

### Abstract

Warfarin is the most widely used oral blood anticoagulant agent worldwide, with over 30 million prescriptions in the United States alone in 2004. However, prescribing the appropriate dosage of Warfarin for each patient can be challenging due to significant individual variability. Incorrect dosages can result in severe consequences, including dangerous bleeding or inadequate prevention of blood clots. While various approaches, such as pharmacogenetic and clinical dosing algorithms, have been developed to determine the initial dosage, they still rely on a trial-and-error procedure, which can lead to adverse effects.

To address this challenge, this project aims to apply multi-armed bandit algorithms to predict the correct dosage of Warfarin without relying on a trial-and-error procedure. Specifically, the project investigates the performance of linear bandit algorithms, demonstrating that LinUCB outperforms both the fixed and clinical dosing baselines. We also investigate alternative model formulations and assess performance.

### Models

#### Benchmark Models

We utilize two baselines to evaluate bandit algorithms. The first is a fixed-dose baseline, which prescribes a constant dose. The second is a clinical dosing baseline, which uses a linear equation that calculates the dose based on patient characteristics.

$$\begin{aligned} & \sqrt{\text{weekly dose}} \\ &= 4.0376 - 0.2546 * \text{Age} \\ & \quad - 0.6752 * \text{Asian} + 0.001 * \text{Enzyme inducer status} \end{aligned}$$

#### LinUCB with Dose

Our primary bandit model is LinUCB, which uses a linear model to predict the correct dosage of Warfarin. The model is trained on a dataset of patient dosages and outcomes. The patient dose is selected based on the model's prediction, and the patient's characteristics are used to update the model. The model is evaluated using a sufficient exploration matrix of past outcomes.

# Default Final Project Estimation of the Warfarin Dose

Eva Lestant, Amar Venugopal

[elestant@stanford.edu](mailto:elestant@stanford.edu), [amarvenu@stanford.edu](mailto:amarvenu@stanford.edu)

## Models

benchmark models against which we compare performance. The first is the fixed dose model, which simply assigns every patient a medium dose. The second is the clinical dosing algorithm, a simple linear model that calculates the square root of weekly dose as a function of patient covariates:

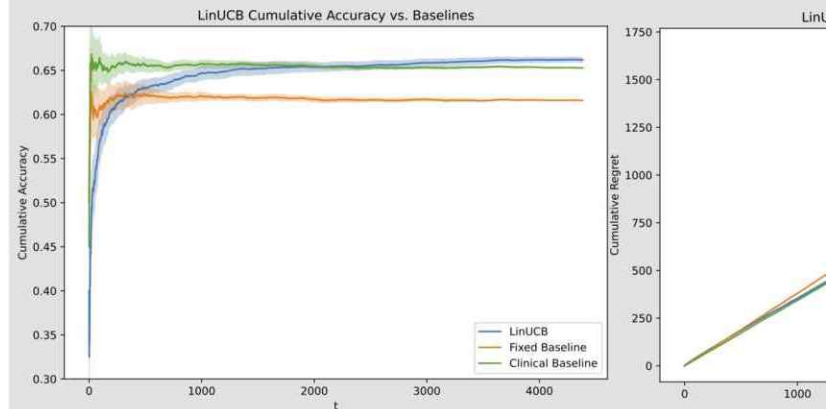
$$\text{Dose} = 0.0118 * \text{Age} + 0.0118 * \text{Height (cm)} + 0.0134 * \text{Weight (kg)} + 0.4060 * \text{Black} + 0.0443 * \text{Missing or mixed race} + 1.2799 * \text{Amiodarone status} - 0.5695 * \text{Amiodarone status}$$

## Disjoint Linear Models

The default algorithm is LinUCB with disjoint linear models (Li et al. (2010)). In this model, the optimal dose is selected based on a linear function of the patient characteristics plus an optimism term to ensure exploration. In the below equations,  $D$  refers to the dose assigned to patients assigned to arm  $a$ .

