

Bias-variance trade-off with infinite-dimensional nuisance parameters

Anders Munch

September 24, 2021

Outline

Setting and motivation

Bias/variance trade-off for plug-in problems

Functional derivatives

Exercises

References

A statistical estimation problem

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

A statistical estimation problem

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

Example (Average treatment effect)

We are given n iid. samples of $O \sim P$, with $P \in \mathcal{P}$ 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

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

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

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

A statistical estimation problem

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

Example (Average treatment effect)

We are given n iid. samples of $O \sim P$, with $P \in \mathcal{P}$ 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

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

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

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

Target and nuisance parameters

Target and nuisance parameters

Target parameter is low-dimensional, scientifically meaningful.

Target and nuisance parameters

Target parameter is low-dimensional, scientifically meaningful.

Nuisance parameters are needed to express the target parameter.

Target and nuisance parameters

Target parameter is low-dimensional, scientifically meaningful.

Nuisance parameters are needed to express the target parameter.

Example (Average treatment effect)

We can express the ATE as

$$\Psi(P) = \tilde{\Psi}(f, \mu_X) := \int_{\mathbb{R}^p} [f(1, x) - f(0, x)] d\mu_X(x).$$

The nuisance parameters are f and μ_X .

Target and nuisance parameters

Target parameter is low-dimensional, scientifically meaningful.

Nuisance parameters are needed to express the target parameter.

Example (Average treatment effect)

We can express the ATE as

$$\Psi(P) = \tilde{\Psi}(f, \mu_X) := \int_{\mathbb{R}^p} [f(1, x) - f(0, x)] d\mu_X(x).$$

The nuisance parameters are f and μ_X . This immediately suggests the target estimator $\hat{\psi}_n^{\text{g-formula}} = \Psi(\hat{f}_n, \hat{\mu}_X)$; for instance, if we use $\hat{\mu}_X = \hat{\mathbb{P}}_n$ we have

$$\hat{\psi}_n^{\text{g-formula}} = \tilde{\Psi}(\hat{f}_n, \hat{\mathbb{P}}_n) = \frac{1}{n} \sum_{i=1}^n \left\{ \hat{f}_n(1, X_i) - \hat{f}_n(0, X_i) \right\}.$$

Hence we just have to select a nuisance estimator \hat{f}_n .

Nuisance parameter estimation and plug-in strategy

In many interesting cases the nuisance parameter is a conditional mean

Nuisance parameter estimation and plug-in strategy

In many interesting cases the nuisance parameter is a conditional mean

- Estimation of the nuisance parameter reduces to the standard statistical / machine learning task of fitting a prediction model.

Nuisance parameter estimation and plug-in strategy

In many interesting cases the nuisance parameter is a conditional mean

- Estimation of the nuisance parameter reduces to the standard statistical / machine learning task of fitting a prediction model.
- This is a very well studied problem, and a zoo of estimators exists for these kinds of problems (most machine learning problems are of this type)

Nuisance parameter estimation and plug-in strategy

In many interesting cases the nuisance parameter is a conditional mean

- Estimation of the nuisance parameter reduces to the standard statistical / machine learning task of fitting a prediction model.
- This is a very well studied problem, and a zoo of estimators exists for these kinds of problems (most machine learning problems are of this type)
- Provides estimators of the target parameter by plugging in a nuisance estimator

Nuisance parameter estimation and plug-in strategy

In many interesting cases the nuisance parameter is a conditional mean

- Estimation of the nuisance parameter reduces to the standard statistical / machine learning task of fitting a prediction model.
- This is a very well studied problem, and a zoo of estimators exists for these kinds of problems (most machine learning problems are of this type)
- Provides estimators of the target parameter by plugging in a nuisance estimator

Such a plug-in strategy can perform poorly, **even when the nuisance estimator is optimized for the nuisance problem.**

Nuisance parameter estimation and plug-in strategy

In many interesting cases the nuisance parameter is a conditional mean

- Estimation of the nuisance parameter reduces to the standard statistical / machine learning task of fitting a prediction model.
- This is a very well studied problem, and a zoo of estimators exists for these kinds of problems (most machine learning problems are of this type)
- Provides estimators of the target parameter by plugging in a nuisance estimator

Such a plug-in strategy can perform poorly, **even when the nuisance estimator is optimized for the nuisance problem.**

We will demonstrate this with the following a toy example.

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 $F_P(x) = P(X \leq x)$ 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

$F_P(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$. Our target parameter is then $\Psi(P) = F_P(x)$ which we can express as

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

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

$F_P(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$. Our target parameter is then $\Psi(P) = F_P(x)$ which we can express as

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

because of our assumption about \mathcal{P} . If we wanted to use “machine learning” for this problem, we could use a kernel estimator

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

to estimate the density f , where $h \in \mathbb{R}_+$ is the bandwidth (a tuning parameter).

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

$F_P(x) = P(X \leq x)$ for unknown $P \in \mathcal{P}$. Our target parameter is then $\Psi(P) = F_P(x)$ which we can express as

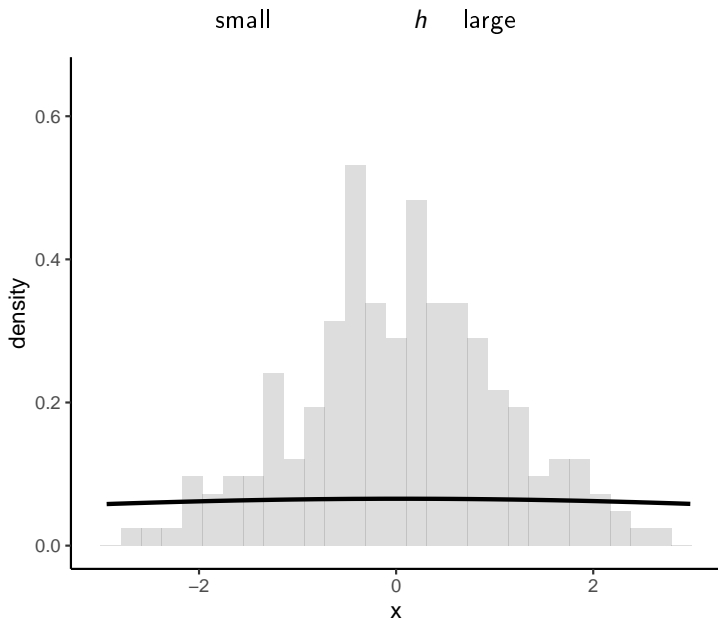
$$\Psi(P) = \tilde{\Psi}(f) := \int_{-\infty}^x f(z) \, dz, \quad \text{for } P = f \cdot \lambda,$$

because of our assumption about \mathcal{P} . If we wanted to use “machine learning” for this problem, we could use a kernel estimator

