
Learning for Dose Allocation in Adaptive Clinical Trials with Safety Constraints

Cong Shen¹ Zhiyang Wang² Sofia S. Villar³ Mihaela van der Schaar^{3,4}

Abstract

Phase I dose-finding trials are increasingly challenging as the relationship between efficacy and toxicity of new compounds (or combination of them) becomes more complex. Despite this, most commonly used methods in practice focus on identifying a Maximum Tolerated Dose (MTD) by learning only from toxicity events. We present a novel adaptive clinical trial methodology, called Safe Efficacy Exploration Dose Allocation (SEEDA), that aims at maximizing the cumulative efficacies while satisfying the toxicity safety constraint with high probability. We evaluate performance objectives that have operational meanings in practical clinical trials, including cumulative efficacy, recommendation/allocation success probabilities, toxicity violation probability, and sample efficiency. An extended SEEDA-Plateau algorithm that is tailored for the increase-then-plateau efficacy behavior of molecularly targeted agents (MTA) is also presented. Through numerical experiments using both synthetic and real-world datasets, we show that SEEDA outperforms state-of-the-art clinical trial designs by finding the optimal dose with higher success rate and fewer patients.

1. Introduction

An adaptive clinical trial utilizes the accumulated results to dynamically modify its future trajectory for better efficiency and ethics, while preserving the integrity and validity of the study. Studies such as the phase I trial in Acute Myeloid Leukaemia in (Yap et al., 2013) and Cancer Research UK study CR0720-11 in (Whitehead et al., 2012) have suggested

that even some simple forms of adaptive design lead to better usage of resources and require fewer participants. These promising results have spawned the interest in developing adaptive clinical trial methodologies in recent years (Villar et al., 2015a; Pallmann et al., 2018; Atan et al., 2019; Lee et al., 2020), which is of great importance because running an actual clinical trial on human subjects is expensive and ethically sensitive. A well-designed trial methodology with thorough theoretical and simulated investigation is widely acknowledged as a crucial first step.

Traditionally, the goal of phase I clinical trials is to identify the Maximum Tolerated Dose (MTD) of a cytotoxic (CTX) or therapeutic agent, which is then used for subsequent studies (Storer, 1989). However, modern cancer phase I trials test antineoplastic agents in patients with advanced cancer stages, who have often exhausted all other available treatment options (Roberts et al., 2004). These participants usually expect therapeutic benefit from participating in the trial, which has motivated the trial design to *include efficacy as a co-primary end point of phase I dose-finding studies* (Yan et al., 2017; Paoletti & Postel-Vinay, 2018). In addition, the monotonic assumption for the dose-efficacy relationship is widely adopted in state of the art designs, which is reasonable for cytotoxic agents but may not apply to the new molecularly targeted agents (MTA) such as monoclonal antibodies (see (Postel-Vinay et al., 2009) for an exemplary trial that illustrates this issue). Designing adaptive clinical trials that can properly address the intrinsic conflict between learning and treatment effectiveness for general dose-response models has become an important task for phase I clinical trials.

In addition to the well-known 3+3 design (Storer, 1989) and continual reassessment method (CRM) (O’Quigley et al., 1990) (and its many variants), Bayesian approaches such as Thompson Sampling (TS) (Aziz et al., 2019) and Gittins index (Villar et al., 2015a;b) have been proposed in the literature for dose-finding studies. However, these methods were originally designed for simplified models that do not capture some of the unique characteristics of clinical trials, often leading to lack of randomization (Villar et al., 2015b), inefficient use of side information (Villar & Rosenberger, 2018), and reduced power levels and estimation issues. Notably, for cases where the best dose for combination therapies

¹University of Virginia, USA ²University of Pennsylvania, USA ³University of Cambridge, United Kingdom ⁴University of California, Los Angeles, USA. Correspondence to: Cong Shen <cong@virginia.edu>.

Table 1: Representative adaptive clinical trial studies

Study	Treatment	Category	Methodology	Evaluation
(Tighiouart et al., 2014)	Veliparib	CTX	EWOC-PH	simulated trial
(Whitehead et al., 2012)	MK-0752	CTX	joint phase I and II design	simulated trial
(Lee et al., 2017)	Erlotinib	MTA	extended TITE-CRM	simulated trial
(Thiessen et al., 2010)	Lapatinib	MTA	escalation to DLT	real-world trial data

is to be found, unknown synergistic/antagonist effects are likely to exist and naive designs will fail to identify them. For MTA, the existence of a plateau of efficacy has been discussed in (Zang et al., 2014) and (Riviere et al., 2018), which indicates that the toxicity constraint must be jointly studied with the dose-efficacy relationship for certain new compounds. This is also confirmed by the real-world trial result; see (Tighiouart et al., 2014). Last but not the least, safety constraints such as minimizing the adverse events (AE) (Petroni et al., 2017) have not been properly evaluated with theoretical guarantees. Table 1 summarizes some representative studies in this direction.

