

Topic 4: Making decisions

Eric B. Laber

Department of Statistical Science, Duke University

Statistics 561



On choices

*Dad always thought laughter was the best medicine, which
I guess is why several of us died of tuberculosis.*
—Robert Koch



On decisions

When you're ten years old, and a car drives by and splashes a puddle of water all over you, it's hard to decide if you should go to school like that or try to go home and change and probably be late. So while he was trying to decide, I drove by and splashed him again..

—Bertrand Russell



Warm-up (5 minutes)

- ▶ Explain to your group
 - ▶ What is a randomized clinical trial? Why do we randomize?
 - ▶ What is confounding?
 - ▶ What is a one-armed bandit?
- ▶ True or false
 - ▶ Regression + randomization = causality
 - ▶ Bandit problems were invented by computer scientists
 - ▶ Laber has had food poisoning from Pizza Hut twice

Decision problems

- ▶ Nearly all statistical analyses drive decisions
 - ▶ Estimate treatment effect \Rightarrow treatment recommendations
 - ▶ Identify gene associated with disease \Rightarrow follow-up study
 - ▶ Model click-through-rate as function of customer + ad attributes \Rightarrow ad selection for website
 - ▶ Forecast product demand \Rightarrow manufacturing decisions
 - ▶ Model wins-above-replacement \Rightarrow contract decisions
 - ▶ ...

Roadmap

- ▶ One-stage decision problems
- ▶ K-stage decision problems
- ▶ Contextual bandits
- ▶ Markov Decision Problems
- ▶ Freedom!



One-stage setup

- ▶ Observe $\{(\mathbf{X}_i, A_i, Y_i)\}_{i=1}^n$ iid from P
 - ▶ $\mathbf{X} \in \mathcal{X} \subseteq \mathbb{R}^p$ covariates (decision context)
 - ▶ $A \in \mathcal{A} = \{-1, 1\}$ action (treatment, intervention, decision, etc.)
 - ▶ $Y \in \mathbb{R}$ utility (outcome, output, reward, etc.) higher is better
- ▶ Goal: select actions to maximize expected utility

Policies

- ▶ $\psi : \mathcal{X} \rightarrow 2^{\mathcal{A}}$ is set of allowable actions, i.e., $\psi(\mathbf{x}) \subseteq \mathcal{A} \setminus \emptyset$
- ▶ Policy $\pi : \mathcal{X} \rightarrow \mathcal{A}$ such that $\pi(\mathbf{x}) \in \psi(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$
 - ▶ Under π decision maker will select action $\pi(\mathbf{x})$ in context \mathbf{x}
 - ▶ Define $V(\pi) \triangleq \mathbb{E}^{\pi} Y$ to be expected utility if actions are selected according the policy π
 - ▶ Optimal policy satisfies $V(\pi^{\text{opt}}) \geq V(\pi)$ for all π

Formalizing the optimal policy

- ▶ Potential outcome $Y^*(a)$ under action a , i.e., the outcome under action a (which may be contrary to what was observed)
 - ▶ Imagine each individual has two potential outcomes $Y^*(1)$ and $Y^*(-1)$, one associated with each action
 - ▶ The potential outcome under policy π is

$$Y^*(\pi) = Y^*(1)1_{\pi(\mathbf{x})=1} + Y^*(-1)1_{\pi(\mathbf{x})=-1}$$

formally, the value of a policy is $V(\pi) = \mathbb{E}Y^*(\pi)$

laber draws a table illustrating potential outcomes



Identifying the optimal policy

- ▶ Optimal policy defined in terms of potential outcomes, need to link to data-generating model
- ▶ Standard causal assumptions
 - ▶ No unmeasured confounders: $\{Y^*(1), Y^*(-1)\} \perp A | \mathbf{X}$
 - ▶ Consistency: $Y = Y^*(A)$, i.e., outcome is potential outcome under action taken
 - ▶ Positivity: there exists $\epsilon > 0$ such that $P(A = a | \mathbf{X} = \mathbf{x}) \geq \epsilon$ for (almost) all $\mathbf{x} \in \mathcal{X}$

No interference: the action assigned to one unit do affect outcomes of others

No unmeasured confounders

- ▶ Actions may be selected according to X and the perceived impact on the outcome Y , e.g., clinical decisions
 - ▶ No unmeasured confounders says we captured all factors affecting action selection and the outcome
 - ▶ One minute: construct example in which this assumption is violated

Spillover effects

- ▶ One minute: generate three examples where this is violated



A sad fact about your life

- ▶ Requisite causal assumption are not testable using observed data
 - ▶ There is no test, procedure, etc. that can be applied to the observed data (no matter how much there is) to evaluated needed causal conditions
 - ▶ Randomization ensures no unmeasured confounders by construction
 - ▶ Must use external information: i.e., knowledge of underlying science, richness of the features \mathbf{X} , etc.