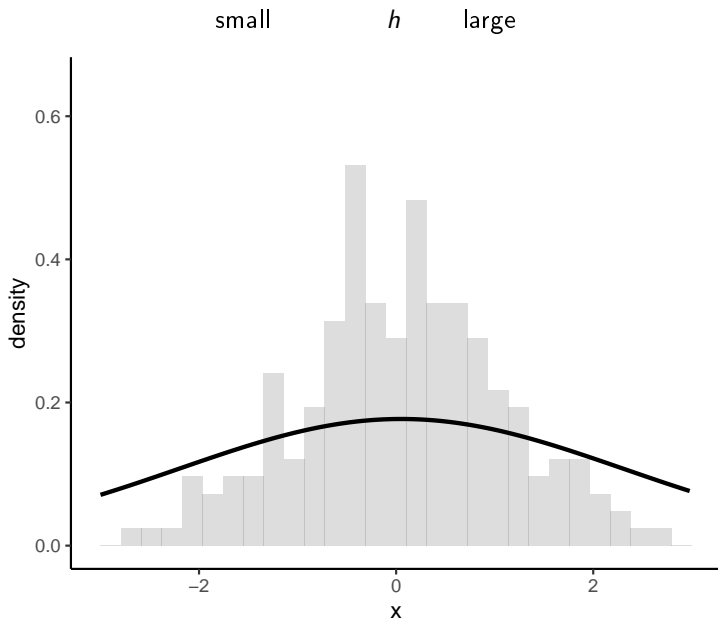
$$\hat{f}_n(x) = \hat{\mathbb{P}}_n[k_h(X, x)] = \frac{1}{n} \sum_{i=1}^n k_h(X_i, x),$$

to estimate the density f , where $h \in \mathbb{R}_+$ is the bandwidth (a tuning parameter). We could then obtain the target estimator $\hat{\psi}_n = \tilde{\Psi}(\hat{f}_n)$.

