

Constrained Estimation of Single-stage Optimal Treatment Regimes

Shuping Ruan, Eric Laber

Department of Statistics, North Carolina State University

March 23, 2018

1 Introduction

Precision medicine aims to accommodate interventions to individual patient attributes. For chronic illnesses, such as cancer, diabetes and so on, clinicians often need to make sequences of treatment decisions based on the evolving status of the patient’s condition. To personalize this multi-stage intervention process, researchers have been developing dynamic treatment regimes (DTRs) to inform clinicians of treatment decisions adaptively. DTRs are a sequential decision making process, of which each decision is made based on the evolving patient status with the goal of maximizing the overall long-term treatment efficacy. It is well-studied in the statistical and biomedical literature [19, 18, 21, 26, 10, 13]. From the standpoint of Markov decision process, reinforcement learning algorithms, such as Q-learning [21], A-learning [3], V-learning [17] etc., are developed to estimate optimal treatment regimes.

Most previous methods for construct optimal dynamic treatment regimes have focused on optimizing a scalar measurement of the long-term efficacy over a fixed time period (finite stage). However, in practice, the clinical situations are more complex. First, they often require consideration of the trade-off among multiple objectives, e.g., effectiveness, side-effect, cost, and so on. The preference of those objectives varies among people and changes over time, thus a single scalar reward or value can not represent the quality of a policy well enough. Thus, considering multiple rewards are necessary. Previous works by Lizotte et al. learn the value function and optimal policy for all preferences, i.e., all the possible convex combination of all the rewards [14, 15]. More recent works adopt multi-objective Markov decision processes (MOMDPs) framework with finite stage. Practical domination is proposed for flexibility based on Pareto domination, and a set of policies that are maximal based on the partial order are treated as indistinguishably optimal [9, 16]. Secondly, patients with chronic diseases are often monitored and treated throughout their life. It often requires taking real-time actions and has no a-priori fixed end point (infinite stage),

and progress is made in infinite stage reinforcement learning for health applications [5, 20, 17].

To deal with the trade-off between multiple objectives, we take a different perspective and adopt the constrained Markov decision processes (CMDPs) framework with infinite horizon. CMDPs are a well-studied framework for reinforcement learning under constraints [2]. The goal is to find the optimal policy, while satisfying constraints on expectations of secondary costs. For many applications of reinforcement learning, the constrained approach is more intuitive and more practical than eliciting a single reward function in order to achieve desirable results. For instance, satisfying safety constraints is necessary for systems that physically interact with humans. Previously, linear programming is used to seek 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 without modeling the underlying dynamics. Recently, Achiam et al [1] proposed constrained policy optimization, a general-purpose policy search algorithm for constrained reinforcement learning guaranteeing near-constraint satisfaction at each iteration. Taking into consideration properties of clinical applications, such as data scarcity and off-policy learning, we develop an algorithm by taking advantage of least-squares policy evaluation and interior-point methods for estimating constrained optimal dynamic treatment regimes. Our method is applied to a simulated cancer trial dataset based on a chemotherapy mathematical model.

2 Methodology

2.1 Set-up

Observed Data

We use dataset observed over a finite length of time steps to construct a regime in the setting of infinite horizon Markov decision process. The structure of the available data is $\mathbf{D} = \left\{ (\mathbf{S}_0^i, A_0^i, \mathbf{R}_0^i, \mathbf{S}_1^i, \dots, \mathbf{S}_{T_i-1}^i, A_{T_i-1}^i, \mathbf{R}_{T_i-1}^i, \mathbf{S}_{T_i}^i) \right\}_{i=1}^n$, a set of n independent, identically distributed trajectories of $(\mathbf{S}_0, A_0, \mathbf{R}_0, \mathbf{S}_1, \dots, \mathbf{S}_{T-1}, A_{T-1}, \mathbf{R}_{T-1}, \mathbf{S}_T)$. Note $T \in \mathbb{N}$ denotes the total number of follow-up time steps for a patient. For each patient, his or 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 , referred as *state* in the reinforcement learning vocabulary. $\mathcal{S} \subseteq \mathbb{R}^m$ denotes the support for state variable. Adopting the time homogeneous Markov decision process model, we assume \mathcal{S} is the same across all time points t . $A_t \in \mathcal{A}$ denotes the treatment assignment at time point t after measuring \mathbf{S}_t , referred as *action* in reinforcement learning. \mathcal{A} denotes the support for treatment assignment, a finite set of all possible treatment options, and is assumed to be the same across all time

points t . $\mathbf{R}_t \in \mathbb{R}^J$ is the reward obtained after treatment A_t is assigned. We assume the reward, possibly defined based on domain expertise, is a known vector-valued function $\mathbf{r}: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^J$, so that $\mathbf{R}_t = \mathbf{r}(\mathbf{S}_t, A_t, \mathbf{S}_{t+1})$. The vector-valued reward function is the same across all the time points t as well. Moreover, if a patient dies during the follow-up, say at decision point t , we set $\mathbf{S}_t = \emptyset$, referred to as the absorbing state in reinforcement learning. Then, the patient's treatment assignment at time t is $A_t = \emptyset$, and his/her length of follow-up $T = t$.

Potential outcomes

In reality, a patient can only be assigned to one of the sequences of treatment assignments. Hence, we can only observe the consequence of that treatment sequence, while the others remain unobserved. To identify the average causal effect of a certain regime, we adopt the counter-factual or potential outcomes framework, established by Neyman, Rubin and Robins for assessment of the time-dependent treatment effect from either randomized or observational studies [8, 22, 23]. 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 , $t \geq 0$, where $\bar{\mathcal{A}}_t = \mathcal{A} \times \dots \times \mathcal{A}$ is the set of all possible the treatment assignment sequences up to time point t . Let $\bar{\mathbf{s}}_t = (s_0, s_1, \dots, s_t) \in \bar{\mathcal{S}}_t$ be a possible state sequences up to time point t , $t \geq 0$, where $\bar{\mathcal{S}}_t = \mathcal{S} \times \dots \times \mathcal{S}$ is the set of all possible state sequences up to time point t . The set of potential outcomes is $\mathbf{W}^* = \{\mathbf{S}_1^*(a_0), \mathbf{S}_2^*(\bar{\mathbf{a}}_1), \dots, \mathbf{S}_{t+1}^*(\bar{\mathbf{a}}_t), \dots\}$, for all $\bar{\mathbf{a}}_\infty \in \bar{\mathcal{A}}_\infty$, 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 $\bar{\mathbf{a}}_t$, $t \geq 0$. Moreover, if $\mathbf{S}_{t+1}^*(\bar{\mathbf{a}}_t) = \emptyset$ happens at time point $t+1$, then $\mathbf{S}_{t+2}^*(\bar{\mathbf{a}}_{t+1}) = \mathbf{S}_{t+3}^*(\bar{\mathbf{a}}_{t+2}) = \dots = \emptyset$, which indicates the patient, if followed treatment assignment sequence $\bar{\mathbf{a}}_t$, would have died at time point $t+1$ in the counter-factual world. Rewards with respect to potential states are $\mathbf{R}_t^* = \mathbf{r}(\mathbf{S}_t^*, A_t, \mathbf{S}_{t+1}^*)$. The following assumptions are made in the potential outcome framework [8, 24, 23, 22, 7].

