

# A Framework For Considering Costs and Benefits of Additional Complexities in Personalized Medicine

Working Title Need a New One

Demetri Pananos, Dan Lizotte

## 1 Introduction

Personalized medicine has four goals: 1) to identify drugs for which between-subject variability in effectiveness or toxicity is a key issue for effective treatment, 2) to identify predictors which may explain this variability, 3) to decide on the right dose of the right drug by considering these factors, and 4) prevent adverse reactions to drugs [1]. Progress in all four goals has accelerated within the last decade. As an example, recent studies on DPYD genotype testing prior to starting fluoropyrimidine-based chemotherapy showed promise in preventing adverse events, making good arguments for integration of DPYD genotype testing into standard of care practices [2]. Despite this progress, personalized medicine still faces several barriers to widespread adoption, including economic burden, patient burden, and expertise burden required for new methods of personalization.

Personalized medicine reduces costs to the healthcare system by identifying patients who are at greater risk for adverse events or dose adjustments, thereby optimizing safety. If the patient does not undergo an adverse event, then this ultimately saves the healthcare system the cost of the hospital stay [3]. More ambitiously, personalized medicine has the potential to save the healthcare system costs by more effectively using resources [4]. The cost of instruments, technicians, and leadership required to operate a personalized medicine clinic are high burden, and it is not yet clear if personalized medicine is sufficiently cost effective to offset operating costs in all circumstances [5]. In their 2019 scoping review of personalized medicine cost effectiveness, Kasztura et. al [5] found that willingness-to-pay thresholds vary wildly from country to country (citing that cost per quality adjusted life year for some modes of personalized medicine range from \$20, 000 USD per quality adjusted life year in for studies in Europe and the United Kingdom to \$200,000 USD per quality adjusted life year for studies in the United States). This high variability in cost effectiveness means the burden required for start up may result in a positive return on investment in some areas but not others. This variability should prompt would be adopters to more closely examine if taking on the initial burden is worth the result.

The dominant perspective on personalized medicine focuses on the use of clinical and physiological information (including biomarkers, genotyping, and diagnostic tests) as a means of optimizing treatments, but largely ignore needs, constraints, and utilities of the patient [6, 7]. Patients can be burdened by frequent followup for clinical measurement (as in the case with Warfarin), be burdened by costly expenses related to obtaining care, or may be more risk adverse/tolerant than the “typical” patient. As an example, transportation has been found to be a large financial burden for patients receiving cancer treatment [8], and continues to burden patients, with a 2020 study finding that the cost of parking alone can climb as high as \$1600 over the course of treatment in

the United States [9]. Additional visits to a clinic have the potential to further burden patients by requiring them to miss a day of work, and find means of childcare during their absence (if necessary). Incorporating patient preferences and reducing the burden of personalization on the patient can result in sustained adherence [10], thereby increasing effectiveness and further preventing adverse events.

An additional expertise burden is added as machine learning (used interchangeably with the term “artificial intelligence”) is adopted into personalized medicine initiatives. Cutting edge machine learning models for prediction or decision making can be prohibitively burdensome to implement effectively. Failure to do so may result in pernicious bias inadvertently affecting subpopulations, as was found to be the case in algorithms for credit scoring [11], crime prediction [12], and hiring [13]. A 2019 study found an instance of this bias in a widely used risk scoring algorithm in healthcare [14], demonstrating that despite the intention of those involved, the use of a model does not guarantee optimized healthcare. Implementation of new approaches requires the partnership of experts in data science, computer science, statistics, and engineering. Collaboration between experts in these disciplines and physicians (among other stakeholders) is crucial to make effective use of data, rigorously internally validate models, report them appropriately and transparently, and temper expectations which may be skewed from hype surrounding these algorithms [15].

These costs may be payable for some, but the question then turns to if the result is worth the expense. Answering that question is difficult without an idea how the additional cost of collecting data, or implementing new algorithms, will benefit the clinic or the patient subject to inherent constraints.

In this study, we present a new framework for helping practitioners interested in implementing personalized medicine to answer a) If the burden of personalization could result in more favorable outcomes for their population of patients, b) how working around patient burden (such as coming into the clinic for fewer measurements) may change effectiveness, and c) if more complex models for personalization can lead to better effectiveness. For our case study, we fit a Bayesian model to existing data on the pharmacokinetics of apixaban. The resulting Bayesian model is used to generate synthetic pharmacokinetic data for use in experiments to compare different forms of personalization. Treating personalization as a dynamic treatment regime, we propose six policies, each increasing in complexity and clinic/patient burden, for personalizing doses of apixaban with the goal of keeping blood serum concentrations within a desired range for as long as possible. Under the assumption that the fitted Bayesian model can produce similar data to what might be observed in the future from new patients, we can make inferences as to how different policies for personalizing doses may improve upon one another, and compare if the additional burden of implementing a more complex or costly form of personalization can generate a more desirable outcome for the patient or healthcare provider.

We begin with an overview of dynamic treatment regimes. We then describe how to estimate an optimal dynamic treatment regime by combining Bayesian pharmacokinetic modelling with Q-learning. We then present our case study, beginning with the details of the Bayesian model we use to fit the real pharmacokinetic data, and present model fit diagnostics to argue that our model is satisfactory for generating synthetic data for use in our simulations. We then present and discuss the results of our simulation in light of the burden presented to a clinic to implement personalized medicine, and how this framework can be integrated to answer questions such as if personalization can produce a positive return on investment, how burdening the patient with an additional clinic visit will improve the clinic’s understanding of how to treat that patient, and if using a more advanced mode of personalization will result in more favorable health outcomes.

## 2 Dynamic Treatment Regimes

In this section, we discuss the theory of dynamic treatment regimes and how personalization can be thought of as a dynamic treatment regime. We describe our experiments in the context of dynamic treatment regimes and introduce our reward function.

### 2.1 Trajectories

Our goal is to find the dose or sequences of doses for a subject to keep their blood serum concentration within a desired range for as long as possible given the constraints: a) subject's blood serum concentrations cannot be measured very frequently, and b) we are limited to pre-dose clinical measurements to make our initial dosing decision. The theory of dynamic treatment regimes and statistical reinforcement learning offers a framework through which to understand our problem and construct one possible solution.