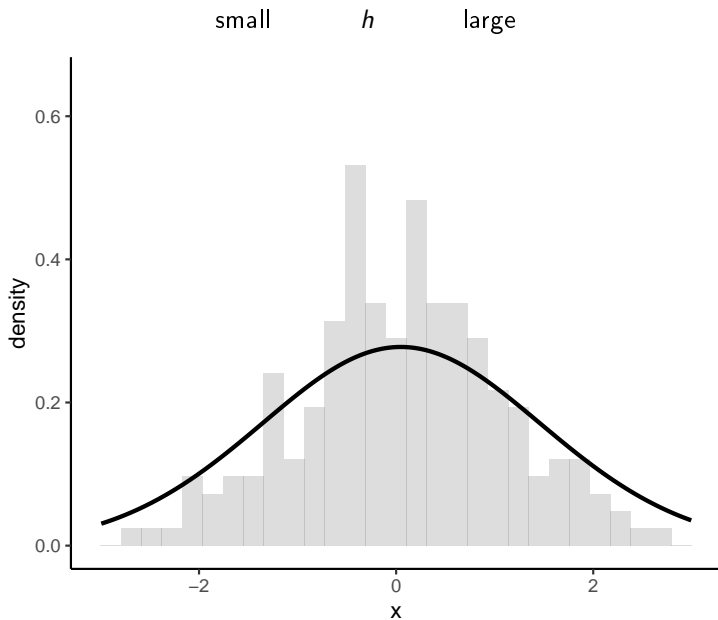
Visualize kernel estimator



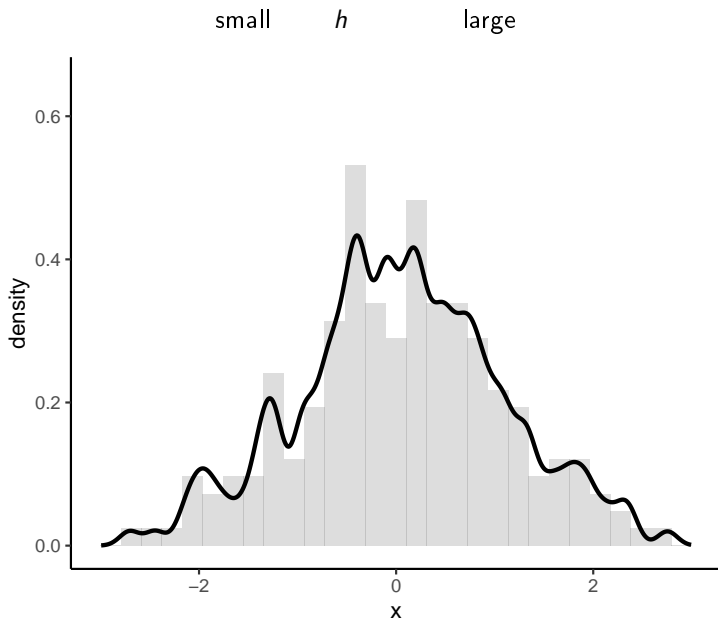
Visualize kernel estimator



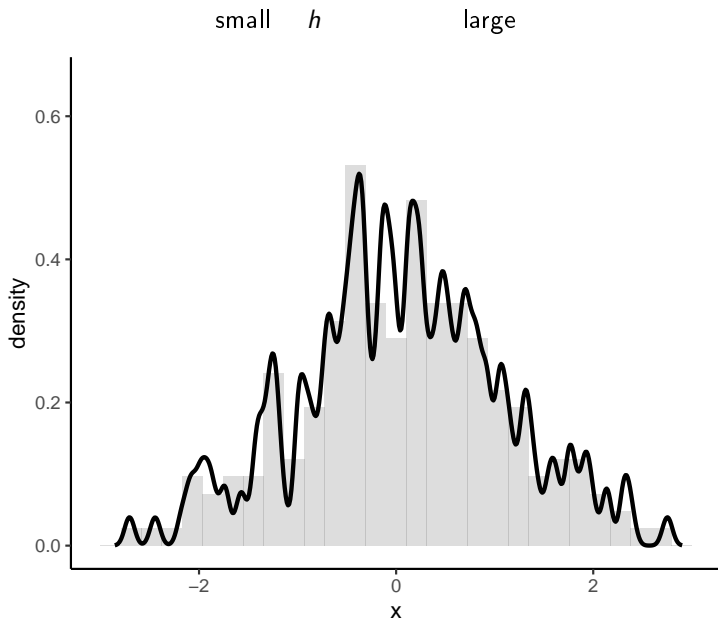
Visualize kernel estimator



Visualize kernel estimator

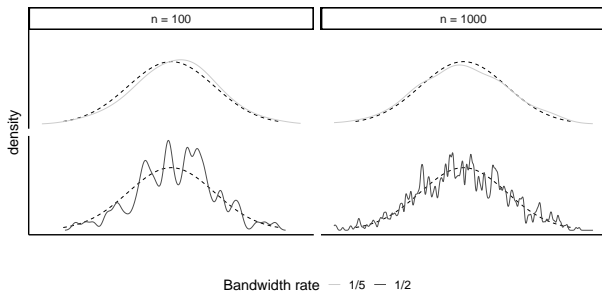


Visualize kernel estimator

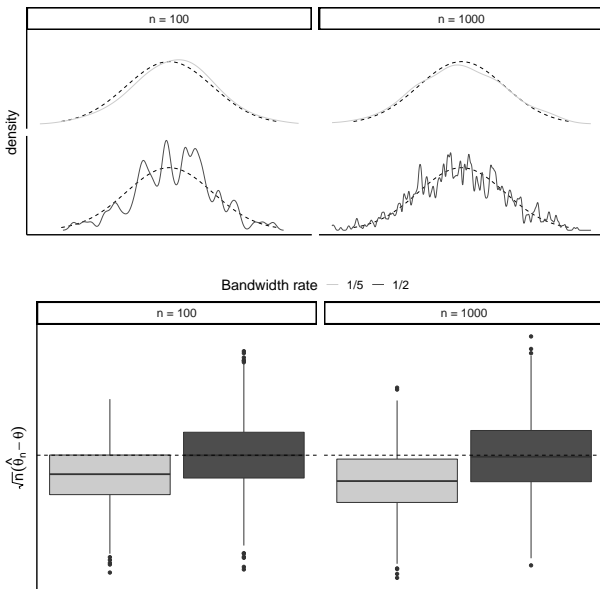


How does this work in practice?

How does this work in practice?



How does this work in practice?



What happened?

The bias-variance trade-off for the nuisance parameter f is

$$\text{MSE}(\hat{f}_n) = C_1 h^4 + C_2 (nh)^{-1} + o(h^2) + o(n^{-1}),$$

where $n \rightarrow \infty$ and $h \rightarrow 0$. This implies that the optimal value for the bandwidth h is $h \asymp n^{-1/5}$ [van der Vaart, 2000, chp. 24].

What happened?

The bias-variance trade-off for the nuisance parameter f is

$$\text{MSE}(\hat{f}_n) = C_1 h^4 + C_2 (nh)^{-1} + o(h^2) + o(n^{-1}),$$

where $n \rightarrow \infty$ and $h \rightarrow 0$. This implies that the optimal value for the bandwidth h is $h \asymp n^{-1/5}$ [van der Vaart, 2000, chp. 24].

For the target parameter $\int_{-\infty}^x f(z) d\lambda(z)$ the bias-variance “trade-off” becomes

$$\text{MSE}(\hat{\theta}_n) = o(n^{-1}) + o(hn^{-1}) + o(h^4).$$

The optimal value of h is now found by picking h as small as possible.

What happened?

The bias-variance trade-off for the nuisance parameter f is

$$\text{MSE}(\hat{f}_n) = C_1 h^4 + C_2 (nh)^{-1} + o(h^2) + o(n^{-1}),$$

where $n \rightarrow \infty$ and $h \rightarrow 0$. This implies that the optimal value for the bandwidth h is $h \asymp n^{-1/5}$ [van der Vaart, 2000, chp. 24].

For the target parameter $\int_{-\infty}^x f(z) d\lambda(z)$ the bias-variance “trade-off” becomes

$$\text{MSE}(\hat{\theta}_n) = o(n^{-1}) + o(hn^{-1}) + o(h^4).$$

The optimal value of h is now found by picking h as small as possible.

Using $h = 0$ can be interpreted as just using the empirical distribution function \hat{F}_n , i.e.,

$$\int_{-\infty}^x \hat{f}_n(z) dz = \hat{\mathbb{P}}_n \left[\int_{-\infty}^x k_h(X_i, z) dz \right] \longrightarrow \hat{\mathbb{P}}_n[\mathbb{1}(X_i \leq x)] =: \hat{F}_n(x).$$

Finite- versus infinite-dimensional nuisance parameters

This phenomena only occurs when there are infinite-dimensional nuisance parameters (e.g., smooth functions, shape-constrained densities, etc.) – or more precisely, when we do not have parametric ($n^{-1/2}$) rate of convergence for the nuisance parameter estimator.

Finite- versus infinite-dimensional nuisance parameters

This phenomena only occurs when there are infinite-dimensional nuisance parameters (e.g., smooth functions, shape-constrained densities, etc.) – or more precisely, when we do not have parametric ($n^{-1/2}$) rate of convergence for the nuisance parameter estimator. Under suitable regularity conditions, if

$$\sqrt{n}(\hat{\nu}_n - \nu) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

then, by the delta method, also

$$\sqrt{n}(\Psi(\hat{\nu}_n) - \Psi(\nu)) \rightsquigarrow \mathcal{N}(0, \dot{\Psi}\sigma^2).$$

Finite- versus infinite-dimensional nuisance parameters

This phenomena only occurs when there are infinite-dimensional nuisance parameters (e.g., smooth functions, shape-constrained densities, etc.) – or more precisely, when we do not have parametric ($n^{-1/2}$) rate of convergence for the nuisance parameter estimator. Under suitable regularity conditions, if

$$\sqrt{n}(\hat{\nu}_n - \nu) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

then, by the delta method, also

$$\sqrt{n}(\Psi(\hat{\nu}_n) - \Psi(\nu)) \rightsquigarrow \mathcal{N}(0, \dot{\Psi}\sigma^2).$$

Even more, if $\hat{\nu}_n$ is efficient $\Psi(\hat{\nu}_n)$ will also be efficient [van der Vaart, 2000, Bickel et al., 2003].

Finite- versus infinite-dimensional nuisance parameters

This phenomena only occurs when there are infinite-dimensional nuisance parameters (e.g., smooth functions, shape-constrained densities, etc.) – or more precisely, when we do not have parametric ($n^{-1/2}$) rate of convergence for the nuisance parameter estimator. Under suitable regularity conditions, if

$$\sqrt{n}(\hat{\nu}_n - \nu) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

then, by the delta method, also

$$\sqrt{n}(\Psi(\hat{\nu}_n) - \Psi(\nu)) \rightsquigarrow \mathcal{N}(0, \dot{\Psi}\sigma^2).$$

Even more, if $\hat{\nu}_n$ is efficient $\Psi(\hat{\nu}_n)$ will also be efficient [van der Vaart, 2000, Bickel et al., 2003].