Regression-based characterization of optimal regime

- ▶ Define $Q(\mathbf{x}, a) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x}, A = a)$
- ▶ Under standard causal assumptions

$$V(\pi) = \mathbb{E} Q \{ \mathbf{X}, \pi(\mathbf{X}) \}$$

given this expression suggest an estimator of π^{opt}

Derivation of regression-based estimator



Q-learning: part I

- Bound: for any policy π it follows that

$$V(\pi) = \mathbb{E} Q\{\mathbf{X}, \pi(\mathbf{X})\} \leq \mathbb{E} \sup_{a \in \psi(\mathbf{X})} Q(\mathbf{X}, a)$$

- Note that the policy

$$\pi^{\text{opt}}(\mathbf{x}) = \arg \max_{a \in \psi(\mathbf{x})} Q(\mathbf{x}, a)$$

attains this bound and is thus optimal

Q-learning: part I cont'd

- ▶ Idea: estimate $Q(\mathbf{x}, a)$ by regressing Y on \mathbf{X}, A to obtain $\hat{Q}_n(\mathbf{x}, a)$ and subsequently $\hat{\pi}_n(\mathbf{x}) = \arg \max_{a \in \psi(\mathbf{x})} \hat{Q}_n(\mathbf{x}, a)$
- ▶ Ex. suppose $\psi(\mathbf{x}) = \{-1, 1\}$ and posit linear model $Q(\mathbf{x}, a; \beta) = \mathbf{x}_0^\top \beta_0 + a \mathbf{x}_1^\top \beta_1$ where $\beta = (\beta_0^\top, \beta_1^\top)^\top$ and $\mathbf{x}_0, \mathbf{x}_1$ features of \mathbf{x}
 - ▶ $\hat{\beta}_n \triangleq \arg \min_{\beta} \mathbb{P}_n \{Y - Q(\mathbf{X}, A; \beta)\}^2$
 - ▶ $\hat{\pi}_n(\mathbf{x}) = \arg \max_a Q(\mathbf{x}, a; \hat{\beta}_n) = \text{sign}(\mathbf{x}_1^\top \hat{\beta}_{1,n})$

Q-learning: flexible models

- ▶ No need to restrict to linear models
- ▶ Construct estimator $\hat{Q}_n(\mathbf{x}, a)$ of $Q(\mathbf{x}, a) = \mathbb{E}(Y|\mathbf{X} = \mathbf{x}, A = a)$ using random forest and then take $\hat{\pi}_n(\mathbf{x}) = \arg \max_{a \in \psi(\mathbf{x})} \hat{Q}_n(\mathbf{x}, a)$
- ▶ Quantum theory of decision rules: if a regression method exists, someone, somewhere has published a paper applying it in Q-learning
 - ▶ Boosting
 - ▶ Neural nets
 - ▶ Nearest neighbors
 - ▶ Gaussian processes
 - ▶ NP-Bayes
 - ▶ ...

Advantage learning v1

- Note that we can always write

$$Q(\mathbf{x}, a) = \mu(\mathbf{x}) + a\Delta(\mathbf{x})$$

where

$$\mu(\mathbf{x}) = \frac{Q(\mathbf{x}, 1) + Q(\mathbf{x}, -1)}{2}, \text{ and } \Delta(\mathbf{x}) = \frac{Q(\mathbf{x}, 1) - Q(\mathbf{x}, -1)}{2}$$

only need to estimate Δ to identify optimal policy

Advantage learning v1 cont'd

- ▶ Write $Q(\mathbf{x}, a) = \tilde{\mu}(\mathbf{x}) + \{(1 + a)/2 - q(\mathbf{x})\} \tilde{\Delta}(\mathbf{x})$, where $\tilde{\mu}(\mathbf{x}) = \mu(\mathbf{x})(1/2 - q(\mathbf{x}))\Delta(\mathbf{x})$, $q(\mathbf{x}) = P(A = 1|\mathbf{X} = \mathbf{x})$, and $\tilde{\Delta}(\mathbf{x}) = 2\Delta(\mathbf{x})$
- ▶ Advantage learning solves

$$\hat{\mu}_n, \hat{\Delta}_n = \arg \min_{\mu, \Delta} \mathbb{P}_n [Y - \mu(\mathbf{X}) + \{(A + 1)/2 - q(\mathbf{X})\} \Delta(\mathbf{X})]^2$$

so that the estimated optimal policy is given by

$$\hat{\pi}_n(\mathbf{x}) = \text{sign} \left\{ \hat{\Delta}_n(\mathbf{x}) \right\} \text{ if } \psi(\mathbf{x}) = \{-1, 1\}$$

Why A-learning v1 works



Why A-learning v1 works



Advantage learning v2

- ▶ Let $A \in \mathcal{A} \subset \mathbb{R}$ be more general action space and for simplicity assume that $\psi(\mathbf{x}) = \mathcal{A}$ for all \mathbf{x}
- ▶ Define $\Gamma(\mathbf{x}, a) = Q(\mathbf{x}, a) - \max_a Q(\mathbf{x}, a)$ so that $\Gamma(\mathbf{x}, a) \leq 0$ and $\Gamma(\mathbf{x}, a) = 0$ if $a = \pi^{\text{opt}}(\mathbf{x})$ then

$$Q(\mathbf{x}, a) = \omega(\mathbf{x}) + \Gamma(\mathbf{x}, a),$$

where $\omega(\mathbf{x}) = \max_a Q(\mathbf{x}, a)$

- ▶ $\Gamma(\mathbf{x}, a)$ is the advantage of selecting action a in context \mathbf{x} and $\pi^{\text{opt}}(\mathbf{x}) = \arg \max_a \Gamma(\mathbf{x}, a)$

Advantage learning v2

- ▶ Goal: estimate Γ without estimating ω
- ▶ Claim: Γ satisfies

$$\Gamma = \arg \min_{\gamma} \mathbb{E} \left\{ Y - \gamma(\mathbf{X}, A) + \int \gamma(\mathbf{X}, a) p(a|\mathbf{X}) d\eta(a) \right\}^2,$$

where η is a dominating measure

- ▶ Note there's no ω in the above expression!

