

動的治療計画のためのポリシー・ロバストなQ学習

江口 真透 医療健康データ科学研究センター 特任教授

概要: 動的治療計画においてダブルロバストな価値関数の推定やSVMを使ったアウトカム重み付け学習などが提案されて以来、強化学習の内容で生物統計において積極的に被験者の個人の履歴データに基づく、より良い治療の計画を提供する斬新な方法論が開発されつつある。治療計画の主要な目的の一つは履歴に基づく最適な治療を導くことである。最適治療とはQ関数に個人の共変量履歴を代入して、治療に関するQ関数の最大化によって定義される。しかし被験者の治療の感受性に異質性があるときには、尤度に基づくQ関数の推定は問題を含む。この問題に対して**治療の割付確率とアウトカムを観測した事後割付確率の比**で定められる重み付け尤度方程式による推定を提案する。従来のロバスト回帰法と比較して、特に**ポリシー・ロバスト性**の観点から提案手法の良さが主張される。

問題設定と提案推定量

単一ステージの場合について考察する。 (X, A, Y) を共変量ベクトル、治療ラベル、アウトカムの組とする。Q関数 $Q_0(x, a) = \mathbb{E}[Y|x, a]$ が既知ならば最適ポリシーは

$$d_0(x) = \operatorname{argmax}_a Q_0(x, a)$$

と得られる。Q関数が未知ならば、例えば、**一般化線形モデル**

$$p_\beta(y|x, a) = \exp\left\{\frac{y\nu_\beta(x, a) - \psi(\nu_\beta(x, a))}{\phi} + c(x, a, \phi)\right\}. \quad (1)$$

を仮定すると、Q関数は $Q_\beta(x, a) = \frac{\partial \psi}{\partial \nu}(\nu_\beta(x, a))$ と書ける。これより、

パラメータ β の推定がキーとなる。標準的な推定は全尤度を使ったスコア関数

$$\mathcal{E}_F(\beta, x, a, y) = \frac{y - Q_\beta(x, a)}{\phi} \frac{\partial}{\partial \beta} \nu_\beta(x, a) \quad (2)$$

が考えられ、最尤推定量 $\hat{\beta}_F$ はデータ $\{(x_i, a_i, y_i) : i=1, \dots, n\}$ に基づき

$$\sum_{i=1}^n \mathcal{E}_F(\beta; x_i, a_i, y_i) = 0$$

の解として求められる。これより最適ポリシーは

$$\hat{d}_F(x) = \operatorname{argmax}_a Q_{\hat{\beta}_F}(x, a)$$

と推定される。しかし、あとで導入するポリシー・ロバスト性の観点から深刻な問題が懸念される。データの条件付き分布が**モデルと乖離**して

$$p_0(y|x, a) = (1 - \epsilon)p_\beta(y|x, a) + \epsilon p_h(y|x, a) \quad (3)$$

と書けたとする。例えば、主要なグループのYの条件付き分布はモデル(1)に従うが、少数のグループが存在して $p_h(y|x, a)$ に従っていると考えられる。

このように真のQ関数 $Q_0(x, a) = \mathbb{E}[Y|x, a]$ はモデルQ関数 $Q_\beta(x, a)$ と異なってしまう。とくに $p_h(y|x, a)$ と $p_\beta(y|x, a)$ の平均がかけ離れると最尤推定量に大きなバイアスが生じるので推定されたポリシー $\hat{d}_F(x)$ は判断を大きく誤る可能性がある。最適ポリシー $d_0(x) = \operatorname{argmax}_a Q_0(x, a)$ に対して*i*番目の観測値 (x_i, a_i, y_i) が $a_i = d_0(x_i)$ となるときアウトカム y_i が想定よりも低い値であった時、**尤度方程式(2)にダメージを与える**。逆に y_i が想定よりも大きな値であればQ関数の治療に関する最大値を過大推定するだけなのでポリシーの推定には影響しないことになる。

このような考察から、既存のロバスト法ではなく、**重み付け推定関数**

$$\mathcal{E}_W(\beta, x, a, y) = \frac{p_\beta(a|x, y)}{p(a|x)} \mathcal{E}_F(\beta, x, a, y) - B(\beta, x, a) \quad (4)$$

を使った推定量を提案する。ここで

$$B(\beta, x, a) = \mathbb{E}\left[\frac{p_\beta(a|x, Y)}{p(a|x)} \mathcal{E}_F(\beta, x, a, Y)\right]$$

提案推定量 $\hat{\beta}_W$ は、**推定方程式**

$$\sum_{i=1}^n \mathcal{E}_W(\beta, x_i, a_i, y_i) = 0$$

の解と定義される。

ポリシー・ロバスト性

つぎに、最適ポリシーの推定量のロバスト性について考察する。簡単のため、治療空間を $\{1, -1\}$ とする。このとき最適ポリシーは

$$d_0(x) = \operatorname{sign}\{Q_0(x, 1) - Q_0(x, -1)\}$$

と書ける。データの分布が(3)のように誤特定されるときでも、推定されたポリシー $d_{\hat{\beta}}(x)$ が最適ポリシーと一致する確率が高いとき、**ポリシー・ロバスト性**を持つという。

