

# BANDIT THEORY: APPLICATIONS TO LEARNING HEALTHCARE SYSTEMS AND CLINICAL TRIALS

Michael Sklar<sup>1</sup>, Mei-Chiung Shih<sup>1,2</sup> and Philip Lavori<sup>1</sup>

<sup>1</sup>*Stanford University and*

<sup>2</sup>*VA Palo Alto Cooperative Studies Program Coordinating Center*

*Abstract:* In recent years, statisticians and clinical scientists have defined two new approaches for studying the effects of medical practice, extending the “gold standard” classical randomized clinical trial to remedy some of its defects, improve its fit to clinical practice, and conform more closely to ethical principles. The contextual multi-armed bandit provides a natural statistical structure for a learning health-care system, allowing the optimization of patient outcomes by adaptively assigning treatments, while building in experimental strength for accuracy in learning. The sequential multiple assignment randomized trial has become the standard for comparing entire dynamic treatment strategies for the management of chronic disease, which more closely matches the goals and practice of clinicians. The theory and methods developed by Professor Tze Leung Lai over the course of his career are of central importance in bringing these two apparently different approaches to bear in efforts to improve clinical practice. We review these methods in this article.

*Key words and phrases:* Clinical trials, medical and pharmaceutical statistics, sequential analysis and optimal stopping.

## 1. Introduction

Throughout his career, Professor Tze Leung Lai has made major contributions to the design of randomized clinical trials. We note his creative work with many students and colleagues in group sequential stopping (Lai and Shih (2004)) and adaptive trials (Lai, Lavori and Liao (2014); Bartroff, Lai and Shih (2012); Lai and Liao (2012)), which have advanced the field of clinical trials. Here, we focus on two areas of Lai’s work that may have an even greater impact on the future of clinical research, as the medical community grapples with the challenges of generating and applying knowledge at point of care in fulfillment of the concept of the “learning healthcare system” (LHS) (Chamberlayne et al. (1998)). “A learning healthcare system is one that is designed to generate and apply the best evidence for the collaborative healthcare choices of each patient and provider; to

---

Corresponding author: Michael Sklar, Department of Statistics, Stanford University, Stanford, CA 94305, USA. E-mail: [sklarm@stanford.edu](mailto:sklarm@stanford.edu).

drive the process of discovery as a natural outgrowth of patient care; and to ensure innovation, quality, safety, and value in health care” (Olsen, Aisner and McGinnis (2007)). The first branch of Lai’s work discussed below deals with methods for incorporating true experimental strength into efforts to explore the comparative effects of different treatments, while exploiting what is learned to improve outcomes in patients. The underlying idea of the “multi-armed bandit” (MAB) for clinical decision-making goes back almost a century (Thompson (1933)), and Lai made his foundational contributions to that subject at the start of his career (Lai and Robbins (1985)). Recent work on the contextual generalization of the MAB (CMAB) has brought the idea back to the fore, as a sound theoretical basis for the LHS enterprise. We surveyed the current literature on the use of the CMAB in clinical medicine as part of a demonstration project currently underway in the US Department of Veterans Affairs (VA) Cooperative Studies Program (CSP), and below we discuss the current “state of the art.”

The other branch of Lai’s work discussed below has been used to further the recent development of a natural framework for defining and comparing dynamic treatment regimes (DTRs), also known as adaptive treatment strategies (ATS), in the management of chronic disease using variants of the sequential multiple assignment randomized trial (SMART) (Lavori and Dawson (2000, 2004); Murphy (2005)). This is again an attempt to bring experimental rigor to the study of clinical decision-making, taking full account of the inherently dynamic nature of ongoing clinical management of chronic disease.

We hope that the review and discussion in this celebratory paper will illustrate two of the less well known, but no less promising and valuable contributions made by Professor Lai in the course of his career, and perhaps encourage others to take up his ideas and carry them forward, as we have done.

## **2. The Multi-Armed Bandit Problem**

The name “multi-arm bandit” suggests a row of slot machines, which in the 1930s were nicknamed “one-armed bandits.” (Presumably the name is inspired by their pull-to-play levers and the often large house edge.) For a gambler in an unfamiliar casino, the “multi-arm bandit problem” would refer to a particular challenge: to maximize the expected winnings over a total of  $T$  plays, moving between machines as desired. The distribution of payouts from pulling each arm may be unknown and different for each machine. How should the gambler

play? Research into the MAB problem and its variants has led to foundational insights for problems in sequential sampling, sequential decision-making, and reinforcement learning.

Mathematical analysis of the MAB problem has been motivated by medical applications since Thompson (1933), with different medical treatments playing the role of bandit machines. Subsequent theory has found wide application across disciplines including finance, recommender systems, and telecommunications (Bouneffouf and Rish (2019)). According to Whittle (1979), the bandit problem was considered by Allied scientists in World War II, but it “so sapped [their minds] that the suggestion was made that the problem be dropped over Germany, as the ultimate instrument of intellectual sabotage.” It was Lai and Robbins (1985) who gave the first tractable asymptotically efficient solution.

Given a set of arms  $k \in 1, \dots, K$ , Lai and Robbins frame the question: How should we sample  $y_1, y_2, \dots$  sequentially from the  $K$  arms in order to achieve the greatest possible expected value of the sum  $S_T = y_1 + \dots + y_T$  as  $T \rightarrow \infty$ ? They model each sample from arm  $k$  as an independent draw from a population  $\Pi_k$  from a family of densities  $f_{\theta_k}$  indexed by parameter  $\theta_k$ . Then, they formalize the space of (possibly random) strategies  $\phi \in \Phi$ , defining  $\phi$  to be an *adaptive allocation rule* if it is a collection of random variables that makes the arm selection at each timestep,  $\phi := (\phi_1, \phi_2, \dots, \phi_T)$ . Thus, each  $\phi_t$  is a random variable on  $\{1, \dots, K\}$ , where the event  $\{\phi_t = k\}$  (“arm  $k$  is chosen at time  $t$ ”) belongs to the  $\sigma$ -field generated by prior decisions and observations  $(\phi_1, x_1, \phi_2, x_2, \dots, \phi_{t-1}, x_{t-1})$ . In this framework, Lai and Robbins (1985) define the *cumulative regret* of an adaptive allocation, rule which measures the strategy’s expected performance against the best arm, equivalent to