In this paper, we address these challenges by developing new dose-finding methods that explicitly impose safety constraints to the allocation and recommendation of dose levels in a phase I clinical trial. Through the lens of multi-armed bandits (MAB), we propose the *Safe Efficacy Exploration Dose Allocation (SEEDA)* algorithm that adaptively updates the admissible set of dose levels satisfying the safety constraints, thus limiting the exploration of doses with harmful effect. Performance analysis for SEEDA is carried out with respect to several measures that have operational meanings in clinical trials, including the probability of safety constraints violation, the average efficacy for patients, and the recommendation and allocation probabilities. Noting that SEEDA only leverages the dose-toxicity logistic model and makes no assumptions on the efficacy, we then show that, by considering the increasing-then-plateau feature of the dose-efficacy relationship for MTA, *SEEDA-Plateau* leads to better performance by leveraging the unimodal structure. Experiments on simulated datasets as well as clinical trials built from real-world datasets show that the proposed methods are capable of finding the optimal dose with higher success rate and fewer patients in most cases, compared to other state-of-the-art designs.

2. Model and problem formulation

2.1. The dose-finding model

In a phase I dose-finding clinical trial, a total of K doses are given where the k -th dose is denoted as $d_k \in \mathcal{D}$, $k \in \mathcal{K} = \{1, 2, \dots, K\}$. The performance is characterized by both *efficacy* and *toxicity*. We model the efficacy X and toxicity Y for dose d_k as Bernoulli random variables with unknown probabilities q_k and p_k , respectively, where $X = 1$ ($X = 0$)

indicates that the dose level is effective (not effective), and $Y = 1$ ($Y = 0$) suggests that the dose is harmful (not harmful) to the patient¹.

We consider adaptive clinical trials where information learned from previous trial patients can be used in allocating doses to subsequent patients (Atan et al., 2019; Villar et al., 2015a; Aziz et al., 2019). For the t -th patient, dose $I(t)$ is selected based on a policy that uses past observations, and administrated to the patient. The efficacy outcome X_t and toxicity response Y_t are realized based on their distributions $X_t \sim \text{Ber}(q_{I(t)})$ and $Y_t \sim \text{Ber}(p_{I(t)})$, and observed by the trialist.

We adopt a well-known dose-toxicity logistic model proposed by in (O’Quigley et al., 1990) to describe the toxicity probability for different dose levels:

$$p_k(a) = \left(\frac{\tanh d_k + 1}{2} \right)^a, \quad (1)$$

where a is a global parameter for all the dose levels. It can be verified that Eqn. (1) satisfies the assumption that the toxicity monotonically increases with dose d_k . The unsafe dose levels are defined as those whose toxicity probabilities p_k ’s are above a pre-determined target toxicity probability θ , which is referred as the MTD threshold. Hence the toxicities of all doses can be written as $p_1 \leq p_2 \leq \dots \leq p_M < \theta < p_{M+1} \leq \dots \leq p_K$ where the (unknown) M denotes the number of safe doses. The efficacy-dose relationship is not modeled to allow for the development of a general algorithm. The specific increase-then-plateau efficacy behavior of MTA will be exploited in Section 4.

2.2. Problem formulation

Several objectives are often desired for a successful dose-finding study, which are summarized as follows.

- **Successful recommendation.** At the end of the trial (n patients) a *dose recommendation* \hat{k}_n is made, which is desired to match the optimal dose k^* that is the lowest safe dose that achieves the highest efficacy (Zang et al., 2014): $k^* = \min\{k : q_k = \max_{l: l \in \mathcal{K}, p_l \leq \theta} q_l\}$.
- **Effective treatment.** The cumulative treatment for trial participants $\sum_{i=1}^n X_t$ is desired to be maximized.

¹This is typically measured by the presence of absence of a dose-limiting toxicity (DLT) reported in a fixed evaluation window after administrating the drug.

- **Minimal violation of the safety constraint.** There are different formulations for the safety constraint. One is to minimize $\mathbb{E}[\sum_{k \in \mathcal{K}, p_k > \theta} N_k(n)/n]$ where $N_k(t)$ denotes the number of times dose k is allocated to the first t patients. Another formulation is to minimize the probability that the average toxicity exceeds the MTD threshold.
- **Small sample size.** Most phase I trials have a pre-determined n which is decided as the minimum number of trial participants to achieve a pre-defined confidence level of successful recommendation. It is desirable to have a small n for cost and efficiency considerations.

