

Optimal Treatment Regimes under Constraints

Ph.D. Thesis Presentation

Shuping Ruan

Supervised by Dr. Eric Laber
Department of Statistics
North Carolina State University

December 19, 2017

Precision medicine

- Categorize individuals into subpopulations based on
 - Demographics
 - Genetic information
 - Results of diagnostic test
 - Susceptibility to a certain disease
 - Response to a specific treatment
- Target therapeutic or preventive interventions to individuals
 - Provide more effective treatments
 - Avoid unnecessary side effects and costs

Dynamic treatment regimes

- Definition: a set of decision rules, one function for each decision point, that map patient characteristics to recommended treatments
- Opposite of the acute care model: active but short-term treatment for an urgent medical condition
- The goal is to optimize the long-term cumulative clinical outcome for chronic care
- Analogous to a policy in reinforcement learning, or a controller in control theory
- Apply to time-varying policies in other fields, such as education, marketing, and economics

Optimal treatment regimes

- Most of DTRs focus on a single scalar outcome
- Indirect methods: Q-learning, A-learning, etc.
- Policy search: outcome weighted learning, doubly robust estimators, etc.
- Balancing multiple competing outcomes
 - Treatment effectiveness
 - Side-effect burden
 - Cost and so on
- Current work
 - Composite outcomes
 - Set-valued treatment regimes

Constrained optimal treatment regimes

- Propose a new framework for estimating optimal treatment regimes under constraints
- Optimize the primary outcome of interest, subject to constraints on the secondary outcomes
 - Single-stage setting
 - Multi-stage setting
 - Infinite-stage setting

Single-stage

Dataset

- Dataset

$$\left\{ (\mathbf{X}^i, A^i, \mathbf{Y}^i) \right\}_{i=1}^n$$

- n i.i.d patient trajectories
- $\mathbf{X} \in \mathcal{X}$: the patient information collected up to the decision point
- $A \in \mathcal{A}$: the treatment assignment
- $\mathbf{Y} \in \mathbb{R}^J$: the vector of outcomes of interest
 - $\mathbf{Y} = (Y_1, Y_2, \dots, Y_J)^\top$
 - Y_1 , the primary outcome of interest, the higher the better
 - Y_2, \dots, Y_J , the secondary outcomes of interest, the lower the better

Potential outcome framework

- The counter-factual framework for identifying the causal effect of a certain regime from randomized or observational studies.
- The set of potential outcomes: $\mathcal{W}^* = \{\mathbf{Y}^*(a), \text{ for all } a \in \mathcal{A}\}$, where $\mathbf{Y}^*(a)$ is the vector-valued outcome that would have been observed if the subject was assigned treatment a .
- Three essential assumptions guarantee that the value for a regime can be estimated using the observed data.

Potential outcome framework

- *Consistency*

$$\mathbf{Y} = \mathbf{Y}^*(A).$$

The actual observed outcome vector \mathbf{Y} for an individual who received treatment A is the same as the potential outcome for that individual assigned with the same treatment, regardless of the experimental conditions used to assign treatment. It also implies that there is no interference among individuals.

Potential outcome framework

- *No unmeasured confounders*

$$\mathbf{W}^* \perp\!\!\!\perp A \mid \mathbf{X}.$$

The set of potential outcomes, $\mathbf{W}^* = \{ \mathbf{Y}^*(a) , \text{for all } a \in \mathcal{A} \}$, are conditionally independent of treatment assignment A given patient information \mathbf{X} . In randomized study, this condition is satisfied by construction. However, it can not be verified in observational studies.

Potential outcome framework

- *Positivity assumption:*

$$\exists \epsilon > 0, \text{s.t. } \Pr(A = a | \mathbf{X}) > \epsilon, \forall a \in \mathcal{A} \text{ w/ prob 1.}$$

This ensures there is a positive probability of receiving every possible treatment assignment for every value of patient covariates in the population. This assumption is satisfied in well-designed randomized studies. It can also be empirically verified in observational studies. Yet, if it is violated, estimating of regimes for certain subsets of patients can be impossible.