- *A1. Consistency:* $\mathbf{S}_{t+1} = \mathbf{S}_{t+1}^*(\bar{\mathbf{A}}_t)$, for all $t \geq 0$.
- *A2. Sequential randomization assumption:* $A_{t+1} \perp \mathbf{W}^* \mid \bar{\mathbf{S}}_{t+1}, \bar{\mathbf{A}}_t$, for all $t \geq 0$.
- *A3. Positivity:* there exists $\epsilon_0 > 0$, so that $P(A_{t+1} = a_{t+1} \mid \bar{\mathbf{S}}_{t+1} = \bar{\mathbf{s}}_{t+1}, \bar{\mathbf{A}}_t = \bar{\mathbf{a}}_t) > \epsilon_0$, 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 $t \geq 0$.

These assumptions link the potential outcome and the observed data, and are guaranteed in well-designed Sequential, Multiple Assignment, Randomized Trials (SMARTs). Therefore, the observed data from those trials are used to infer the average causal effect of a regime of interest.

Markov Decision Processes

To construct a regime in the infinite-horizon setting using a dataset observed over a finite number of time steps, we assume that the underlying dynamics

is a time homogeneous Markov Decision Processes (MDPs). In infinite-horizon setting, MDP is considered as a 5-tuple of $(\mathcal{S}, \mathcal{A}, \mathbb{P}, \mathbf{R}, \gamma)$, where \mathcal{S} , \mathcal{A} and \mathbf{R} is as described above. Additionally, \mathbb{P} is a markovian transition model in which $p(\mathbf{s}' | \mathbf{s}, a)$ denotes the probability density of a transition to state \mathbf{s}' when taking action a in state \mathbf{s} . A discount factor for future reward, $\gamma \in [0, 1)$, is also introduced to form total discounted rewards, which are the value functions to operate constrained optimization on. The following assumptions are made for infinite-stage time homogeneous Markov decision process.

- *A4. Markov assumption:* $\mathbf{S}_{t+1} \perp\!\!\!\perp (\bar{\mathbf{A}}_{t-1}, \bar{\mathbf{S}}_{t-1}) | (A_t, \mathbf{S}_t)$, for all $t \geq 1$.
- *A5. Time homogeneity:* the conditional density $P_t(\mathbf{S}_{t+1} = \mathbf{s}' | A_t = a, \mathbf{S}_t = \mathbf{s}) = p(\mathbf{s}' | a, \mathbf{s})$ 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.

As the time homogeneity is assumed in infinite-stage(3. setting, time step subscripts maybe dropped for simplicity.

Values of dynamic treatment regimes

A dynamic treatment regime is equivalent to a *policy* in reinforcement learning vocabulary, which is mostly to be considered deterministic. Thus, a dynamic treatment regime, $\pi : \mathcal{S} \rightarrow \mathcal{A}$, is defined as a function which maps the support of the state variable to the set of the possible treatment assignments. As time homogeneity is assumed, we consider only stationary deterministic regimes, where this mapping function $\pi(\mathbf{s})$ does not change over time. Hence, a patient with state $\mathbf{S}_t = \mathbf{s}$ at time point t will be assigned with treatment $A_t = \pi(\mathbf{s})$ for all t . The value function $\mathbf{V}^\pi(\mathbf{s})$ 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 . Mathematically, $\mathbf{V}^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{s}}^\pi \sum_{t=0}^{\infty} \gamma^t \mathbf{r}(\mathbf{s}_t, a_t, \mathbf{s}_{t+1})$, where $\mathbb{E}_{\mathbf{s}}^\pi$ is the expectation when the initial state is \mathbf{s} and a policy π is followed. The value function can also be defined recursively via the bellman equation, $\mathbf{V}^\pi(\mathbf{s}) = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, \pi(\mathbf{s})) (\mathbf{R}(\mathbf{s}, \pi(\mathbf{s}), \mathbf{s}') + \gamma \mathbf{V}^\pi(\mathbf{s}'))$. Here, $\mathbf{V}^\pi(\mathbf{s}) \in \mathbb{R}^J$ has the same dimensionality as the reward vector \mathbf{R} , as we are considering multiple reward functions instead of a scalar reward function. Moreover, the state-action value function under policy π , $\mathbf{Q}^\pi(\mathbf{s}, a) = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) (\mathbf{R}(\mathbf{s}, a, \mathbf{s}') + \gamma \mathbf{Q}^\pi(\mathbf{s}', \pi(\mathbf{s}')))$, is defined similar but the first step takes action a and $\mathbf{Q}^\pi(\mathbf{s}, a) \in \mathbb{R}^J$. In clinical cases, the transition model \mathbb{P} is unknown, optimal regimes must be learn from observed dataset. In infinite horizon setting, as time steps are dropped, we break n observed trajectories into 4-tuple of $(\mathbf{s}, a, \mathbf{r}, \mathbf{s}')$ for estimating value functions. The counter-factual assumptions (A1-A3) rule out the confounding phenomena and guarantee the identifiability of the average causal effect of a regime.

Define infinite-stage constrained optimal dynamic treatment regimes

Our strategy to cope with multiple competing outcomes is constrained optimization. We optimize the primary outcome of interest, subject to the constraints on the secondary outcomes, over the space of all the possible regimes under consideration, Π . Here, the average of value functions is referred as competing outcomes, denoted as $\mathbf{V}(\pi) = \mathbb{E}\mathbf{V}^\pi(\mathbf{s})$. The space of regimes under consideration may be crafted by experts with domain knowledge via policy function approximation. As we have $\mathbf{V}^\pi(\mathbf{s}) = \left(V_1(\pi), V_2^\pi(\mathbf{s}), \dots, V_q^\pi(\mathbf{s})\right)^\top$, Let $V_1(\pi) = \mathbb{E}V_1(\pi)$ be the primary outcome of interest, and $V_j(\pi) = \mathbb{E}V_j^\pi(\mathbf{s})$, $j = 2, 3, \dots, J$ be the secondary outcomes. Mathematically,

$$\begin{aligned} & \max_{\pi \in \Pi} V_1(\pi), \\ & \text{subject to } V_j(\pi) \leq \nu_{j-1}, \end{aligned} \quad (1)$$

where $j = 2, \dots, J$. The constraint upper-bounds $\boldsymbol{\nu} = (\nu_1, \nu_2, \dots, \nu_{J-1})^\top$ can be specified based on patient preference and/or expert domain knowledge. Therefore, we define an infinite-stage constrained optimal regime as $\pi_{\boldsymbol{\nu}}^* = \operatorname{argmax}_{\pi \in \Pi} V_1(\pi)$, subject to $V_j(\pi) - \nu_{j-1} \leq 0$, where $j = 2, \dots, J$. Denote the feasible policy space, which is the set of all policy satisfying the constraints, as $\mathcal{F}(\Pi)$. For all $\pi \in \mathcal{F}(\Pi)$ and all $j = 2, \dots, J$, $V_j(\pi) \leq \nu_{j-1}$. Then, an infinite-stage constrained optimal regime can also be written as $\pi_{\boldsymbol{\nu}}^* = \operatorname{argmax}_{\pi \in \mathcal{F}(\Pi)} V_1(\pi)$. To search over a feasible policy space with manageable computation complexity, we use policy function approximation such that $\pi(\mathbf{s}) = \pi(\mathbf{s}; \boldsymbol{\theta})$, where $\boldsymbol{\theta} \in \mathbb{R}^q$ is the indexing parameter for policies. Hence, we use $V_j(\pi)$ and $V_j(\boldsymbol{\theta})$ interchangeably, for all j . The search space is reduced from the set of all feasible policies to the feasible space of the indexing parameter $\boldsymbol{\theta}$, denoted as $\mathcal{F}(\boldsymbol{\Theta}) = \{\boldsymbol{\theta} \in \mathbb{R}^q : V_j(\boldsymbol{\theta}) \leq \nu_{j-1}, j = 2, \dots, J\}$.

