

Contents:

§1 Policy Learning

1. Policy 的定义
2. 无限制时的 optimal policy

§2 Policy Evaluation

1. Policy Learning 的流程
2. Policy evaluation

§3 Empirical-Welfare Maximization

1. regret 的定义
2. Empirical-welfare maximization
3. Empirical-welfare maximization 的 bound

14 Policy learning

So far, we have focused on methods for estimating treatments effects. In many application areas, however, the fundamental goal of performing a causal analysis is not to estimate treatment effects, but rather to guide decision making: We want to understand treatment effects so that we can effectively prescribe treatment and allocate limited resources. In other words, for a set of new units X_{n+1}, \dots , we want to assign it to either placebo or treatment $\pi(X_i) \in \{0, 1\}$ such that the average outcome is optimal.

The problem of learning optimal treatment assignment policies is closely related to but subtly different from the problem of estimating treatment heterogeneity. On one hand, policy learning appears easier: All we care about is assigning people to treatment or to control, and we don't care about accurately estimating treatment effects beyond that. On the other hand, when learning policies, we need to account for considerations that were not present when simply estimating treatment effects: Any policy we actually want to use must be simple enough we can actually deploy it, cannot discriminate on protected

Policy 的定义: 一个 function $\pi: \mathcal{R}^p \mapsto \{0, 1\}$ s.t. we assign treatment when $\pi(X_i) = 1$

$$\Rightarrow Y(\pi(X_i)) = \pi(X_i) \cdot Y(1) + [1 - \pi(X_i)] \cdot Y(0)$$

eg. ① $\pi(X) = 1 \{ \sum_{i=1}^p a_i X_i > 1 \}$ ② binary decision tree

characteristics, should not rely on game-able features, etc. In other words, when we do policy learning, we always have sorts of constraints.

1. Policy

Definition 14.1. Suppose the features $X \in \mathcal{R}^d$, a policy $\pi(\cdot)$ is a function $\mathcal{R}^d \rightarrow \{0, 1\}$ such that an individual with feature $X = x$ is treated if and only if $\pi(x) = 1$. We define the **Value of the policy** π as

$$V(\pi) = \mathbf{E} Y_i(\pi(X_i)). \quad (\text{对 } X_i \text{ 取期望}) \quad (4)$$

Notice that

$$Y_i(\pi(X_i)) = \begin{cases} Y_i(1) & \text{if } \pi(X_i) = 1, \\ Y_i(0) & \text{if } \pi(X_i) = 0. \end{cases}$$

Therefore,

$$Y_i(\pi(X_i)) = Y_i(1)\pi(X_i) + Y_i(0)(1 - \pi(X_i)).$$

In such case, for any $\pi(\cdot)$,

$$\begin{aligned} V(\pi) &= \mathbf{E} Y_i(0) + \mathbf{E} (Y_i(1) - Y_i(0)) \times \pi(X_i). \\ &= \mathbf{E} [\underbrace{\pi(X_i)}_{\text{CATE}} \cdot \underbrace{Y_i(1) - Y_i(0)}_{\text{CATE}}] \\ &= \mathbf{E} [Y_i(1) - Y_i(0) | X_i] \end{aligned}$$

The problem we are interested in is to find

$$\hat{\pi} = \arg \max_{\pi \in \Pi} V(\pi). \quad (5)$$

2. 无限制时的 optimal policy

If there is no constraints, finding the optimal policy is simple, which is demonstrated as follows:

Theorem 14.1. Suppose we do not have any constraints, then the optimal policy is given by

$$\hat{\pi}(x) = \mathbf{1}_{\tau(x) \geq 0}, \text{ where } \tau(x) = \mathbf{E}(Y_i(1) - Y_i(0)) | X_i = x. \quad (6)$$