Proposing a learning model that explicitly guarantees all of the above objectives is elusive and non-constructive in developing the dose-allocation policy. We thus formulate dose-finding clinical trials as an *online efficacy learning problem with explicit safety constraint*, and subsequently provide performance analysis on the metrics of interest. Specifically, we aim at maximizing the cumulative efficacy over a finite number of patients n while simultaneously guaranteeing that the average toxicity observed from the n dose allocations is kept under the probability threshold θ with high probability. This can be written as:

$$\begin{aligned} & \text{maximize} && \mathbb{E} \left[\sum_{t=1}^n X_t \right] \\ & \text{subject to} && \mathbb{P} \left[\frac{1}{n} \sum_{t=1}^n Y_t \leq \theta \right] \geq 1 - \delta. \end{aligned} \quad (2)$$

Essentially, problem formulation (2) focuses on safe exploration among all the dose levels to maximize cumulative efficacies. Clinical trial designs for (2) thus need to pursue both objectives of toxicity and efficacy.

3. The SEEDA algorithm

3.1. Algorithm description

The proposed Safe Efficacy Exploration Dose Allocation (SEEDA) design is completely described in Algorithm 1. In particular, $\hat{p}_k(t)$ and $\hat{q}_k(t)$ are the estimated toxicity and efficacy, respectively, after administrating the t -th patient. The principle of dose selection is to first dynamically construct the admissible set $\mathcal{D}_1(t)$ using the Upper Confidence Bound (UCB) principle (Auer et al., 2002), where the confidence interval $\alpha(t)$ is constructed as

$$\alpha(t) = \bar{C}_1 K \left(\frac{\log \frac{2K}{\delta}}{2t} \right)^{\frac{\bar{\gamma}_1}{2}}, \quad (3)$$

where \bar{C}_1 and $\bar{\gamma}_1$ are algorithm parameters². Note that the admissible set consists of doses that, with high confidence, satisfy the toxicity constraint.

²See Section B in the supplementary material for a discussion on how to select these algorithm parameters.

Then, limiting to those in the admissible set $\mathcal{D}_1(t)$, the algorithm again applies the UCB principle (UCB-1 from (Auer et al., 2002)) to select a dose with the largest $F(p, s, n)$ for the efficacy estimate:

$$F(p, s, n) = p + \sqrt{\frac{c \log(n)}{s}}, \quad (4)$$

with c denoting the UCB-1 coefficient. It should be noted that (4) can be replaced by other UCB principles, e.g., KL-UCB (Garivier & Cappè, 2011).

Algorithm 1 The Safe Efficacy Exploration Dose Allocation (SEEDA) Algorithm

Input: $p_k(a)$ for each $k \in \mathcal{K}$; MTD threshold θ ; total number of patients n .
Initialize: $N_k(1) = 0, \hat{p}_k(1) = 0, \hat{q}_k(1) = 0, \forall k \in \mathcal{K}$; Sample each dose once and set: $I(t) = t, \hat{q}_{I(t)}(K) = X_t, \hat{p}_{I(t)}(K) = Y_t, N_{I(t)}(K) = 1$, for $t = 1$ to K ; $t = K + 1$.
 1: **while** $t \leq n$ **do**
 2: Compute the estimated parameter: $\hat{a}(t) = \sum_{k=1}^K w_k(t-1) \hat{a}_k(t-1)$;
 3: Set the admissible set: $\mathcal{D}_1(t) = \{d_k \in \mathcal{D} : p_k(\hat{a}(t) + \alpha(t)) \leq \theta\}$;
 4: Select dose: $I(t) = \arg \max_{d_k \in \mathcal{D}_1(t)} F(\hat{q}_k(t), N_k(t), t)$;
 5: Observe the revealed outcomes X_t and Y_t ;
 6: Update estimations: $\hat{q}_{I(t)}(t) = \frac{\hat{q}_{I(t)}(t-1)N_{I(t)}(t-1) + X_t}{N_{I(t)}(t-1) + 1}$, $\hat{p}_{I(t)}(t) = \frac{\hat{p}_{I(t)}(t-1)N_{I(t)}(t-1) + Y_t}{N_{I(t)}(t-1) + 1}$, $N_{I(t)}(t) = N_{I(t)}(t-1) + 1$;
 7: Update parameter estimation: $\hat{a}_{I(t)}(t) = \arg \min_{a \in \mathcal{A}} |p_{I(t)}(a) - \hat{p}_{I(t)}(t)|$;
 8: Update weights: $w_k(t) = N_k(t)/t, \forall d_k \in \mathcal{D}$;
 9: $t = t + 1$.
 10: **end while**
Output: $\hat{d}(n) = \arg \max_{d_k: p_k(\hat{a}(n)) \leq \theta} p_k(\hat{a}(n))$.