A dynamic treatment regime (DTR) is a sequence of decision rules for adapting a treatment plan to the time-varying state of an individual subject [16]. In DTRs, and their cousin topic in computer science *reinforcement learning*, an agent (often thought of as a robot in reinforcement learning, but within medicine sometimes thought of as a physician's computerized decision support system) interacts with a system for a number of stages. At each stage, the agent receives an *observation* of the system and then decides which *action* to take. This action will result in an observed *reward* which is followed by a new observation of the system after it has been impacted by the action. This cycle of observation, action, reward then repeats, with the agent aiming to take actions which yield the largest total reward.

Key to our DTR is the concept of a *trajectory*. Define a stage to be a triple containing an observation, chosen action, and resulting reward. Let  $O_i$  denote an observation at the  $i$ th stage,  $A_i$  be the action at the  $i^{th}$  stage, and  $Y_i$  denote the reward at the  $i^{th}$  stage, denoted in capital letters when considering the observation, action, and reward as random variables. A trajectory is then the tuple  $(O_1, A_1, O_2, A_2, \dots, O_K, A_K, O_{K+1})$ . Following notation by Chakraborty and Moodie [16], we will denote a system's history at stage  $j$  as  $H_j = (O_1, A_1, O_2, A_2, \dots, O_{j-1}, A_{j-1}, O_j)$ . The reward at stage  $j$  is then a function of the system's history, the action taken, and the next observation  $Y_j = Y_j(H_j, A_j, O_{j+1})$ .

### 2.2 Policies, Value Functions, and Q-Learning

A deterministic policy  $d = (d_1, \dots, d_k)$  is a vector of decision rules each which take as input the system's history and output an action to take. The stage  $j$  value function for a policy  $d$  is the expected reward the agent would receive starting from history  $h_j$  (here in lower case since it is an observed quantity) and then choose actions according to  $d$  for every action thereafter. The value function is written as

$$V_j^d(h_j) = E_d \left[ \sum_{k=j}^K Y_k(H_k, A_k, O_{k+1}) \middle| H_j = h_j \right]. \quad (1)$$

Here, the expectation is computed over the distribution of trajectories. Importantly, the stage  $j$  value function can be decomposed into the expectation of reward at stage  $j$  plus the stage  $j + 1$  value function [16]

$$V_j^d(h_j) = E_d [Y_j(H_j, A_j, O_{j+1}) + V_{j+1}^d(H_{j+1}) | H_j = h_j]. \quad (2)$$

The optimal stage  $j$  value function is the value function under a policy which yields maximal value. Mathematically,

$$V_j^{opt}(h_j) = \max_{d \in \mathcal{D}} V_j^d(h_j) \quad (3)$$

A natural question is “what policy maximizes the value?”. Estimating such a policy can be achieved by estimating the optimal Q function [16]. The optimal Q function at stage  $j$  is a function of the system’s history  $h_j$  and a proposed action  $a_j$ ,

$$Q_j^{opt} = E [Y_j(H_j, A_j, O_{j+1}) + V_{j+1}^{opt}(H_{j+1}) | H_j = h_j, A_j = a_j] \quad (4)$$

Note that the optimal Q function has similar form and interpretation to the optimal value function (namely, it is the expected reward starting at stage  $j$  but with the added condition that we take action  $a_j$  and then follow the optimal policy thereafter). Due to the decomposition of the reward function at stage  $j$ , estimation of the optimal Q function can be performed by choosing the action which yields the largest reward at each stage assuming we act optimally in the future. Below, we describe how we choose optimal actions using a posterior distribution of a subject’s pharmacokinetics.

### 2.3 Experimental Design In Terms of Stages of a DTR

In our experiments, we develop a DTR for selecting the best dose for keeping a patient’s blood plasma concentration within a desired range. Here, we present details of the experimental design in the DTR framework, leaving simulation details (including how the data were simulated) for our methods section.

Our experiment consists of 1000 simulated subjects taking a dose of apixaban once every 12 hours with perfect adherence for a total of 10 days. Sometime in the second 12 hour period on the fourth day (between 108 and 120 hours after the initial dose), we have the opportunity to measure the simulated subject’s blood concentration, should our policy allow for it. At the start of the fifth day, the dose is adjusted based on all the pre-dose clinical measurements plus the observed concentration. The dose will be adjusted so as to attempt to maximize the time spent between 0.1 mg/L and 0.3 mg/L. Thus, our DTR consists of two stages (the first five days, and the latter five days), however the size of the range may be adapted for different scenarios. We choose this range as it is not so narrow that even optimal doses perform poorly, but not so wide that any dose can achieve high reward.

In terms of the DTR, the system is the patient for whom a dose is selected, the actions correspond to selection of dose sizes, and the reward is the proportion of time spent within the desired concentration range. The trajectories we will use to estimate the optimal Q functions are of the form

$$O_1, A_1, Y_1, O_2, A_2, Y_2, O_3 \quad (5)$$

The interpretation of a given trajectory is:

- $O_1$  is any pre-dose clinical measurements of the subject. In our experiments, we consider age in years, renal function (as measured by serum creatinine in mMol/L), weight in kilograms, and dichotomous biological sex (dummy coded so that male=1 and female=0). We choose these variables as they are known to affect the pharmacokinetics of apixaban [17].

- $A_1$  is dual action of initial dose to provide the subject plus a time in the future at which to measure the subject's blood serum concentration.
- $Y_1$  is the proportion of time spent within the concentration range in the first five days.
- $O_2$  is the pre clinical measurements of the subject plus the observed concentration made on the fourth day.
- $A_2$  is the dose adjustment
- $Y_2$  is the proportion of time spent within the concentration range in the last five days after the dose adjustment.
- $O_3$  would be pre-dose clinical measurements, the observed concentration made on the fourth day, and the next concentration measurement, were it to be made. As we examine just the two actions  $A_1$  and  $A_2$ , we do not make use of  $O_3$  but include it here to adhere with our definition of trajectories above.