$$R_T(\phi, \theta) := \sum_{t=1}^T \mu^*(\theta) - \mathbb{E}[\mu(\theta_{\phi_t})],$$

where  $\mu(\theta_k)$  is the expected value of arm  $k$ , and  $\mu^*(\theta) := \max_k \{\mu(\theta_k)\}$ . Lai and Robbins (1985) give a strategy that achieves an expected cumulative regret of order  $O(\log T)$ , and provide a matching lower bound to show it is nearly optimal. This strategy creates an *upper confidence bound* (UCB) for each arm, where the estimated return is given a bonus for uncertainty. A simple example of a UCB is the UCB1 of Auer, Cesa-Bianchi and Fischer (2002), which at round  $t$ , picks the arm maximizing

$$\bar{y}_{k,t} + \sqrt{\frac{2 \ln(t)}{n_{k,t}}},$$

where the rewards  $y_t$  are in  $[0, 1]$ ,  $\bar{y}_{k,t}$  is the average of the observed rewards from arm  $k$ , and  $n_{k,t}$  is the number of samples observed from arm  $k$ . Typically, UCBs are designed so that inferior arm(s) are discarded with minimal investment, and the best arm(s) are guaranteed to remain in play; a key contribution of Lai and Robbins (1985) was to show how such statements can be quantified using Chernoff bounds (or other concentration inequality arguments), and then converted into an upper bound on the cumulative regret. Their approach has been generalized and extended to yield algorithms and regret guarantees across a variety of applications, with UCBs acting as a guiding design principle.

The richness of the bandit problem has generated a multitude of other approaches. By adding to the above model a prior distribution for the arm parameters  $\theta$ , the bandit problem can be framed as a Bayesian optimization over  $\phi$  to find the allocation strategy that minimizes the expected regret  $\int R_T(\theta, \phi) d p \theta$ . This optimization can, in principle, be solved with dynamic programming (as in Cheng and Berry (2007)); however, dynamic programming does not scale well to large or complicated experiments, because the number of possible states explodes. Using results from Whittle (1980), Villar, Bowden and Wason (2015) show how the computation can be reduced considerably by framing the optimal solution as an index policy.

When solving for the optimal strategy is not feasible, the heuristic solution of Thompson sampling is a popular choice, with good practical and theoretical performance (Chapelle and Li (2011); Kaufmann, Korda and Munos (2012); Russo and Van Roy (2016)). The decision rule proposed by Thompson (1933) is an adaptive allocation rule, where  $\phi_t$ , given all data observed prior to time  $t$ , is non-deterministic and chooses arm  $k$  with probability equal to its posterior chance of being the best arm. That is,  $\phi_t = k$  with probability  $p_{k,t} := P_{\mathbb{F}_t} \{k^* = k\}$ , where  $P_{\mathbb{F}_t}$  is the posterior probability distribution given  $(\phi_1, x_1, \dots, \phi_{t-1}, x_{t-1})$ , and  $k^* := \operatorname{argmax}_k (\mu(\theta_k))$  is the index of the best arm (which is a random variable). If the best arm is not unique, the tie should be broken to ensure the uniqueness of  $k^*$ . In fact, a Thompson allocation can be performed with just one sample from the posterior  $\mathbb{F}_t$ , as shown in the following workflow:

---

**Algorithm 1:** Bayesian Workflow with Thompson Sampling

---

- 1 Assume a likelihood model parametrized by  $\theta$ , such that  $\theta$  determines the arm means by  $\mu(\theta) = (\mu_1(\theta), \dots, \mu_k(\theta))$ ;
- 2 Assume a prior  $\mathbb{F}_1$ ;
- 3 **for** each sample  $t \in \{1, \dots, T\}$  **do**
- 4     Draw from the posterior a sample of the vector of arm means *sample*  $\theta' \sim \mathbb{F}_t$ ; set  $\mu' := (\mu_1(\theta'), \dots, \mu_k(\theta'))$ ;
- 5     Allocate to the arm corresponding to the largest entry of  $\mu'$ :  
       set  $\phi_t := \operatorname{argmax}_k \{\mu'_k\}$  (breaking ties at random);
- 6     Receive from arm  $\phi_t$  the next payoff  $x_t$ ;
- 7     Given the new observation, update posterior to  $\mathbb{F}_{t+1}$
- 8 **end**

---

Exact sampling from the posterior is not always tractable. A popular technique for sampling the posterior approximately is the Markov Chain Monte Carlo (MCMC) method. The convergence properties of MCMC to the posterior distribution, and in particular the number of steps that must be run to achieve accurate sampling, are well understood only in special cases (Diaconis (2009); Dwivedi et al. (2018)). Where theory falls short, practitioners may appeal to a variety of diagnostics tools to provide evidence of convergence to the posterior (Roy (2020)).

There are many other approaches to the bandit problem, including epsilon-greedy (Sutton and Barto (1998)), knowledge gradient (Ryzhov, Powell and Frazier (2012)), and information-directed sampling (Russo and Van Roy (2014a)).

### 3. Adaptive Randomization in an LHS

In an LHS, the arms of an MAB are treatments and the rewards are patient outcomes. Thus, minimizing the cumulative regret corresponds to maximizing patients' measured quality of care, a primary function of the LHS. However, typically, there is a secondary goal of learning from a trial: useful takeaways may include confidence intervals for the treatment effects, developing a treatment guide, or making recommendations for non-participating patients in parallel with the trial.