3.2. Performance analysis

The SEEDA algorithm is developed with the aim to solve problem (2). It is thus important to analyze (a) whether the cumulative efficacy is maximized, and (b) how often the toxicity constraint is violated. For metric (a), it can be equivalently formulated as regret minimization, i.e., the cumulative efficacy difference between the oracle policy with full information and that of the learning algorithm. Formally, the efficacy regret is defined as

$$R(n) = q^* n - \mathbb{E} \left[\sum_{t=1}^n q_{I(t)} \right], \quad (5)$$

where $q^* = q_{k^*}$ denotes the efficacy associated with the optimal dose defined in Section 2.2, and a^* denotes the true

parameter in (1). As for metric (b), we need to evaluate

$$e(n) = \mathbb{P} \left[\frac{1}{n} \sum_{t=1}^n p_{I(t)}(a^*) > \theta \right],$$

in conjunction with (5), i.e., whether the proposed *SEEDA* algorithm minimizes $R(n)$ and satisfies $e(n) \leq \delta$ at the same time. In addition, other performance measures such as successful recommendation probability and sample efficiency are of practical interest, and we provide theoretical guarantees for them as well. Due to space limitations, all proofs are provided in the supplementary material.

3.2.1. CUMULATIVE EFFICACY

We start the theoretical analysis by showing that for each patient t in *SEEDA*, the dose levels whose toxicities are below the MTD threshold are included in the admissible set with high probability. This corresponds to the *type I* error event that is of interest in clinical trials.

Lemma 1 $\mathbb{P}[p_k(\hat{a}(t) + \alpha(t)) > \theta] \leq \delta, \forall p_k(a^*) \leq \theta$.

Next we prove that with sufficient patients, the dose levels exceeding the toxicity threshold are excluded from the admissible set with high probability. This corresponds to the *type II* error event in clinical trials.

Lemma 2 If $t > t_1 = \frac{1}{2} \left(\frac{\bar{C}_1 K}{|\Delta - \epsilon|} \right)^{\frac{2}{\gamma_1}} \log \frac{2K}{\delta}$, $\Delta = \min_{k \in \mathcal{K}} \Delta_k$, where $\Delta_k = |a^* - p_k^{-1}(\theta)|$ represents the gap between a^* and the parameter when the toxicity is at θ , then:

$$\mathbb{P}[p_k(\hat{a}(t) + \alpha(t)) \leq \theta] \leq \exp(-2t\epsilon^2), \forall p_k(a^*) > \theta. \quad (6)$$

Combining Lemmas 1 and 2 leads to the main result on cumulative efficacy regret.

Theorem 1 With t_1 defined in Lemma 2, the regret of *SEEDA* can be upper bounded as:

$$R(n) \leq \sum_{d_k: p_k(a^*) \leq \theta} \frac{c \log(n)}{q^* - q_k} + \left(n\delta Q + \frac{1}{2}t_1 + \frac{K - M}{2\epsilon^2} \right) \quad (7)$$

where $Q = \max_{i \in \mathcal{K}} |q_i - q_k^*|$ denotes the maximal single-step regret, and $\epsilon > 0$ is a constant. Furthermore, if $\delta = O(\frac{1}{n})$, we have that $R(n) \leq O(\log n)$.

Theorem 1 indicates that the efficacy regret is bounded by $O(\log n)$. A closer look at this scaling reveals that it consists of two parts. The first is due to the structureless model for efficacy – we impose no assumption on the efficacy of different dose levels. The second part, which is reflected through

t_1 , is determined by the structured model for toxicity, which affects the admissible set. As will be shown in Section 4, with the increase-then-plateau efficacy assumption, the first $\log n$ component can be further improved.

3.2.2. SAFETY CONSTRAINT VIOLATION

We now move on to analyzing the safety constraint violation. The first result is to verify whether the *SEEDA* algorithm indeed satisfies the safety constraint in problem (2).

Theorem 2 For any given n , the average toxicity observed from the *SEEDA* algorithm satisfies

$$\mathbb{P} \left[\frac{1}{n} \sum_{t=1}^n p_{I(t)} - \theta \leq C_2 \epsilon^{\gamma_2} \right] \geq 1 - \delta,$$

for an arbitrary $\epsilon > 0$. C_2 and γ_2 are problem-dependent parameters defined in Section A of the supplementary material.

The safety constraint in problem (2) is formulated based on the average toxicity exceeding the MTD threshold. In practice, we are often interested in minimizing the number of patients that have been exposed to unsafe dose levels, $\mathbb{E}[\sum_{k \in \mathcal{K}, p_k > \theta} N_k(n)/n]$. Corollary 1 analyzes this metric.