The *functional* delta method gives r_n -rate convergence of the plug-in estimator $\Psi(\hat{\nu}_n)$ when $\hat{\nu}_n$ converges at rate r_n ;

Finite- versus infinite-dimensional nuisance parameters

This phenomena only occurs when there are infinite-dimensional nuisance parameters (e.g., smooth functions, shape-constrained densities, etc.) – or more precisely, when we do not have parametric ($n^{-1/2}$) rate of convergence for the nuisance parameter estimator. Under suitable regularity conditions, if

$$\sqrt{n}(\hat{\nu}_n - \nu) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

then, by the delta method, also

$$\sqrt{n}(\Psi(\hat{\nu}_n) - \Psi(\nu)) \rightsquigarrow \mathcal{N}(0, \dot{\Psi}\sigma^2).$$

Even more, if $\hat{\nu}_n$ is efficient $\Psi(\hat{\nu}_n)$ will also be efficient [van der Vaart, 2000, Bickel et al., 2003].

The *functional* delta method gives r_n -rate convergence of the plug-in estimator $\Psi(\hat{\nu}_n)$ when $\hat{\nu}_n$ converges at rate r_n ; but it does not tell us how to get $n^{-1/2}$ rate convergence of an estimator $\hat{\psi}_n$ using $\hat{\nu}_n$, when $\hat{\nu}_n$ converges at a slower rate.

The bias-variance tradeoff revisited

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

The bias-variance tradeoff revisited

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

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

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

The bias-variance tradeoff revisited

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

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

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

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

The bias-variance tradeoff revisited

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

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

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

$\mathbb{G}_n[\varphi(O, \hat{\nu}_n)]$ determines the (main) variance – often
 $\mathbb{G}_n[\varphi(O, \hat{\nu}_n)] = \mathbb{G}_n[\varphi(O, \nu)] + o_P(1)$

The bias-variance tradeoff revisited

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

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

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

$\mathbb{G}_n[\varphi(O, \hat{\nu}_n)]$ determines the (main) variance – often
 $\mathbb{G}_n[\varphi(O, \hat{\nu}_n)] = \mathbb{G}_n[\varphi(O, \nu)] + o_P(1)$

$\tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu)$ is bias

The bias-variance tradeoff revisited

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

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

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

$\mathbb{G}_n[\varphi(O, \hat{\nu}_n)]$ determines the (main) variance – often
 $\mathbb{G}_n[\varphi(O, \hat{\nu}_n)] = \mathbb{G}_n[\varphi(O, \nu)] + o_P(1)$

$\tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu)$ is bias \rightarrow This is what ruined the naive plug-in strategy!

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

What is a derivative?

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

What is a derivative?

A linear approximation $\dot{\Psi}_P$ of the map Ψ at $P \in \mathcal{P}$, i.e.,

$$\left\| \Psi(P + \varepsilon_n g_n) - \Psi(P) - \dot{\Psi}_P(\varepsilon_n g_n) \right\| = o(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

What is a derivative?

A linear approximation $\dot{\Psi}_P$ of the map Ψ at $P \in \mathcal{P}$, i.e.,

$$\left\| \Psi(P + \varepsilon_n g_n) - \Psi(P) - \dot{\Psi}_P(\varepsilon_n g_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

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

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

What is a derivative?

A linear approximation $\dot{\Psi}_P$ of the map Ψ at $P \in \mathcal{P}$, i.e.,

$$\left\| \Psi(P + \varepsilon_n g_n) - \Psi(P) - \dot{\Psi}_P(\varepsilon_n g_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

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

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

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

What is a derivative?

A linear approximation $\dot{\Psi}_P$ of the map Ψ at $P \in \mathcal{P}$, i.e.,

$$\left\| \Psi(P + \varepsilon_n g_n) - \Psi(P) - \dot{\Psi}_P(\varepsilon_n g_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

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

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

Defining a functional derivative

Functional derivatives and von Mises expansions are useful for analyzing and handling the issues we have encountered [Serfling, 1980].

What is a derivative?

A linear approximation $\dot{\Psi}_P$ of the map Ψ at $P \in \mathcal{P}$, i.e.,

$$\left\| \Psi(P + \varepsilon_n g_n) - \Psi(P) - \dot{\Psi}_P(\varepsilon_n g_n) \right\| = \mathcal{O}(\varepsilon_n),$$

when $\varepsilon_n \rightarrow 0$.

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

- ▶ For which g_n should this hold? Along “lines”, “paths”, or “uniformly” (g_n fixed, converging, or bounded)?
- ▶ Which norm on \mathcal{P} should we use?
- ▶ Is there a natural space \mathcal{M} in which to embed \mathcal{P} ?

Gâteaux and Hadamard differentiability

The weakest kind of differentiability is *Gâteaux differentiability*. When $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ the Gâteaux derivative $\dot{\Psi}_P$ is the directional derivative

$$\dot{\Psi}_P(g) = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon g). \quad (1)$$

Gâteaux and Hadamard differentiability

The weakest kind of differentiability is *Gâteaux differentiability*. When $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ the Gâteaux derivative $\dot{\Psi}_P$ is the directional derivative

$$\dot{\Psi}_P(g) = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon g). \quad (1)$$

A stronger type is *Hadamard differentiability* which demands that the linear approximation is valid along any converging sequence $g_n \rightarrow g$.

Gâteaux and Hadamard differentiability

The weakest kind of differentiability is *Gâteaux differentiability*. When $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ the Gâteaux derivative $\dot{\Psi}_P$ is the directional derivative

$$\dot{\Psi}_P(g) = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon g). \quad (1)$$

A stronger type is *Hadamard differentiability* which demands that the linear approximation is valid along any converging sequence $g_n \rightarrow g$. A Hadamard differentiable map is also called *compactly differentiable*, because this requirement is equivalent to a uniformly valid linear approximation over compact subsets [Reeds, 1976].

Gâteaux and Hadamard differentiability

The weakest kind of differentiability is *Gâteaux differentiability*. When $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ the Gâteaux derivative $\dot{\Psi}_P$ is the directional derivative

$$\dot{\Psi}_P(g) = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon g). \quad (1)$$

A stronger type is *Hadamard differentiability* which demands that the linear approximation is valid along any converging sequence $g_n \rightarrow g$. A Hadamard differentiable map is also called *compactly differentiable*, because this requirement is equivalent to a uniformly valid linear approximation over compact subsets [Reeds, 1976]. This makes Hadamard differentiability well-suited for statistical problems [Gill et al., 1989, van der Vaart, 2000].