## Results

We find that LinUCB significantly outperforms the baseline, both in terms of both cumulative accuracy and cumulative expected regret. The resulting performance follows, with the shaded regions corresponding to the confidence intervals over the 20 different random patient sets.

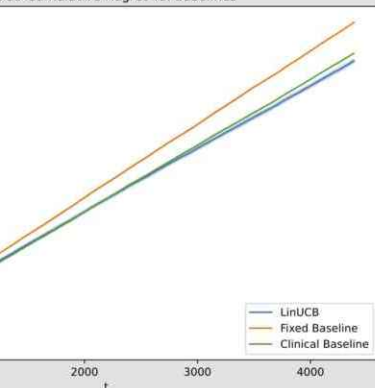


We also go on to find that our extension

e

forms the fixed dose  
ative accuracy and  
formance plots are as  
ng to 95% confidence  
orderings tested:

LinUCB Cumulative Regret vs. Baselines



on models, namely

## Discussion & Conclusion

Our results show clear performance improvements when comparing bandit methods to both the fixed dose baseline and, in the case of LinUCB, to the clinical dosing algorithm as well. This clearly demonstrates the power of bandit methods for Warfarin dose assignment and suggests that such methods should be further explored to achieve better patient outcomes.

However, there is one key caveat to these results: while LinUCB outperforms the clinical dosing algorithm by the end of the “online” training period, it takes nearly 2000 patients before its performance truly matches that of the clinical dosing algorithm. As a result, if this were a truly online setting, several of those patients may have received a lower quality of care than if the clinical dosing algorithm had been applied. This raises ethical considerations regarding potential harm to patients used in training bandit methods, particularly for healthcare applications. This suggests that some degree of offline learning, such as what



This project leverages a publicly available patient dataset collected by PharmGKB. Comprising of 5528 patients drawn from studies across 9 countries, this dataset includes optimal patient-specific warfarin doses, as well as patient features such as gender, race, height, weight, medical history, genotypes, and phenotypes.

In order to handle these categorical genotypic features, we employ one-hot encoding, constructing a dummy variable for each class membership. We then drop one such dummy variable from each group and include an intercept in our feature space, in keeping with standard practice for linear models. The end result is a dataset with 4386 observations and 23 features.

The label, which we try to predict, is given by the correct therapeutic dose of warfarin, which we discretize into 3 bins corresponding to “*low*”, “*medium*”, and “*high*”.

Where  $A_a = D_a^\top D_a$

Our second band to contextual bandits (Lattimore & Goyal (2013)). There, we sample  $\theta_t$  continuously upon each round, and the prescribed  $\theta_t$  is a separate set of parameters for each setting similar to

$$\hat{\mu}_g(t)$$

Our final algorithm is a linear model based on the information, this is better than all others tested. The performance of linear

$a_t$

...ates for patients assigned to arm  $a$ .

$$\arg \max_{a \in \mathcal{A}_t} \left( x_{t,a}^\top \hat{\theta}_a + \alpha \sqrt{x_{t,a}^\top A_a^{-1} x_{t,a}} \right) + I_d$$

## Sampling

This algorithm is Thompson Sampling, applied to bandits with linear payoffs, based on Agrawal. This algorithm works by posterior sampling, i.e. sample the linear parameter vector from a Gaussian posterior distribution. We deviate from the standard implementation by giving each arm a set of parameters, thereby making it a disjoint LinUCB.