The reward function we use depends on the subject's true latent concentration. Let  $c_j, j = 1 \dots K$  be the  $j^{th}$  latent concentration value at time  $t_j$ . The reward function is

$$Y(c_1, c_2, \dots, c_k) = \frac{1}{k} \sum_{j=1}^K \mathbb{I}(0.1 < c_j < 0.3) \quad (6)$$

Here,  $\mathbb{I}$  is an indicator function returning 1 if the argument is true, and 0 else. To leverage off-the-shelf optimization tools, we approximate this reward function with a continuously differentiable function, namely

$$Y(c_1, c_2, \dots, c_k) = \frac{1}{k} \sum_{j=1}^K \exp \left( - \left[ \frac{c_j - 0.15}{0.05} \right]^{2\beta} \right) \quad (7)$$

Here,  $\beta$  is a positive integer. For sufficiently large beta, our approximation becomes arbitrarily close to our intended reward function. In practice we set  $\beta=5$  to balance between good approximation of our intended reward and vanishing gradients impeding our optimization. We suppress the dependence on the history in the definition of the reward as the reliance on the history is implicit. The reward depends on the latent concentrations which depend on previous doses (actions) and potentially on the previous dose measurements (observations of the system).

Our stage 2 optimal Q function is then

$$Q_2^{opt}(H_2, A_2) = E \left[ Y(c_{j+1}, c_{j+2}, \dots, c_{j+n}) \middle| H_2, A_2 \right], \quad (8)$$

and our stage 1 optimal Q function is

$$Q_1^{opt}(H_1, A_1) = E \left[ Y(c_1, c_2, \dots, c_j) + \max_{a_2 \in \mathcal{A}} Q_2^{opt}(H_2, a_2) \middle| H_1, A_1 \right] \quad (9)$$

We seek to maximize the stage 1 optimal Q function to learn the optimal policy for dosing subjects under the constraint we can measure them at most once and are limited to the aforementioned pre-dose clinical variables. The interpretation of stage 1 optimal Q function is as follows: *Given the pre-dose clinical variables of the subject and a proposed initial dose and measurement time, the stage 1 optimal Q function gives the proportion of time the subject's blood serum concentration is between 0.1mg/L and 0.3mg/L assuming that we provide the subject*

with the best dose possible at the start of the 5<sup>th</sup> day. The actions  $A_1$  and  $A_2$  which maximize these functions constitute the optimal policy.

The concentration values  $c_j$  in the optimal Q functions are latent, meaning we have no direct access to them in practice. Furthermore, obtaining measurements with high enough frequency so that the reward is faithfully estimated would be too burdensome on the patient. What is left to explain is how these concentrations are computed. In the next section, we describe how we use a Bayesian model to obtain latent concentration predictions and compute the required expectations.

### 3 Bayesian Models of Pharmacokinetics

In order to estimate the optimal Q functions, we need to be able to predict how a subject’s concentrations respond to any given dose. These predictions can be passed to the reward function to compute the values of  $Q_1$  and  $Q_2$ . Bayesian models of pharmacokinetics are capable of providing distributions of concentrations and can account for pre-dose clinical variables, as we will describe. The model we construct takes as inputs pre-dose clinical variables (age, sex, weight, and creatinine) and provides a distribution of plausible concentrations prior to seeing data. We can pass these prior predictions to our Q function and determine what dose keeps the subject in range for the longest time prior to seeing any data. Once we observe a concentration measurement from the subject, we can condition our model on that observation, and perform the optimization again to adjust the subject’s dose in light of new information. For our experiments, we first fit a Bayesian pharmacokinetic model to real pharmacokinetic data collected by our colleagues in clinical pharmacology [18] and then use the fitted model to simulate patients and make predictions for our experiments.

We extend a previously proposed one compartment Bayesian pharmacokinetic model [19] to include fixed effects of covariates on pharmacokinetic parameters. The model presented in [19] is a hierarchical Bayesian model of apixaban pharmacokinetics, in which the clearance rate (L/hour), time to max concentration (hours), absorption time delay (hours), and ratio between the elimination and absorption rate constants (called alpha, a unitless parameter) are hierarchically modelled. We extend that model by regressing the latent pharmacokinetic parameters on the aforementioned pre-dose clinical variables.

The model fit to real data, which we refer to as  $\mathcal{M}_1$ , is

$$y_{i,j} \sim \text{Lognormal}(C_i(t_j), \sigma_y^2) \quad (10)$$

$$\sigma^2 \sim \text{Lognormal}(0.1, 0.2) \quad (11)$$

$$C_i(t_j) = \begin{cases} \frac{D_i \cdot F}{Cl_i} \cdot \frac{k_{e,i} \cdot k_{a,i}}{k_{e,i} - k_{a,i}} \left( e^{-k_{a,i}(t_j - \delta_i)} - e^{-k_{e,i}(t_j - \delta_i)} \right) & t_j > \delta_i \\ 0 & \text{else} \end{cases} \quad (12)$$

$$k_{e,i} = \alpha_i \cdot k_{a,i} \quad (13)$$

$$k_{a,i} = \frac{\log(\alpha_i)}{t_{max,i} \cdot (\alpha_i - 1)} \quad (14)$$

$$\delta_i \sim \text{Beta}(\phi, \kappa) \quad (15)$$

$$\text{logit}(\alpha_i) | \beta_\alpha, \sigma_\alpha^2 \sim \text{Normal}(\mu_\alpha + \mathbf{x}_i^T \beta_\alpha, \sigma_\alpha^2) \quad (16)$$

$$\log(t_{max,i}) | \beta_{t_{max}}, \sigma_{t_{max}}^2 \sim \text{Normal}(\mu_{t_{max}} + \mathbf{x}_i^T \beta_{t_{max}}, \sigma_{t_{max}}^2) \quad (17)$$

$$\log(Cl_i) | \beta_{Cl}, \sigma_{Cl}^2 \sim \text{Normal}(\mu_{Cl} + \mathbf{x}_i^T \beta_{Cl}, \sigma_{Cl}^2) \quad (18)$$

