

動的治療計画と強化学習

江口 真透

医療健康データ科学研究センター 教授

1. 研究の流れ

Logistic Regression, Cox (1958)

Perceptron, Rosenblatt (1958)

Dynamic programming, Bellman (1957)

Bellman equation (principle of optimality)

Q-learning, Watkins (1989)

Deep Learning, LeCun, Bengio, Hinton (2015)

Deep Q-learning, Silver, Hassabis,...

Dynamic treatment regime, Murphy (2003), Zhao, et al (2012)

2. ランダム臨床試験と動的治療計画

Random Clinical Trial

Treatment = Intervention

Random assignment for treatments

Dynamic treatment regimes

Treatment = Action

Adaptive biased coin, cf. Lavori-Dawson (2002)

Multiple Assignment Randomized Trial (SMART)

Classification-based approach, cf. Zhao (2012, 2015)

Chakraborty & Moodie (2013). Statistical methods for dynamic treatment regimes, cf. Robins (1986); Murphy (2003).

3. 確率フレームワーク

- (\mathbf{X}, A, Y)

$$\begin{cases} \text{state } \mathbf{X} \in \mathcal{X} \subseteq \mathbb{R}^p \\ \text{action } A \in \mathcal{A} = \{1, \dots, M\} \\ \text{reward } Y \in \mathcal{Y} = \{y \in \mathbb{R} : y \geq 0\} \end{cases}$$

$$(\mathbf{X}, A, Y) \sim p(\mathbf{x}, a, y) = p(y|\mathbf{x}, a)p(a|\mathbf{x})p(\mathbf{x})$$

- Deterministic policy $d: \mathcal{X} \rightarrow \mathcal{A}$ has a value function

$$V_d = \mathbb{E}_d[Y] = \mathbb{E} \left[\frac{\mathbb{I}(d(\mathbf{X}) = A)}{p(A|\mathbf{X})} Y \right]$$

- Optimal policy $d^{\text{opt}} = \operatorname{argmax}_{d \in \mathcal{D}} V_d$ where \mathcal{D} is the space of all policy functions

Cf. Supervised learning: A feature vector \mathbf{X} predicts an outcome Y

$$p(\mathbf{x}, y) = p(y|\mathbf{x})p(\mathbf{x})$$

4. Q-関数

- Q-function $Q(\mathbf{x}, a) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}, A = a]$
- Optimal policy $d^{\text{opt}}(\mathbf{x}) = \operatorname{argmax}_{a \in \mathcal{A}} Q(\mathbf{x}, a)$ (Sutton and Barto 1998)

- Decomposition of Q-function $Q(\mathbf{x}, a) = \pi(a|\mathbf{x})\eta(\mathbf{x})$
where $\pi(a|\mathbf{x}) = \frac{Q(\mathbf{x}, a)}{\sum_{b \in \mathcal{A}} Q(\mathbf{x}, b)}$ and $\eta(\mathbf{x}) = \sum_{a \in \mathcal{A}} Q(\mathbf{x}, a)$

Note: $d^{\text{opt}}(\mathbf{x}) = \operatorname{argmax}_{a \in \mathcal{A}} \pi(a|\mathbf{x})$

We consider a parametric model of Q-function

$$Q(\mathbf{x}, a) = \pi(a|\mathbf{x}, \boldsymbol{\theta})\eta(\mathbf{x})$$

Here we call $\eta(\mathbf{x})$ nuisance function.

5. ガンマ・ベキダイバージェンス

Cf. Fujisawa-Eguchi (2008)

$$D_\gamma(Q_1, Q_2) = \int_{\mathcal{X}} \frac{\sum_{a \in \mathcal{A}} Q_1(\mathbf{x}, a) Q_2(\mathbf{x}, a)^\gamma}{\{\sum_{b \in \mathcal{A}} Q_2(\mathbf{x}, b)^{\gamma+1}\}^{\frac{\gamma}{\gamma+1}}} p(\mathbf{x}) d\mathbf{x} - \int_{\mathcal{X}} \left\{ \sum_{a \in \mathcal{A}} Q_1(\mathbf{x}, a)^{\gamma+1} \right\}^{\frac{1}{\gamma+1}} p(\mathbf{x}) d\mathbf{x}$$

We model a Q-function as $Q_\theta(\mathbf{x}, a) = \pi(a|\mathbf{x}, \boldsymbol{\theta})\eta(\mathbf{x})$. Then

$$D_\gamma(Q_0, Q_\theta) = \sum_{a \in \mathcal{A}} \int_{\mathcal{X}} Q_0(\mathbf{x}, a) \frac{\pi(a|\mathbf{x}, \boldsymbol{\theta})^\gamma}{\{\sum_{b \in \mathcal{A}} \pi(b|\mathbf{x}, \boldsymbol{\theta})^{\gamma+1}\}^{\frac{\gamma}{\gamma+1}}} p(\mathbf{x}) d\mathbf{x} + \text{const.}$$