このとき、重み付け推定関数(4)の重み関数は**ロジスティック関数形**

$$p_\beta(a|x, y) = \begin{cases} \frac{1}{1 + \exp(-\alpha_1(x)y - \alpha_0(x))} & \text{if } a = 1 \\ \frac{1}{1 + \exp(\alpha_1(x)y + \alpha_0(x))} & \text{if } a = -1 \end{cases}$$

で表せる。ここで、

$$\alpha_1(x) = \frac{\nu_\beta(x, 1) - \nu_\beta(x, -1)}{\phi}$$

$$\alpha_0(x) = \frac{\psi(\nu_\beta(x, -1)) - \psi(\nu_\beta(x, 1))}{\phi} + \log \frac{p(1|x)}{p(-1|x)}$$

従って、この表現と一般化線形モデルの基本性質から

$$\alpha_1(x) > 0 \text{ if and only if } Q(x, 1) > Q(x, -1)$$

が成立することが分かるので $Q(x, 1) > Q(x, -1)$ ならばYが大きな値を持つときは**重み付き推定関数のノルム**は大きくなるがYが小さな値を持つときは有界なノルムを持つことが分かる。具体的には以下の例題で示される。

例題. Yの条件付き分布が正規分布 $\mathcal{N}(Q_\beta(x, a), \sigma^2)$ として、線形モデル

$$Q_\beta(x, a) = \beta_0 + a\beta_1 + \mathbf{x}^\top \beta_2 + a\mathbf{x}^\top \beta_3$$

を仮定する。重み付き推定関数(3)は

$$\lim_{y \rightarrow -\infty} \|\mathcal{E}_W(\beta, x, a, y)\| = 0, \quad \lim_{y \rightarrow \infty} \|\mathcal{E}_W(\beta, x, a, y)\| = \infty$$

となり、切片項に対する成分は下の図のような挙動を示す。ガンマ・べき推定関数はYの両側で有界になり、極限で0になるが、この重み付き推定関数は片側だけが有界になる。このように**ポリシーが影響を受けるYの外れ値**のみにロバスト性があることが分かる。

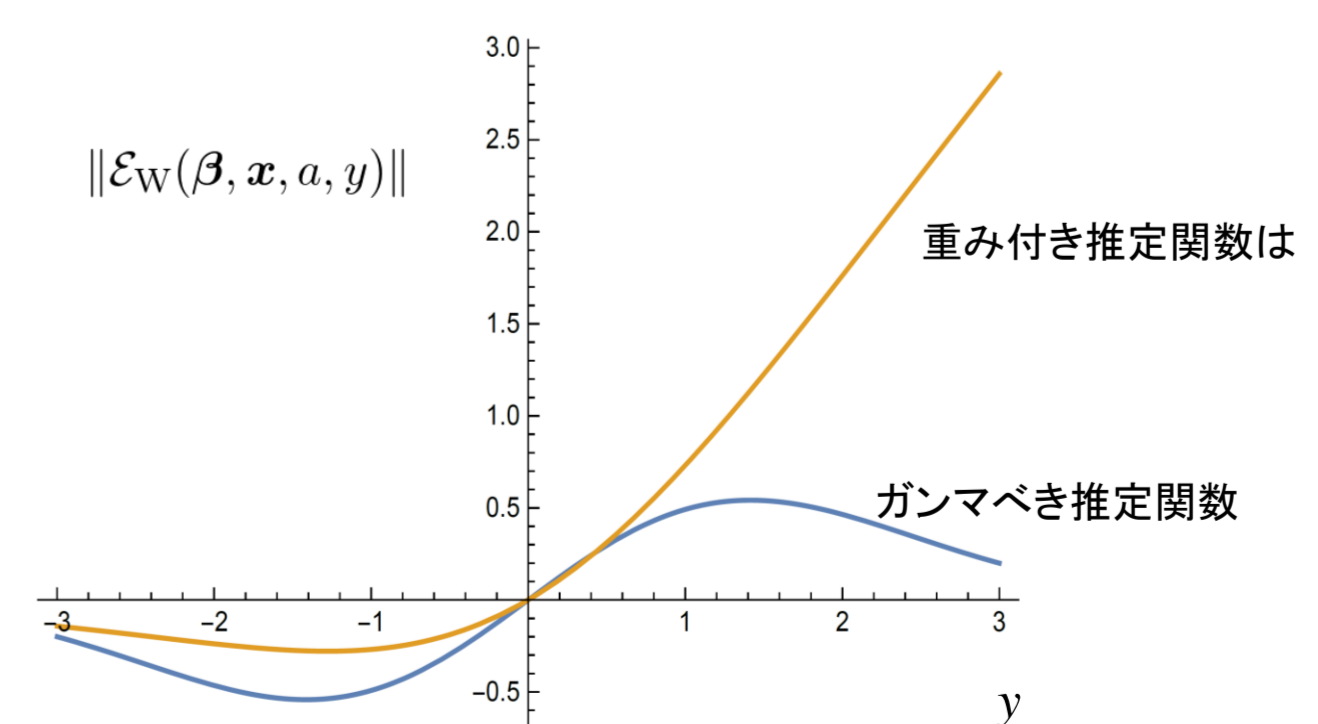


図1. 最適ポリシー $d_0(x) = 1$ のときの推定関数グラフ