CATE > 0 时施加 treatment

Proof Notice that

$$\begin{aligned} V(\pi) &= \mathbf{E} Y_i(0) + \mathbf{E} [\pi(X_i) \times \mathbf{E}(Y_i(1) - Y_i(0)) | X_i] \\ &= \mathbf{E} Y_i(0) + \mathbf{E} [\pi(X_i) \times \tau(X_i)] \leq \mathbf{E} Y_i(0) + \mathbf{E} \tau(X_i) \times \mathbf{1}_{\tau(X_i) \geq 0}, \\ &= \mathbf{E} [\underbrace{\pi(X_i)}_{\leq 1} \cdot \underbrace{\tau(X_i)}_{\substack{\geq 0 \\ \text{CATE} > 0}} \cdot \mathbf{1}_{\tau(X_i) \geq 0}] + \mathbf{E} [\underbrace{\pi(X_i)}_{\substack{\leq 1 \\ \text{CATE} < 0}} \cdot \underbrace{\tau(X_i)}_{< 0} \cdot \mathbf{1}_{\tau(X_i) < 0}] \end{aligned}$$

which reaches maximum when $\pi(x) = \mathbf{1}_{\tau(x) \geq 0}$. □

This approach may be reasonable in some applications, but may result in policies that are hard to interpret or may not respect other practical constraints that are called for in the application. The focus of this chapter will be on developing methods for learning policies that do respect such constraints.

注: 现实中(为了 policy 的易懂性等因素)会存在 constraints, 如:

$$\Pi = \{1\} \{ \sum_{i=1}^p a_i X_i \geq 0 \} \mid a_i \in \mathcal{R}, \forall i = 1, \dots, p \}$$

$$\Pi = \{2\text{-layer decision tree}\}$$

Example 7 In Kitagawa and Tetenov (2018), the authors considered linear treatment rules of the form

$$\pi(x) = \mathbf{1}_{\text{prior earining} \times \alpha_1 + \text{education} \times \alpha_2 > c}.$$

15 Policy evaluation

1. PL的流程 The workflow of policy learning has 3 steps:

1. Collect data with random or quasi-random treatment assignments W_i to learn a (optimal) policy $\hat{\pi}$.

2. In a second (optional) phase, we may want to evaluate the quality of the learned policy, i.e., estimate $V(\hat{\pi})$. In this stage, we always resort to a out-of-the-sample test set which are independent of training set.

3. Finally, once were done learning, we enter the last phase where we may choose to deploy the learned policy, i.e., to set $W_i = \hat{\pi}(X_i)$ for future observations.

If we know the value $V(\cdot)$, evaluating the value of a policy is simple. Unfortunately, in most of the time, we do not know the value function. In such case, with a (suppose estimated) given policy $\hat{\pi}$, and a new out-of-the-sample test set, we want to estimate $V(\hat{\pi})$. Notably, because the test set and training sets are independent of each other, this task is equivalent to using the test set to estimate $V(\pi)$ for an arbitrary fixed policy π .

2. Policy evaluation **Theorem 15.1.** Suppose we have n observations $(Y_i, W_i, X_i) \in \mathbf{R} \times \{0, 1\} \times \mathbf{R}^d$, where the unconfoundedness condition holds true. Suppose the propensity score $e(x) = \mathbf{E}W_i|X_i = x$ is known, then the Inverse-propensity weighting estimator

$$\hat{V}_{IPW}(\pi) = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{1}_{W_i=\pi(X_i)} Y_i}{\text{Prob}(W_i = \pi(X_i)|X_i)}, \quad (7)$$

is an unbiased estimator for $V(\pi)$. $= \frac{1}{n} \left(\sum_{\pi(X_i)=1} \frac{\mathbf{1}_{W_i=1} Y_i}{e(X_i)} + \sum_{\pi(X_i)=0} \frac{\mathbf{1}_{W_i=0} Y_i}{1-e(X_i)} \right)$

Notice that

$$\text{Prob}(W_i = \pi(X_i)|X_i) = \begin{cases} \text{Prob}(W_i = 1|X_i) = e(X_i) & \text{if } \pi(X_i) = 1, \\ \text{Prob}(W_i = 0|X_i) = 1 - e(X_i) & \text{if } \pi(X_i) = 0. \end{cases}$$