To carry out policy search, we need to solve the constrained optimization problem (3.1). This is done using interior-point methods, which are constrained nonlinear optimization methods for finding local optimums, implemented in Matlab `fmincon` [25, 4]. We also need a method to estimate the value functions using observed dataset. This is done by least-square policy evaluation (LSQ), a part of the least-squares policy iteration (LSPI) algorithm. [12, 11].

2.2 Interior point method for constrained optimization

To solve our constrained optimization problem (3.1) above, interior point algorithm is used. As the optimization softwares often implemented as minimization instead of maximization, we denote $v_1(\boldsymbol{\theta}) = -V_1(\boldsymbol{\theta})$ and $v_j(\boldsymbol{\theta}) = V_j(\boldsymbol{\theta}) - \nu_{j-1}$, for $j = 2, \dots, J$. Hence, problem (3.1) notation is simplified as

$$\begin{aligned} & \min_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} v_1(\boldsymbol{\theta}) \\ & \text{subject to } v_j(\boldsymbol{\theta}) \leq 0, \end{aligned} \quad (2)$$

where $j = 2, \dots, J$. The interior point method solves a following sequence of approximate minimization problem (2), where ρ is always positive and approaches to zero in the limit. For each $\rho > 0$, the approximate problem is

$$\min_{\boldsymbol{\theta}, \mathbf{z}} f_{\rho}(\boldsymbol{\theta}, \mathbf{z}) = \min_{\boldsymbol{\theta}} v_1(\boldsymbol{\theta}) - \rho \sum_{j=2}^J \ln(z_j), \text{ subject to } v_j(\boldsymbol{\theta}) + z_j = 0, \quad (3)$$

where $j = 2, \dots, J$. There are as many slack variables z_j as there are inequality constraints v_j . The z_j are restricted to be positive to keep $\ln(z_j)$ bounded. As ρ decreases to zero, the minimum of f_{ρ} should approach the minimum of v_1 . The added logarithmic term is called a barrier function [25, 4].

2.3 Least-squares policy evaluation

Least-squares policy evaluation, is adopted to approximate the state-action value function of a fixed regime/policy. As it is the state-action value function being approximated, instead of the state value function, changing policy is allowed without a model for the underlying dynamics. The exact \mathbf{Q}^{π} values for all state-action pairs can be found by solving the linear system of the Bellman equations,

$$\mathbf{Q}^{\pi}(\mathbf{s}, a) = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) \left\{ \mathbf{R}(\mathbf{s}, a, \mathbf{s}') + \gamma \sum_{a' \in \mathcal{A}} \pi(a' | \mathbf{s}') \mathbf{Q}^{\pi}(\mathbf{s}', a') \right\},$$

for any $\mathbf{s} \in \mathcal{S}$ and $a' \in \mathcal{A}$. Thus, 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

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

Linear approximation is a common way for estimating value functions, so that each component of the vector $\mathbf{Q}^{\pi}(\mathbf{s}, a; w)$ are approximated by a linear parametric combination of K basis functions (features). As $\mathbf{Q}(\mathbf{s}, a) = (Q_1(\mathbf{s}, a), Q_2(\mathbf{s}, a), \dots, Q_J(\mathbf{s}, a))^{\top}$, the approximation for each component is $Q_j^{\pi}(\mathbf{s}, a; w) = \sum_{k=1}^K \phi_{j,k}(\mathbf{s}, a) w_{j,k} = \boldsymbol{\phi}_j^{\top}(\mathbf{s}, a) \mathbf{w}_j$, where $\mathbf{w}_j = (w_{j,1}, \dots, w_{j,K})^{\top}$ are the parameters to estimate. Moreover, the basis functions $\boldsymbol{\phi}_j(\mathbf{s}, a) = (\phi_{j,1}(\mathbf{s}, a), \dots, \phi_{j,K}(\mathbf{s}, a))^{\top}$ are arbitrary and fixed, which are often non-linear functions of \mathbf{s} and a . It is also required that the basis functions $\phi_{j,k}$ are linearly independent to ensure that there are no redundant parameters and that the matrices involved in the computations are full rank.

Substituting each component of the Q function vector with the linear approximator, we get, for $j = 1, \dots, J$,

$$\boldsymbol{\phi}_j^{\top}(\mathbf{s}, a) \mathbf{w}_j = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) R_j(\mathbf{s}, a, \mathbf{s}') + \gamma \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) \sum_{a' \in \mathcal{A}} \pi(a' | \mathbf{s}') \boldsymbol{\phi}_j^{\top}(\mathbf{s}', a') \mathbf{w}_j.$$

This fixed point equation then can be rearranged as

$$\left\{ \phi_j^\top(\mathbf{s}, a) - \gamma \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) \sum_{a' \in \mathcal{A}} \pi(a' | \mathbf{s}') \phi_j^\top(\mathbf{s}', a') \right\} \mathbf{w}_j = \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}, a) R_j(\mathbf{s}, a, \mathbf{s}'). \quad (4)$$

Given a sample set of 4-tuples $\mathcal{D} = \{\mathbf{s}^i, a^i, \mathbf{s}^{i'}, \mathbf{r}^i\}_{i=1}^N$, the equation (3.4) above becomes a over-constrained/overdetermined linear system over the parameter vector \mathbf{w}_j . The linear system can be written as

$$\mathbf{B}_j \mathbf{w}_j = \mathbf{b}_j,$$

where $\mathbf{B}_j = \Phi_j^\top (\Phi_j - \gamma \mathbf{P}^\pi \Phi_j)$ and $\mathbf{b}_j = \Phi_j^\top \mathbf{R}_j$. Moreover,

$$\Phi_j = \begin{pmatrix} \phi_j^\top(\mathbf{s}^1, a^1) \\ \phi_j^\top(\mathbf{s}^2, a^2) \\ \dots \\ \phi_j^\top(\mathbf{s}^n, a^n) \end{pmatrix}, \quad \mathbf{P}^\pi \Phi_j = \begin{pmatrix} \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}^1, a^1) \phi_j^\top(\mathbf{s}', \pi(\mathbf{s}')) \\ \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}^2, a^2) \phi_j^\top(\mathbf{s}', \pi(\mathbf{s}')) \\ \dots \\ \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}^n, a^n) \phi_j^\top(\mathbf{s}', \pi(\mathbf{s}')) \end{pmatrix},$$