Thus, the expected/empirical loss is

$$L_\gamma(\boldsymbol{\theta}) = - \sum_{a \in \mathcal{A}} \int_{\mathcal{X}} Q_0(\mathbf{x}, a) \frac{\pi(a|\mathbf{x}, \boldsymbol{\theta})^\gamma}{\{\sum_{b \in \mathcal{A}} \pi(b|\mathbf{x}, \boldsymbol{\theta})^{\gamma+1}\}^{\frac{\gamma}{\gamma+1}}} p(\mathbf{x}) d\mathbf{x}$$

$$L_\gamma^{\text{emp}}(\boldsymbol{\theta}, \mathcal{D}) = - \sum_{i=1}^n \frac{y_i}{p(a_i|\mathbf{x}_i)} \frac{\pi(a_i|\mathbf{x}_i, \boldsymbol{\theta})^\gamma}{\{\sum_{b \in \mathcal{A}} \pi(b|\mathbf{x}_i, \boldsymbol{\theta})^{\gamma+1}\}^{\frac{\gamma}{\gamma+1}}} \quad \text{where } \mathcal{D} = \{(\mathbf{x}_i, a_i, y_i)_{i=1}^n\}$$

Assume that the true Q-function is $Q_{\theta_0}(\mathbf{x}, a)$ and let $\hat{\boldsymbol{\theta}}_\gamma = \operatorname{argmin}_{\boldsymbol{\theta}} L_\gamma^{\text{emp}}(\boldsymbol{\theta}, \mathcal{D})$. Then $\hat{\boldsymbol{\theta}}_\gamma$ is consistent with $\boldsymbol{\theta}_0$.

We note: Only the class of γ -estimators is free from η -dependence.

6. 決定関数

- Q-function $Q(x, a) = \mathbb{E}[Y|X = a, X = x]$

$$d^{\text{opt}}(x) = \operatorname{argmax}_{a \in \mathcal{A}} Q(x, a)$$

- Decision function $f(x, a)$ is defined to satisfy

$$\sum_{a \in \mathcal{A}} f(x, a) = 0 \quad (\forall x \in \mathcal{X})$$

Definition Decision function $f(x, a)$ is said to be **D-consistent** if

$$d^{\text{opt}}(x) = \operatorname{argmax}_{a \in \mathcal{A}} f(x, a)$$

NB1: Q-learning aims to estimate the optimal Q-function.

NB2: $Q(x, a)$ vs $f(x, a)$; regression vs prediction (Zhao, 2012, 2015)

7. 決定一貫性

- Let Ψ be a strictly decreasing and convex function on \mathbb{R}

$$\Psi\text{-loss function } L_\Psi(f) = \mathbb{E} \left[\frac{Y}{p(A|\mathbf{X})} \Psi(f(\mathbf{X}, A)) \right]$$

Let \mathcal{D} be an empirical dataset, $\mathcal{D} = \{(\mathbf{x}_i, a_i, y_i)_{i=1}^n\}$

$$\Psi\text{-empirical loss function } L_\Psi(f, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \frac{y_i}{p(a_i|\mathbf{x}_i)} \Psi(f(\mathbf{x}_i, a_i))$$

Proposition

Let $f^* = \operatorname{argmin}_{f \in \mathcal{F}} L_\Psi(f)$. Then, $f^*(\mathbf{x}, a)$ is D-consistent.

NB3: Any $L_\Psi(f)$ leads to the optimal policy d^{opt}

証明の概略

- It suffices to show

$$d^{\text{opt}}(\mathbf{x}) = \operatorname{argmax}_{a \in \mathcal{A}} f^*(\mathbf{x}, a) \quad \text{where } d^{\text{opt}}(\mathbf{x}) = \operatorname{argmax}_{a \in \mathcal{A}} Q(\mathbf{x}, a)$$

- By definition, $L_\Psi(f) = \mathbb{E} \left[\frac{Y}{p(A|\mathbf{X})} \Psi(f(\mathbf{X}, A)) \right]$

$$= \sum_{a \in \mathcal{A}} \int_{\mathcal{X}} \int_{\mathcal{Y}} \frac{y}{p(a|\mathbf{x})} \Psi(f(\mathbf{x}, a)) p(y|\mathbf{x}, a) p(a|\mathbf{x}) p(\mathbf{x}) dy d\mathbf{x}$$

$$= \sum_{a \in \mathcal{A}} \int_{\mathcal{X}} \Psi(f(\mathbf{x}, a)) Q(\mathbf{x}, a) d\mathbf{x}$$

Let $\mathcal{L}(f) = L_\Psi(f) - \int_{\mathcal{X}} \lambda(\mathbf{x}) \sum_{a \in \mathcal{A}} f(\mathbf{x}, a) d\mathbf{x}$

- Then, we find the Euler-Lagrange equilibrium condition

$$\Psi'(f^*(\mathbf{x}, a)) Q(\mathbf{x}, a) = \lambda(\mathbf{x}) \quad (\forall (\mathbf{x}, a) \in \mathcal{X} \times \mathcal{A})$$