Corollary 1 The number of unsafe dose allocations from *SEEDA*, i.e., the selected dose levels exceed the MTD threshold, can be bounded as:

$$\mathbb{E} \left[\sum_{d_k: p_k > \theta} N_k(n) \right] \leq t_1 + \frac{K - M}{2\epsilon^2}.$$

Interestingly, Corollary 1 indicates that unsafe dose allocations in *SEEDA* are upper bounded by a constant, which is linear in the number of unsafe doses $K - M$ regardless of the number of trial participants n .

3.2.3. RECOMMENDATION ACCURACY

Finally, we analyze the recommendation accuracy of *SEEDA* at the end of the n -th dose allocation.

Corollary 2 The probability that *SEEDA* recommends the MTD satisfies:

$$\mathbb{P} \left[\hat{d}(n) = \arg \max_{d_k: p_k \leq \theta} p_k \right] \geq 1 - \delta_1, \quad (8)$$

where $\delta_1 = 2K \exp \left(-2 \left(\frac{\Delta_M}{\bar{C}_1 K} \right)^{2\gamma_1} n \right)$.

Corollary 2 guarantees the finding of the MTD with high probability. The recommendation error rate decays *exponentially* with the number of trial participants, which is a

nice property. It is worth noting that a lower bound of the minimal number of trial participants for a given accuracy requirement can be inferred from the upper bound of recommendation error rate (8). This is a practically important result, as sample efficiency directly relates to the cost and ethical constraints of a trial. This is further illustrated in the numerical experiments in Section 5.1.3.

4. Extension to the increase-then-plateau efficacy model

Algorithm 2 The SEEDA-Plateau Algorithm

Input: $p_k(a)$ for each $k \in \mathcal{K}$; MTD threshold θ ; total number of patients n .

Initialize: $N_k(1) = 0, \hat{p}_k(1) = 0, \hat{q}_k(1) = 0, \forall k \in \mathcal{K}$; $L(1) = K; \eta = 2; l_k = 0, \forall k \in \mathcal{K}$; Sample each dose once and set: $I(t) = t, \hat{q}_{I(t)}(K) = X_t, \hat{p}_{I(t)}(K) = Y_t, N_{I(t)}(K) = 1$, for $t = 1$ to K ; $t = K + 1$.

- 1: **while** $t \leq n$ **do**
- 2: Compute the estimated parameter: $\hat{a}(t) = \sum_{k=1}^K w_k(t-1)\hat{a}_k(t-1)$;
- 3: Set the admissible set: $\mathcal{D}_1(t) = \{d_k \in \mathcal{D} : p_k(\hat{a}(t) + \alpha(t)) \leq \theta\}$;
- 4: Set $L(t) = \arg \max_{d_k \in \mathcal{D}_1(t)} \hat{q}_k(t)$ and increase $l_{L(t)}$ by 1;
- 5: If $\frac{l_{L(t)}-1}{\eta+1} \in \mathbb{N}$, $I(t) = L(t)$; Otherwise $I(t) = \arg \max_{\substack{\{L(t)-1, L(t), L(t)+1\} \\ \cap \mathcal{D}_1(t)}} F(\hat{q}_k(t), N_k(t), t)$;
- 6: Observe the revealed outcomes X_t and Y_t ;
- 7: Update estimations: $\hat{q}_{I(t)}(t) = \frac{\hat{q}_{I(t)}(t-1)N_{I(t)}(t-1)+X_t}{N_{I(t)}(t-1)+1}$, $\hat{p}_{I(t)}(t) = \frac{\hat{p}_{I(t)}(t-1)N_{I(t)}(t-1)+Y_t}{N_{I(t)}(t-1)+1}$, $N_{I(t)}(t) = N_{I(t)}(t-1) + 1$;
- 8: Update parameter estimation: $\hat{a}_{I(t)}(t) = \arg \min |p_{I(t)}(a) - \hat{p}_{I(t)}(t)|$;
- 9: Update weights: $w_k(t) = N_k(t)/t, \forall d_k \in \mathcal{D}$;
- 10: $t = t + 1$.

11: **end while**

12: Estimate the turning point of efficacy as:

$$L_1(n) = \min_{k: d_k \in \mathcal{D}_1(n)} \left\{ m \geq k : |\hat{q}_m(n) - \hat{q}_{m+1}(n)| \leq \sqrt{\frac{c \log(n)}{N_m(n)}} + \sqrt{\frac{c \log(n)}{N_{m+1}(n)}}, \hat{q}_m(n) \leq \hat{q}_{m+1}(n) \right\},$$

$$L_2(n) = \arg \max_{d_k: p_k(\hat{a}(n)) \leq \theta} p_k(\hat{a}(n)).$$

Output: $\hat{d}(n) = \min\{L_1(n), L_2(n)\}$.

The proposed SEEDA dose allocation policy is general in the sense that no efficacy model is assumed. In practice, however, efficacy often exhibits certain structure which, if utilized correctly, may further improve the performance.