A single-stage treatment regime

- A single-stage treatment regime $\pi : \mathcal{X} \rightarrow \mathcal{A}$ is a function that maps the support of patient information \mathbf{X} to the set of all possible treatments.
- Under a regime π , a patient with $\mathbf{X} = \mathbf{x}$ is recommended to receive treatment $\pi(\mathbf{x})$.
- The vector-valued potential outcome of the regime π is $\mathbf{Y}^*(\pi) = \sum_{a \in \mathcal{A}} \mathbf{Y}^*(a) \mathbb{I}\{\pi(\mathbf{X}) = a\}$.

Values of a regime

- The value functions of a regime π is defined as

$$\mathbf{V}(\pi) = \mathbb{E} \mathbf{Y}^*(\pi),$$

of which each component is $V_j(\pi) = \mathbb{E} Y_j^*(\pi)$, $j = 1, \dots, J$. That is the expected outcome if every patient in the population of interest is assigned treatment according to π ,

- The Q functions are defined as

$$\mathbf{Q}(\mathbf{X}, A) = \mathbb{E}(\mathbf{Y} | \mathbf{X}, A),$$

of which each component is $Q_j(\mathbf{X}, A) = \mathbb{E}(Y_j | \mathbf{X}, A)$, $j = 1, \dots, J$.

- Under the three assumptions, we have

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

Constrained optimal treatment regimes

- A single-stage constrained optimal regime is

$$\max_{\pi \in \Pi} V_1(\pi)$$

$$\text{subject to } V_j(\pi) \leq \nu_{j-1},$$

where $j = 2, \dots, J$, and $\nu = (\nu_1, \nu_2, \dots, \nu_{J-1})^\top$ are the constraints.

Constrained optimal treatment regimes

- Linear decision rules, $\pi(\mathbf{x}; \boldsymbol{\theta}) = \text{sgn}(\mathbf{x}^\top \boldsymbol{\theta})$, $\|\boldsymbol{\theta}\|_2^2 = \boldsymbol{\theta}^\top \boldsymbol{\theta} = 1$
- Further simplify the notation

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^q} v_1(\boldsymbol{\theta})$$

subject to $v_j(\boldsymbol{\theta}) \leq 0$, $h(\boldsymbol{\theta}) = 0$,

where $j = 2, \dots, J$, $v_1(\boldsymbol{\theta}) = -V_1(\boldsymbol{\theta})$, $v_j(\boldsymbol{\theta}) = V_j(\boldsymbol{\theta}) - v_j$, and $h(\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \boldsymbol{\theta} - 1$.

Interior-point method

- The interior-point method solves a sequence of the following approximate minimization problem,

$$\min_{\theta, \mathbf{z}} \phi_\mu(\theta, \mathbf{z}) = v_1(\theta) - \mu \sum_{j=2}^J \ln z_j,$$

subject to $v_j(\theta) + z_j = 0$, and $h(\theta) = 0$,

where μ is always positive and approaches to zero in the limit.

- As μ decreases to zero, the minimums of ϕ_μ form a trajectory path that approaches the feasible optimum of $v_1(\theta)$, θ_v^* , in the limit.

Estimation of the value functions

- $V_j(\theta) = \mathbb{E} (Q_j(\mathbf{X}, \pi(\mathbf{X}, \theta))),$ for $j = 1, \dots, J.$
- A linear model $Q_j(\mathbf{X}, A) = \mathbf{X}^\top \alpha_j + A \cdot \mathbf{X}^\top \beta_j.$
- A regime is approximated by $\pi(\mathbf{X}) = \text{sgn}(\mathbf{X}^\top \theta).$
- Let $(z_1, z_2) = (\mathbf{x}^\top \theta, \mathbf{x}_1^\top \beta_j),$ and $f_{\beta_j}(z_1, z_2; \theta)$ be the joint dist'n.
- The estimated values of a regime π are

$$\hat{V}_j(\theta) = \mathbf{x}^\top \hat{\alpha}_j + \iint \text{sgn}(z_1) z_2 \hat{f}_{\beta_j}(z_1, z_2; \theta) dz_1 dz_2,$$

where $\hat{\alpha}_j$ and $\hat{\beta}_j$ are the least-squares estimators; $\hat{f}_{\beta_j}(z_1, z_2; \theta)$ is a kernel density estimator (KDE).