$$\text{and } \mathbf{R}_j = \begin{pmatrix} \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}^1, a^1) R_j(\mathbf{s}^1, a^1, \mathbf{s}') \\ \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}^2, a^2) R_j(\mathbf{s}^2, a^2, \mathbf{s}') \\ \dots \\ \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}' | \mathbf{s}^n, a^n) R_j(\mathbf{s}^n, a^n, \mathbf{s}') \end{pmatrix},$$

where $\phi_j^\top(\mathbf{s}', \pi(\mathbf{s}')) = \sum_{a' \in \mathcal{A}} \pi(a' | \mathbf{s}') \phi_j^\top(\mathbf{s}', a')$. Since the transition probability function $p(\mathbf{s}' | \mathbf{s}, a)$ and reward functions $R_j(\mathbf{s}, a, \mathbf{s}')$ may be unknown, we can construct approximators for \mathbf{B} and \mathbf{b} using samples. More precisely, we have approximated versions of Φ_j , $\mathbf{P}^\pi \Phi_j$ and \mathbf{R}_j based on the sample set as follows:

$$\widehat{\Phi}_j = \begin{pmatrix} \phi_j^\top(\mathbf{s}^1, a^1) \\ \phi_j^\top(\mathbf{s}^2, a^2) \\ \dots \\ \phi_j^\top(\mathbf{s}^n, a^n) \end{pmatrix}, \quad \widehat{\mathbf{P}^\pi \Phi}_j = \begin{pmatrix} \phi_j^\top(\mathbf{s}^{1'}, \pi(\mathbf{s}^{1'})) \\ \phi_j^\top(\mathbf{s}^{2'}, \pi(\mathbf{s}^{2'})) \\ \dots \\ \phi_j^\top(\mathbf{s}^{n'}, \pi(\mathbf{s}^{n'})) \end{pmatrix}, \quad \text{and } \widehat{\mathbf{R}}_j = \begin{pmatrix} r_1 \\ r_2 \\ \dots \\ r_n \end{pmatrix}.$$

Given $\widehat{\Phi}_j$, $\widehat{\mathbf{P}^\pi \Phi}_j$, and $\widehat{\mathbf{R}}_j$, \mathbf{B}_j and \mathbf{b}_j can be approximated as $\widehat{\mathbf{B}}_j = n^{-1} \widehat{\Phi}_j^\top (\widehat{\Phi}_j - \gamma \widehat{\mathbf{P}^\pi \Phi}_j) = n^{-1} \sum_{i=1}^n \phi_j(\mathbf{s}^i, a^i) \left(\phi_j(\mathbf{s}^i, a^i) - \gamma \phi_j(\mathbf{s}^{i'}, \pi(\mathbf{s}^{i'})) \right)$ and $\widehat{\mathbf{b}}_j = n^{-1} \widehat{\Phi}_j^\top \widehat{\mathbf{R}}_j = n^{-1} \sum_{i=1}^n \phi_j(\mathbf{s}^i, a^i) r_j^i$. It is shown in the least-squares policy iteration paper that $\lim_{n \rightarrow \infty} \widehat{\mathbf{B}}_j = \mathbf{B}_j$ and $\lim_{n \rightarrow \infty} \widehat{\mathbf{b}}_j = \mathbf{b}_j$, if the samples are uniformly distributed over the state space. Moreover, the Markov property ensures that the solution $\widehat{\mathbf{w}}^\pi = \widehat{\mathbf{B}}^{-1} \widehat{\mathbf{b}}$ will converge to the true solution \mathbf{w}^π for sufficiently large n whenever \mathbf{w}^π exists [12, 11]. The least-squares policy evaluation algorithm is listed in Algorithm 2 below.

Equipped with least-squares policy evaluation, we can hence calculate the values of any arbitrary regime/policy π . For $j = 1, \dots, J$, $V(\pi)$ is estimated by $n^{-1} \sum_{i=1}^n \widehat{Q}_j^\pi(\mathbf{s}^i, \pi(\mathbf{s}^i)) = n^{-1} \sum_{i=1}^n \phi_j^\top(\mathbf{s}^i, \pi(\mathbf{s}^i)) \widehat{\mathbf{w}}_j$.

2.3.1 Algorithm

Putting together in the following box, our algorithm uses interior point method for policy search in terms policy indexing parameters θ , and least-squares policy evaluation for policy evaluation.

Algorithm 1: Constrained optimal regime with least-squares policy evaluation [12, 11] for policy evaluation and interior-point method [25, 4] for policy search.

Input : A sample set D of 4 tuples (s', a, s, r)
Output: Estimated constrained optimal regime $\hat{\pi}_\nu$ indexed by $\hat{\theta}_\nu$

$\pi \leftarrow$ random initialization
until converge
 $\hat{Q}^\pi(s, a) \leftarrow$ least-squares policy evaluation (D, π)
 $\hat{V}(\pi) \leftarrow \frac{1}{n} \sum_{i=1}^n \hat{Q}^\pi(s^i, \pi(s^i))$
 $\pi \leftarrow \operatorname{argmax}_{\pi \in \mathcal{F}(\Pi)} V(\pi)$
end
 $\hat{\pi}_\nu \leftarrow \pi$
return $\hat{\pi}_\nu$

Algorithm 2: Least-squares policy evaluation (LSQ). [12, 11]

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

$\hat{B} \leftarrow 0$ // $(k \times k)$ matrix
 $\hat{b} \leftarrow 0$ // $(k \times 1)$ vector
for each $(s, a, r, s') \in D$
 $\hat{B} \leftarrow \hat{B} + \phi(s, a) (\phi(s, a) - \gamma \phi(s', \pi(s')))^T$
 $\hat{b} \leftarrow \hat{b} + \phi(s, a)r$
 $\hat{w}^\pi \leftarrow \hat{B}^{-1} \hat{b}$
return \hat{w}^π

3 Simulation

3.1 Chemotherapy mathematical model

The chemotherapy mathematical model, a system of ordinary differential equations (ODE), proposed by Zhao et al [27]. is modified and used to generate a hypothetical clinical trial data. Their 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, $\mathcal{A} = \{0.00, 0.25, 0.50, 0.75, 1.00\}$. Two state variables W_t and M_t are considered. W_t denotes the patient negative wellness (toxicity) at time point t . M_t denotes the tumor size observed at time point t . A_t denotes chemotherapy agent dose at time point t . The ODE system is modeled as [27]

$$\begin{aligned}\dot{W}_t &= b_1 \max(M_t, M_0) + c_1(A_t - d_1), \\ \dot{M}_t &= (b_2 \max(W_t, W_0) - c_2(A_t - d_2)) \times \mathbb{I}(M_t > 0),\end{aligned}$$

where decision points are $t = 0, 1, \dots, T - 1$. Moreover, \dot{W}_t and \dot{M}_t are the transition functions. The indicator function term $\mathbb{I}(M_t > 0)$ means tumor size is absorbed at 0, the patient has been cured and no future tumor recurrence considered. These changing rate yields a piece-wise linear model over time. Constants value are set as $b_1 = 0.1, b_2 = 0.15, c_1 = 1.2, c_2 = 1.2, d_1 = 0.5$ and $d_2 = 0.5$. The initial states are draw as $M_0 \sim \text{Uniform}(0, 2)$ and $W_0 \sim \text{Uniform}(0, 2)$. The initial dose level assignment is draw as $A_0 \sim \text{Discrete Uniform}(0.25, 0.50, 0.75, 1.00)$. The state variables for the next time point can be obtained via $W_{t+1} = \max(W_t + \dot{W}_t, 0)$, $M_{t+1} = \max(M_t + \dot{M}_t, 0)$. The dose level assignment is drawn as $A_t \sim \text{Discrete Uniform}(0.00, 0.25, 0.50, 0.75, 1.00)$, for $t = 1, \dots, T - 1$.