Gâteaux and Hadamard differentiability

The weakest kind of differentiability is *Gâteaux differentiability*. When $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ the Gâteaux derivative $\dot{\Psi}_P$ is the directional derivative

$$\dot{\Psi}_P(g) = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon g). \quad (1)$$

A stronger type is *Hadamard differentiability* which demands that the linear approximation is valid along any converging sequence $g_n \rightarrow g$. A Hadamard differentiable map is also called *compactly differentiable*, because this requirement is equivalent to a uniformly valid linear approximation over compact subsets [Reeds, 1976]. This makes Hadamard differentiability well-suited for statistical problems [Gill et al., 1989, van der Vaart, 2000].

If Ψ is Hadamard differentiable, Ψ is also Gâteaux differentiable, and in that case the Hadamard and Gâteaux derivative are identical.

Hadamard differentiability and the canonical gradient

The “gradient” of $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is called the *canonical gradient* or *efficient influence function* of a statistical estimation problem (\mathcal{P}, Ψ) , and it is a fundamental object for semi-parametric efficiency theory – we will see why in a moment.

Hadamard differentiability and the canonical gradient

The “gradient” of $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is called the *canonical gradient* or *efficient influence function* of a statistical estimation problem (\mathcal{P}, Ψ) , and it is a fundamental object for semi-parametric efficiency theory – we will see why in a moment. More formally we define:

Definition (Tangent space)

The *tangent space* $\dot{\mathcal{P}}_P$ for the model \mathcal{P} at $P \in \mathcal{P}$ is the (closed linear span of the) collection of (Hadamard) derivatives \dot{P}_ε for all one-dimensional parametric submodel $P_\varepsilon \subset \mathcal{P}$ with $P_0 = P$.

Hadamard differentiability and the canonical gradient

The “gradient” of $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is called the *canonical gradient* or *efficient influence function* of a statistical estimation problem (\mathcal{P}, Ψ) , and it is a fundamental object for semi-parametric efficiency theory – we will see why in a moment. More formally we define:

Definition (Tangent space)

The *tangent space* $\dot{\mathcal{P}}_P$ for the model \mathcal{P} at $P \in \mathcal{P}$ is the (closed linear span of the) collection of (Hadamard) derivatives \dot{P}_ε for all one-dimensional parametric submodel $P_\varepsilon \subset \mathcal{P}$ with $P_0 = P$.

Definition (Canonical gradient)

If $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is Hadamard differentiable at P tangential to the tangent space $\dot{\mathcal{P}}_P$,¹ we refer to the Hadamard derivative $\dot{\Psi}_P$ as the *canonical gradient of the statistical estimation problem*.

Hadamard differentiability and the canonical gradient

The “gradient” of $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is called the *canonical gradient* or *efficient influence function* of a statistical estimation problem (\mathcal{P}, Ψ) , and it is a fundamental object for semi-parametric efficiency theory – we will see why in a moment. More formally we define:

Definition (Tangent space)

The *tangent space* $\dot{\mathcal{P}}_P$ for the model \mathcal{P} at $P \in \mathcal{P}$ is the (closed linear span of the) collection of (Hadamard) derivatives \dot{P}_ε for all one-dimensional parametric submodel $P_\varepsilon \subset \mathcal{P}$ with $P_0 = P$.

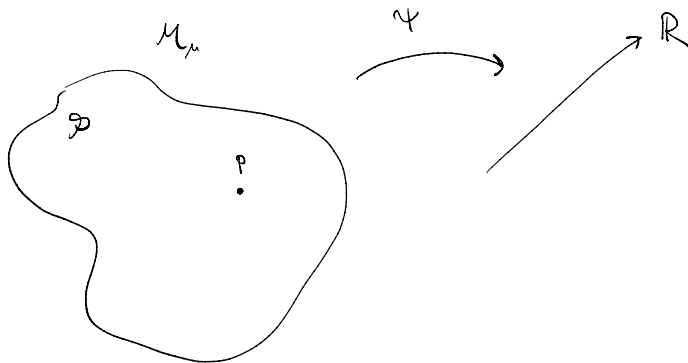
Definition (Canonical gradient)

If $\Psi: \mathcal{P} \rightarrow \mathbb{R}$ is Hadamard differentiable at P tangential to the tangent space $\dot{\mathcal{P}}_P$,¹ we refer to the Hadamard derivative $\dot{\Psi}_P$ as the *canonical gradient of the statistical estimation problem*.

¹This just means that the Ψ and $\dot{\Psi}_P$ need only be defined on the subsets $\mathcal{P} \subset \mathcal{M}$ and $\dot{\mathcal{P}}_P \subset \mathcal{M}$, respectively.

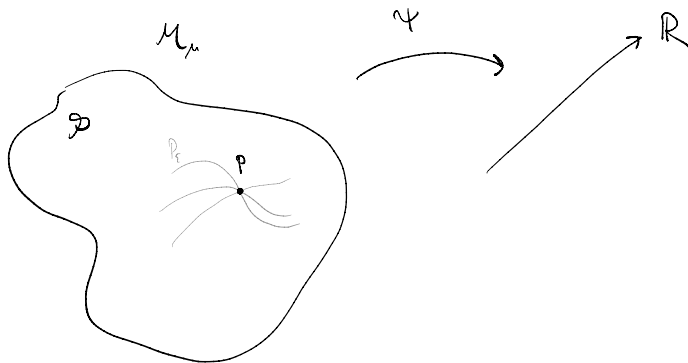
Visualization

Think of the gradient of a function defined on manifold (e.g., a surface).



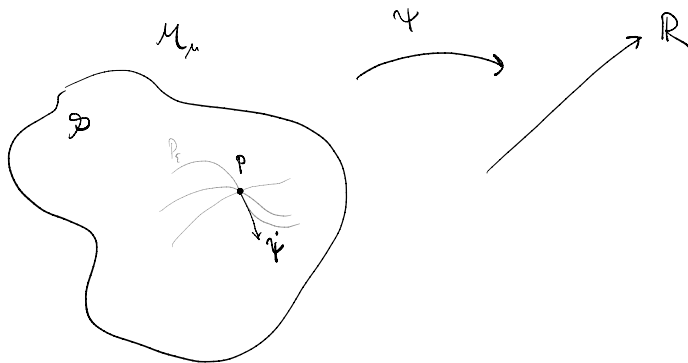
Visualization

Think of the gradient of a function defined on manifold (e.g., a surface).



Visualization

Think of the gradient of a function defined on manifold (e.g., a surface).