For conventional cytotoxic agents, efficacy monotonically increases with dose levels. The same is not true for MTAs, for which the dose-efficacy curve increases initially and then plateaus after reaching the level of saturation (Zang et al., 2014; Riviere et al., 2018). In this section, we modify the SEEDA algorithm to handle the increase-then-plateau efficacy model, and analyze its performance.

Formally, we introduce the following increase-then-plateau efficacy assumption, which holds for MTA.

Assumption 1 $q_k, k \in \mathcal{K}$ satisfies $q_1 \leq q_2 \leq q_3 \leq \dots \leq q_N = q_{N+1} = \dots = q_K$.

The *SEEDA-Plateau* algorithm is given in Algorithm 2. With Assumption 1, the efficacy has an inherent non-decreasing structure. The key idea is to combine the selection rule of OSUB in (Combes & Proutière, 2014) and reform step 4 in Algorithm 1. Note that step 4 calculates $L(t)$ as the estimated dose level with the optimal efficacy and safe toxicity at t . Algorithm 2 not only selects this dose level frequently enough, but also keeps exploring its neighboring dose levels.

We now analyze the regret of SEEDA-Plateau and present the result in Theorem 3. Compared to Theorem 1 for SEEDA without the increase-then-plateau efficacy model, one can see that the first $\log(n)$ coefficient improves from $c \sum_{d_k: p_k(a^*) \leq \theta} (q^* - q_k)^{-1}$ to $c(q^* - q_{N-1})^{-1}$. This gain comes precisely from the increase-then-plateau efficacy model, as the unimodal structure that exploits this structure leads to $\log(n)$ regret only from the neighboring arm.

Theorem 3 *The regret of SEEDA-Plateau satisfies:*

$$R(n) \leq \frac{c \log(n)}{q^* - q_{N-1}} + O(\log \log(n)) + (n\delta Q + t_1 + \frac{K-M}{2\epsilon^2}). \quad (9)$$

Furthermore, if $\delta = O(\frac{1}{n})$, we have that $R(n) \leq O(\log n)$.

The optimal dose level we have defined before can be rewritten as $k^* = \min\{M, N\}$, the recommendation accuracy of SEEDA-Plateau is given in Theorem 4.

Theorem 4 *With c set as $2 < c < \frac{5}{2}$, the probability that SEEDA-Plateau fails to recommend the optimal dose can be bounded as:*

$$\mathbb{P}[\hat{d}(n) \neq k^*] \leq \frac{3}{n^c} + \delta_1. \quad (10)$$

Compared to Corollary 2, the error probability of SEEDA-Plateau is increased by $\frac{3}{n^c}$. This is due to the ambiguity of the efficacy-optimal dose and the toxicity-optimal one, which leads to the two candidate doses $L_1(n)$ and $L_2(n)$. In practice, however, this ambiguity can be eliminated via preliminary experiments.

5. Experiments

5.1. Synthetic dataset

To investigate the operational characteristics and evaluate the performance of the proposed adaptive designs, we present an experimental study with $K = 6$ dose levels and $n = 300$ trial cohorts, with each cohort consists of 3 patients. The estimation is updated after observing all individual outcomes from a cohort. All experiment results are obtained with 1000 trial repetitions. The MTD threshold is set as $\theta = 0.35$.

The trial setup is the same as (Riviere et al., 2018) and (Zang et al., 2014), and we have simulated eight different efficacy and toxicity scenarios³. Due to the space limitation, we only report the results of the first scenario, where efficacy reaching the maximal value (the optimal dose) before toxicity hits MTD threshold. Additional results for this setting as well as the other seven scenarios are reported in Section K to M in the supplementary material.

The following baseline designs are used for comparison (whenever appropriate), whose details can be found in the supplementary material: 3+3, CRM, MCRM, Independent TS, KL-UCB, UCB-1, and multi-objective bandits. Note that MTA-RA and other TS variants in (Riviere et al., 2018) are not included because they assume a different truncated efficacy model, which needs to be perfectly known to the algorithm. For algorithms that require prior information of toxicity and efficacy, they are set as $[0.02, 0.06, 0.12, 0.20, 0.30, 0.40]$ and $[0.12, 0.20, 0.30, 0.40, 0.50, 0.59]$, respectively.

5.1.1. RECOMMENDATION AND ALLOCATION ACCURACY

We report the allocation and recommendation percentages of each dose for all considered designs in Table 2. Dose 3 (in bold font) is the optimal biological dose for this scenario. However, we comment that dose 4 also satisfies the optimality condition without violating the safety constraint. Nevertheless, it has a higher toxicity probability (although still below MTD) without increasing efficacy; thus less preferable to Dose 3. We note that for all the considered designs, the recommendation rule is $\hat{d}(n) = \arg \max_{k: \hat{p}_k(n) \leq \theta} \hat{q}_k(n)$, where $\hat{q}_k(n)$ and $\hat{p}_k(n)$ are the final estimations of toxicity and efficacy for dose level d_k , respectively. This suggests that safety constraint is considered in recommendation.