The goals of regret minimization and knowledge generation, often framed as "exploitation vs. exploration," are indeed in fundamental conflict: Bubeck, Munos and Stoltz (2011) formalized a notion of exploration-based experiments, where recommendations are made outside the trial. They define the *simple regret*

to be

$$r_T = \mu^* - \mu_{c,T},$$

where  $\mu_{c,T}$  is the expectation of the recommended arm after round  $T$ , and  $\mu^*$  is the expectation of the best arm. Bubeck, Munos, and Stoltz show that upper bounds on the cumulative regret  $R_T$  lead to lower bounds on  $r_T$ , and vice versa. In this sense, algorithms that minimize the cumulative regret occupy an extreme point of a design space: they maximize the welfare of trial patients, but sacrifice knowledge about inferior treatments. At the other extreme point of the design space, an ideal trial for knowledge generation, with two arms of equal variance, will split the sample sizes equally, consigning half of the patients to the inferior treatment.

Most practical implementations of adaptive randomization in clinical trials use modified bandit algorithms. A common prescription is to lead with a first phase of equal randomization. Or, allocation probabilities may be shrunk toward  $1/K$  in some fashion. Wathen and Thall (2017) discuss the design options of restricting allocations to  $[0.1, 0.9]$ , leading with a period of equal randomization to prevent the algorithm from “getting stuck” on a worse arm, and altering the Thompson sampling to allocate with probability proportional to  $p_{k,t}^c$ , for  $c \in (0, 1]$ . Villar, Bowden and Wason (2015) consider forced sampling of the control arm every  $1/K$  patients. Kasy and Sautmann (2019) modify the Thompson sampling to tamp down selection of the best arm(s), asymptotically leading to equal randomization between the best candidates. Lai, Liao and Kim (2013) give a design that maintains a preferred set of arms, randomizing equally between them, and adaptively drops arms from this set at interim analyses. These various design choices and algorithmic tweaks are typically investigated and tuned by simulation. Even without explicit modification to the standard bandit approach, most medical applications will have a *delay* between the treatment assignment and the observation of an outcome; the resulting reduction in available information leads to more exploration for most algorithms.

There are many benefits to using nearer-to-equal randomization probabilities. First, balancing sample sizes between a pair of arms serves inference goals such as increased power of hypothesis tests, shorter confidence intervals, and more accurate future recommendations. Second, closer-to-equal randomization may improve the information for interim decisions such as early stopping and sample size re-estimation. Third, without tuning, there may be an unacceptably high chance of sending a majority of patients to the wrong arm (Thall, Fox and Wathen (2015)). Fourth, more equal randomization can help detect violated as-

sumptions, such as time trends or a model misspecification. Fifth, the possibility of violated assumptions suggests treating data as slightly less informative. Finally, probabilities nearer  $1/2$  are helpful for inverse-probability weighting and randomization tests.

On the other hand, when a treatment is strongly disfavored for a patient, ethical health care requires setting its randomization chance to zero. This may be achieved by thresholding allocation probabilities according to some rule, or suspending or dropping treatment arms at interim analyses. Furthermore, more equal randomization comes at an opportunity cost to the welfare of trial participants. Practical trial design in an LHS must seek a balance between these competing objectives of knowledge generation and participant welfare.

#### 4. Inference for MABs in an LHS

The LHS may desire several forms of knowledge from an adaptive randomization trial, including confidence intervals for the outcomes of arms (and their differences), guarantees about selecting arms correctly, and recommendations for treatments in non-participating patients.

Frequentist inference under adaptive randomization designs can be challenging. Owing to adaptive sampling, the distribution of standard estimates for the mean of an arm is typically nonGaussian, and not pivotal with respect to the treatment effect. Concentration techniques for UCBs, such as Chernoff bounds, can be applied for confidence bounds that may hold uniformly over possible stopping times (Jamieson and Nowak (2014); Zhao et al. (2016); Karnin, Koren and Somekh (2013)). The concentration approach has been extended to FDR control with the always-valid p-values framework (Johari, Pekelis and Walsh (2015); Yang et al. (2017)). Furthermore, self-normalization techniques from de la Peña, Lai and Shao (2008) permit extensions to large classes of distributions. However, confidence intervals from concentration bounds may be conservative, slack by a constant or logarithmic factor of width.

In confirmatory trial design, adaptivity may be managed by dividing the trial into segments, each having constant randomization probabilities so that Gaussian theory can be used (with numerical integration for stopping boundaries to compute the type-I error and power at fixed alternatives). Lai, Liao and Kim (2013) and Shih and Lavori (2013) show how to do this for their MAB-inspired designs. Alternatively, Korn and Freidlin (2011) suggest block-randomization and block-stratified analysis. Compared to the constantly changing allocation strategies of the standard bandit algorithms, discretization of strategy can come

at a moderate or minimal cost, depending on the design and goals.

For analyzing MAB designs with a constantly updating allocation strategy, a key idea for constructing valid frequentist p-values is the randomization test. The randomization test assumes the sharp null hypothesis that the treatment has exactly zero effect, and relies on probabilistic randomization in the allocation algorithm to generate power. In exchange, with other minimal assumptions, it grants valid p-values, even in the presence of time trends and other confounders in the patient population (Simon and Simon (2011)). To form confidence intervals, a sharp additive model for the treatment effect may be considered. Confidence bounds then follow by inverting the randomization test, as in Ernst (2004).

Another tool for constructing confidence intervals is hybrid resampling, by Lai and Li (2006). This procedure considers families of different shifts and scales of the observed data, and simulates via resampling to infer which distributions are consistent with the observed treatment effects. Lai and Li show that for group sequential trials, confidence intervals from hybrid resampling can have more accurate coverage than that of standard normal approximations.

Hadad et al. (2019) suggest a double-robust estimation approach. In addition to using an augmented inverse-probability weight (AIPW) model, they propose further adaptively re-weighting the data to force the treatment effect estimate into an asymptotically Gaussian distribution. Double-robust estimation may help to correct for time trends or other confounding. However, data re-weighting comes at a cost to efficiency, as pointed out by Tsiatis and Mehta (2003).