Alternative characterization

The tangent space can be represented as

$$\dot{\mathcal{P}}_P = \overline{\text{span}}\{\dot{\ell}_0\} \subset \mathcal{L}_0^2(P), \quad \text{where} \quad \dot{\ell}_0 = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(p_\varepsilon),$$

and $\mathcal{L}_0^2(P)$ is the space of all functions $g \in \mathcal{L}^2(P)$ and $P[g] = 0$.

Alternative characterization

The tangent space can be represented as

$$\dot{\mathcal{P}}_P = \overline{\text{span}}\{\dot{\ell}_0\} \subset \mathcal{L}_0^2(P), \quad \text{where} \quad \dot{\ell}_0 = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(p_\varepsilon),$$

and $\mathcal{L}_0^2(P)$ is the space of all functions $g \in \mathcal{L}^2(P)$ and $P[g] = 0$.

By Riesz' representation theorem, and the (Hadamard) chain rule, this implies that the canonical gradient (efficient influence function) $\dot{\Psi}_P$ can be identified with an element $\varphi_P \in \dot{\mathcal{P}}_P \subset \mathcal{L}_0^2(P)$ such that

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P_\varepsilon) = \dot{\Psi}_P \left(\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} P_\varepsilon \right) = \dot{\Psi}_P(\dot{\ell}_0) = \langle \varphi_P, \dot{\ell}_0 \rangle_P, \quad (2)$$

holds for any differentiable submodel P_ε with score function $\dot{\ell}_0$, where $= \langle \varphi_P, \dot{\ell}_0 \rangle_P = \int \varphi_P(o) \dot{\ell}_0(o) dP(o)$.

Alternative characterization

The tangent space can be represented as

$$\dot{\mathcal{P}}_P = \overline{\text{span}}\{\dot{\ell}_0\} \subset \mathcal{L}_0^2(P), \quad \text{where} \quad \dot{\ell}_0 = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(p_\varepsilon),$$

and $\mathcal{L}_0^2(P)$ is the space of all functions $g \in \mathcal{L}^2(P)$ and $P[g] = 0$.

By Riesz' representation theorem, and the (Hadamard) chain rule, this implies that the canonical gradient (efficient influence function) $\dot{\Psi}_P$ can be identified with an element $\varphi_P \in \dot{\mathcal{P}}_P \subset \mathcal{L}_0^2(P)$ such that

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P_\varepsilon) = \dot{\Psi}_P \left(\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} P_\varepsilon \right) = \dot{\Psi}_P(\dot{\ell}_0) = \langle \varphi_P, \dot{\ell}_0 \rangle_P, \quad (2)$$

holds for any differentiable submodel P_ε with score function $\dot{\ell}_0$, where $= \langle \varphi_P, \dot{\ell}_0 \rangle_P = \int \varphi_P(o) \dot{\ell}_0(o) dP(o)$.

Note the function φ_P is determined *both* by the map Ψ *and* the model \mathcal{P} : Many functions $\tilde{\varphi}_P \in \mathcal{L}_0^2(P)$ might fulfill (2), but only one will be contained in the tangent space.

Alternative characterization

The tangent space can be represented as

$$\dot{\mathcal{P}}_P = \overline{\text{span}}\{\dot{\ell}_0\} \subset \mathcal{L}_0^2(P), \quad \text{where} \quad \dot{\ell}_0 = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(p_\varepsilon),$$

and $\mathcal{L}_0^2(P)$ is the space of all functions $g \in \mathcal{L}^2(P)$ and $P[g] = 0$.

By Riesz' representation theorem, and the (Hadamard) chain rule, this implies that the canonical gradient (efficient influence function) $\dot{\Psi}_P$ can be identified with an element $\varphi_P \in \dot{\mathcal{P}}_P \subset \mathcal{L}_0^2(P)$ such that

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P_\varepsilon) = \dot{\Psi}_P \left(\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} P_\varepsilon \right) = \dot{\Psi}_P(\dot{\ell}_0) = \langle \varphi_P, \dot{\ell}_0 \rangle_P, \quad (2)$$

holds for any differentiable submodel P_ε with score function $\dot{\ell}_0$, where $= \langle \varphi_P, \dot{\ell}_0 \rangle_P = \int \varphi_P(o) \dot{\ell}_0(o) dP(o)$.

Note the function φ_P is determined *both* by the map Ψ *and* the model \mathcal{P} : Many functions $\tilde{\varphi}_P \in \mathcal{L}_0^2(P)$ might fulfill (2), but only one will be contained in the tangent space.

In the important special case of a *fully non-parametric model* \mathcal{P} , $\dot{\mathcal{P}}_P = \mathcal{L}_0^2(P)$ and thus (2) alone uniquely identifies the function φ_P .

Efficient estimation and the canonical gradient

Efficient estimation and the canonical gradient

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* $\text{IF}(\cdot, P) \in \mathcal{L}_P^2$, if $P[\text{IF}(O, P)] = 0$ for all $P \in \mathcal{P}$, and

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

Efficient estimation and the canonical gradient

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* $\text{IF}(\cdot, P) \in \mathcal{L}_P^2$, if $P[\text{IF}(O, P)] = 0$ for all $P \in \mathcal{P}$, and

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

Theorem (Efficient influence function)

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

Efficient estimation and the canonical gradient

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* $\text{IF}(\cdot, P) \in \mathcal{L}_P^2$, if $P[\text{IF}(O, P)] = 0$ for all $P \in \mathcal{P}$, and

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

Theorem (Efficient influence function)

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

Information bound

The *information bound* for estimating Ψ in the model \mathcal{P} is

$$\mathcal{I}(\mathcal{P}, \Psi) := \inf_{P_\varepsilon} \{ \mathcal{I}(P_\varepsilon, \Psi) \}, \quad \text{with} \quad \mathcal{I}(P_\varepsilon, \Psi) := \frac{P[\dot{\ell}_0^2]}{(\partial_0 \Psi(P_\varepsilon))^2}.$$

It holds that $\mathcal{I}(\mathcal{P}, \Psi)^{-1} = P[\varphi^2]$.