Why A-learning v2 works



Why A-learning v2 works



Classification-based representation: quiz

- ▶ Warm-up: discuss with your stats group (3.25 minutes)
 - ▶ What is sampling bias? When does it occur?
 - ▶ What is the Horvitz-Thompson estimator?
 - ▶ What is cost-sensitive classification?
- ▶ True or false
 - ▶ Double sampling is when you use the same spoon twice in the same trough of bean dip (I'm looking at you Chad)
 - ▶ Survey sampling is mostly relegated to the census and marketing
 - ▶ Laber discovered 'Cart Narcs' on YouTube at 2AM and binged every episode could find (let's keep the tough questions to a minimum)



Classification-based representation overview

- ▶ Define the propensity score $P(A = 1|\mathbf{X} = \mathbf{x})$
- ▶ Under our standard causal conditions

$$V(\pi) = P \left\{ \frac{Y 1_{A=\pi(\mathbf{X})}}{P(A|\mathbf{X})} \right\}$$

this is the classic Horvitz-Thompson (HT) estimator from survey sampling!

- ▶ HT more commonly known as inverse probability weighted (IPW) representation; What's the intuition behind this estimator?

HT representation details



IPWE

- ▶ Inverse probability weighted estimator (IPWE) of $V(\pi)$

$$\hat{V}_n^{\text{IPWE}}(\pi) = \mathbb{P}_n \left\{ \frac{Y 1_{A=\pi(\mathbf{X})}}{\hat{P}_n(A|\mathbf{X})} \right\},$$

where $\hat{P}_n(a|\mathbf{x})$ is the estimated propensity score, e.g., estimated using logistic regression, nnet, etc.

- ▶ Estimated optimal decision rule

$$\hat{\pi}_n = \arg \max_{\pi \in \Pi} \mathbb{P}_n \left\{ \frac{Y 1_{A=\pi(\mathbf{X})}}{P(A|\mathbf{X})} \right\}$$

feasible for some classes Π if n isn't too large but generally not tractable



Linking IPWE with cost-sensitive classification

- ▶ A tale of brilliance and the despair of being too late

$$\begin{aligned}\arg \max_{\pi} \widehat{V}_n^{\text{IPWE}}(\pi) &= \arg \max_{\pi} \mathbb{P}_n \left\{ \frac{Y 1_{A\pi(\mathbf{X}) \geq 0}}{\widehat{P}_n(A|\mathbf{X})} \right\} \\ &= \arg \min_{\pi} \mathbb{P}_n \left\{ \frac{Y 1_{A\pi(\mathbf{X}) < 0}}{\widehat{P}_n(A|\mathbf{X})} \right\} \\ &= \arg \min_{\pi} \mathbb{P}_n \widehat{W}_n 1_{YA\pi(\mathbf{X}) < 0},\end{aligned}$$

where $\widehat{W}_n = |Y|/\widehat{P}_n(A\mathbf{X})$

- ▶ Egad! This looks like a weighted classification problem!

Derivation of weighted classification



Convex surrogates and optimal decisions

- ▶ Consider decision rules $\pi(\mathbf{x}) = \text{sign}\{f(\mathbf{x})\}$ where $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is a (generally smooth) fn, e.g., $f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}$
- ▶ Let $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ be one of our convex surrogates from classification, e.g., hinge loss, logistic loss, exp loss, etc.
- ▶ Let \mathcal{F} be a class of functions from \mathbb{R}^p into \mathbb{R} , the outcome weighted estimator (OWL) is given by $\hat{\pi}_n(\mathbf{x}) = \text{sign}\{\hat{f}_n(\mathbf{x})\}$ where

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \mathbb{P}_n \widehat{W}_n \phi \{Y A f(\mathbf{X})\}$$

OWL example

- ▶ OWL of linear decision rule is given by

$$\hat{\beta}_n = \arg \min_{\beta} \mathbb{P}_n \widehat{W}_n \phi(YA\mathbf{X}^T\beta < 0) + \lambda \|\beta\|^2$$

so that $\hat{\pi}_n^{\text{OWL}}(\mathbf{x}) = \text{sign}(\mathbf{x}^T \hat{\beta}_n)$

- ▶ Same theory applies as in classification!

Fact: IPWE is terrible

- ▶ IPWE is highly unstable b/c small propensities inflate variance and only a fraction of data is used, e.g., in a randomized trial only 1/2 the data appear (on average) in the weighted sum¹
- ▶ Better: AIPWE which is given by

$$\hat{V}_n^{\text{AIPWE}}(\pi) = \mathbb{P}_n \left[\frac{Y 1_{A=\pi(\mathbf{X})}}{\hat{P}_n(A|\mathbf{X})} - \frac{1_{A=\pi(\mathbf{X})} - \hat{P}_N(A|\mathbf{X})}{\hat{P}_n(A|\mathbf{X})} \hat{Q}_n\{\mathbf{X}, \pi(\mathbf{X})\} \right]$$

¹All the data are used in the estimation of the propensity score.

Why AIPWE work? Double robustness.

Why AIPWE work? Double robustness. cont'd

Roadmap

- ▶ One-stage decision problems
- ▶ **K-stage decision problems**
- ▶ Contextual bandits
- ▶ Markov Decision Problems
- ▶ Freedom!



Sequential decision problems

- ▶ Multi-stage treatment strategies
- ▶ Planning a transition to carbon-neutrality
- ▶ Navigation
- ▶ Optimizing vaccine distribution over the next three months
- ▶ ...