Finally, if one assumes a prior and enters the Bayesian framework, posterior inference is a highly flexible approach to analysis. Because Bayes' rule decouples the experimenter's allocation decisions from the rest of the likelihood, the standard Bayesian workflow can be applied to the data without concern for the adaptivity of the design (Berger and Wolpert (1988)). Subject to typical caveats on prior selection and accurate posterior sampling, posterior inference can yield Bayes factors for testing, credible intervals for treatment effects, and decision analysis for treatment recommendations.

## 5. Contextual MABs and Personalized Medicine

For an LHS that continuously seeks to improve and personalize treatment, the important question is not *which* treatment is best, but *for whom* each treatment is best. To address this question, one must augment the bandit model with information about each patient. Calling this side information "covariates" or "*contexts*," one arrives at the CMAB problem.



CMABs have found great success in the internet domain for problems such as serving ads, presenting search results, and testing website features. In contrast, applications in medicine have lagged (with the prominent exception of mobile health (Greenewald et al. (2017); Xia (2018))). The design of trials in an LHS brings new challenges to the CMAB framework, such as ethical requirements, small sample sizes (roughly  $10^2 - 10^4$  patients, in comparison to  $10^4 - 10^9$  clicks for internet applications), requirements for medical professionals to inspect and understand processes, feedback times, and demand for generalizable conclusions.

In the following section, we focus on correctly specified linear models. This assumption derives some justification from the features of an LHS: assuming that covariates are continuous and low dimensional, the patient population of greatest interest is expected to occupy a small region of the covariate domain, owing to the systematic filtering of equipoise requirements and further shrinking of the population under experimental focus as “exploiting” increases. Additionally, the conditional expectation of the response is typically a smooth function of the covariates. Therefore, assuming both smoothness of conditional expectation and locality of the studied population, Taylor’s theorem implies approximate correctness of the linear model. Similar arguments can be applied to logistic models and other smooth model classes.

### 5.1. Linear models for the reward

If at step  $t$  we observe a context vector  $x_t$  of length  $d$ , sample from arm  $\phi_t = k$ , and receive reward  $y_t$ , we may consider the following simple linear model for the expected reward:

$$E[y_t | x_t, \phi_t = k] = x_t^T \theta_k^*,$$

where  $\theta_k^*$  is an unknown parameter vector of length  $d$ . The LinUCB algorithm of Li et al. (2010) brings the UCB of Lai and Robbins (1985) to this linear model. Assuming the linear model parameters are not shared between arms and that contexts do not depend on the arm chosen (see Li et al. (2010) for the general case,) they suggest estimating  $\theta_k^*$  for each arm using a ridge regression  $\hat{\theta}_k$ . That is, if  $X_{k,t}$  is a design matrix whose rows are the contexts of the individuals previously assigned to arm  $k$  before time  $t$  and  $Y_{k,t}$  is a vector of their rewards, the ridge estimator with tuning parameter  $\lambda$  is

$$\hat{\theta}_{k,t} = (X_{k,t}^T X_{k,t} + \lambda I_d)^{-1} X_{k,t}^T Y_{k,t}.$$

Next, Li et al. (2010) construct a UCB for the expected reward around the

ridge regression prediction, suggesting the confidence interval

$$|x_t^T \hat{\theta}_{k,t} - x_t^T \theta_k^*| \leq \alpha \sqrt{x_t^T (X_{k,t}^T X_{k,t} + \lambda I_d)^{-1} x_t}$$

where  $\lambda$  is set to one and  $\alpha$  is a tuning parameter. This confidence interval implicitly assumes a correctly specified linear model and independence of  $Y_{k,t}$  given  $X_{k,t}^T$ , an assumption which is typically broken by the allocation mechanism unless  $(x_t, y_t)$  is independent and identically distributed (i.i.d.) for all  $t$ . Nevertheless, analogously to the basic UCB algorithm, they propose the LinUCB algorithm, which chooses the arm with the highest UCB,

$$\phi_t^{UCB} := \operatorname{argmax}_k \left\{ x_t^T \hat{\theta}_{k,t} + \alpha \sqrt{x_t^T (X_{k,t}^T X_{k,t} + \lambda I_d)^{-1} x_t} \right\}.$$

LinUCB is easy to implement and has proven popular in applications, inspiring further improvements and competitors. Chu et al. (2011) analyze a theoretical fix to LinUCB and give a regret analysis for a modified algorithm of order  $O(\sqrt{Td \ln^3(KT \ln(T)/\delta)})$ . They also give a nearly-matching general lower bound for the problem of order  $\Omega(\sqrt{KT})$ .

Alternatively, Abbasi-Yadkori, Pál and Szepesvári (2011), working within a more general framework called “linear bandits” or “linear stochastic bandits,” construct self-normalized confidence sets for the arm parameters. In the linear bandit, rather than choosing among a discrete set of arms, one chooses the context  $x_t$  from a set  $D_t$ , and the rewards are modeled as  $y_t = x_t^T \theta^* + \eta_t$ . Note that model (5.1) can be embedded within the linear bandit by sufficiently increasing the dimensions of  $x_t$  and  $\theta^*$  and taking  $D_t$  as an appropriate finite set of  $K$  vectors. Abbasi-Yadkori, Pál and Szepesvári (2011) assume that, conditioned on data prior to time  $t$ ,  $\eta_t$  is mean-zero and  $R$ -sub-Gaussian for some  $R \geq 0$ . Further, it is assumed that  $\|\theta^*\|_2 \leq S$ , for some  $S \geq 0$ . Then, defining  $X_t$  as a  $(t-1) \times d$  matrix whose rows consist of the contexts  $x_s^T$ , for  $s = 1, \dots, t-1$ , defining the reward vector  $Y_t$  as a vector of length  $(t-1)$  of the corresponding rewards  $y_s$ , for  $s = 1, \dots, t-1$ , and denoting  $\bar{V}_t := \lambda I_d + X_t^T X_t$ , for all  $t \geq 1$ , one may write the ridge estimator as