Simulation design

- The generative model for simulation is

$$\mathbf{X} \sim \text{MVN}(\mathbf{0}, \mathbf{I}),$$

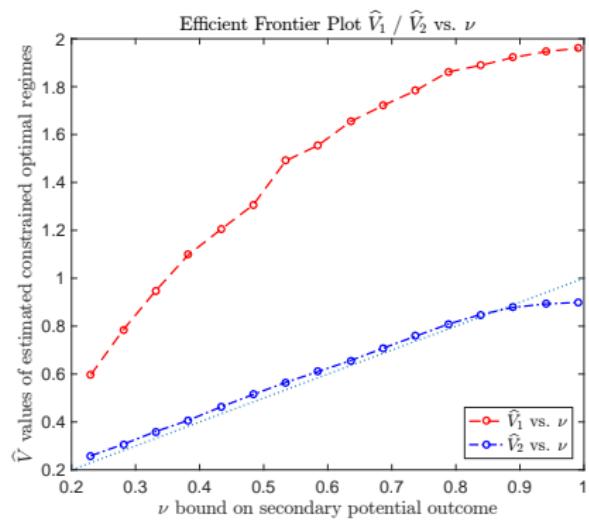
$$A \sim \text{Uniform}\{-1, 1\},$$

$$Y_1 = \bar{\mathbf{X}}^\top \boldsymbol{\alpha}_1 + A \cdot (\bar{\mathbf{X}}^\top \boldsymbol{\beta}_1) + \epsilon_1, \quad \epsilon_1 \sim N(0, \sigma_1^2),$$

$$Y_2 = \bar{\mathbf{X}}^\top \boldsymbol{\alpha}_2 + A \cdot (\bar{\mathbf{X}}^\top \boldsymbol{\beta}_1) + \epsilon_2, \quad \epsilon_2 \sim N(0, \sigma_2^2),$$

where \mathbf{I} is a 2×2 identity matrix and $\bar{\mathbf{X}}^\top = (1, \mathbf{X}^\top)$. For simplicity, we consider two competing outcomes, i.e., $J = 2$. Also, let $\mathbf{X}_0 = \mathbf{X}_1 = \mathbf{X}$.

- $M_{MC} = 200$, $N_{train} = 1000$, $N_{test} = 10000$.



Multi-stage

Data

- The dataset

$$\{(\mathbf{X}_1^i, A_1^i, \mathbf{X}_2^i, A_2^i, \dots, \mathbf{X}_T^i, A_T^i, \mathbf{Y}^i)\}_{i=1}^n,$$

- n i.i.d. patient trajectories
- \mathbf{X}_t , the patient covariate collected between the (t-1)-th decision point and the t-th decision point
- A_t , the t th-stage treatment variable
- \mathbf{Y} , the final outcome vector. Y_1 , the primary outcome of interest, the higher the better; Y_2, \dots, Y_J , the secondary outcomes of interest, the lower the better.
- Let \mathbf{H}_t denote the patient history information up to the decision point t , i.e., $\mathbf{H}_t^\top = (\mathbf{H}_{t-1}^\top, A_{t-1}, \mathbf{X}_t^\top), \dots$
- $\bar{\mathbf{A}}_t = (A_1, A_2, \dots, A_t)$ denotes a sequence of treatment history up to time point t