Setup

- ▶ Observe $\{\mathbf{X}_{1,i}, A_{1,i}, \mathbf{X}_{2,i}, A_{2,i}, \dots, \mathbf{X}_{T,i}, A_{T,i}, Y_i\}_{i=1}^n$ comprising n i.i.d. trajectories drawn from unknown distn P
 - ▶ $\mathbf{X}_t \in \mathcal{X} \subseteq \mathbb{R}^p$ measurements at time t
 - ▶ $A_t \in \mathcal{A} = \{-1, 1\}$ action/decision/txt at time t
 - ▶ $Y \in \mathbb{R}$ outcome/utility coded so that higher is better
- ▶ Define $\mathbf{H}_1 = \mathbf{X}_1$ and $\mathbf{H}_t = (\mathbf{H}_{t-1}, A_{t-1}, \mathbf{X}_t)$ to be history, i.e., info available to decision maker before decision at time t

Multi-stage policy

- ▶ Policy $\pi = (\pi_1, \dots, \pi_T)$ so that $\pi_t : \text{supp } \mathbf{H}_t \rightarrow \mathcal{A}$
 - ▶ Given info $\mathbf{H}_t = \mathbf{h}_t$, select decision $\pi_t(\mathbf{h}_t)$
 - ▶ Indices $t = 1, \dots, T$ need not correspond to fixed calendar times
- ▶ Optimal regime π^{opt} maximizes mean utility to select decisions at each decision point; formalize using potential outcomes

Ex. Treatment regime: mHealth for PTSD in cancer patients (PI S. Smith)

First stage decision rule

If distress ≥ 3 **then:** Cancer Distress Coach (CDC)

Else if PTSD symptom score ≥ 20 **then:** CDC

Else: usual care

Second stage decision rule

If responder **then:** continue first stage treatment

Else if using CDC and PSTD change ≥ 3 **then:** add mCoaching

Else if using CDC and distress ≥ 4 **then:** add FaceTime CBT

Else FaceTime CBT only

Characterizing the optimal policy

- ▶ Let $\bar{\mathbf{a}}_t = (a_1, \dots, a_t)$ denote decision sequence
 - ▶ $H_t^*(\bar{\mathbf{a}}_{t-1})$ potential history under $\bar{\mathbf{a}}_{t-1}$
 - ▶ $Y^*(\bar{\mathbf{a}}_t)$ potential outcome under $\bar{\mathbf{a}}_T$
 - ▶ Potential outcome under policy π is thus

$$Y^*(\pi) = \sum_{\bar{\mathbf{a}}_T} Y^*(\bar{\mathbf{a}}_T) \prod_{v=1}^T 1_{\pi_v \{H_v^*(\bar{\mathbf{a}}_{v-1})\} = a_v}$$

- ▶ Define $V(\pi) \triangleq \mathbb{E} Y^*(\pi)$ opt policy satisfies $V(\pi^{\text{opt}}) \geq V(\pi)$ for all other policies π

Standard causal assumptions (generalized)

- ▶ No unmeasured confounders (sequential ignorability)

$$\{Y^*(\bar{\mathbf{a}}_T), \mathbf{H}_t^*(\bar{\mathbf{a}}_{t-1}), t = 1, \dots, T\} \perp A_t | \mathbf{H}_t$$

for all $t = 1, \dots, T$

- ▶ Consistency, $Y = Y^*(\bar{\mathbf{A}}_T)$ and $\mathbf{H}_t = \mathbf{H}_t^*(\bar{\mathbf{A}}_{t-1})$ for all $t = 1, \dots, T$
- ▶ Positivity, $P(A_t = a | \mathbf{H}_t) > 0$ with probability one for all $a \in \mathcal{A}$

Regression-based characterization of optimal regime

- Define Q-functions

$$Q_T(\mathbf{h}_T, a_T) \triangleq \mathbb{E}(Y | \mathbf{H}_T = \mathbf{h}_T, A_T = a_T)$$

$$Q_t(\mathbf{h}_t, a_t) \triangleq \mathbb{E} \left\{ \max_{a_{t+1}} Q_{t+1}(\mathbf{H}_{t+1}, a_{t+1}) | \mathbf{H}_t = \mathbf{h}_t, A_t = a_t \right\}$$

for $t = T - 1, \dots, 1$

- Dynamic programming (Bellman 1957)

$$\pi_t^{\text{opt}}(\mathbf{h}_t) = \arg \max_{a_t} Q_t(\mathbf{h}_t, a_t)$$

for $t = 1, \dots, T$

Q-learning (again)

- ▶ Obvious idea: estimate Q -functions using regression
- ▶ E.g., posit linear models $Q_t(\mathbf{h}_t, a_t; \boldsymbol{\beta}_t) = \mathbf{h}_{t,0}^\top \boldsymbol{\beta}_{t,0} + a_t \mathbf{h}_{t,1}^\top \boldsymbol{\beta}_{t,1}$ where $\mathbf{h}_{t,0}$ and $\mathbf{h}_{t,1}$ are features of \mathbf{h}_t and $\boldsymbol{\beta}_t = (\boldsymbol{\beta}_{t,0}^\top, \boldsymbol{\beta}_{t,1}^\top)^\top$ are unknown coeffs

- ▶ Compute $\hat{\boldsymbol{\beta}}_{T,n} = \arg \min_{\boldsymbol{\beta}} \mathbb{P}_n \{Y - Q(\mathbf{H}_T, A_T; \boldsymbol{\beta}_T)\}^2$