The survival status of the patient, denoted by F_t , is also modeled. We assume everyone is alive at the initial decision point $t = 0$, that is $p_0 = 0$ and $F_0 = 0$. Death events occur during time interval $(t - 1, t]$, $t = 1, 2, \dots, 6$, and are recorded at the end of each interval as variable F_t , $t = 1, 2, \dots, 6$. Assume that survival status depends on both toxicity and tumor size. For each time interval $(t - 1, t]$, define the hazard function as $\lambda(t)$, where $\log \lambda(t) = \mu_0 + \mu_1 W_t + \mu_2 M_t$, $\mu_1 = \mu_2 = 1$ and $\mu_0 = -8.5$. This again is a piece-wise linear approximation with $\lambda(t) = \exp(\mu_0 + \mu_1 W_t + \mu_2 M_t)$. Then, the cumulative hazard function during time interval $(t - 1, t]$ is $\Delta\Lambda(t) = \sum_{s=t-1}^t \lambda(s) ds = \sum_{s=t-1}^t \exp(\mu_0 + \mu_1 W_t + \mu_2 M_t) ds = \exp(\mu_0 + \mu_1 W_t + \mu_2 M_t)$. The survival function is $\Delta F(t) = \exp(-\Delta\Lambda(t)) = \exp(-\exp(\mu_0 + \mu_1 W_t + \mu_2 M_t))$. The random event of death during time interval $(t - 1, t]$ is drawn as $F_t \sim \text{Bernoulli}(p_t)$, where $p_t = 1 - \exp(-\exp(\mu_0 + \mu_1 W_t + \mu_2 M_t))$. If $F_{t-1} = 1$, then $F_t = 1$. Also, as long as death occurred, all the other state variables at the following decision points are all set to null.

The reward functions here is a bivariate vector, consisting of positive reward and negative reward, denoted as $\mathbf{R}_t = (R_t^+, R_t^-)^\top$. The positive reward function is used to assess tumor size reduction, while the negative reward to assess the increase of patient negative wellness (toxicity). Specifically, the reward functions

are defined as follow.

$$\begin{aligned}
R^+(\cdot, t) = & -15 \times \mathbb{I}(F_{t+1} = 1) \\
& + 5 \times \mathbb{I}(F_{t+1} \neq 1, M_{t+1} = 0) \\
& + 5 \times |M_{t+1} - M_t| \times \mathbb{I}(F_{t+1} \neq 1, M_{t+1} - M_t \leq 0) \\
& + 5 \times |M_{t+1} - M_t| \times \mathbb{I}(F_{t+1} \neq 1, M_{t+1} - M_t \leq -0.5),
\end{aligned} \tag{5}$$

$$\begin{aligned}
R_t^- = & 5 \times |W_{t+1} - W_t| \times \mathbb{I}(W_{t+1} - W_t) \geq -0.5) + \\
& 5 \times |W_{t+1} - W_t| \times \mathbb{I}(W_{t+1} - W_t \geq 0.5).
\end{aligned} \tag{6}$$

To sum up, the trajectories / training data generated according to the ODE model, where with $N = 1000$ and $T = 6$, are as follow

$$\begin{array}{ccccccccc}
\mathbf{S}_0 & \xrightarrow{A_0} & \mathbf{S}_1 & \xrightarrow{A_1} & \mathbf{S}_2 & \xrightarrow{A_2} & \mathbf{S}_3 & \xrightarrow{A_3} & \mathbf{S}_4 & \xrightarrow{A_4} & \mathbf{S}_5 & \xrightarrow{A_5} & \mathbf{S}_6, \\
& & \widehat{\mathbf{R}}_0 & & \widehat{\mathbf{R}}_1 & & \widehat{\mathbf{R}}_2 & & \widehat{\mathbf{R}}_3 & & \widehat{\mathbf{R}}_4 & & \widehat{\mathbf{R}}_5
\end{array}$$

where $\mathbf{S}_t = (M_t, W_t, F_t)$, $t = 0, 1, \dots, 6$. Moreover, $\mathbf{R}_t = (R_t^+, R_t^-)$, $t = 0, 1, \dots, 5$. There are 7 decision points. The last decision point $T = 6$ has only states $\mathbf{S}_6 = (M_6, W_6, F_6)$, without following action nor reward. The trajectories is then broken down into 4-tuples of $(\mathbf{s}, a, \mathbf{s}', \mathbf{r})$ with the time stamps dropped.

3.2 Function approximation

3.2.1 Q function approximation

To construct linear approximators for Q functions [6], we use $K = 4$ Gaussian radial basis functions and an intercept of one. The Gaussian radial basis function has the form of $\phi(x) = \exp(-\|x - \mu\|^2 / 2\sigma^2)$, where μ and σ^2 are the parameters to be specified. Denote the Q function for positive rewards as $Q^+(\mathbf{s}, a)$, and the one for negative rewards as $Q^-(\mathbf{s}, a)$. As the positive reward function is a function of M_t , we only incorporate M_t in the basis functions for estimating $Q^+(\mathbf{s}, a)$. Hence, we can rewrite $\hat{Q}^+(\mathbf{s}, a)$ as $\hat{Q}^+(m, a)$. Specifically, $\hat{Q}^+(m, a) = \hat{w}_0^+ + \sum_{k=1}^4 \hat{w}_k^+ \exp(-\|m - \mu_k^+\|^2 / 2(\sigma_k^+)^2)$, where \hat{w}_k^+ are the weights to be estimated via least-squares policy evaluation. μ_k^+ 's are set as the 20, 40, 60, 80 percentiles of the states, and σ_k^+ 's the average distance of the percentiles to all the sample points. $\hat{Q}^-(\mathbf{s}, a)$ is constructed similarly.

Policy function approximation

For policy function approximation [6], we focus on simple decision rule to reduce the search space. Let $\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5)^\top$ be the index parameters for a policy. The policy function is defined as $\pi(\mathbf{s}; \boldsymbol{\theta}) = 0.00 \times \mathbb{I}(f(\mathbf{s}, \boldsymbol{\theta}) < 1) + 0.25 \times \mathbb{I}(1 \leq f(\mathbf{s}, \boldsymbol{\theta}) < 2) + 0.50 \times \mathbb{I}(2 \leq f(\mathbf{s}, \boldsymbol{\theta}) < 3) + 0.75 \times \mathbb{I}(3 \leq f(\mathbf{s}, \boldsymbol{\theta}) < 4) + 1 \times \mathbb{I}(f(\mathbf{s}, \boldsymbol{\theta}) > 4)$, where $f(\mathbf{s}, \boldsymbol{\theta}) = \theta_0 + \theta_1 m + \theta_2 m^2 + \theta_3 w + \theta_4 w^2 + \theta_5 mw$.