Potential outcome framework

- The set of potential outcomes is
 $\mathbf{W}^* = \{\mathbf{X}_2^*(a_1), \mathbf{X}_3^*(\bar{a}_2), \dots, \mathbf{X}_T^*(\bar{a}_{T-1}), \mathbf{Y}_T^*(\bar{a}_T), \text{for all } \bar{a}_t \in \bar{\mathcal{A}}_t, t = 1, 2, \dots, T\}$
- Assumptions to connect observed data w/ potential outcomes
 - Consistency:* $\mathbf{Y} = \mathbf{Y}^*(\bar{A}_T)$, and $\mathbf{X}_t = \mathbf{X}_t^*(\bar{A}_{t-1})$, $t = 2, \dots, T$.
 - Sequential randomization assumption:* $A_t \perp\!\!\!\perp \mathbf{W}^* \mid \mathbf{H}_t$, for $t = 1, 2, \dots, T$.
 - Positivity assumption:* $\exists \epsilon_t > 0$, such that $\Pr(A_t = a_t \mid \mathbf{H}_t) > \epsilon_t$, w/ prob 1, for all $a_t \in \mathcal{A}_t$, $t = 1, 2, \dots, T$.

Values of a regime

- The values of a dynamic treatment regime, $V(\pi) = \mathbb{E} Y^*(\pi)$, is defined as the expected final outcome if each patient in the population of interest is treated according to π .
- Values of a regime are estimated using G-computation and the mean-variance modeling techniques by Linn et al.

Constrained optimal treatment regimes

- A multi-stage constrained optimal treatment regime

$$\max_{\pi \in \Pi} V_1(\pi)$$

$$\text{subject to } V_j(\pi) \leq v_{j-1},$$

where $j = 2, 3, \dots, J$ and Π is the class of dynamic treatment regimes under consideration.

Constrained optimal treatment regimes

- Linear decision rules for each stage, $\pi_t(\mathbf{h}_t) = \text{sgn}(\mathbf{h}_t^\top \boldsymbol{\theta}_t)$, and $\boldsymbol{\theta}_t^\top \boldsymbol{\theta}_t = 1$.
- Interior-point method

$$\min_{\boldsymbol{\theta}, \mathbf{z}} \phi_\mu(\boldsymbol{\theta}, \mathbf{z}) = v_1(\boldsymbol{\theta}) - \mu \sum_{j=2}^J \ln z_j,$$

subject to $v_j(\boldsymbol{\theta}) + z_j = 0$, $h_t(\boldsymbol{\theta}_t) = 0$,

where $v_1(\boldsymbol{\theta}) = -V_1(\boldsymbol{\theta})$, $v_j(\boldsymbol{\theta}) = V_j(\boldsymbol{\theta}) - \nu_j$ for $j = 2, \dots, J$, and $h_t(\boldsymbol{\theta}_t) = \boldsymbol{\theta}_t^\top \boldsymbol{\theta}_t - 1$ for $t = 1, \dots, T$.

Simulation design

This model is a simple representation of the data from a two-stage sequential randomized trials.

$$X_1 \sim \text{Normal}(1, 1), \mathbf{H}_1 = (1, X_1)^\top,$$

$$A_1 \sim \text{Uniform} \{-1, 1\},$$

$$X_2 = \mathbf{H}_1^\top \boldsymbol{\beta}_{1,0} + A_1 \mathbf{H}_1^\top \boldsymbol{\beta}_{1,1} + \epsilon, \epsilon \sim \text{Normal}(0, 1),$$