$$b(t) = I_d + \sum_{\tau=1}^{t-1} b(\tau) b(\tau)^\top \mathbf{1}\{a_\tau = a\}$$

$$= B_a(t)^{-1} \left( \sum_{\tau=1}^{t-1} b(\tau) r(\tau) \mathbf{1}\{a_\tau = a\} \right)$$

$$\tilde{\mu}_a(t) \sim \mathcal{N}(\hat{\mu}_a(t), v^2 B_a(t)^{-1})$$

$$a_t = \arg \max_a b^\top(t) \tilde{\mu}_a(t)$$

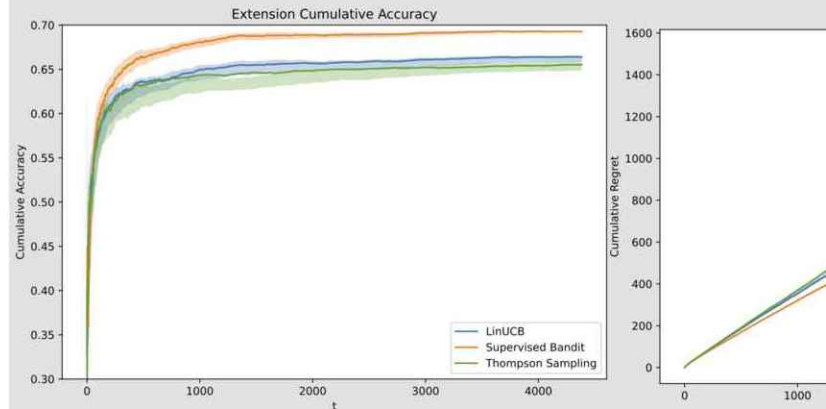
## Supervised Linear Bandit

This is a supervised linear bandit, which trains a model on all prior observed true Warfarin doses, and a standard binary reward. By utilizing this extra information, the model should achieve better performance than standard bandit models and will set a kind of upper bound on the performance of standard bandit models.

$$\hat{\theta}_t = (X_t^\top X_t)^{-1} X_t^\top Y_t$$

$$a_{t+1} = \arg \max_{a \in \{0,1,2\}} \|a - x_{t+1}^\top \hat{\theta}_t\|$$

...supervised linear regression “bandit” and Thompson Sampling achieve good performance:



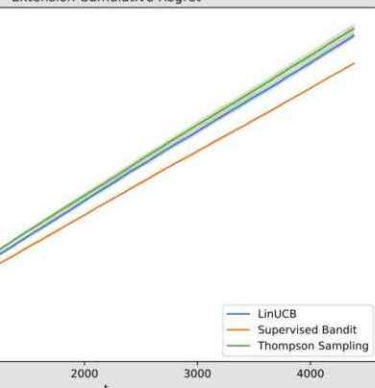
We can see that Thompson Sampling appears to have similar performance to LinUCB (within error). As expected, the supervised linear “bandit” achieves significantly higher cumulative regret compared to both models. This is unsurprising since the supervised model has access to the true treatment applied to past patients, which amounts to just a binary reward, which amounts to just a binary correct/incorrect treatment.

## References

1. I. W. P. Consortium. Estimation of the warfarin dose: a pharmacogenomic study. *Journal of Clinical Pharmacology*, 360(8):753–764, 2009.
2. L. Li, W. Chu, J. Langford, and R. E. Schapire. *Proceedings of the 19th international conference on machine learning*, 2006.
3. S. Agrawal and N. Goyal. Thompson sampling. *Learning*, pages 127–135, 2013.

Thompson Sampling,

Extension Cumulative Regret



to achieve similar  
ected, the supervised  
accuracy and lower  
prising, as it is given  
ients, rather than the  
nary indicator for

we have performed in this project, is necessary prior to the rollout of a bandit algorithm in such contexts.

## Future

The key direction of future research in this project is to investigate further implementations of linear bandit models to attempt to close the performance gap between our current optimal implementation, LinUCB, and the supervised linear bandit. While the latter's performance presents an upper bound on the potential performance of linear bandit models and likely cannot be replicated by a model that observes only binary rewards, the existing gap suggests that other model formulations may be able to provide better performance.

Possible candidates that we have not tested in this project include robust algorithms or regularized linear models.

irin dose with clinical and pharmacogenetic data. *New England Journal of Medicine*,

e. A contextual-bandit approach to personalized news article recommendation. In  
*nal conference on World wide web*, pages 661–670, 2010.

ng for contextual bandits with linear payoffs. In *International Conference on Machine*