- ▶ Recursively for $t = T - 1, \dots, 1$ compute

$$\hat{\boldsymbol{\beta}}_{t,n} = \arg \min_{\boldsymbol{\beta}_t} \mathbb{P}_n \left\{ \max_{a_{t+1}} Q(\mathbf{H}_{t+1}, a_{t+1}; \hat{\boldsymbol{\beta}}_{t+1,n}) - Q_t(\mathbf{H}_t, A_t; \boldsymbol{\beta}_t) \right\}^2$$

- ▶ $\hat{\pi}_t(\mathbf{h}_t) = \arg \max_{a_t} Q_t(\mathbf{h}_t, a_t; \hat{\boldsymbol{\beta}}_{t,n})$

Q-learning cont'd

- ▶ No need to stick to linear models can use ML/nonpar models
- ▶ Highly extensible, e.g., if regression works, Q-learning works (mostly)
 - ▶ Censoring
 - ▶ High-dimensional inputs
 - ▶ Image, text, other complex data structures
 - ▶ Functional inputs
 - ▶ ...

Nonregularity: the talk you parents never gave you

- ▶ Inference for Q -functions notoriously difficult
 - ▶ Max operator \Rightarrow non-smooth functional of distn
 - ▶ No asymptotically unbiased estimator exists
 - ▶ No regular estimator exists
 - ▶ Bootstrap, normal approximations, and credible regions can perform poorly

laber derives limiting distribution in Q-learning linear models

laber derives limiting distribution in Q-learning

linear models cont'd

Simple treatment selection problem

- ▶ Let P be fixed but unknown distn on \mathbb{R}^p ($p \geq 2$) and define:

$$\mu_0 \triangleq P\mathbf{X} = \int \mathbf{x} dP(\mathbf{x})$$

$$\theta_0 \triangleq \bigvee_{j=1}^p \mu_{0,j} = \max(\mu_{0,1}, \mu_{0,2}, \dots, \mu_{0,p})$$

- ▶ Observe $\mathbf{X}_1, \dots, \mathbf{X}_n$ drawn i.i.d. from P and define:

$$\hat{\mu}_n \triangleq \mathbb{P}_n \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

$$\hat{\theta}_n \triangleq \bigvee_{j=1}^p \hat{\mu}_{n,j} = \max(\hat{\mu}_{n,1}, \hat{\mu}_{n,2}, \dots, \hat{\mu}_{n,p})$$

- ▶ Inference for θ_0 based on limiting distn of $\sqrt{n}(\hat{\theta}_n - \theta_0)$

Fixed parameter asymptotics for $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- For any $\nu \in \mathbb{R}^p$ define $\mathfrak{U}(\nu) = \arg \max_j \nu_j$ which may contain more than one element if no unique maximizer

Fixed parameter asymptotics for $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- For any $\nu \in \mathbb{R}^p$ define $\mathfrak{U}(\nu) = \arg \max_j \nu_j$ which may contain more than one element if no unique maximizer

Lemma

Assume regularity conditions ensuring that $\sqrt{n}(\mathbb{P}_n - P)\mathbf{X} \rightsquigarrow N(0, \Sigma)$. Then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow \bigvee_{j \in \mathfrak{U}(\mu_0)} Z_j,$$

where $\mathbf{Z} \sim N(0, \Sigma)$.

Proof overview

Proof

Let $\mathcal{E}_n = 1 \left\{ \max_{k \notin \mathfrak{L}(\mu_0)} \hat{\mu}_{n,k} \geq \max_{k \in \mathfrak{L}(\mu_0)} \hat{\mu}_{n,k} \right\}$ then $\mathcal{E}_n = o_P(1)$. Write

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_0) &= \bigvee_{j \in \mathfrak{L}(\mu_0)} \sqrt{n}(\hat{\mu}_{n,j} - \theta_0) \\ &\quad + \left\{ \bigvee_{k \notin \mathfrak{L}(\mu_0)} \sqrt{n}(\hat{\mu}_{n,k} - \theta_0) - \bigvee_{j \in \mathfrak{L}(\mu_0)} \sqrt{n}(\hat{\mu}_{n,j} - \theta_0) \right\} \mathcal{E}_n \\ &= \bigvee_{j \in \mathfrak{L}(\mu_0)} \sqrt{n}(\hat{\mu}_{n,j} - \mu_{0,j}) + r_n, \end{aligned}$$

where we have used $\mu_{0,j} = \theta_0$ for all $j \in \mathfrak{L}(\mu_0)$. The result follows from the CLT, continuous mapping theorem, and showing that $r_n = o_P(1)$.