$$p(\phi) \sim \text{Beta}(20, 20) \quad (19)$$

$$p(\kappa) \sim \text{Beta}(20, 20) \quad (20)$$

$$p(\mu_{Cl}) \sim \text{Normal}(\log(3.3), 0.15^2) \quad (21)$$

$$p(\mu_{t_{max}}) \sim \text{Normal}(\log(3.3), 0.1^2) \quad (22)$$

$$p(\mu_\alpha) \sim \text{Normal}(-0.25, 0.5^2) \quad (23)$$

$$p(\sigma_y) \sim \text{Lognormal}(\log(0.1), 0.2^2) \quad (24)$$

$$p(\sigma_{CL}) \sim \text{Gamma}(15, 100) \quad (25)$$

$$p(\sigma_{t_{max}}) \sim \text{Gamma}(5, 100) \quad (26)$$

$$p(\sigma_\alpha) \sim \text{Gamma}(10, 100) \quad (27)$$

$$p(\beta_{Cl,k}) \sim \text{Normal}(0, 0.25^2) \quad k = 1 \dots 4 \quad (28)$$

$$p(\beta_{t_{max},k}) \sim \text{Normal}(0, 0.25^2) \quad k = 1 \dots 4 \quad (29)$$

$$p(\beta_{\alpha,k}) \sim \text{Normal}(0, 0.25^2) \quad k = 1 \dots 4 \quad (30)$$

Here, normal distributions are parameterized by their mean and variance, lognormal distributions are parameterized by the mean and variance of the random variable on the log scale, and gamma distributions are parameterized by their shape and rate. The  $\mu$  in the model above represent population means on either the log or logit scale, the  $\beta$  are regression coefficients for the indicated pharmacokinetic parameter, the sigmas are the population level standard deviations on the log or logit scale,  $\delta$  is a parameter which relaxes the assumption that the dose is absorbed into the blood immediately upon ingestion,  $F$  is the bioavailability of apixiban (which we fix to 0.5 [17]) and  $D$  is the size of the dose in milligrams. All continuous variables were standardized using the sample mean and standard deviation prior to being passed to the model.

Once fit,  $\mathcal{M}_1$  can be used to predict the pharmacokinetics of new patients, using the patient's covariates as predictors. To do so, the marginal posterior distributions for  $\mu_{Cl}$ ,  $\mu_{t_{max}}$ ,  $\mu_\alpha$ ,  $\beta_{Cl}$ ,  $\beta_{t_{max}}$ ,  $\beta_\alpha$ ,  $\sigma_{Cl}$ ,  $\sigma_{t_{max}}$ ,  $\sigma_\alpha$ , and  $\sigma_y$  must be summarized. We use maximum likelihood on the posterior samples to summarize the marginal

posterior distributions. We model the population means and regression coefficients as normal, and the standard deviations as gamma. The maximum likelihood estimates are used to construct priors for a new model, which we call  $\mathcal{M}_2$ . We construct  $\mathcal{M}_2$  so as to be able to predict plasma concentration after multiple doses (of potentially different sizes) administered over time, and remove the time delay ( $\delta$ ) to simplify our simulations. Model priors for  $\mathcal{M}_2$  are then

$$p(\mu_{Cl}) \sim \text{Normal}(0.5, 0.04) \quad (31)$$

$$p(\mu_{t_{max}}) \sim \text{Normal}(0.93, 0.05) \quad (32)$$

$$p(\mu_{\alpha}) \sim \text{Normal}(-1.35, 0.13) \quad (33)$$

$$p(\sigma_{Cl}) \sim \text{Gamma}(69.15, 338.31) \quad (34)$$

$$p(\sigma_{t_{max}}) \sim \text{Gamma}(74.96, 349.56) \quad (35)$$

$$p(\sigma_{\alpha}) \sim \text{Gamma}(10.1, 102.07) \quad (36)$$

$$p(\beta_{Cl,1}) \sim \text{Normal}(0.39, 0.08^2) \quad (37)$$

$$p(\beta_{Cl,2}) \sim \text{Normal}(0.19, 0.04^2) \quad (38)$$

$$p(\beta_{Cl,3}) \sim \text{Normal}(0.02, 0.04^2) \quad (39)$$

$$p(\beta_{Cl,4}) \sim \text{Normal}(0.01, 0.04^2) \quad (40)$$

$$p(\beta_{t_{max},1}) \sim \text{Normal}(-0.01, 0.08^2) \quad (41)$$

$$p(\beta_{t_{max},2}) \sim \text{Normal}(0.09, 0.05^2) \quad (42)$$

$$p(\beta_{t_{max},3}) \sim \text{Normal}(-0.05, 0.04^2) \quad (43)$$

$$p(\beta_{t_{max},4}) \sim \text{Normal}(-0.01, 0.04^2) \quad (44)$$

$$p(\beta_{\alpha,1}) \sim \text{Normal}(-0.19, 0.17^2) \quad (45)$$

$$p(\beta_{\alpha,2}) \sim \text{Normal}(0.33, 0.11^2) \quad (46)$$

$$p(\beta_{\alpha,3}) \sim \text{Normal}(-0.06, 0.1^2) \quad (47)$$

$$p(\beta_{\alpha,4}) \sim \text{Normal}(-0.09, 0.1^2) \quad (48)$$

$$(49)$$

For our experiments, we generate the pharmacokinetic parameters of 1000 subjects from the prior predictive model of  $\mathcal{M}_2$ . Bayesian models are generative models, meaning they can generate pseudodata by drawing random variables according to the model specification going from top (model priors) to bottom (model likelihood). To do so, we begin by resampling 1000 tuples of age, sex, weight, and creatinine from the dataset used to fit  $\mathcal{M}_{\infty}$ . We sample one draw of  $\mu_{Cl}$ ,  $\mu_{t_{max}}$ ,  $\mu_{\alpha}$ ,  $\beta_{Cl}$ ,  $\beta_{t_{max}}$ , and  $\beta_{\alpha}$  from their respective prior distributions in  $\mathcal{M}_2$ . The values of these parameters remained fixed for all 1000 subjects. Conditioned on the values of these mus and betas, we compute the expectation of the population distribution for each pharmacokinetic parameter