3.3 Simulation results

The goal here is to maximize $V^+(\pi)$, subject to $V^-(\pi) \leq \nu$, where ν is the bound on the secondary outcome. We applied our method to the simulated dataset. Table 3.1. shows the values of primary and secondary outcomes of the estimated constrained optimal regimes, along with their standard deviations. Table 3.2. shows the estimated indexing parameters of the estimated regimes, along with their standard deviations. Figure 3.1. shows the values of the primary objective(red) / secondary objective (blue) vs. constraint ν . Figure 3.2-3.6 shows the actions of the estimated regime for each state under different constraint values. As the ν increases, we start to observe more higher dosages being assigned to patients. Higer dosage leads to better treatment effect (more reduced tumor size), but more toxicity on patient's wellness.

Table 1: Values of estimated optimal regimes under different constraint bounds.

ν	\widehat{V}^+	std^+	\widehat{V}^-	std^-
5.49	0.39	0.36	5.42	0.16
6.85	1.35	0.29	6.62	0.43
8.21	3.58	1.06	7.66	0.30
9.57	4.07	0.89	7.77	0.35
10.93	4.97	0.56	9.53	1.42
12.29	5.98	1.10	11.43	0.69
13.65	6.13	0.99	11.90	1.46
15.01	7.08	0.93	14.88	0.42
16.37	7.98	1.31	14.69	2.15
17.73	8.89	0.61	16.84	0.57
19.09	9.15	0.30	16.93	0.71
20.45	9.85	1.16	17.86	2.33
21.81	9.76	1.01	18.18	2.88
23.18	9.76	1.01	18.18	2.88
24.54	11.62	1.38	21.88	0.63
25.89	12.02	1.17	23.45	2.47
27.25	12.04	1.17	23.57	2.69
28.61	12.86	0.54	28.27	1.06
29.98	13.69	0.57	30.25	1.04
31.34	14.53	0.97	31.11	1.10

The constraint bounds are denoted by ν . \widehat{V}^+ denotes the primary outcome values of the estimated regimes. \widehat{V}^- denotes the secondary outcome values of the estimated regimes. Standard deviations of those estimated regime values are reported as well.

Table 2: The estimated indexing parameters of estimated regimes under different constraint bounds.

ν	$\hat{\theta}_{\nu,1}$	std_1	$\hat{\theta}_{\nu,2}$	std_2	$\hat{\theta}_{\nu,3}$	std_3	$\hat{\theta}_{\nu,4}$	std_4	$\hat{\theta}_{\nu,5}$	std_5	$\hat{\theta}_{\nu,6}$	std_6
5.49	0.36	0.79	-0.38	0.76	0.13	0.48	0.03	0.84	-0.29	0.71	-0.06	0.76
6.85	0.11	1.47	-1.55	1.10	0.29	0.67	0.35	1.44	-1.06	1.19	-0.20	1.41
8.21	-0.00	1.66	-2.02	1.09	0.58	0.80	0.27	1.48	-1.26	1.40	-0.41	1.53
9.57	0.29	1.65	-2.38	1.16	0.58	0.84	0.32	1.71	-1.43	1.37	0.02	1.66
10.93	0.49	1.64	-2.61	1.13	0.67	0.90	0.18	1.72	-1.47	1.57	0.23	1.64
12.29	0.46	1.70	-2.90	1.16	0.71	0.89	0.08	1.90	-1.27	1.58	0.26	1.59
13.65	0.56	1.76	-3.20	1.12	0.70	0.92	0.01	1.70	-1.06	1.66	0.32	1.58
15.01	0.67	1.74	-3.46	1.08	0.65	0.90	-0.02	1.87	-0.83	1.77	0.14	1.55
16.37	0.89	1.81	-3.49	1.10	0.48	0.95	-0.05	1.95	-0.43	1.81	0.16	1.61
17.73	0.88	1.92	-3.54	1.17	0.39	0.99	-0.21	2.03	-0.07	1.70	0.05	1.56
19.09	1.20	2.08	-3.36	1.54	0.47	1.17	-0.41	2.14	0.11	1.56	0.15	1.52
20.45	1.69	2.18	-3.06	1.79	0.72	1.44	-0.84	2.28	0.13	1.52	0.16	1.61
21.81	1.95	2.20	-2.78	2.06	1.05	1.64	-0.88	2.40	0.12	1.48	0.27	1.88
23.18	2.21	2.21	-2.64	2.19	1.30	1.76	-1.09	2.32	0.15	1.47	0.03	1.93
24.54	2.46	2.27	-2.46	2.23	1.46	1.86	-1.32	2.31	0.36	1.43	-0.14	2.13
25.89	2.84	2.14	-2.14	2.32	1.70	2.03	-1.57	2.12	0.61	1.45	-0.19	2.31
27.25	3.25	1.82	-1.95	2.26	1.66	2.14	-1.60	2.05	1.01	1.60	-0.12	2.41
28.61	3.41	1.73	-1.86	2.26	1.58	2.08	-1.10	2.06	1.46	1.65	-0.05	2.53
29.98	3.65	1.42	-1.53	2.27	2.00	2.06	-0.98	2.04	1.99	1.68	-0.07	2.65
31.34	3.89	1.14	-1.07	2.18	2.23	1.91	-0.60	2.15	2.64	1.61	-0.13	2.70

Here, ν denotes the constraint bounds. $\hat{\theta}_{\nu,j}$ denotes the j -th component of the estimated parameter vector $\hat{\theta}_\nu$ of the estimated regimes. Standard deviations of those estimated regime values are reported as well.

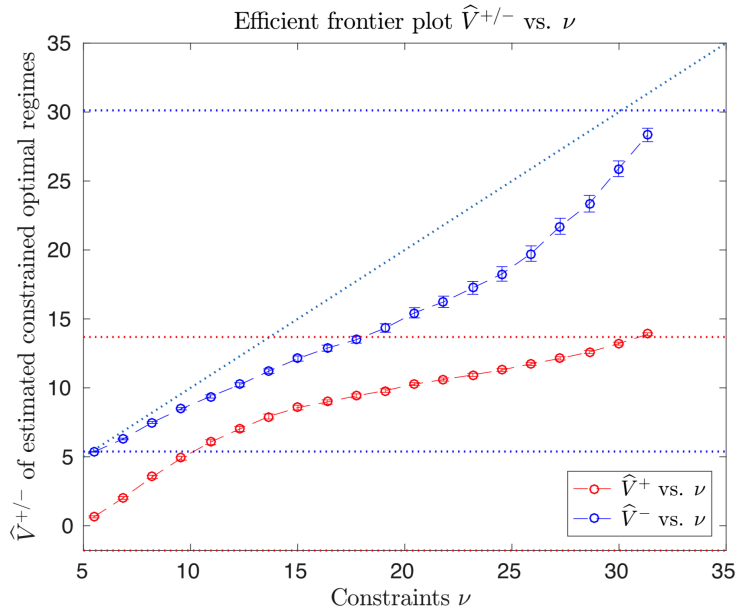


Figure 1: Efficient frontier for estimated constrained optimal regimes (infinite-stage).

The red dashed line is for the primary outcome to maximized. The blue dashed line is for the secondary outcome to be constrained. The red dotted lines are the minima and maxima for unconstrained optimization of the primary objective. The blue dotted lines are the minimal and maximal for unconstrained optimization of the secondary objective.