Looking closer at $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- ▶ Lemma states $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow \bigvee_{j \in \mathfrak{A}(\mu_0)} Z_j$, where $Z_j \sim N(0, \Sigma)$
- ▶ Limiting distribution of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ depends abruptly on μ_0
 - ▶ Ex., if $p = 2$ and $\mu_0 = (0, 0)^\top$ and $\Sigma = I_2$ then the limit is the max of two independent normals but if $\mu_0 = (0, \epsilon)^\top$ for $\epsilon > 0$, say $\epsilon = 10^{-20}$, the limit is standard normal!
 - ▶ Asymptotic distn depends on maximizers of μ_0 and associated submatrix of Σ but not gaps in μ_0 this is problematic as the finite sample distn can depend critically on these gaps

Finite sample behavior of $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \sim N(\mu_0, \mathbf{I}_p)$ and μ_0 has a unique maximizer $\mu_1 > \mu_j$ for $j = 2, \dots, p$, then

$$P \left\{ \sqrt{n} (\hat{\theta}_n - \theta_0) \leq t \right\} = \Phi(t) \prod_{j=2}^p \Phi \left\{ t + \sqrt{n}(\theta_0 - \mu_{0,j}) \right\}$$

which can be far from asymptotic limit, $\Phi(t)$, when gaps $\theta_0 - \mu_{0,j}$ are small relative to \sqrt{n}

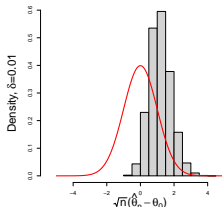
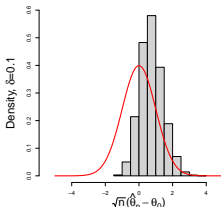
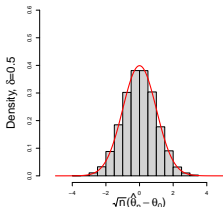
Finite sample behavior of $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \sim N(\mu_0, \mathbf{I}_p)$ and μ_0 has a unique maximizer $\mu_1 > \mu_j$ for $j = 2, \dots, p$, then

$$P\left\{\sqrt{n}(\hat{\theta}_n - \theta_0) \leq t\right\} = \Phi(t) \prod_{j=2}^p \Phi\left\{t + \sqrt{n}(\theta_0 - \mu_{0,j})\right\}$$

which can be far from asymptotic limit, $\Phi(t)$, when gaps $\theta_0 - \mu_{0,j}$ are small relative to \sqrt{n}

- Ex. set $p = 6$, $n = 100$, and $\mu_{0,j} = \mu_{0,1} - \delta$ for $j \geq 2$



Problems with fixed parameter asymptotics

- ▶ Asymptotic approximation to distn of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ poor as it fails to account for uncertainty in $\mathfrak{L}(\mu_0)$
 - ▶ No uncertainty about $\mathfrak{L}(\mu_0)$ as $n \rightarrow \infty$, i.e., power of $H_{0,j} : \mu_{0,j} = \theta_0$ goes to one
 - ▶ Mismatch of asymptotic framework and finite sample behavior
- ▶ Need asymptotic framework that retains uncertainty (instability) about $\mathfrak{L}(\mu_0)$ as n diverges

Moving parameter asymptotics

- ▶ Idea: consider sequence of generative models so that key features of sampling distribution retained asymptotically
 - ▶ Commonly used for comparing power of hypothesis tests
 - ▶ $1/\sqrt{n}$ perturbations of generative model²
- ▶ Consider $\mathbf{X}_{n,1}, \dots, \mathbf{X}_{n,n}$ drawn i.i.d from P_n

$$\mu_{0,n} \triangleq P_n \mathbf{X} = \mu_0 + \mathbf{s}/\sqrt{n}$$
$$\theta_{0,n} \triangleq \bigvee_{j=1}^p \mu_{0,n,j} = \max(\mu_{0,n,1}, \dots, \mu_{0,n,p}),$$

where $\mathbf{s} \in \mathbb{R}^p$ is a called local parameter

²It can also be useful to think of local alternatives as a kind of derivative, as we'll note later.

Moving parameter asymptotics cont'd

- ▶ Moving parameter asymptotic allow for gaps $\theta_{0,n} - \mu_{0,n,j}$ that shrink with sample size \Rightarrow uncertainty about $\mathfrak{L}(\mu_0)$ persists

Lemma

Let $\mathbf{s} \in \mathbb{R}^P$ be fixed. Assume for each n that $\{\mathbf{X}_{n,i}\}_{i=1}^n$ are drawn i.i.d. from P_n , which satisfy (i) $P_n \mathbf{X} = \mu_0 + \mathbf{s}/\sqrt{n}$, and (ii) $\sqrt{n}(\mathbb{P}_n - P_n)\mathbf{X} \rightsquigarrow N(0, \Sigma)$. Then, under P_n ,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow \bigvee_{j \in \mathfrak{L}(\mu_0)} (Z_j + s_j) - \bigvee_{j \in \mathfrak{L}(\mu_0)} s_j,$$

where $\mathbf{Z} \sim N(0, \Sigma)$.

Moving parameter limit of $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- ▶ Under P_n : $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow \bigvee_{j \in \mathfrak{L}(\mu_0)} (Z_j + s_j) - \bigvee_{j \in \mathfrak{L}(\mu_0)} s_j$
- ▶ Limit depends on local parameter $\mathbf{s} \Rightarrow$ non-regular
 - ▶ Anticipates poor performance of asymptotic approximation
 - ▶ Bootstrap and series approximations invalid w/o modification
 - ▶ Non-regularity is ubiquitous in decision problems (think max operator in Bellman optimality eqns)
- ▶ Want inference procedures that are valid under moving parameter asymptotics

Valid inference via projection

- ▶ Building block idea: projection confidence set
 - ▶ μ_0 is regular
 - ▶ $\theta_0 = \theta_0(\mu_0)$
- ▶ Projection confidence set for max of means
 - ▶ Construct $(1 - \alpha)$ confidence region for μ_0

$$\zeta_{n,1-\alpha} = \{\mu \in \mathbb{R}^p : n(\hat{\mu}_n - \mu)^\top \hat{\Sigma}_n^{-1} (\hat{\mu}_n - \mu) \leq \chi_{p,1-\alpha}^2\}$$

then $P(\mu_{0,n} \in \zeta_{n,1-\alpha}) \geq 1 - \alpha + o_{P_n}(1)$

- ▶ Projection confidence set for θ_0

$$\Gamma_{n,1-\alpha} = \left\{ \theta \in \mathbb{R} : \theta = \bigvee_{j=1}^p \mu_j \text{ for some } \mu \in \zeta_{n,1-\alpha} \right\}$$