by computing  $\mu_{Cl} + \mathbf{x}^T \beta_{Cl}$ ,  $\mu_{t_{max}} + \mathbf{x}^T \beta_{t_{max}}$ ,  $\mu_{\alpha} + \mathbf{x}^T \beta_{\alpha}$ , where  $\mathbf{x}^T$  is the resampled tuple. From the prior distribution of M2, we sample one draw of  $\sigma_{Cl}$ ,  $\sigma_{t_{max}}$ ,  $\sigma_{\alpha}$ , and  $\sigma_y$ . These remained fixed for all 1000 subjects. Using the previously computed expectations and  $\sigma$ , we sample 1000 tuples of pharmacokinetic parameters, one for each of the simulated subjects. The clearance rate and time to max concentration were sampled assuming a lognormal distribution. Alpha was sampled using a logitnormal distribution. The pharmacokinetics can then be determined conditional on the pharmacokinetic parameters. Each of simulated subjects' pharmacokinetic parameters remained fixed through the experiments. We simulate the latent concentration using  $C(t)$  as written in  $\mathcal{M}_2$ , and can simulate observed concentrations by drawing a sample from a lognormal distribution with mean  $\ln(C(t))$  and standard deviation  $\sigma_y$ .

We use Stan, an open source probabilistic programming language, for fitting our Bayesian models via Hamiltonian Monte Carlo (a Markov Chain Monte Carlo technique) and computing markov chain diagnostics. Twelve chains are initialized and run for 2000 iterations each (1000 for warmup allowing the Markov chain the opportunity to find the correct target distribution and 1000 to use as samples from the posterior).

## 4 Modes of Personalization

Thus far, we have motivated personalization of dose sizes through Q learning. Q learning is the most complex solution that could be implemented at this time, and its implementation in practice would be burdensome due to this complexity. This makes implementation of Q learning a tall order, especially considering alternative DTRs exist which are not as costly to implement. The cost of implementing Q learning may be worth paying if the benefit of Q learning over these other DTRs is substantial, but we need a framework in which to estimate the size of this benefit.

Our study considers the following modes of personalization. Each one has different requirements in terms of computation, data needs, clinical overhead, and patient burden. Our study aims to understand, in this particular setting (apixaban with potential PK monitoring) what the relative benefits of these different modes might be in practice. It also presents a general framework for evaluating these different modes of personalization in other settings. The six modes of personalization we consider are:

- 1) Dose selection using a hierarchical Bayesian model which does not incorporate subject covariates. This model was presented in Pananos & Lizotte [19]. We refer to this mode as the “No Covariate Model”.
- 2) 1) and conditioning the model on a single sample from the subject taken sometime in the final 12 hours before the half way point. At the start of the fifth day, a new dose is selected and used for the remaining time. We refer to this mode as “No Covariate + 1 Sample”.
- 3) Dose selection from M2. A single dose is selected at the start of the regiment and is used throughout the 10 simulated days. We refer to this mode as “Covariate Model”.
- 4) 3) and conditioning the model on a single sample from the subject taken sometime in the final 12 hours before the half way point. At the start of the fifth day, a new dose is selected and used for the remaining time. We refer to this mode as “Covariate model + 1 Sample”.
- 5) A two stage DTR, however the initial dose is the result of the procedure in 3). The best time to sample the patient is then determined via Q learning. We refer to this mode as “Optimal Sampling Time”.

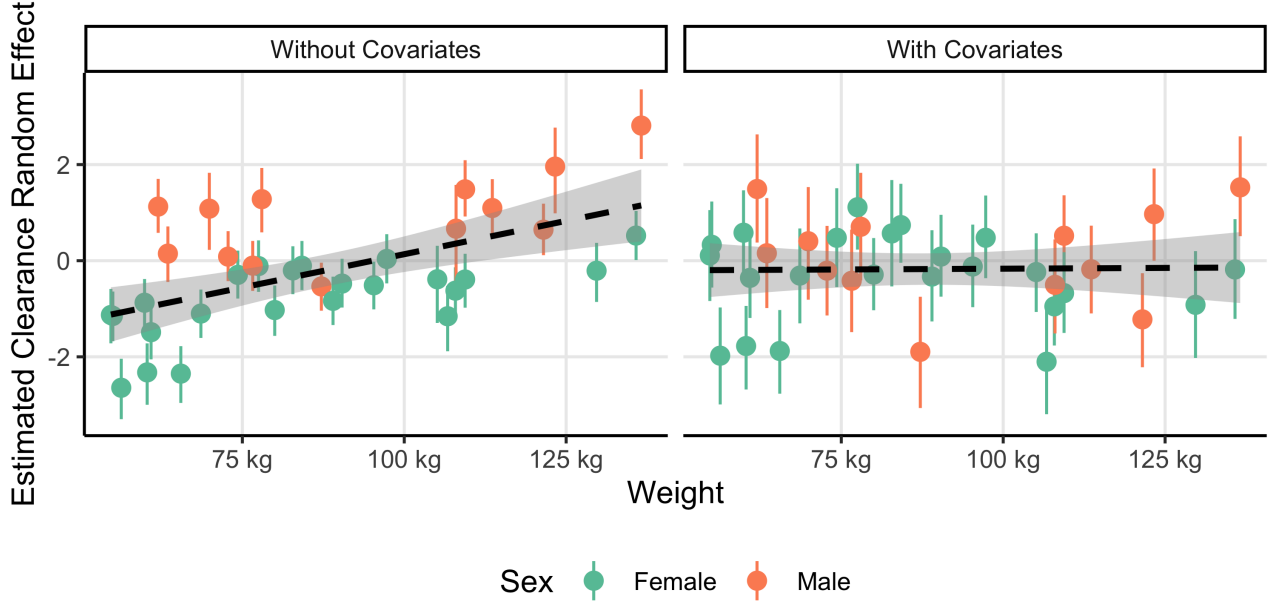


Figure 1: Random effects estimates for clearance rate  $Cl_i$  and 95% credible intervals (left). Random effects estimates are colored by patient sex. Prior to adjusting for covariates, a general trend in weight can be seen in the random effects. Subjects who are heavier tend to have larger random effect, and males tend to have larger random effects than females of the same weight. Patterns such as these indicate that weight and sex can be used to explain variation in the random effects. After adjusting for sex and weight (right), the random effects have no discernable pattern.