$$\mathbf{H}_2 = (1, X_2)^\top,$$

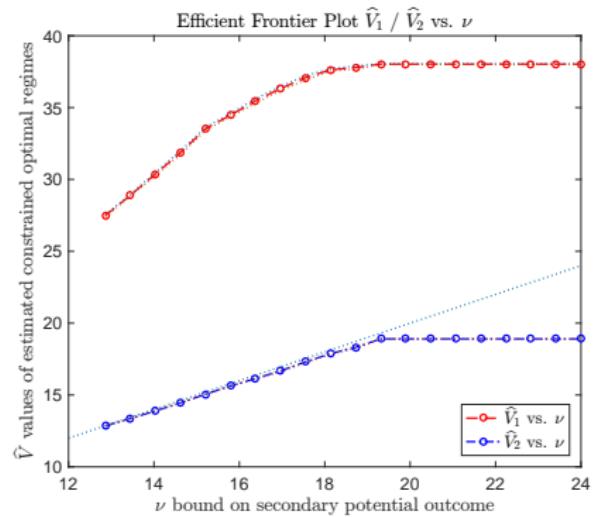
$$A_2 \sim \text{Uniform} \{-1, 1\},$$

$$Y = \mathbf{H}_2^\top \boldsymbol{\beta}_{2,0,Y} + A_2 \mathbf{H}_2^\top \boldsymbol{\beta}_{2,1,Y} + \epsilon_Y$$

$$Z = \mathbf{H}_2^\top \boldsymbol{\beta}_{2,0,Z} + A_2 \mathbf{H}_2^\top \boldsymbol{\beta}_{2,1,Z} + \epsilon_Z$$

$$(\epsilon_Y, \epsilon_Z)^\top \sim \text{Normal}(\mathbf{0}_2, \Sigma_{Y,Z})$$

The class of regimes under consideration is $\pi_1 = \text{sgn}(\mathbf{h}_1^\top \boldsymbol{\theta}_1)$ and $\pi_2 = \text{sgn}(\mathbf{h}_2^\top \boldsymbol{\theta}_2)$.



Infinite-stage

Constrained reinforcement learning

- Patients with chronic diseases are often monitored and treated throughout their life.
 - Real-time actions
 - No a-priori fixed end point
 - Infinite stage reinforcement learning
- Adopt the constrained Markov Decision Processes (cMDPs) framework for infinite horizon.
- Previously, linear programming is used to sought out constrained optimal policies in the setting of finite cMDPs with known models.
- However, few methods have been proposed for high-dimensional constrained reinforcement learning problems w/o modeling the underlying dynamics.

Data

- The structure of the available data is

$$\mathcal{D} = \left\{ (\mathbf{S}_0^i, A_0^i, \mathbf{R}_0^i, \mathbf{S}_1^i, \dots, \mathbf{S}_{T_i-1}^i, A_{T_i-1}^i, \mathbf{R}_{T_i-1}^i, \mathbf{S}_{T_i}^i) \right\}_{i=1}^n,$$
- n i.i.d. patient trajectories.
- For each patient, his/her follow up time length T_i may be different.
- $\mathbf{S}_t \in \mathcal{S}$ denotes a vector of patient clinical information recorded up to and including time point t
- $A_t \in \mathcal{A}$ denotes the treatment assignment at time point t after measuring \mathbf{S}_t ,
- $\mathbf{R}_t \in \mathbb{R}^J$ is the reward obtained after treatment A_t is assigned.

Potential outcome frameworks

- Let $\bar{\mathbf{a}}_t = (a_0, a_1, \dots, a_t) \in \bar{\mathcal{A}}_t$ be a possible treatment assignment sequence up to time point t .
- Let $\bar{\mathbf{s}}_t = (\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_t) \in \bar{\mathcal{S}}_t$ be a possible state sequences up to time point t , $t \geq 0$.
- The set of potential outcomes is

$$\mathbf{W}^* = \left\{ \mathbf{S}_1^*(a_0), \mathbf{S}_2^*(\bar{\mathbf{a}}_1), \dots, \mathbf{S}_{t+1}^*(\bar{\mathbf{a}}_t), \dots, \text{for all } \bar{\mathbf{a}}_\infty \in \bar{\mathcal{A}}_\infty \right\},$$
where $\mathbf{S}_{t+1}^*(\bar{\mathbf{a}}_t)$ is the potential state at $(t + 1)$ -th time point that would have been observed if the individual had been assigned the treatment sequence $\bar{\mathbf{a}}_t$, $t \geq 0$.