$$\hat{\theta}_t := \bar{V}_t^{-1} X_t^T Y_t.$$

Abbasi-Yadkori, Pál and Szepesvári (2011) then derive the confidence set

$$C_t := \left\{ \|\hat{\theta}_t - \theta^*\|_{\bar{V}_t} \leq R \sqrt{2 \log \left( \frac{\det(\bar{V}_t)^{1/2} \det(\lambda I_d)^{-1/2}}{\delta} \right) + \lambda^{1/2} S} \right\}$$

where  $\|\cdot\|_{\bar{V}_t}$  is a matrix weighted 2-norm. The collection of these sets,  $\mathbb{C} := \bigcap_{t \geq 1} C_t$ , provides  $1 - \delta$  uniform confidence that  $\theta^* \in \mathbb{C}$ , regardless of an adaptive mechanism for the context choice. Abbasi-Yadkori, Pál and Szepesvári (2011) leverage this confidence approach into a strategy that generalizes the UCB. They follow the underlying principle of “optimism in the face of uncertainty” to select the context

$$x_t := \operatorname{argmax}_{x \in D_t} \max_{\theta \in C_t} x^T \theta,$$

and prove regret guarantees for the linear bandit with this algorithm. For a  $K$ -arm trial designer, a key takeaway is that uniform confidence sets offer an approach to model inference (noting that practical use requires strong modeling assumptions, a choice of  $\lambda$ , and bounds for the unknown parameters  $R$  and  $S$ ).

A different approach to the CMAB problem is to generalize the  $\epsilon$ -greedy algorithm: periodic exploration can be used to estimate a model, and to verify that estimates based on adaptive data collection are not far off. Under the simple linear model (5.1), Goldenshluger and Zeevi (2013) propose maintaining two sets of linear model estimates:  $\hat{\theta}_k^*$ , estimated on a small amount of equal-randomized data, and  $\tilde{\theta}_k^*$ , based on all of the (adaptively allocated) data. If the estimated rewards from equal randomization  $x_t^T \hat{\theta}_k^*$  are well separated, the arm with the largest estimate is chosen. Else, the arm with the largest value of  $x_t^T \tilde{\theta}_k^*$  is chosen. Under strong assumptions including  $K = 2$  arms, i.i.d samples, and a *margin condition* that ensures that the decision boundary between the arms is sharp, that is,

$$\mathbb{P} \{ |(\theta_1^* - \theta_2^*)^T X_t| \leq \rho \} \leq L\rho, \quad \forall \rho \in (0, \rho_0],$$

they derive a cumulative regret bounded by  $O(d^3 \log T)$ . Bastani and Bayati (2019) improve these bounds and extend this approach to high-dimensional sparse linear models using  $L^1$  penalization. Bastani, Bayati and Khosravi (2020) also show that under certain conditions, a pure greedy approach can yield rate-optimal regret.

## 5.2. More general models for the reward

The Bayesian workflow for the MAB naturally extends to linear models and beyond. Russo and Van Roy (2014b) show that for several classes of well-specified Bayesian problems with contexts, Thompson sampling achieves near-optimal performance and behaves like a problem-adaptive UCB. A variety of competitive risk bounds have been proven for Thompson sampling (Agrawal and Goyal (2012, 2013); Kaufmann, Korda and Munos (2012); Korda, Kaufmann and Munos

(2013)). In empirical studies, Thompson sampling often outperforms competitors by a small margin (Scott (2010); Chapelle and Li (2011); Dimakopoulou et al. (2017)).

An alternative for the nonBayesian is what we call “pseudo-Thompson bootstrapping.” Given a black box algorithm that models the outcomes under each arm, the idea is to bootstrap-resample data to generate variation in the model’s estimates. Pretending that this resampling distribution is a posterior, one can drop the estimated “probabilities” of arm superiority into the Thompson rule and hope to recover its performance advantages. While this technique approximates Thompson sampling for some known cases (Eckles and Kaptein (2014)), its general theoretical properties remain unclear. The main appeal of the approach is to offer a wrapper for popular estimation algorithms for large data sets, including regression trees, random forests, and neural networks (Elmachtoub et al. (2017); Osband et al. (2016)).

Vaswani et al. (2019) propose the RandUCB algorithm, which gives LinUCB nondeterministic allocation probabilities by perturbing the confidence bound randomly in a way that somewhat resembles bootstrapping. For the linear model, RandUCB can be viewed as a generalization of Thompson sampling under a Gaussian model. Vaswani et al. also prove competitive regret guarantees for RandUCB.

Finally, there are nonparametric methods that leverage the smoothness of the expected response. Rigollet and Zeevi (2010) discretize space into buckets, and run MABs on each of them independently. Lu, Pál and Pál (2010) give a contextual bandit that clusters data adaptively and provides guarantees under Lipschitz assumptions. Kim, Lai and Xu (2020) perform a local linear regression and pair it with  $\epsilon$ -greedy randomization and arm elimination, meeting minimax lower bounds on regret under certain regularity conditions.

## 6. Dynamic Treatment Regimes

An LHS bears responsibility for patients over time as their clinical status and treatment needs evolve. Formalizing the notion of a complete care strategy, a DTR is a set of rules that dictates treatment decisions, given a patient’s history of covariates and prior treatment (Lavori and Dawson (2004)).

Thus, if a patient is observed at time-points  $\tau_i$  when observations  $x_i$  are recorded and treatment action  $a_i$  is taken, a DTR is a function that maps  $(\tau_i, x_{1:i}, a_{1:i-1})$  to  $a_i$ .

DTRs may be studied using a SMART, which begins with an initial treatment

randomization and at each subsequent decision point, re-randomizes patients among further treatment options. A SMART culminates in an outcome  $Y$  for each individual (which may be a function of  $(x_1, \dots, x_I)$ ), by which the treatments will be assessed. Lavori and Dawson (2007, 2008) construct confidence intervals for comparing DTRs based on their expected outcomes.