6) The two stage DTR we describe in the previous sections, estimated via Q learning. We refer to this mode as “Q Learning”.

We compare all methods on their achieved reward as well as the difference between the achieved reward and theoretically largest reward (the reward we would achieve if we knew the pharmacokinetic parameters exactly). Because we know the true latent pharmacokinetic parameters of the simulated subjects, we can optimize the reward with the known pharmacokinetics of the subject, thereby yielding the largest reward possible.

## 5 Results

### 5.1 Bayesian Model

In this section we present results of two different analyses. We first assess the fit of our Bayesian model, as the model is crucial for our decision making experiments. We then present the results of our decision making experiments.

We fit M1 to real pharmacokinetic data using Stan. Stan monitors several markov chain diagnostics none of which detected problematic markov chain behavior, which indicates that Stan’s sampling algorithm was able to converge to the target distribution (0 divergences, all all Gelman-Rubin diagnostics  $< 1.01$ , all effective sample size ratios  $> 22\%$ ).

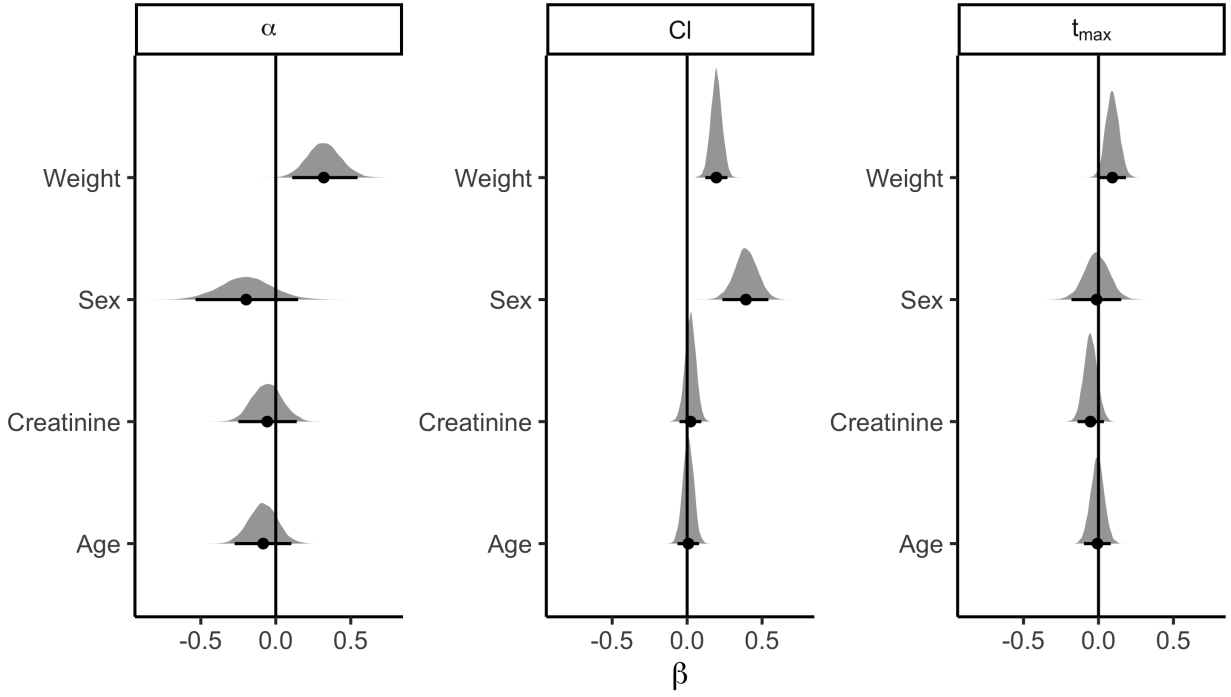


Figure 2: Posterior distributions of regression coefficients. Expectations are shown as black dots, 95% credible intervals are shown as horizontal black lines. Solid black vertical line is  $\beta = 0$  for reference. Note, regression coefficients for  $Cl$  and  $t_{max}$  act multiplicatively (a one unit increase in weight leads to a change in  $Cl$  of  $\exp(\beta)$ ), while regression coefficients for  $\alpha$  are interpreted on the log odds scale.

The inclusion of covariates in the model results in a better fit than excluding them. Shown in figure 1 are the estimated random effects for the clearance pharmacokinetic parameter of each subject as a function of weight. Subject sex is indicated by color, the overall trend is shown in the black dashed line. Failing to include subject sex and weight results in males having on average a larger random effect than females of the same weight, and heavier subjects having a larger random effect than lighter subjects. When covariates are added into the model, the variation in the random effects attenuates, resulting in closer alignment to model assumptions. A better fit to the data means data generated from the model may be closer aligned with the true data generating process.

Examining the posterior distributions of the regression coefficients provides further insight. Subject weight increases the expected value of alpha (which is used to compute the elimination and absorption rates in the first order one compartment PK model. The parameter  $\alpha$  is the ratio of how fast the drug exits the central compartment to how fast the drug enters the central compartment) as well as the time to max concentration. There appears to be an effect of sex on  $\alpha$  (males have smaller alpha than females, meaning the drug leaves their central compartment slower or enters the central compartment quicker), however the uncertainty is large (estimated effect -0.2 on the logit scale, 95% credible interval -0.54 to 0.15). See supplementary table X for a full summary of the regression coefficients.

Model training error sees a very small improvement. Including subject covariates decreases model training error from 6.89 ng/ml to 6.84 ng/ml. Estimates of concentration uncertainty remain similar between the two models as well. We conclude the inclusion of covariates in the model improves model inferences but does not improve the fit of the model to the data in any substantial way. Either model would require additional validation prior to using in a predictive capacity.