Debiasing and the canonical gradient

Recall the statistical estimation problem $\Psi(P) = \tilde{\Psi}(P, \nu) = P[\varphi(O, \nu)]$ and the decomposition

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n} \left\{ \tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) \right\}.$$

Debiasing and the canonical gradient

Recall the statistical estimation problem $\Psi(P) = \tilde{\Psi}(P, \nu) = P[\varphi(O, \nu)]$ and the decomposition

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n} \left\{ \tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) \right\}.$$

If φ is the canonical gradient then the Gâteaux derivative of $\nu \mapsto \tilde{\Psi}(P, \nu)$ is zero.²

Debiasing and the canonical gradient

Recall the statistical estimation problem $\Psi(P) = \tilde{\Psi}(P, \nu) = P[\varphi(O, \nu)]$ and the decomposition

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n} \left\{ \tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) \right\}.$$

If φ is the canonical gradient then the Gâteaux derivative of $\nu \mapsto \tilde{\Psi}(P, \nu)$ is zero.² Hence a von Mises expansion suggests that

$$\tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) = 0 + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|^2).$$

Debiasing and the canonical gradient

Recall the statistical estimation problem $\Psi(P) = \tilde{\Psi}(P, \nu) = P[\varphi(O, \nu)]$ and the decomposition

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n} \left\{ \tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) \right\}.$$

If φ is the canonical gradient then the Gâteaux derivative of $\nu \mapsto \tilde{\Psi}(P, \nu)$ is zero.² Hence a von Mises expansion suggests that

$$\tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) = 0 + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|^2).$$

In particular, we can obtain $n^{-1/2}$ convergence of $\hat{\theta}_n$ even though $\hat{\nu}_n$ only converges at rate $n^{-1/4}$.

Debiasing and the canonical gradient

Recall the statistical estimation problem $\Psi(P) = \tilde{\Psi}(P, \nu) = P[\varphi(O, \nu)]$ and the decomposition

$$\sqrt{n}(\hat{\theta}_n - \theta) = \mathbb{G}_n[\varphi(O, \hat{\nu}_n)] + \sqrt{n} \left\{ \tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) \right\}.$$

If φ is the canonical gradient then the Gâteaux derivative of $\nu \mapsto \tilde{\Psi}(P, \nu)$ is zero.² Hence a von Mises expansion suggests that

$$\tilde{\Psi}(P, \hat{\nu}_n) - \tilde{\Psi}(P, \nu) = 0 + \mathcal{O}_P(\|\hat{\nu}_n - \nu\|^2).$$

In particular, we can obtain $n^{-1/2}$ convergence of $\hat{\theta}_n$ even though $\hat{\nu}_n$ only converges at rate $n^{-1/4}$.

²This property is referred to as *Neyman orthogonality*, and is a central component of “debiased machine learning”; see Chernozhukov et al. [2018] for details.

Finding the canonical gradient for a given Ψ

We can (informally and heuristically³) find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O , where δ_O is the Dirac measure in the point O :

³To make the arguments more precise, see Ichimura and Newey [2015].

Finding the canonical gradient for a given Ψ

We can (informally and heuristically³) find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O , where δ_O is the Dirac measure in the point O :

For any $h \in \mathbb{R}_+$, define the sub-model $P_\varepsilon^h := P + \varepsilon K_O^h$, where $K_O^h \rightarrow \delta_O$, $h \rightarrow 0$.

³To make the arguments more precise, see Ichimura and Newey [2015].

Finding the canonical gradient for a given Ψ

We can (informally and heuristically³) find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O , where δ_O is the Dirac measure in the point O :

For any $h \in \mathbb{R}_+$, define the sub-model $P_\varepsilon^h := P + \varepsilon K_O^h$, where $K_O^h \rightarrow \delta_O$, $h \rightarrow 0$. The score function of this model is

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(dP + \varepsilon dK_O^h) = \frac{dK_O^h}{dP}.$$

³To make the arguments more precise, see Ichimura and Newey [2015].

Finding the canonical gradient for a given Ψ

We can (informally and heuristically³) find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O , where δ_O is the Dirac measure in the point O :

For any $h \in \mathbb{R}_+$, define the sub-model $P_\varepsilon^h := P + \varepsilon K_O^h$, where $K_O^h \rightarrow \delta_O$, $h \rightarrow 0$. The score function of this model is

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(dP + \varepsilon dK_O^h) = \frac{dK_O^h}{dP}.$$

Thus, if Ψ is Hadamard differentiable, we can use (2) to write

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P_\varepsilon^h) = \left\langle \varphi_P, \frac{dK_O^h}{dP} \right\rangle_P = \int \varphi_P(o) \frac{dK_O^h(o)}{dP(o)} dP(o) = \int \varphi_P(o) dK_O^h(o).$$

³To make the arguments more precise, see Ichimura and Newey [2015].

Finding the canonical gradient for a given Ψ

We can (informally and heuristically³) find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O , where δ_O is the Dirac measure in the point O :

For any $h \in \mathbb{R}_+$, define the sub-model $P_\varepsilon^h := P + \varepsilon K_O^h$, where $K_O^h \rightarrow \delta_O$, $h \rightarrow 0$. The score function of this model is

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(dP + \varepsilon dK_O^h) = \frac{dK_O^h}{dP}.$$

Thus, if Ψ is Hadamard differentiable, we can use (2) to write

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P_\varepsilon^h) = \left\langle \varphi_P, \frac{dK_O^h}{dP} \right\rangle_P = \int \varphi_P(o) \frac{dK_O^h(o)}{dP(o)} dP(o) = \int \varphi_P(o) dK_O^h(o).$$

Letting $h \rightarrow 0$, we get a candidate for the efficient influence function:

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon \delta_O) = \varphi_P(O).^4$$

³To make the arguments more precise, see Ichimura and Newey [2015].

Finding the canonical gradient for a given Ψ

We can (informally and heuristically³) find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O , where δ_O is the Dirac measure in the point O :

For any $h \in \mathbb{R}_+$, define the sub-model $P_\varepsilon^h := P + \varepsilon K_O^h$, where $K_O^h \rightarrow \delta_O$, $h \rightarrow 0$. The score function of this model is

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \log(dP + \varepsilon dK_O^h) = \frac{dK_O^h}{dP}.$$

Thus, if Ψ is Hadamard differentiable, we can use (2) to write

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P_\varepsilon^h) = \left\langle \varphi_P, \frac{dK_O^h}{dP} \right\rangle_P = \int \varphi_P(o) \frac{dK_O^h(o)}{dP(o)} dP(o) = \int \varphi_P(o) dK_O^h(o).$$

Letting $h \rightarrow 0$, we get a candidate for the efficient influence function:

$$\left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \Psi(P + \varepsilon \delta_O) = \varphi_P(O).^4$$

³To make the arguments more precise, see Ichimura and Newey [2015].

⁴One should then verify that the found the candidate φ_P fulfills (2) for *all* parametric sub-model. In addition, if we impose restrictions on \mathcal{P} such that $\dot{\mathcal{P}}_P$ is a proper subset of $\mathcal{L}_0^2(P)$ we also need to check that $\varphi_P \in \dot{\mathcal{P}}_P$.