In the group sequential clinical trial setting, Zhong (2018) demonstrates asymptotic multivariate normal approximations to estimated outcomes and transition probabilities. Zhong proposes likelihood-based Wald tests with simultaneous coverage for comparing DTRs, and demonstrates the approach on adaptive play-the-winner designs.

Key challenges in the design and analysis of SMARTs include incorporating patient covariates and handling estimations as the length of the decision tree grows, because the number of treatment strategies and possible patient histories explodes rapidly. Most SMARTS do not consider more than two decision points per patient.

One approach for handling patient covariates is Q-learning (Sutton and Barto (1998); Murphy (2005)). Q-learning seeks to model the patient's expected final outcome, conditional on taking action  $a_i$  and given the history  $(\tau_i, x_{1:i}, a_{1:i-1})$ , and assuming optimal decision making thereafter. This model is thus a function,  $E[Y|\tau_i, x_{1:i}, a_{1:i}] \sim Q(\tau_i, x_{1:i}, a_{1:i})$ , called the *Q-function*. In order to estimate a Q-function, Q-learning alternates between model estimation of expected values of states and backward induction to select optimal actions, using a modified version of Bellman's inequality. Q-learning is therefore an approximate dynamic programming technique. It has been combined with a variety of modeling approaches including linear models (Murphy (2005)), regression trees (Ernst, Geurts and Wehenkel (2005)), and kernels (Ormonoit and Sen (2002)). Chakraborty and Murphy (2014) discuss nonregular asymptotics for Q-learning with the linear model.

## 7. Recent Advances for Learning about DTRs

Intuitions and approaches from MAB and CMAB theory, including Lai and Robbins (1985), are proving fruitful in constructing algorithms with theoretical guarantees and good performance for analyzing DTRs.

In the general DTR setting, Zhang and Bareinboim (2019) use the UCB approach to motivate a reinforcement learning algorithm and derive regret guarantees. Following the techniques of Lai and Robbins (1985) and Auer, Jaksch and Ortner (2009), at each time  $t$ , they construct a uniform confidence set  $\mathbb{M}_t$

with two-sided bounds on the final payouts and transitions, and then use the Bellman equation recursively to find an optimal DTR in  $\mathbb{M}_t$ . Note that this approach also permits confidence bands for the values of individual DTRs. Zhang and Bareinboim derive regret guarantees, and further show that weak evidence from observational data collection can be used to narrow the range of possible transitions, thus narrowing  $\mathbb{M}_t$  and improving performance.

Hu and Kallus (2020) analyze a two-stage DTR model. Assuming a linear model for Q-functions, they extend the contextual bandit approaches of Golden-shluger and Zeevi (2013) and Bastani, Bayati and Khosravi (2020) to the two-stage two-treatment DTR setting, using a combination of unbiased estimates from a small sample and biased estimates from the full sample. They derive regret bounds under several margin conditions on the Q-functions, notably showing under a sharp margin condition a regret bound of order  $O(d(\log d)^{2/3} \log T + (d \log d)^2)$ . They demonstrate that their bounds have optimal dependence on  $T$  by applying lower bounds from contextual bandits.

Wang and Powell (2016) demonstrate an important connection between DTRs and contextual bandits in a Bayesian framework. They model binary outcomes using Bayesian generalized linear models and handle posterior computations using Laplace approximations. With quick recursive computation of the value function, they show how to collapse the DTR problem into a CMAB problem, where each decision point becomes a bandit sample, and payoffs are given by the change in the posterior expected value. This formulation naturally permits them to use Bayesian CMAB algorithms, including the knowledge gradient, Thompson sampling, and greedy Bayes algorithms, for learning and executing DTRs.

## 8. Conclusion

In this paper, we have reviewed the current literature on the use of adaptive randomization, as represented by the CMAB, as a natural model for an LHS that seeks to add experimental strength to its portfolio of learning methodologies. Lai's contributions to the theory and practice of bandits goes back almost 40 years, and his basic insights have led to innovations in diverse areas, such as finance, internet commerce, and medical drug development.

This past year, the world has been engulfed in a global pandemic, and many of the ideas reviewed here have taken on much greater significance. Governments are learning how to control the outbreak using a combination of interventions, including masks and "social distancing," rapid learning by intensivists faced with a virus having pleiotropic clinical effects, development of new therapeutics and

the repurposing of existing drugs, and ambitious programs of vaccine development. They are also dealing with the economic disruption, both directly from the pandemic and indirectly from the control efforts.

The limitations of observational, non-experimental approaches and conventional randomized clinical trials have been cast in sharp relief. Each scientific specialty has begun to propose ways to make its responses better and faster “next time around.” Virologists propose beginning therapeutic drug and vaccine development, even in advance of knowing the identity of the new pandemic agent. Ecologists and wildlife conservationists urge a greatly expanded global project to survey likely animal sources of the next spillover event, and to target those agents that are likely to pose a substantial global threat. Clinical scientists and trialists seek to create pre-formed platforms for rapid testing of the drugs, non-pharmacologic interventions, and vaccines that will be proposed. In this area, innovative experimental design will be critical, and as statisticians are recruited to help prepare for the next emergency, they will find, as we have, that the work of Tze Leung Lai will provide a sturdy basis and flexible, but reliable framework for their efforts.

## Acknowledgments

Michael Sklar acknowledges the support of the National Science Foundation under DMS-1811818.