Potential outcome frameworks

- *Consistency*: $\mathbf{S}_{t+1} = \mathbf{S}_{t+1}^*(\bar{\mathbf{A}}_t)$, for all $t \geq 0$.
- *Sequential randomization assumption*: $A_{t+1} \perp\!\!\!\perp \mathbf{W}^* \mid \bar{\mathbf{S}}_{t+1}, \bar{\mathbf{A}}_t$, for all $t \geq 0$.
- *Positivity*: there exists $\epsilon_0 > 0$, so that
 $P\left(A_{t+1} = a_{t+1} \mid \bar{\mathbf{S}}_{t+1} = \bar{\mathbf{s}}_{t+1}, \bar{\mathbf{A}}_t = \bar{\mathbf{a}}_t\right) > \epsilon_0$, for all $a_{t+1} \in \mathcal{A}$,
 $\bar{\mathbf{a}}_t \in \bar{\mathcal{A}}_t$ and $\bar{\mathbf{s}}_{t+1} \in \bar{\mathcal{S}}_{t+1}$, and all $t \geq 0$.

Markov Decision Process

- *Markov assumption:* $\mathbf{S}_{t+1} \perp\!\!\!\perp (\bar{\mathbf{A}}_{t-1}, \bar{\mathbf{S}}_{t-1}) \mid (A_t, \mathbf{S}_t)$, for all $t \geq 1$.
- *Time homogeneity:* the conditional density
 $P_t(\mathbf{S}_{t+1} = \mathbf{s}' \mid A_t = a, \mathbf{S}_t = \mathbf{s}) = P(\mathbf{s}' \mid a, \mathbf{s})$ for all $\mathbf{s}, \mathbf{s}' \in \mathcal{S}$ and $a \in \mathcal{A}$ and $t \geq 0$, where \mathbf{s} and \mathbf{s}' denote the current state and the next state, respectively.

Value functions of dynamic treatment regimes

- A dynamic treatment regime, $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- The value function $V^\pi(\mathbf{s})$ of a state under a certain policy π is defined as the expected total discounted rewards when the process begins in state \mathbf{s} and all decisions are made according to policy π .
- $V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{s}}^\pi \sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, a_t, \mathbf{s}_{t+1})$.
- Bellman equation,

$$V^\pi(\mathbf{s}) = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, \pi(\mathbf{s})) (R(\mathbf{s}, \pi(\mathbf{s}), \mathbf{s}') + \gamma V^\pi(\mathbf{s}')).$$

Value functions of dynamic treatment regimes

- Moreover, the state-action value function under policy π ,
$$Q^\pi(\mathbf{s}, a) = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) \left(R(\mathbf{s}, a, \mathbf{s}') + \gamma Q^\pi(\mathbf{s}', \pi(\mathbf{s}')) \right),$$
 is defined similar but the first step takes action a .
- In clinical cases, the transition model \mathbb{P} is unknown, optimal regimes must be learned from observed dataset.
- In infinite horizon setting, as time steps are dropped, we break n observed trajectories into 4-tuple of $(\mathbf{s}, a, r, \mathbf{s}')$ for estimating value functions.

Least-squares policy evaluation

- The state-action value function \mathbf{Q}^π is considered the fixed point of the Bellman operator: $\mathbf{Q}^\pi = T^\pi \mathbf{Q}^\pi$, where the Bellman operator defined as

$$T^\pi \mathbf{Q}^\pi = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) \left(\mathbf{R}(\mathbf{s}, a, \mathbf{s}') + \gamma \sum_{a' \in \mathcal{A}} \pi(a' | \mathbf{s}') \mathbf{Q}^\pi(\mathbf{s}', a') \right).$$

- The approximation for each component is

$$Q_j^\pi(\mathbf{s}, a; \mathbf{w}) = \sum_{k=1}^K \phi_{j,k}(\mathbf{s}, a) w_{j,k} = \phi_j^\top(\mathbf{s}, a) \mathbf{w}_j$$

- $\mathbf{w}_j = (w_{j,1}, \dots, w_{j,K})^\top$ are the parameters to estimate.
- The basis functions $\phi_j(\mathbf{s}, a) = (\phi_{j,1}(\mathbf{s}, a), \dots, \phi_{j,K}(\mathbf{s}, a))^\top$ are arbitrary and fixed, which are often non-linear functions of \mathbf{s} and a .