Exercise 1 – efficient influence function of the toy example

Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_X for $\Psi(P) = F_P(x)$ where

$$F_P(x) = P(X \leq x) = \int_{-\infty}^x dP(z)$$

is the cumulative distribution function at x .

Exercise 1 – efficient influence function of the toy example

Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_X for $\Psi(P) = F_P(x)$ where

$$F_P(x) = P(X \leq x) = \int_{-\infty}^x dP(z)$$

is the cumulative distribution function at x .

Strategy

1. Write the target parameter $\Psi(P)$ explicitly as a functional of P .

Exercise 1 – efficient influence function of the toy example

Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_X for $\Psi(P) = F_P(x)$ where

$$F_P(x) = P(X \leq x) = \int_{-\infty}^x dP(z)$$

is the cumulative distribution function at x .

Strategy

1. Write the target parameter $\Psi(P)$ explicitly as a functional of P .
2. Calculate the ordinary derivative of

$$\varepsilon \longmapsto \Psi(P + \varepsilon \delta_X).$$

and evaluate at $\varepsilon = 0$.

Exercise 2 – the treatment-specific mean response

Let $O = (Y, A, X) \sim P$. Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O with

$$\Psi(P) = \mathbb{E}_P[f_P(1, X)] = \int f_P(1, x) \, d\mu_P(x),$$

with $f_P(a, x) = \mathbb{E}_P[Y \mid A = a, X = x]$ and μ_P the marginal distribution of X .

Exercise 2 – the treatment-specific mean response

Let $O = (Y, A, X) \sim P$. Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O with

$$\Psi(P) = \mathbb{E}_P[f_P(1, X)] = \int f_P(1, x) \, d\mu_P(x),$$

with $f_P(a, x) = \mathbb{E}_P[Y \mid A = a, X = x]$ and μ_P the marginal distribution of X .

Strategy

1. Write the target parameter $\Psi(P)$ explicitly as a functional of P .

Exercise 2 – the treatment-specific mean response

Let $O = (Y, A, X) \sim P$. Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O with

$$\Psi(P) = \mathbb{E}_P[f_P(1, X)] = \int f_P(1, x) \, d\mu_P(x),$$

with $f_P(a, x) = \mathbb{E}_P[Y \mid A = a, X = x]$ and μ_P the marginal distribution of X .

Strategy

1. Write the target parameter $\Psi(P)$ explicitly as a functional of P .
 - 1.1 Write f_P and μ_P explicitly as functions of P .

Exercise 2 – the treatment-specific mean response

Let $O = (Y, A, X) \sim P$. Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O with

$$\Psi(P) = \mathbb{E}_P[f_P(1, X)] = \int f_P(1, x) \, d\mu_P(x),$$

with $f_P(a, x) = \mathbb{E}_P[Y \mid A = a, X = x]$ and μ_P the marginal distribution of X .

Strategy

1. Write the target parameter $\Psi(P)$ explicitly as a functional of P .
 - 1.1 Write f_P and μ_P explicitly as functions of P .
 - 1.2 Insert the these into the definition of Ψ above.

Exercise 2 – the treatment-specific mean response

Let $O = (Y, A, X) \sim P$. Find a candidate for the efficient influence function by calculating the Gâteaux derivative of Ψ at δ_O with

$$\Psi(P) = \mathbb{E}_P[f_P(1, X)] = \int f_P(1, x) d\mu_P(x),$$

with $f_P(a, x) = \mathbb{E}_P[Y \mid A = a, X = x]$ and μ_P the marginal distribution of X .

Strategy

1. Write the target parameter $\Psi(P)$ explicitly as a functional of P .
 - 1.1 Write f_P and μ_P explicitly as functions of P .
 - 1.2 Insert the these into the definition of Ψ above.
2. Calculate the ordinary derivative of

$$\varepsilon \mapsto \Psi(P + \varepsilon \delta_O).$$

and evaluate at $\varepsilon = 0$.

Exercise 2 – first step

Express $\Psi(P) = \int_{\mathcal{X}} f(1, x) \mu_P(dx)$ explicitly as a functional of P :

Exercise 2 – first step

Express $\Psi(P) = \int_{\mathcal{X}} f(1, x) \mu_P(dx)$ explicitly as a functional of P : We can write

$$\mu_P(x) = \sum_{a=0}^1 \int_{\mathcal{Y}} P(dy, a, x),$$

Exercise 2 – first step

Express $\Psi(P) = \int_{\mathcal{X}} f(1, x) \mu_P(dx)$ explicitly as a functional of P : We can write

$$\mu_P(x) = \sum_{a=0}^1 \int_{\mathcal{Y}} P(dy, a, x),$$

and

$$f(1, x) = \int_{\mathcal{Y}} y \frac{P(dy, 1, x)}{\int_{\mathcal{Y}} P(dy, 1, x)},$$

Exercise 2 – first step

Express $\Psi(P) = \int_{\mathcal{X}} f(1, x) \mu_P(dx)$ explicitly as a functional of P : We can write

$$\mu_P(x) = \sum_{a=0}^1 \int_{\mathcal{Y}} P(dy, a, x),$$

and

$$f(1, x) = \int_{\mathcal{Y}} y \frac{P(dy, 1, x)}{\int_{\mathcal{Y}} P(dy, 1, x)},$$

so

$$\Psi(P_\varepsilon) = \int_{\mathcal{X}} \left\{ \frac{\int_{\mathcal{Y}} y P_\varepsilon(dy, 1, x)}{\int_{\mathcal{Y}} P_\varepsilon(dy, 1, x)} \left(\sum_{a=0}^1 \int_{\mathcal{Y}} P_\varepsilon(dy, a, dx) \right) \right\}.$$

References

- P. J. Bickel, Y. Ritov, et al. Nonparametric estimators which can be "plugged-in". *The Annals of Statistics*, 31(4):1033–1053, 2003.
- 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.
- R. D. Gill, J. A. Wellner, and J. Præstgaard. Non-and semi-parametric maximum likelihood estimators and the von mises method (part 1)[with discussion and reply]. *Scandinavian Journal of Statistics*, pages 97–128, 1989.
- H. Ichimura and W. K. Newey. The influence function of semiparametric estimators. *arXiv preprint arXiv:1508.01378*, 2015.
- J. Reeds. On the definition of von mises functionals. *Ph. D. dissertation, Harvard Univ.*, 1976.
- R. J. Serfling. *Approximation theorems of mathematical statistics*, volume 162. John Wiley & Sons, 1980.
- A. W. van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.