

CS 234 Winter 2022  
HW 4  
Due: Feb 27 at 11:59 pm (PST)

For submission instructions please refer to [website](#). For all problems, if you use an existing result from either the literature or a textbook to solve the exercise, you need to cite the source.

This time, there is no submission script for you to run in the starter code of Problem 1. Please only zip “main.py” into a zip file when you upload to Gradescope.

## 1 Estimation of the Warfarin Dose [60pts]

### 1.1 Introduction

**Warfarin** is the most widely used oral blood anticoagulant agent worldwide; with more than 30 million prescriptions for this drug in the United States in 2004. The appropriate dose of warfarin is difficult to establish because it can vary substantially among patients, and the consequences of taking an incorrect dose can be severe. If a patient receives a dosage that is too high, they may experience excessive anti-coagulation (which can lead to dangerous bleeding), and if a patient receives a dosage which is too low, they may experience inadequate anti-coagulation (which can mean that it is not helping to prevent blood clots). Because incorrect doses contribute to a high rate of adverse effects, there is interest in developing improved strategies for determining the appropriate dose ([Consortium, 2009](#)).

Commonly used approaches to prescribe the initial warfarin dosage are the *pharmacogenetic algorithm* developed by the IWPC (International Warfarin Pharmacogenetics Consortium), the *clinical algorithm* and a *fixed-dose* approach.

In practice a patient is typically prescribed an initial dose, the doctor then monitors how the patient responds to the dosage, and then adjusts the patient’s dosage. This interaction can proceed for several rounds before the best dosage is identified. However, it is best if the correct dosage can be initially prescribed.

This question is motivated by the challenge of Warfarin dosing, and considers a simplification of this important problem, using real data. The goal of this question is to explore the performance of multi-armed bandit algorithms to best predict the correct dosage of Warfarin for a patient *without* a trial-and-error procedure as typically employed.

**Problem setting** Let  $T$  be the number of time steps. At each time step  $t$ , a new patient arrives and we observe its individual feature vector  $X_t \in \mathbb{R}^d$ : this represents the available knowledge about the patient (e.g., gender, age, ...). The decision-maker (your algorithm) has access to  $K$  arms,

where the arm represents the warfarin dosage to provide to the patient. For simplicity, we discretize the actions into  $K = 3$

- Low warfarin dose: under 21mg/week
- Medium warfarin dose: 21-49 mg/week
- High warfarin dose: above 49mg/week

If the algorithm identifies the correct dosage for the patient, the reward is 0, otherwise a reward of  $-1$  is received.

Lattimore and Szepesvári have a nice series of blog posts that provide a good introduction to bandit algorithms, available here: [BanditAlgs.com](http://BanditAlgs.com). The [Introduction](#) and the [Linear Bandit](#) posts may be particularly of interest. For more details of the available Bandit literature you can check out the [Bandit Algorithms Book](#) by the same authors.

## 1.2 Dataset

We use a publicly available patient dataset that was collected by staff at the Pharmacogenetics and Pharmacogenomics Knowledge Base (PharmGKB) for 5700 patients who were treated with warfarin from 21 research groups spanning 9 countries and 4 continents. You can find the data in `warfarin.csv` and metadata containing a description of each column in `metadata.xls`. Features of each patient in this dataset includes, demographics (gender, race, ...), background (height, weight, medical history, ...), phenotypes and genotypes.

Importantly, this data contains the true patient-specific optimal warfarin doses (which are initially unknown but are eventually found through the physician-guided dose adjustment process over the course of a few weeks) for 5528 patients. You may find this data in mg/week in **Therapeutic Dose of Warfarin**<sup>1</sup> column in `warfarin.csv`. There are in total 5528 patient with known therapeutic dose of warfarin in the dataset (you may drop and ignore the remaining 173 patients for the purpose of this question). Given this data you can classify the right dosage for each patient as *low*: less than 21 mg/week, *medium*: 21-49 mg/week and *high*: more than 49 mg/week, as defined in [Consortium \(2009\)](#) and [Introduction](#).

**The data processing is already implemented for you**

## 1.3 Implementing Baselines [10pts]

Please implement the following two baselines in `main.py`

1. *Fixed-dose*: This approach will assign 35mg/week (medium) dose to all patients.
2. *Warfarin Clinical Dosing Algorithm*: This method is a linear model based on age, height, weight, race and medications that patient is taking. You can find the exact model is section S1f of [appx.pdf](#).

Run the fixed dosing algorithm and clinical dosing algorithm with the following command:

```
python main.py --run-fixed --run-clinical
```

---

<sup>1</sup>You cannot use **Therapeutic Dose of Warfarin** data as an input to your algorithm.

You should see the `total_fraction_correct` to be fixed at about 0.61 for fixed dose and 0.64 for clinical dose algorithm. You can run them individually as well. Just use one of the command line arguments instead.