Besides,

$$\begin{aligned} \mathbf{1}_{W_i=\pi(X_i)} Y_i &= \mathbf{1}_{W_i=\pi(X_i)} \times (Y_i(1) \times \mathbf{1}_{W_i=1} + Y_i(0) \times \mathbf{1}_{W_i=0}) \\ &= Y_i(1) \times \mathbf{1}_{W_i=\pi(X_i)=1} + Y_i(0) \times \mathbf{1}_{W_i=\pi(X_i)=0}. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbf{E} \frac{\mathbf{1}_{W_i=\pi(X_i)} Y_i}{\text{Prob}(W_i = \pi(X_i) | X_i)} | X_i &= \frac{1}{\text{Prob}(W_i = \pi(X_i) | X_i)} \times (\mathbf{E} Y_i(1) \times \mathbf{1}_{W_i=\pi(X_i)=1} | X_i + \mathbf{E} Y_i(0) \times \mathbf{1}_{W_i=\pi(X_i)=0} | X_i) \\ &= \frac{1}{\text{Prob}(W_i = \pi(X_i) | X_i)} \times (\pi(X_i) \mathbf{E} Y_i(1) | X_i \times e(X_i) + (1 - \pi(X_i)) \mathbf{E} Y_i(0) | X_i \times (1 - e(X_i))) \\ &= \begin{cases} \pi(X_i) \times \mathbf{E} Y_i(1) | X_i & \text{if } \pi(X_i) = 1, \\ (1 - \pi(X_i)) \mathbf{E} Y_i(0) | X_i & \text{if } \pi(X_i) = 0, \end{cases} = \pi(X_i) \times \mathbf{E} Y_i(1) | X_i + (1 - \pi(X_i)) \mathbf{E} Y_i(0) | X_i, \end{aligned}$$

which proves the result.

If we want to incorporate a parametric model, we can consider a double-robust type estimator

$$\hat{V}_{AIPW}(\pi) = \frac{1}{n} \sum_{i=1}^n \hat{\mu}(X_i, \pi(X_i)) + \frac{\mathbf{1}_{W_i=\pi(X_i)}}{\text{Prob}(W_i = \pi(X_i) | X_i)} (Y_i - \hat{\mu}(X_i, \pi(X_i))).$$

16 Empirical-welfare maximization

1. regret

We now return to the task of learning a policy, i.e., using experimental or quasi- experimental data to choose a good treatment assignment rule $\hat{\pi}$. Throughout, we assume that the policymaker is constrained to choose a policy π belonging to some class Π of acceptable policies. For example, one may consider $\pi(x) = \mathbf{1}_{a^\top x \geq c}$ for some vector a and threshold c , or some fixed-depth decision trees.

In such setting, we define the optimal policy as

$$\pi^* = \arg \max V(\pi) \text{ such that } \pi \in \Pi.$$

We define the regret

$$R(\pi) = \sup_{\zeta \in \Pi} V(\zeta) - V(\pi). \quad = V(\pi^*)$$

Our goal is to learn a policy with guaranteed worst-case bounds on the regret. We refer this task as a learning (rather than estimation) task because the performance of π is only assessed in terms of its regret. No requirements will be made on π converging to the optimal policy π^* in terms of its functional form (and in fact no assumption is made that there is a unique optimal policy).

2. Empirical-welfare maximization

If the optimal policy π^* maximizes the true value function, then it is natural to attempt learn $\hat{\pi}$ by maximizing an estimated value function

$$\hat{\pi} = \arg \max \{ \hat{V}(\pi) : \pi \in \Pi \}.$$

This approach is called empirical-welfare maximization by Kitagawa and Tetenov(2018). A plausible

estimator is the IPW estimator $\hat{V}_{IPW}(\pi)$, in such case we define

$$\hat{\pi}_{IPW} = \arg \max \{ \hat{V}_{IPW}(\pi) : \pi \in \Pi \}.$$

3. Empirical-welfare
maximization 的
bound

To bound the regret,

$$\begin{aligned} R(\hat{\pi}) &= V(\pi^*) - V(\hat{\pi}) = V(\pi^*) - \underbrace{\hat{V}(\pi^*) + \hat{V}(\pi^*) - V(\hat{\pi})}_{\leq |V(\pi^*) - \hat{V}(\pi^*)| + |V(\hat{\pi}) - \hat{V}(\pi^*)|} \\ &\leq |V(\pi^*) - \hat{V}(\pi^*)| + |V(\hat{\pi}) - \hat{V}(\pi^*)| \leq 2 \sup_{\pi \in \Pi} |V(\pi) - \hat{V}(\pi)|. \end{aligned}$$

In other words, bounding the estimator of the value function is core for the learning of a policy.

References

Ding, Peng (2023). *A First Course in Causal Inference*.