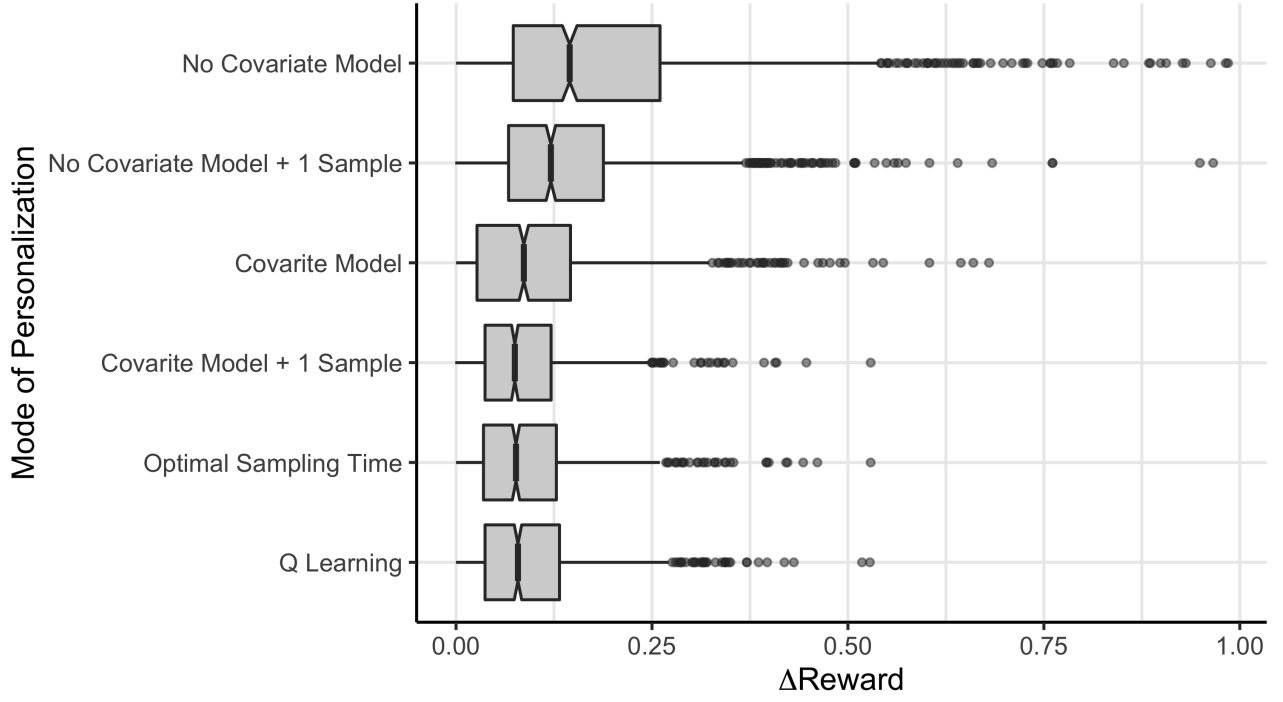


Figure 3: Boxplots of the difference between theoretically largest reward and achieved reward for each of the 1000 simulated subjects. Subjects who achieve a reward close to their maximum reward have a difference on 0, subjects who achieve a reward less than their maximum have larger differences, with the largest difference being 1.

## 5.2 Simulation Results

We consider 6 modes of personalization which range in the amount of information used in the decision process as well as burden placed on the patient and clinic, and burden of implementation. We present the results of our simulation in figure X below in terms of difference between theoretically largest reward and reward achieved by the mode of personalization. The results are ordered from least amount of information and burden (top) to most amount of information and burden (bottom).

Modes of personalization which use less information have a larger difference (i.e. yield smaller reward on average than what is theoretically possible). The no covariate model (which uses no information about the subject) performs worst with a median difference of 0.19. The distribution of differences for this mode is right skewed with some differences exceeding 0.95, meaning the subject could have been in range for nearly the entire time but the mode selected a dose which failed to put the subject in range.

The use of covariates in the model nearly cuts the difference in half, achieving an average difference of 0.1 with smaller right skew. There is a diminishing in the difference in rewards as additional burden is undertaken. Modes which use observed concentration information (Covariate Model + 1 Sample, Optimal Sampling Time, and Q Learning) lead to marginally more reward on average.

## 6 Discussion

Including age, sex, weight, and creatinine in our Bayesian model improved the inferences from that model. Though the predicted concentrations and estimates of uncertainty changed negligibly, the covariates explained residual confounding in the random effects. This results in a model which better explains the observed variation in the data and hence should generate more plausible data for simulation.

That modes of personalization which use more information result in larger reward is ultimately unsurprising. From our perspective, the more pertinent result is the diminishing return on investment observed when using additional information (and consequently, taking on additional implementation burden to effectively use that information). Were doses to be personalized for the simulated population, we would recommend the covariate model be used, as it is easier to implement, puts smaller burden on the clinic, and results in approximately mean/median rewards as compared to other methods. Taking an additional sample doesn't seem to improve the expected mean/median reward appreciably to be worth the burden of having the patient take time to come in, drawing their blood, measuring their concentration in the lab, and reporting results to the decision maker, in which a decision to not adjust the dose might be made anyway. Similar arguments can be made for Q learning, which has even higher implementation burden.

But mean/median reward does not tell the whole story. As noted, the distribution of differences in reward is right skewed. Some subjects have a very large difference, and the possibility of these differences might not be acceptable in different contexts with different drugs. There is a tradeoff between less extreme differences and taking on additional clinic and implementation burden, and that tradeoff should be examined on a case by case basis.

Context is crucial, and how we adapt to that context is perhaps a question in need of closer examination. Traditional methods of personalization include conditioning only on a subject's covariates (not unlike the Covariate model we present here). But of course patients are not their age, sex, weight, and creatinine. Additionally, safety information and best available practices might change in the future as more research on drugs is performed. Were new safety information to be published, one might imagine the reward function might be affected, which may result in a new mode of personalization being more/less preferable or more/less feasible. Any number of factors in flux can change the context in which personalization occurs, and that change in context may prompt for a re-evaluation in how personalization is done.

Thus, our results are not about apixaban per se. We don't offer recommendations on how personalization for apixaban should be done because we can't anticipate the context. What we offer is a framework for developing strategies of personalization and evaluating their performance against their implementation and clinic burden. Context can be changed where needed, either through the reward function, or by adjusting when the clinic is able to take measurements, or by including additional information such as genotype in the Bayesian model. Using this framework, clinics have flexibility to personalize the personalization.