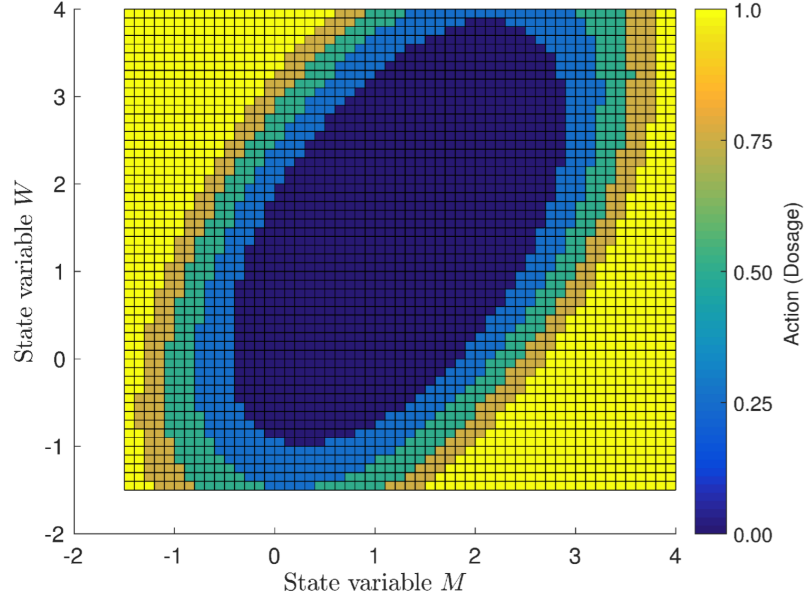


Figure 2: Action for each state under constraint $\nu = 10.93$

Yellow represents high dosage treatment assignment. Blue represents low dosage assignment. As the constraint bound gets loose, more higher dosage treatments are assigned to patients.

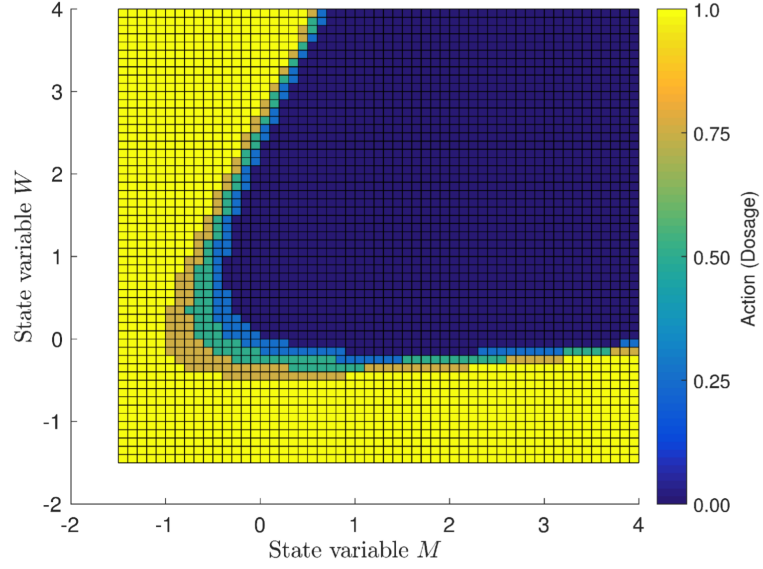


Figure 3: Action for each state under constraint $\nu = 17.73$

Yellow represents high dosage treatment assignment. Blue represents low dosage assignment. As the constraint bound gets loose, more higher dosage treatments are assigned to patients.

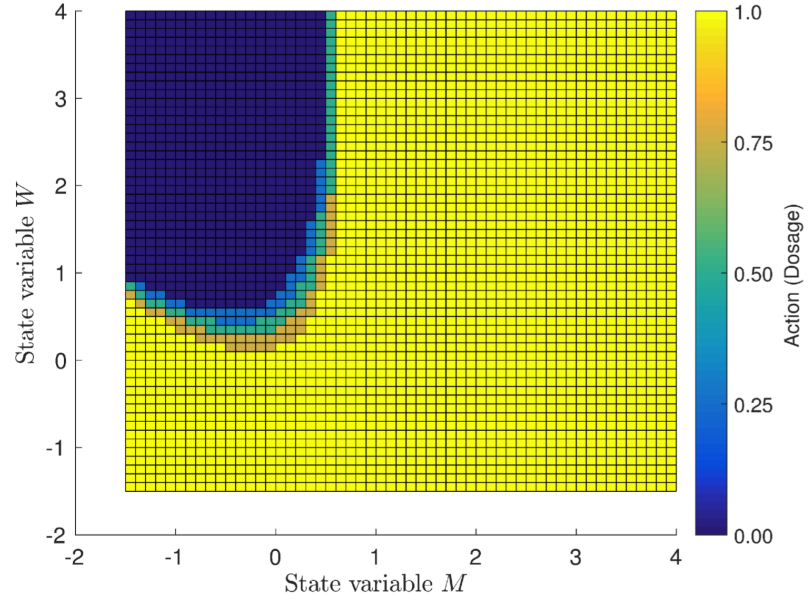


Figure 4: Action for each state under constraint $\nu = 24.54$

Yellow represents high dosage treatment assignment. Blue represents low dosage assignment. As the constraint bound gets loose, more higher dosage treatments are assigned to patients.

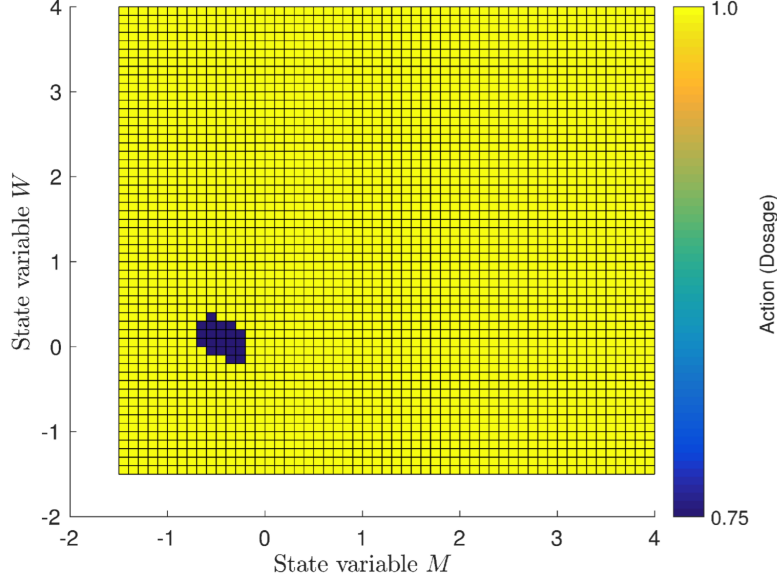


Figure 5: Action for each state under constraint $\nu = 31.34$

Yellow represents high dosage treatment assignment. Blue represents low dosage assignment. As the constraint bound gets loose, more higher dosage treatments are assigned to patients.