Least-squares policy evaluation

Algorithm 1: Least-squares policy evaluation LSQ

Input : A sample set D of 4 tuples (s', a, s, r)

k : Number of basis functions

ϕ : Basis functions

γ : Discount factor

π : policy whose value function is sought

Output: Weights \hat{w}^π

$\hat{A} \leftarrow 0$ // $(k \times k)$ matrix

$\hat{b} \leftarrow 0$ // $(k \times 1)$ vector

for each $(s, a, r, s') \in D$

$\hat{A} \leftarrow \hat{A} + \phi(s, a) (\phi(s, a) - \gamma \phi(s', \pi(s')))^T$

$\hat{b} \leftarrow \hat{b} + \phi(s, a)r$

$\hat{w}^\pi \leftarrow \hat{A}^{-1}\hat{b}$

return \hat{w}^π



Constrained optimal dynamic treatment regimes

- The average of value functions is referred as competing outcomes, denoted as $\mathbf{V}(\pi) = \mathbb{E}\mathbf{V}^\pi(\mathbf{s})$.
- Mathematically,

$$\max_{\pi \in \Pi} V_1(\pi),$$

$$\text{subject to } V_j(\pi) \leq v_{j-1},$$

where $j = 2, \dots, J$.

Chemotherapy mathematical model

- The chemotherapy mathematical model, a system of ordinary differential equations (ODE), proposed by Zhao et al, is modified and used to generate a hypothetical clinical trial data.
- The model reflects the capability of the drug to suppress tumor growth, as well as its negative impact on patient wellness due to the toxicity of chemotherapy.
- The dose assignment is discretized to $L = 5$ levels,
 $\mathcal{A} = \{0.00, 0.25, 0.50, 0.75, 1.00\}$.

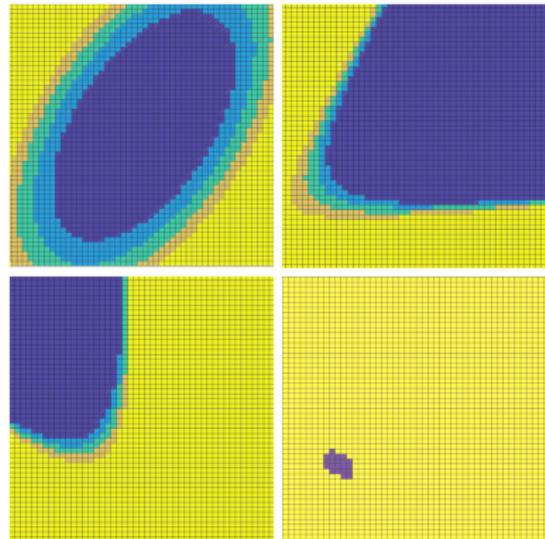
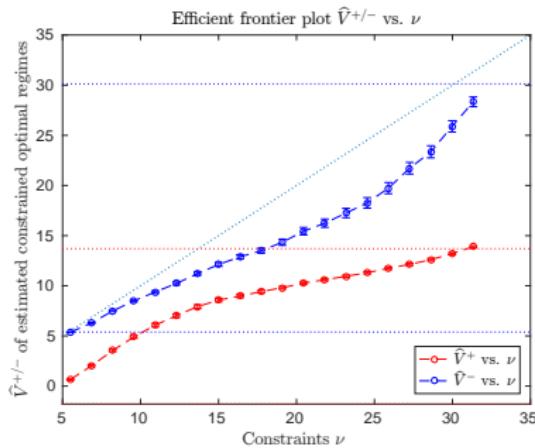


Figure: Treatment assignments under four different constraints

Discussion and Future Work

- Consistency and asymptotic normality were proven for single-stage and finite-stage
- Theoretical work for infinite-stage constrained problem
- Real clinical data application
- Incorporate more complex input data, such as omics data, medical images, etc.
- Learning rewards/objective/constraints from clinical experts

Thank You & Happy Holidays!