## References

- Abbasi-Yadkori, Y., Pál, D. and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems 24* (Edited by J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira and K. Q. Weinberger), 2312–2320. Curran Associates, Inc.
- Agrawal, S. and Goyal, N. (2012). Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, 39–1.
- Agrawal, S. and Goyal, N. (2013). Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning* **28**, 127–135.
- Auer, P., Cesa-Bianchi, N. and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning* **47**, 235–256.
- Auer, P., Jaksch, T. and Ortner, R. (2009). Near-optimal regret bounds for reinforcement learning. In *Advances in Neural Information Processing Systems*, 89–96.
- Bartroff, J., Lai, T. L. and Shih, M. (2012). *Sequential Experimentation in Clinical Trials: Design and Analysis*. Springer Science & Business Media.
- Bastani, H. and Bayati, M. (2019). Online decision making with high-dimensional covariates. *Operations Research*.
- Bastani, H., Bayati, M. and Khosravi, K. (2020). Mostly exploration-free algorithms for contex-

- tual bandits. *Management Science*.
- Berger, J. O. and Wolpert, R. L. (1988). *The Likelihood Principle*. IMS.
- Bouneffouf, D. and Rish, I. (2019). A survey on practical applications of multi-armed and contextual bandits. *arXiv preprint arXiv:1904.10040*.
- Bubeck, S., Munos, R. and Stoltz, G. (2011). Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science* **412**, 1832–1852.
- Chakraborty, B. and Murphy, S. (2014). Dynamic treatment regimes. *Annual Review of Statistics and its Application* **1**, 447–464.
- Chamberlayne, R., Green, B., Barer, M. L., Hertzman, C., Lawrence, W. J. and Sheps, S. B. (1998). Creating a population-based linked health database: A new resource for health services research. *Canadian Journal of Public Health* **89**, 270–273.
- Chapelle, O. and Li, L. (2011). An empirical evaluation of thompson sampling. In *Advances in Neural Information Processing Systems*, 2249–2257.
- Cheng, Y. and Berry, D. (2007). Optimal adaptive randomized designs for clinical trials. *Biometrika* **94**, 673–689.
- Chu, W., Li, L., Reyzin, L. and Schapire, R. (2011). Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 208–214.
- Diaconis, P. (2009). The markov chain monte carlo revolution. *Bulletin of the American Mathematical Society* **46**, 179–205.
- Dimakopoulou, M., Zhou, Z., Athey, S. and Imbens, G. (2017). Estimation considerations in contextual bandits. *arXiv preprint arXiv:1711.07077*.
- Dwivedi, R., Chen, Y., Wainwright, M. and Yu, B. (2018). Log-concave sampling: Metropolis-hastings algorithms are fast! In *Conference on Learning Theory*, 793–797.
- Eckles, D. and Kaptein, M. (2014). Thompson sampling with the online bootstrap. *arXiv preprint arXiv:1410.4009*.
- Elmachtoub, A., McNellis, R., Oh, S. and Petrik, M. (2017). A practical method for solving contextual bandit problems using decision trees. *arXiv preprint arXiv:1706.04687*.
- Ernst, D., Geurts, P. and Wehenkel, L. (2005). Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research* **6**, 503–556.
- Ernst, M. (2004). Permutation methods: A basis for exact inference. *Statistical Science* **19**, 676–685.
- Goldenshluger, A. and Zeevi, A. (2013). A linear response bandit problem. *Stochastic Systems* **3**, 230–261.
- Greenewald, K., Tewari, A., Murphy, S. and Klasnja, P. (2017). Action centered contextual bandits. In *Advances in Neural Information Processing Systems* **30**, 5973–5981.
- Hadad, V., Hirshberg, D., Zhan, R., Wager, S. and Athey, S. (2019). Confidence intervals for policy evaluation in adaptive experiments. *arXiv preprint arXiv:1911.02768*.
- Hu, Y. and Kallus, N. (2020). Dtr bandit: Learning to make response-adaptive decisions with low regret. *arXiv preprint arXiv:2005.02791*.
- Jamieson, K. and Nowak, R. (2014). Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In *2014 48th Annual Conference on Information Sciences and Systems (CISS)*, 1–6. IEEE.
- Johari, R., Pekelis, L. and Walsh, D. J. (2015). Always valid inference: Bringing sequential analysis to a/b testing. *arXiv preprint arXiv:1512.04922*.



- Karnin, Z., Koren, T. and Somekh, O. (2013). Almost optimal exploration in multi-armed bandits. In *International Conference on Machine Learning*, 1238–1246.
- Kasy, M. and Sautmann, A. (2019). Adaptive treatment assignment in experiments for policy choice. *Econometrica* **89**, 113–132.
- Kaufmann, E., Korda, N. and Munos, R. (2012). Thompson sampling: An asymptotically optimal finite-time analysis. In *International Conference on Algorithmic Learning Theory*, 199–213. Springer.
- Kim, D., Lai, T. L. and Xu, H. (2020). Multi-armed bandits with covariates: Theory and applications.
- Korda, N., Kaufmann, E. and Munos, R. (2013). Thompson sampling for 1-dimensional exponential family bandits. In *Advances in Neural Information Processing Systems*, 1448–1456.
- Korn, E. and Freidlin, B. (2011). Outcome-adaptive randomization: is it useful? *Journal of Clinical Oncology* **29**, 771–776.
- Lai, T., Liao, O. and Kim, D. (2013). Group sequential designs for developing and testing biomarker-guided personalized therapies in comparative effectiveness research. *Contemporary Clinical Trials* **36**, 651–663.
- Lai, T. and Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics* **6**, 4–22.
- Lai, T. L., Lavori, P. and Liao, O. (2014). Adaptive choice of patient subgroup for comparing two treatments. *Contemporary Clinical Trials* **39**, 191–200.
- Lai, T. L. and Li, W. (2006). Confidence intervals in group sequential trials with random group sizes and applications to survival analysis. *Biometrika* **93**, 641–654.
- Lai, T. L. and Liao, O. (2012). Efficient adaptive randomization and stopping rules in multi-arm clinical trials for testing a new treatment. *Sequential Analysis* **31**, 441–457.
- Lai, T. L. and Shih, M. (2004). Power, sample size and adaptation considerations in the design of group sequential clinical trials. *Biometrika* **91**, 507–528.
- Lavori, P. and Dawson, R. (2000). A design for testing clinical strategies: Biased adaptive within-subject randomization. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **163**, 29–38.
- Lavori, P. and Dawson, R. (2004). Dynamic treatment regimes: Practical design considerations. *Clinical Trials* **1**, 9–20.
- Lavori, P. and Dawson, R. (2007). Improving the efficiency of estimation in randomized trials of adaptive treatment strategies. *Clinical Trials* **4**, 297–308.
- Lavori, P. and Dawson, R. (2008). Adaptive treatment strategies in chronic disease. *Annu. Rev. Med.* **59**, 443–453.
- Li, L., Chu, W., Langford, J. and Schapire, R. (2010). A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th International Conference on World Wide Web*, 661–670.
- Lu, T., Pál, D. and Pál, M. (2010). Contextual multi-armed bandits. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 485–492.
- Murphy, S. A. (2005). An experimental design for the development of adaptive treatment strategies. *Statistics in Medicine* **24**, 1455–1481.
- Olsen, L., Aisner, D. and McGinnis, J. M. (2007). Institute of medicine roundtable on evidence-based medicine: The learning healthcare system. In *Workshop Summary*.
- Ormonet, D. and Sen, S. (2002). Kernel-based reinforcement learning. *Machine Learning* **49**, 161–178.