We can see from the results that SEEDA almost equally recommends dose 3 and 4 with a total probability of 94.6%.

³We remark that although no real-world trial data is utilized in the experiment, this approach is commonly accepted in clinical trials as the first-step study for a new methodology; see (Whitehead et al., 2012; Yap et al., 2013; Zang et al., 2014; Riviere et al., 2018).

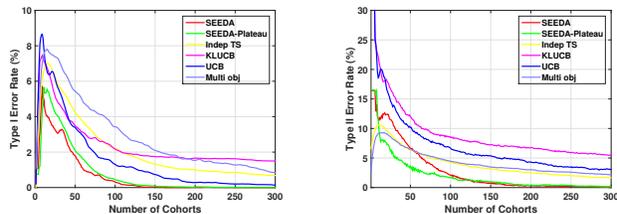


Figure 1: Type I (left) and type II (right) error rates as a function of number of cohorts.

This is because the algorithm cares about maximizing efficacy without violating safety constraint, and both dose 3 and 4 satisfy such conditions. As a result, SEEDA treats both equally as the optimal solution. However, by leveraging the increase-then-plateau model assumption, SEEDA-Plateau can further break the “tie” between dose 3 and 4, and correctly recognize that dose 3 is the optimal biological dose: it chooses dose 3 at 86.6% while dose 4 only 10.4%. We see that the gain of SEEDA-Plateau is significant over all the other designs (even compared to SEEDA). For a more detailed understanding of the recommendation accuracy, the corresponding type I and type II error rates (definitions are given in Section J in the supplementary material) are plotted in Fig. 1, and we observe that both SEEDA and SEEDA-Plateau outperform other baseline methods over the range of cohorts.

As for allocation, we observe that both SEEDA and SEEDA-Plateau concentrate at dose 3 and 4, while spending very little budget on both tail ends of the dosage. In particular, SEEDA-Plateau allocates the fewest percentages (1%) of patients to the most toxic dose 6 among all designs.

5.1.2. CONVERGENCE AND SAFETY VIOLATION

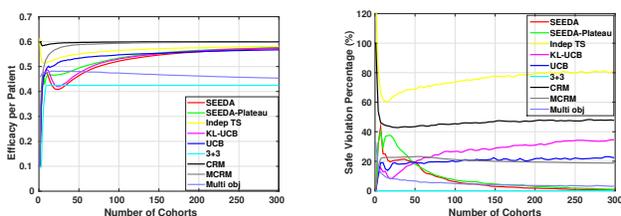


Figure 2: Comparison of efficacy per patient (left) and the safety violation percentage (right).

To have a deeper understanding of the tradeoff between efficacy and toxicity, we plot side-by-side the convergence of efficacy and toxicity as t increases in Fig. 2. KL-UCB, UCB and Independent TS have good convergence but suffer from significant safety violation in the process since they do not consider the safety constraint during exploration. CRM has

Table 2: Recommendation & allocation percentages of different designs. Optimal biological dose is #3. In each cell the first row reports the mean value over 1000 repetitions, and the second row reports the (standard deviation).