Why the projection set works

- ▶ Error of projection set is

$$\begin{aligned} P(\theta_0 \notin \Gamma_{n,1-\alpha}) &= P(\theta_0 \notin \Gamma_{n,1-\alpha}, \mu_0 \notin \zeta_{n,1-\alpha}) \\ &\quad + P(\theta_0 \notin \Gamma_{n,1-\alpha}, \mu_0 \in \zeta_{n,1-\alpha}) \\ &\leq P(\mu_0 \notin \zeta_{n,1-\alpha}) \\ &\leq \alpha + o_P(1), \end{aligned}$$

holds under both fixed and moving parameter frameworks

- ▶ Pros and cons of projection sets
 - ▶ Robust, general, conceptually simple (one-line proof!)
 - ▶ Minimally adaptive \Rightarrow extremely conservative in some settings

Bound-based inference

- ▶ Idea: sandwich non-smooth functional between two smooth bounds and use bootstrap or series approximations on bounds
- ▶ Bounds constructed by taking sup/inf over local perturbations
 - ▶ Valid under moving parameter frameworks by construction
 - ▶ Smallest possible among regular bounds³
 - ▶ Necessarily adaptive because if the functional is smooth (regular) it can serve as its own bound

³In some sense that can be made more precise.

Upper bound for $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- Define conservative (overly inclusive) estimator of $\mathfrak{L}(\mu_0)$ as

$$\hat{\mathfrak{L}}_n(\mu_0) = \left\{ j : \max_k \sqrt{n}(\hat{\mu}_{n,k} - \mu_{n,j}) / \hat{\sigma}_{n,j,k} \leq \tau_n \right\},$$

where $\tau_n \rightarrow \infty$, $\tau_n = o_p(\sqrt{n})$, and $\hat{\sigma}_{n,j,k}$ an estimator of $n\text{Var}(\hat{\mu}_{n,j} - \hat{\mu}_{n,k})$

- Local perturbations

$$\hat{\mathcal{S}}_n(\mu_0) = \left\{ \mathbf{s} \in \mathbb{R}^p : s_j = \mu_{0,j} \text{ if } j \notin \hat{\mathfrak{L}}_n(\mu_0) \right\}$$

- Upper bound on $\sqrt{n}(\hat{\theta}_n - \theta_0)$ is

$$\hat{U}_n = \sup_{\mathbf{s} \in \hat{\mathcal{S}}_n(\mu_0)} \sqrt{n} \left[\bigvee_{j=1}^p \{(\hat{\mu}_{n,j} - \mu_{0,j}) + s_j\} - \bigvee_{j=1}^p s_j \right]$$

lower bound obtained by replace sup with an inf

Bound based confidence interval for θ_0

- Let $\hat{u}_{n,1-\alpha/2}^{(b)}$ be the $(1 - \alpha/2) \times 100$ percentile of \hat{U}_n and $\hat{\ell}_{n,\alpha}^{(b)}$ the $(\alpha/2) \times 100$ percentile of the lower bound \hat{L}_n obtained via bootstrap, then

$$\left[\hat{\theta}_n - \frac{\hat{u}_{n,1-\alpha/2}^{(b)}}{\sqrt{n}}, \hat{\theta}_n - \frac{\hat{\ell}_{n,\alpha/2}^{(b)}}{\sqrt{n}} \right]$$

is a $(1 - \alpha) \times 100\%$ confidence interval for θ_0

Comments on bound-based interval

- ▶ Set $\hat{\mathcal{S}}_n(\mu_0)$ only allows perturbations of means near the max
 - ▶ \sim Pretest used in weakly identifiable problems in econometrics
 - ▶ Local parameter only appears in components that are near maximizers \Rightarrow only want sup / inf to affect these
- ▶ Sup/inf adaptively 'smooths' $\sqrt{n}(\hat{\theta}_n - \theta_0)$ making it regular
- ▶ Alternative to bounding $\sqrt{n}(\hat{\theta}_n - \theta_0)$ and taking quantiles is to bound quantiles directly (may improve precision)

Discussion of max means problem

- ▶ Non-smooth max operator causes instability
 - ▶ Fixed parameter asymptotics provide poor approximation to the finite sample distribution \Rightarrow std inference unreliable
 - ▶ Local (moving) parameter asymptotics faithfully capture instability as sample size grows large \Rightarrow reliable for inference
 - ▶ Projection and bound-based sets two inference procedures that remain valid under local parameter asymptotics

Roadmap

- ▶ One-stage decision problems
- ▶ K-stage decision problems
- ▶ Contextual bandits
- ▶ **Markov Decision Problems**
- ▶ Freedom!



