neyman orthogonality and assumptions we can write

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] \\ &\quad + D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] \\ &\quad + \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_\nu^2)\end{aligned}$$

Constructing estimators: Solve the efficient score equation

Find a parametrization $\Psi(P) = P[\varphi(O, \nu)]$ such that φ is the (canonical) gradient. Then by Neyman orthogonality and assumptions we can write

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] \\ &\quad + D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] \quad = 0 \\ &\quad + \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_\nu^2)\end{aligned}$$

Constructing estimators: Solve the efficient score equation

Find a parametrization $\Psi(P) = P[\varphi(O, \nu)]$ such that φ is the (canonical) gradient. Then by Neyman orthogonality and assumptions we can write

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] &&= \mathbb{G}_n[\varphi(O, \nu)] \\ &+ D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] &&= 0 \\ &+ \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_\nu^2) &&= \mathcal{O}_P(1)\end{aligned}$$

Constructing estimators: Solve the efficient score equation

Find a parametrization $\Psi(P) = P[\varphi(O, \nu)]$ such that φ is the (canonical) gradient. Then by Neyman orthogonality and assumptions we can write

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] &&= \mathbb{G}_n[\varphi(O, \nu)] \\ &\quad + D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] &&= 0 \\ &\quad + \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_\nu^2) &&= \mathcal{O}_P(1) \\ &= \mathbb{G}_n[\varphi(O, \nu)] + \mathcal{O}_P(1).\end{aligned}$$

Constructing estimators: Solve the efficient score equation

Find a parametrization $\Psi(P) = P[\varphi(O, \nu)]$ such that φ is the (canonical) gradient. Then by Neyman orthogonality and assumptions we can write

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] &&= \mathbb{G}_n[\varphi(O, \nu)] \\ &\quad + D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] &&= 0 \\ &\quad + \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_\nu^2) &&= \mathcal{O}_P(1) \\ &= \mathbb{G}_n[\varphi(O, \nu)] + \mathcal{O}_P(1).\end{aligned}$$

Hence $\hat{\theta}_n$ is a RAL estimator, and if $\varphi - P[\varphi]$ is the canonical gradient it will be *asymptotically efficient*.

Constructing estimators: Solve the efficient score equation

Find a parametrization $\Psi(P) = P[\varphi(O, \nu)]$ such that φ is the (canonical) gradient. Then by Neyman orthogonality and assumptions we can write

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta) &= \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] &&= \mathbb{G}_n[\varphi(O, \nu)] \\ &\quad + D_\nu \Psi_0[\sqrt{n}(\hat{\nu}_n - \nu)] &&= 0 \\ &\quad + \mathcal{O}_P(\sqrt{n}\|\hat{\nu}_n - \nu\|_\nu^2) &&= \mathcal{O}_P(1) \\ &= \mathbb{G}_n[\varphi(O, \nu)] + \mathcal{O}_P(1).\end{aligned}$$

Hence $\hat{\theta}_n$ is a RAL estimator, and if $\varphi - P[\varphi]$ is the canonical gradient it will be *asymptotically efficient*.

This is the approach taken in Chernozhukov et al. [2018]. See also Example 4.1 of the note.

- V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins. Double/debiased machine learning for treatment and structural parameters, 2018.
- E. H. Kennedy. Semiparametric theory and empirical processes in causal inference. In *Statistical causal inferences and their applications in public health research*, pages 141–167. Springer, 2016.