- Osband, I., Blundell, C., Pritzel, A. and Van Roy, B. (2016). Deep exploration via bootstrapped DQN. In *Advances in Neural Information Processing Systems*, 4026–4034.
- Peña, V. H., Lai, T. L. and Shao, Q.-M. (2008). *Self-normalized Processes: Limit Theory and Statistical Applications*. Springer Science & Business Media.
- Rigollet, P. and Zeevi, A. (2010). Nonparametric bandits with covariates. In *23rd COLT*, 54–66.
- Roy, V. (2020). Convergence diagnostics for markov chain monte carlo. *Annual Review of Statistics and its Application* **7**, 387–412.
- Russo, D. and Van Roy, B. (2014a). Learning to optimize via information-directed sampling. In *Advances in Neural Information Processing Systems*, 1583–1591.
- Russo, D. and Van Roy, B. (2014b). Learning to optimize via posterior sampling. *Mathematics of Operations Research* **39**, 1221–1243.
- Russo, D. and Van Roy, B. (2016). An information-theoretic analysis of thompson sampling. *The Journal of Machine Learning Research* **17**, 2442–2471.
- Ryzhov, I., Powell, W. and Frazier, P. (2012). The knowledge gradient algorithm for a general class of online learning problems. *Operations Research* **60**, 180–195.
- Scott, S. (2010). A modern bayesian look at the multi-armed bandit. *Applied Stochastic Models in Business and Industry* **26**, 639–658.
- Shih, M. and Lavori, P. (2013). Sequential methods for comparative effectiveness experiments: Point of care clinical trials. *Statistica Sinica* **23**, 1775–1791.
- Simon, R. and Simon, N. (2011). Using randomization tests to preserve type I error with response adaptive and covariate adaptive randomization. *Statistics & Probability Letters* **81**, 767–772.
- Sutton, R. S. and Barto, A. G. (1998). *Introduction to Reinforcement Learning*. MIT Press Cambridge.
- Thall, P., Fox, P. and Wathen, J. (2015). Statistical controversies in clinical research: scientific and ethical problems with adaptive randomization in comparative clinical trials. *Annals of Oncology* **26**, 1621–1628.
- Thompson, W. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* **25**, 285–294.
- Tsiatis, A. A. and Mehta, C. (2003). On the inefficiency of the adaptive design for monitoring clinical trials. *Biometrika* **90**, 367–378.
- Vaswani, S., Mehrabian, A., Durand, A. and Kveton, B. (2019). Old dog learns new tricks: Randomized ucb for bandit problems. *arXiv preprint arXiv:1910.04928*.
- Villar, S., Bowden, J. and Wason, J. (2015). Multi-armed bandit models for the optimal design of clinical trials: benefits and challenges. *Statistical Science: A Review Journal of the Institute of Mathematical Statistics* **30**, 199–215.
- Wang, Y. and Powell, W. (2016). An optimal learning method for developing personalized treatment regimes. *arXiv preprint arXiv:1607.01462*.
- Wathen, J. and Thall, P. (2017). A simulation study of outcome adaptive randomization in multi-arm clinical trials. *Clinical Trials* **14**, 432–440.
- Whittle, P. (1979). Discussion of dr gittins’ paper. *Journal of the Royal Statistical Society: Series B (Methodological)* **41**, 164–177.
- Whittle, P. (1980). Multi-armed bandits and the gittins index. *Journal of the Royal Statistical Society: Series B (Methodological)* **42**, 143–149.
- Xia, I. (2018). *The Price of Personalization: An Application of Contextual Bandits to Mobile*

- Health*. Ph.D. thesis. Harvard University.
- Yang, F., Ramdas, A., Jamieson, K. and Wainwright, M. (2017). A framework for multi-a (rmed)/b (andit) testing with online fdr control. In *Advances in Neural Information Processing Systems*, 5957–5966.
- Zhang, J. and Bareinboim, E. (2019). Near-optimal reinforcement learning in dynamic treatment regimes. In *Advances in Neural Information Processing Systems*, 13401–13411.
- Zhao, S., Zhou, E., Sabharwal, A. and Ermon, S. (2016). Adaptive concentration inequalities for sequential decision problems. In *Advances in Neural Information Processing Systems*, 1343–1351.
- Zhong, X. (2018). *Design and Analysis of Sequential Multiple Assignment Randomized Trial for Comparing Multiple Adaptive Interventions*. Ph.D. thesis. Mailman School of Public Health, Columbia University.

Michael Sklar

Department of Statistics, Stanford University, Stanford, CA 94305, USA.

E-mail: sklarm@stanford.edu

Mei-Chiung Shih

VA Palo Alto Cooperative Studies Program Coordinating Center, Mountain View, CA 94043, USA.

E-mail: mei-chiung.shih@va.gov

Philip Lavori

Department of Statistics, Stanford University, Stanford, CA 94305, USA.

E-mail: lavori@stanford.edu

(Received October 2020; accepted January 2021)