Setup and notation

- ▶ Observe $\left\{ \mathbf{S}_i^1, A_i^1, \dots, \mathbf{S}_i^T, A_i^T, \mathbf{S}^{T+1} \right\}_{i=1}^n$ drawn i.i.d. from P
 - ▶ T observation period, e.g., trial follow-up
 - ▶ $\mathbf{S}^t \in \mathcal{S} \subseteq \mathbb{R}^p$ state at time t
 - ▶ $A \in \mathcal{A} = \{1, \dots, K\}$
- ▶ Utility $U^t = u(\mathbf{S}^t, A^t, \mathbf{S}^{t+1})$, higher is better
- ▶ Assume homogeneous MDP \rightarrow focus on deterministic stationary policies $\pi : \mathcal{S} \rightarrow \mathcal{A}$

What is a homogeneous MDP?

Why is the MPD so powerful?

- ▶ Allows for estimation of optimal policy (to be applied indefinitely) from data collected over a finite horizon
- ▶ Joint distribution determined by
 - ▶ Initial state distn $f(\mathbf{s}^1)$
 - ▶ Transition kernel $f(\mathbf{s}'|\mathbf{s}, a)$, e.g., density of \mathbf{S}_t given \mathbf{S}_t and A_t
 - ▶ If we know these then we know everything!!

Defining an optimal policy

- ▶ Discounted reward $V(\mathbf{s}; \pi) = \mathbb{E}^{\pi} \left(\sum_{v \geq 0} \gamma^v U^{t+v} \mid \mathbf{s}^t = \mathbf{s} \right)$,
optimal regime satisfies $V(\mathbf{s}, \pi^{\text{opt}}) \geq V(\mathbf{s}, \pi)$ for all \mathbf{s} and π
- ▶ Model-based estimation (overview)
 - ▶ Posit model $f(\mathbf{s}^{t+1} \mid \mathbf{s}^t, a; \boldsymbol{\theta})$ indexed by $\boldsymbol{\theta} \in \Theta$
 - ▶ Construct likelihood $\mathcal{L}_n(\boldsymbol{\theta})$ compute $\hat{\boldsymbol{\theta}}_n$ soln to $0 = \nabla_{\boldsymbol{\theta}} \mathcal{L}_n(\boldsymbol{\theta})$
 - ▶ For any π , generate data under $\hat{\boldsymbol{\theta}}_n$ construct estimator $V_{\hat{\boldsymbol{\theta}}_n}(\mathbf{s}; \pi)$ with flexible approximation architecture (e.g., nnets)

Bellman equations

- Warm-up exercise, show:

$$V(\mathbf{S}^t; \pi) = \mathbb{E}^{\pi} \{ U^t + \gamma V(\mathbf{S}^{t+1}; \pi) | \mathbf{S}^t \}$$

and thus for any function $\psi(\mathbf{S}^t)$ it follows that

$$0 = \mathbb{E} \left[\left(\frac{\delta_{A^t}(\mathbf{S}^t)}{\mu^t(A^t | \mathbf{H}^t)} \{ U^t + \gamma V(\mathbf{S}^{t+1}, \pi) - V(\mathbf{S}^t, \pi) \} \right) \psi(\mathbf{S}^t) \right],$$

where δ_w is a point-mass at w , \mathbf{H}^t is the history at time t , and μ^t is the propensity

- Note: term inside parentheses (i.e., omitting ψ) is called importance-weighted TD-error

V-learning

- ▶ Model-free estimation (overview)
 - ▶ Posit model $V(\mathbf{s}, \pi; \boldsymbol{\lambda})$ indexed by $\boldsymbol{\lambda} \in \boldsymbol{\Lambda}$
 - ▶ Solve estimating eqns based on Bellman optimality eqns

$$0 = \mathbb{P}_n \sum_{t=1}^T \left[\frac{\delta_{A^t}(\mathbf{S}^t)}{\mu^t(A^t | \mathbf{H}^t)} \left\{ U^t + \gamma V(\mathbf{S}^{t+1}, \pi; \boldsymbol{\lambda}) - V(\mathbf{S}^t, \pi; \boldsymbol{\lambda}) \right\} \psi(\mathbf{S}^t, \pi; \boldsymbol{\lambda}) \right]$$

- ▶ Optimal weights (per Godambe)

$$\psi(\mathbf{S}^t, \pi; \boldsymbol{\lambda}) = \mathbb{E} \{ \nabla_{\boldsymbol{\lambda}} g(\boldsymbol{\lambda}) | \mathbf{S}^t \} / \mathbb{E} \{ g^2(\boldsymbol{\lambda}) | \mathbf{S}^t \},$$

where $g(\boldsymbol{\lambda})$ is importance-weighted TD-error

- ▶ Est optimal weights, $\hat{\psi}_{\hat{\boldsymbol{\theta}}_n}(\mathbf{S}^t, \pi; \boldsymbol{\lambda})$ using fitted dynamics model

V-learning cont'd

- ▶ Let $\hat{\lambda}^{\pi}$ denote solution to est eqn at π
- ▶ Select reference distribution \mathcal{R} over states⁴
- ▶ Compute estimated optimal policy in class Π as

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \int V(\mathbf{s}, \pi; \hat{\lambda}^{\pi}) d\mathcal{R}(\mathbf{s})$$

⁴How best to select this distribution is a bit of an open question. A typical choice is the estimated distribution of the initial state.



V-learning discussion

- ▶ Does not require correct dynamics model (i.e., model-free)
- ▶ Separate specification of class of policies, Π , which may be chosen to ensure interpretability, satisfaction of logistical/cost constraints, etc.
- ▶ (Blue)devil is in the details, optimization over space of policies can be non-trivial, typically stochastic search methods are used, e.g., SGD over parametric class

Thank you.

`eric.laber@duke.edu`

`laber-labs.com`