#### 1.4 Implementing a Linear Upper Confidence Bandit Algorithm [15pts]

Please implement the Disjoint Linear Upper Confidence Bound (LinUCB) algorithm from [Li et al. \(2010\)](#) in `main.py`. See Algorithm 1 from paper. Please feel free to adjust the `-alpha` argument, but you don't have to. Run the LinUCB algorithm with the following command:

```
python main.py --run-linucb
```

You should see the `total_fraction_correct` to be above 0.64, though the results may vary per run.

#### 1.5 Implementing a Linear eGreedy Bandit Algorithm [5pts]

Is the upper confidence bound making a difference? Please implement the e-Greedy algorithm in `main.py`. Please feel free to adjust the `-ep` argument, but you don't have to. Does eGreedy perform better or worse than Upper Confidence bound? (You do not need to include your answers here) Run the  $\epsilon$ -greedy LinUCB with the following command:

```
python main.py --run-egreedy
```

You should see the `total_fraction_correct` to be above 0.61, though the results may vary per run.

#### 1.6 Implementing a Thompson Sampling Algorithm [20pts]

Please implement the Thompson Sampling for Contextual Bandits from [Agrawal and Goyal \(2013\)](#) in `main.py`. See Algorithm 1 and section 2.2 from paper. Please feel free to adjust the `-v2` argument, but you don't have to. (This actually  $v$  squared from the paper) Run the Thompson Sampling algorithm with the following command:

```
python main.py --run-thompson
```

You should see the `total_fraction_correct` to be **around** 0.64, though the results may vary per run.

#### 1.7 Results [10pts]

At this point, you should see a plot in your results folder titled `"fraction_incorrect.png"`. If not, run the following command to generate the plot:

```
python main.py
```

Include this plot in for this part. Please also comment on your results in a few sentences. How would you compare the algorithms? Which algorithm "did the best" based on your metric?

## 2 Learning a Policy From an Approximated MDP [20pts]

Consider an infinite-horizon, discounted MDP  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$  with  $\gamma \in [0, 1)$  and bounded rewards such that  $R_{\max} = \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{R}(s,a) < \infty$ . In practical settings, we rarely know the true model of the agent-environment interaction. Here, we are interested in the case where the model is estimated from experience data in the real world; scarcity of data then implies that our model will only be approximate.

Recall certainty-equivalence and model-based reinforcement-learning where we attempt to compute an optimal policy by first estimating an approximate MDP  $\widehat{\mathcal{M}} = \langle \mathcal{S}, \mathcal{A}, \widehat{\mathcal{R}}, \widehat{\mathcal{T}}, \gamma \rangle$  from the experience data which is identical to  $\mathcal{M}$  except for the approximate reward function  $\widehat{\mathcal{R}} : \mathcal{S} \times \mathcal{A} \rightarrow [0, R_{\max}]$  and approximate transition function  $\widehat{\mathcal{T}} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ .

Naturally, we are interested in controlling the performance difference between policies computed in  $\widehat{\mathcal{M}}$  but deployed in  $\mathcal{M}$ . Let  $\pi_{\mathcal{M}}^*$  and  $\pi_{\widehat{\mathcal{M}}}^*$  denote the optimal policies of  $\mathcal{M}$  and  $\widehat{\mathcal{M}}$ , respectively. We define the planning loss as follows:

$$\|V_{\mathcal{M}}^{\pi_{\mathcal{M}}^*} - V_{\widehat{\mathcal{M}}}^{\pi_{\widehat{\mathcal{M}}}^*}\|_{\infty}$$

Here,  $V_{\mathcal{M}}^{\pi}$  indicates the value function obtained by running policy  $\pi$  in the MDP  $\mathcal{M}$

If  $\Pi = \{\mathcal{S} \rightarrow \mathcal{A}\}$  denotes the class of all stationary, deterministic policies, let  $\overline{\Pi} \subseteq \Pi$  be a restricted policy class.

(a) [5pts] Prove that

$$\|V_{\mathcal{M}}^{\pi_{\mathcal{M}}^*} - V_{\widehat{\mathcal{M}}}^{\pi_{\widehat{\mathcal{M}}}^*}\|_{\infty} \leq 2 \max_{\pi \in \overline{\Pi}} \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty}.$$

For any function  $f : \mathcal{X} \rightarrow \mathbb{R}$  and any two sets  $A \subseteq B \subseteq \mathcal{X}$ , we have that  $\max_{x \in A} f(x) \leq \max_{x \in B} f(x)$ .