## 7 Limitations

We've examined six modes for making decisions. The next mode improves on a deficiency of the previous mode in a natural manner, and so our experiment constitutes a kind of ablation study. We believe the decision making aspect of our study extracts information in a responsible way and uses the best decision making methodology

available. That being said, the experiment is not without limitations.

The bayesian model of the pharmacokinetics is integral to the methodology we present. Any shortcomings in the model affect the quality of the decision and decision process. Bayesian models are not as ubiquitous as other models in pharmacology, and so particular expertise is required for model development and evaluation. That expertise increases the implementation burden of any decision process involving Bayesian models. However, we demonstrate how one such model can be constructed in a past study [CITE] and include open sourced code and data for practitioners to replicate our model fitting.

Additionally, the data required to construct a high quality Bayesian model of pharmacokinetics require multiple observations of a single patient over an extended time, preferably over multiple well timed doses with near perfect adherence. Obtaining such data requires well organized efforts and is high burden for both investigators and participating subjects. This makes acquiring a robust Bayesian model for use in dose personalization difficult.

## 8 Future Work

Because the data required to build reliable Bayesian pharmacokinetic models are difficult to collect in practice, research into developing these models from observational data may prove fruitful in extending this work. If clinics record data on measured blood concentrations, they may have dozens or hundreds of subjects with only one or two measurements per subject. Moreover, the subjects in question may be on multiple drugs or have comorbidities which may affect the pharmacokinetics of the drug under study. Additional research into constructing Bayesian models which can adjust for polypharmacy and comorbidities while learning an individual's pharmacokinetics from a large but sparse sample would drive this work towards being easier to implement in practice.

## References

- [1] Bridget L Morse and Richard B Kim. Is personalized medicine a dream or a reality? *Critical reviews in clinical laboratory sciences*, 52(1):1–11, 2015.
- [2] Theodore John Wigle, Brandi Povitz, Wendy Teft, Robin Legan, John Gordon Lenehan, Markus Gulilat, Stephanie Nevison, Justin Kritzinger, Veera Punaganty, Denise Keller, et al. Prospective cohort study of the impact of hospital-wide dihydropyrimidine dehydrogenase (dpyd) genotype testing for fluoropyrimidine-based chemotherapy on adverse events and hospital costs. *American Society of Clinical Oncology*, 2019.
- [3] Margot de Looft, Bob Wilffert, Cornelis Boersma, Lieven Annemans, Stefan Vegter, Job FM van Boven, and Maarten J Postma. Economic evaluations of pharmacogenetic and pharmacogenomic screening tests: a systematic review. second update of the literature. *PloS one*, 11(1):e0146262, 2016.
- [4] Fatiha H Shabaruddin, Nigel D Fleeman, and Katherine Payne. Economic evaluations of personalized medicine: existing challenges and current developments. *Pharmacogenomics and personalized medicine*, 8: 115, 2015.
- [5] Miriam Kasztura, Aude Richard, Nefti-Eboni Bempong, Dejan Loncar, and Antoine Flahault. Cost-effectiveness of precision medicine: a scoping review. *International journal of public health*, 64(9):1261–1271, 2019.
- [6] Wolf Rogowski, Katherine Payne, Petra Schnell-Inderst, Andrea Manca, Ursula Rochau, Beate Jahn, Oguzhan Alagoz, Reiner Leidl, and Uwe Siebert. Concepts of ‘personalization’ in personalized medicine: implications for economic evaluation. *Pharmacoeconomics*, 33(1):49–59, 2015.
- [7] Antonello Di Paolo, François Sarkozy, Bettina Ryll, and Uwe Siebert. Personalized medicine in europe: not yet personal enough? *BMC health services research*, 17(1):1–9, 2017.
- [8] Peter S Houts, Allan Lipton, Harold A Harvey, Barbara Martin, Mary A Simmonds, Richard H Dixon, Santo Longo, Thomas Andrews, Robert A Gordon, John Meloy, et al. Nonmedical costs to patients and their families associated with outpatient chemotherapy. *Cancer*, 53(11):2388–2392, 1984.
- [9] Anna Lee, Kanan Shah, and Fumiko Chino. Assessment of parking fees at national cancer institute–designated cancer treatment centers. *JAMA oncology*, 6(8):1295–1297, 2020.
- [10] Rachel A Elliott, Judith A Shinogle, Pamela Peele, Monali Bhosle, and Dyfrig A Hughes. Understanding medication compliance and persistence from an economics perspective. *Value in health*, 11(4):600–610, 2008.
- [11] Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Calif. L. Rev.*, 104:671, 2016.
- [12] Kristian Lum and William Isaac. To predict and serve? *Significance*, 13(5):14–19, 2016.
- [13] Ifeoma Ajunwa. The paradox of automation as anti-bias intervention, 41 cardozo, 1, 2020.
- [14] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.

- [15] Holger Fröhlich, Rudi Balling, Niko Beerenwinkel, Oliver Kohlbacher, Santosh Kumar, Thomas Lengauer, Marloes H Maathuis, Yves Moreau, Susan A Murphy, Teresa M Przytycka, et al. From hype to reality: data science enabling personalized medicine. *BMC medicine*, 16(1):1–15, 2018.
- [16] Bibhas Chakraborty. *Statistical methods for dynamic treatment regimes*. Springer, 2013.
- [17] Wonkyung Byon, Samira Garonzik, Rebecca A Boyd, and Charles E Frost. Apixaban: a clinical pharmacokinetic and pharmacodynamic review. *Clinical pharmacokinetics*, 58(10):1265–1279, 2019.
- [18] Rommel G Tirona, Zahra Kassam, Ruth Strapp, Mala Ramu, Catherine Zhu, Melissa Liu, Ute I Schwarz, Richard B Kim, Bandar Al-Judaibi, and Melanie D Beaton. Apixaban and rosuvastatin pharmacokinetics in nonalcoholic fatty liver disease. *Drug Metabolism and Disposition*, 46(5):485–492, 2018.
- [19] A Demetri Pananos and Daniel J Lizotte. Comparisons between hamiltonian monte carlo and maximum a posteriori for a bayesian model for apixaban induction dose & dose personalization. In *Machine Learning for Healthcare Conference*, pages 397–417. PMLR, 2020.