4 Conclusion and Future

We propose a framework for constrained optimal dynamic treatment regimes to handle the trade-off between the primary objective and all other secondary objectives in infinite stage setting. The simulation results based on the chemotherapy ODE system are presented and visualized. This framework offers an intuitive way for clinicians to exam the trade-off and make treatment decisions based on patient's preference. Different from CPO, our method takes into consideration that clinical data is expensive and scarce. Borrowing strength from least-squares policy evaluation, our method is able to learn from data efficiently. Moreover, least-squares policy evaluation is an iterative method for policy evaluation, and has the advantage of being able to learn both offline and online. Hence, our method can fits in not only the situation where clinical policies needed to be learned after data collect (offline, off-policy batched), but also the situation where online real-time policy learning is needed. Interior-point method is also a well-studies optimization method for constrained estimation. Its theoretical guarantees assure us of good enough optimal solutions. However, it is

obvious that the choice of policy function approximation may have impact on the decision. So does the choice of Q function approximation. Clinical domain expertise may be required. Alternatively, automated feature learning techniques for function approximation from the machine learning community can be incorporated. More complex dataset may be collected, such as text, image, speech and so on, considering the recent technology advancement in mobile devices. How to incorporate those complex information to better describe an individual's state of health is challenging. Rigorous theoretical work for our method is also under investigation.

Besides constraints on the expected value of a policy, we can also consider risk constraints, where the probabilities of adverse events occurring are restricted. Although reinforcement learning is a powerful technique to find optimal treatment regimes for clinical practice, designing appropriate reward functions is crucial for serving the desired clinical purpose, but very difficult. Current approach may not scale well in complex clinical situations or preventive healthcare where multiple subgoals may be involved. How to automatically generate rewards and objectives in complex clinical situations can be an interesting direction for investigation.

Nowadays, many aspects of the clinical practice have been transformed by mobile devices, such as smart phones, tablets, wearable sensors etc. It allows clinicians remotely monitor and intervene patients' chronic conditions in real-time. It also allows for adaptive preventive interventions for motivating and maintaining healthy behaviors, such as physical exercise, diets, and so on. To better understand adaptive interventions, interdisciplinary collaborations become a necessity among clinicians, medical researchers, behavioral scientists, statisticians, and computer scientists.

References

- [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained Policy Optimization. *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [2] E. Altman. *Constrained Markov Decision Processes*. Stochastic Modeling Series. Taylor & Francis, 1999.
- [3] D Blatt, SA Murphy, and J Zhu. A-learning for approximate planning. *Ann Arbor*, 1001:48109–2122, 2004.
- [4] Richard H. Byrd, Mary E. Hribar, and Jorge Nocedal. An Interior Point Algorithm for Large-Scale Nonlinear Programming. *SIAM Journal on Optimization*, 9(4):877–900, 1999.

- [5] Ashkan Ertefaie. Constructing Dynamic Treatment Regimes in Infinite-Horizon Settings. pages 1–39, 2014.
- [6] Alborz Geramifard. A Tutorial on Linear Function Approximators for Dynamic Programming and Reinforcement Learning. *Foundations and Trends® in Machine Learning*, 6(4):375–451, 2013.
- [7] Miguel A Hernan and James M Robins. Estimating causal effects from epidemiological data. *Journal of epidemiology and community health*, (7):578–86, July.
- [8] Jerzy Splawa-Neyman, D. M. Dabrowska and T. P. Speed. On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9. *Statistical Science*, 5(4):465–472, 1990.
- [9] Eric B. Laber, Daniel J. Lizotte, and Bradley Ferguson. Set-valued dynamic treatment regimes for competing outcomes. *Biometrics*, 70(1):53–61, 2014.
- [10] Eric B Laber, Daniel J Lizotte, Min Qian, William E Pelham, and Susan A Murphy. Dynamic treatment regimes : technical challenges and applications. 2014.
- [11] Michail G. Lagoudakis and Ronald Parr. Model-Free Least-Squares Policy Iteration. In *Advances in Neural Information Processing Systems 14 (NIPS 2001)*, pages 1547–1554, 2001.
- [12] Michail G. Lagoudakis and Ronald Parr. Least-squares policy iteration. *The Journal of Machine Learning Research*, 4:1107–1149, 2003.
- [13] Kristin A Linn, Eric B Laber, and Leonard A Stefanski. Interactive Q-learning for Probabilities and Quantiles. 2014.
- [14] Daniel J Lizotte, Michael H. Bowling, and Susan A. Murphy. Efficient reinforcement learning with multiple reward functions for randomized controlled trial analysis. in *Proc. of Int. Conf. on Machine Learning*, pages 695–702, 2010.
- [15] Daniel J. Lizotte, Michael H. Bowling, and Susan A. Murphy. Linear Fitted-Q Iteration with Multiple Reward Functions. *Journal of Machine Learning Research*, 13:3253–3295, 2012.
- [16] Daniel J. Lizotte and Eric B. Laber. Multi-objective markov decision processes for data-driven decision support. *J. Mach. Learn. Res.*, 17(1):7378–7405, January 2016.
- [17] Daniel J. Lockett, Eric B. Laber, Anna R. Kahkoska, David M. Maahs, Elizabeth Mayer-Davis, and Michael R. Kosorok. Estimating Dynamic Treatment Regimes in Mobile Health Using V-learning. 2016.
- [18] Erica Moodie. Dynamic treatment regimes. *Clinical trials (London, England)*, 1(5):471, 2004.

- [19] S. A. Murphy. Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(2):331–355, May 2003.
- [20] S A Murphy, Y Deng, E B Laber, H R Maei, R S Sutton, and K. Witkiewitz. A Batch, Off-Policy, Actor-Critic Algorithm for Optimizing the Average Reward. *arXiv*, pages 1–18, 2016.
- [21] Susan A Murphy. A Generalization Error for Q-Learning. *Journal of machine learning research : JMLR*, 6:1073–1097, July 2005.
- [22] James M. Robins. Causal Inference from Complex Longitudinal Data, 1997.
- [23] D. B. Rubin. Discussion of Randomized analysis of experimental data: The Fisher randomization test by D. Basu. *Journal of the American Statistical Association*, (75):591–593, 1980.
- [24] Donald B. Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.
- [25] R. A. Waltz, J. L. Morales, J. Nocedal, and D. Orban. An interior algorithm for nonlinear optimization that combines line search and trust region steps. *Mathematical Programming*, 107(3):391–408, 2006.
- [26] Baqun Zhang, Anastasios A. Tsiatis, Marie Davidian, Min Zhang, and Eric Laber. Estimating optimal treatment regimes from a classification perspective. *Stat*, 1(1):103–114, 2012.
- [27] Yufan Zhao, Michael R. Kosorok, and Donglin Zeng. Reinforcement learning design for cancer clinical trials. *Statistics in Medicine*, 28(26):3294–3315, 2010.

Appendices

A Least-squares policy evaluation (LSQ) algorithm

Algorithm 3: LSQ

```

 $A \leftarrow 0$ 
 $b \leftarrow 0$ 
for each  $(s, a, s', r)$ 
     $A \leftarrow A + \phi(s, a) (\phi(s, a) - \gamma \phi(s', \pi(s')))^{\top}$ 
     $b \leftarrow b + \phi(s, a)r$ 
 $w \leftarrow A^{-1}b$ 
return  $w$ 

```