This fact immediately implies the following corollary of part (a) whenever  $\pi_{\mathcal{M}}^*, \pi_{\widehat{\mathcal{M}}}^* \in \overline{\Pi}$ :

$$\|V_{\mathcal{M}}^{\pi_{\mathcal{M}}^*} - V_{\widehat{\mathcal{M}}}^{\pi_{\widehat{\mathcal{M}}}^*}\|_{\infty} \leq 2 \max_{\pi \in \overline{\Pi}} \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty}.$$

(b) [7pts] Prove that for any  $\pi \in \Pi$ ,

$$\|Q_{\mathcal{M}}^{\pi} - Q_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \leq \frac{1}{1-\gamma} \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| \widehat{\mathcal{R}}(s,a) + \gamma \mathbb{E}_{s' \sim \widehat{\mathcal{T}}(\cdot|s,a)} [V_{\mathcal{M}}^{\pi}(s')] - Q_{\mathcal{M}}^{\pi}(s,a) \right|.$$

**Hint:** Define the approximate Bellman operator  $\widehat{\mathcal{B}}^{\pi} : \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\} \rightarrow \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$  such that

$$\widehat{\mathcal{B}}^{\pi} Q(s,a) = \widehat{\mathcal{R}}(s,a) + \gamma \mathbb{E}_{s' \sim \widehat{\mathcal{T}}(\cdot|s,a), a' \sim \pi(s')} [Q(s', a')],$$

and consider the sequence of value functions  $\{Q_0, Q_1, \dots, Q_m, \dots, Q_{\infty}\}$  where

$$\begin{cases} Q_0(s,a) = Q_{\mathcal{M}}^{\pi}(s,a) & \forall (s,a) \in \mathcal{S} \times \mathcal{A} \\ V_0(s) = V_{\widehat{\mathcal{M}}}^{\pi} & \forall s \in \mathcal{S} \\ Q_m(s,a) = \widehat{\mathcal{B}}^{\pi} Q_{m-1} & \forall (s,a) \in \mathcal{S} \times \mathcal{A} \\ Q_{\infty}(s,a) = Q_{\widehat{\mathcal{M}}}^{\pi}(s,a) & \forall (s,a) \in \mathcal{S} \times \mathcal{A} \end{cases}.$$

**Hint:** Can you show that  $\widehat{\mathcal{B}}^{\pi}$  is a contraction? You may find the telescoping sum and triangle inequality proof techniques you used in Assignment 1 Q2 to be useful here.

- (c) **[8pts]** Assume that the restricted policy class is finite ( $|\bar{\Pi}| < \infty$ ) and that we have observed  $n$  independent samples of rewards and next-state transitions from each state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Prove for any  $\delta \in [0, 1]$  that the following bound on planning loss holds with probability at least  $1 - \delta$ :

$$\|V_{\mathcal{M}}^{\pi_{\mathcal{M}}} - V_{\mathcal{M}}^{\pi_{\widehat{\mathcal{M}}}}\|_{\infty} \leq \frac{2R_{\max}}{(1-\gamma)^2} \sqrt{\frac{1}{2n} \log \frac{2|\mathcal{S}||\mathcal{A}||\bar{\Pi}|}{\delta}}.$$

**Hint:** We are sampling from  $\widehat{\mathcal{T}}$  in  $\widehat{\mathcal{M}}$  to compute  $Q_{\widehat{M}}^{\pi}$ , while the true mean is  $Q_M^{\pi}$ .

**Hint:** Recall Hoeffding’s inequality – let  $X_1, X_2, \dots, X_n$  be a finite sequence of i.i.d., real-valued random variables such that  $X_i \in [a, b]$  for all  $i \in \{1, \dots, n\}$ . Denote the sample mean and population mean as  $\bar{X} = \frac{1}{n} \sum_i X_i$  and  $\mu = \mathbb{E}[X_1]$ , respectively. Then, the following concentration inequality holds for any  $\varepsilon > 0$ :

$$\mathbb{P}(|\bar{X} - \mu| \leq \varepsilon) \geq 1 - 2 \exp \left( \frac{-2n^2 \varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right).$$

### 3 Analysing features used in Warfarin dosage algorithms [4pts]

1. **[2pts]** As a [recent survey of dosing algorithms](#) suggest, most existing clinical algorithms have not been evaluated for their efficacy in under-served populations. How big of a disparity would warrant intervention? How would you decide on a threshold? In a scenario in which your disparity threshold was met, what constraint(s) ([Thomas et al. \(2017\)](#)) could you introduce to your algorithm?
2. Many current Warfarin dosage algorithms rely on race as a feature. When race is used in medicine, it is often intended as a proxy for genetic difference ([Vyas et. al 2020](#)). However, as ([Goodman & Brett 2021](#)) note, prevalence of genetic markers relevant to a particular disease might vary more within a group than between groups. Observed racial differences in health outcomes are at least as likely to be caused by social and environmental determinants of health (such as differential access to lead-free water) as they are by genetic factors.
  - (a) **[OPTIONAL]** Run at least one of your algorithms again without the “race” or “gender” feature. How does accuracy change?
  - (b) **[2pts]** In addition to knowing whether average accuracy increases or decreases, what else might you measure or investigate in order to determine whether to use race as a feature in your algorithm?

## References

- S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- I. W. P. Consortium. Estimation of the warfarin dose with clinical and pharmacogenetic data. *New England Journal of Medicine*, 360(8):753–764, 2009.

- L. Li, W. Chu, J. Langford, and R. E. Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- P. S. Thomas, B. C. da Silva, A. G. Barto, and E. Brunskill. On ensuring that intelligent machines are well-behaved, 2017.