	Recommended						Allocated					
Toxicity probabilities	0.01	0.05	0.15	0.2	0.45	0.6	0.01	0.05	0.15	0.2	0.45	0.6
Efficacy probabilities	0.1	0.35	0.6	0.6	0.6	0.6	0.1	0.35	0.6	0.6	0.6	0.6
SEEDA	0 (0)	1 (0.71)	47.20 (3.40)	47.40 (3.41)	4.40 (2.46)	0 (0)	11.18 (0.58)	9.18 (1.99)	30.76 (7.76)	31.71 (7.69)	12.06 (3.45)	5.11 (0.62)
SEEDA-Plateau	0.80 (0.32)	2.20 (1.96)	86.60 (8.58)	10.40 (3.65)	0 (0)	0 (0)	7.83 (1.61)	8.98 (4.21)	30.12 (6.01)	37.17 (7.54)	14.91 (3.02)	1.00 (0.61)
Independent TS	2.60 (2.47)	9.40 (3.86)	44.60 (10.25)	35.40 (10.32)	6.60 (2.96)	1.40 (0.69)	3.66 (0.97)	7.26 (3.85)	22.22 (15.47)	21.00 (10.44)	22.26 (10.43)	23.60 (9.22)
KL-UCB	0.20 (0.13)	4.60 (2.71)	48.80 (11.68)	43.60 (11.36)	2.80 (2.64)	0 (0)	10.93 (0.81)	7.16 (0.94)	21.33 (10.52)	20.91 (11.31)	21.21 (10.92)	18.46 (11.10)
UCB	0 (0)	2.40 (2.15)	54.00 (9.92)	40.40 (9.05)	3.20 (3.09)	0 (0)	5.45 (0.49)	9.50 (1.16)	22.13 (2.11)	20.93 (2.24)	20.43 (2.15)	24.57 (2.11)
3+3	0 (0)	2.40 (0.41)	12.00 (4.31)	17.60 (5.31)	45.20 (7.15)	22.80 (4.35)	16.04 (5.12)	17.82 (4.23)	20.19 (10.25)	18.12 (9.15)	16.81 (8.15)	5.82 (4.12)
CRM	0 (0)	0 (0)	0 (0)	33.80 (8.26)	65.80 (10.63)	0.40 (0.40)	0.12 (0.11)	0.35 (0.25)	2.62 (0.32)	33.90 (10.21)	57.69 (11.24)	5.33 (0.23)
MCRM	0 (0)	0 (0)	0.20 (0.15)	61.00 (9.67)	38.80 (8.65)	0 (0)	1.47 (0.24)	1.18 (0.67)	5.64 (3.62)	55.48 (8.63)	34.63 (7.65)	1.60 (0.67)
Multi-obj	0.81 (0.19)	3.23 (0.94)	47.90 (11.86)	41.03 (11.90)	5.88 (1.80)	1.15 (0.36)	18.42 (6.07)	21.92 (5.55)	23.36 (6.68)	18.48 (7.05)	9.92 (5.17)	7.89 (4.65)

higher efficacy at the cost of bad safety constraint violation, while 3+3 performs poorly in efficacy but has the lowest safety probability; this behavior is similarly observed for multi-objective bandits. The SEEDA(-Plateau) algorithm, in comparison, converges to the optimal efficacy at a slower rate, but the exploration process is carefully controlled so that the safety violation is minimized, which is evident from the right subplot of Fig. 2.

5.1.3. SAMPLE EFFICIENCY

Sample efficiency is measured by the minimum number of trial participants to achieve a pre-specified recommendation accuracy (also known as *early stopping* (Montori et al., 2005)). We start the trial with a minimum of 6 patients, and continue recruiting patients until the stopping condition is triggered. Fig. 3 plots the average minimum number of patients to achieve a given recommendation accuracy for different algorithms⁴. We see that SEEDA-Plateau outperforms all other algorithms by a large margin, thanks to the “double dipping” of the model assumptions which gives the most accurate estimation of the optimal dose. In comparison, SEEDA performs similarly to the baseline algorithms. The reason is that the goal of SEEDA is to recommend the efficacy-maximal dose that satisfies the safety constraint. In this particular setting, both dose 3 and 4 satisfy this condition, and SEEDA does not have the mechanism to further minimize toxicity. This leads to a recommendation error that is similar to other baseline designs.

⁴3+3, CRM and MCRM are excluded since they only target finding MTD.

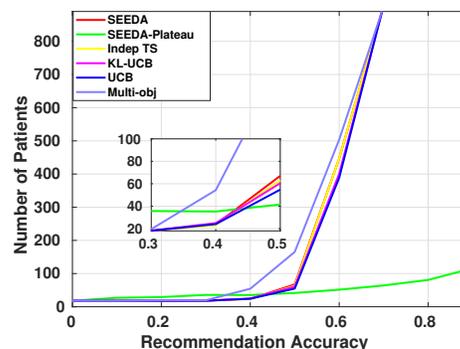


Figure 3: The minimum number of trial participants to achieve a given recommendation accuracy.

The sample efficiency advantage of SEEDA-Plateau is of critical importance in practice, as the significant cost associated with clinical trials is mostly proportional to the number of trial participants. Furthermore, reducing the number of patients while achieving the same level of accuracy minimizes the safety and ethical concern in the trial, which is another important consideration.

5.2. Real-world datasets

We evaluate the SEEDA algorithms in two real-world datasets *neurodeg* and *IBSCovars* based on (Biesheuvel & Hothorn, 2002). We first extract dose and resp variables from the observations reported in the dataset. With these samples, we fit them into a commonly used Emax dose-response model as in (Bornkamp et al., 2011) with an R

Table 3: Recommendation & allocation percentages of the neurodeg dataset. In each cell the first row reports the mean value over 1000 repetitions, and the second row reports the (standard deviation).

	Recommended					Allocated				
Toxicity	0.01	0.08	0.30	0.60	0.80	0.01	0.08	0.30	0.60	0.80
Efficacy	0.01	0.35	0.45	0.52	0.57	0.01	0.35	0.45	0.52	0.57
SEEDA	0.60 (0.40)	32.91 (10.57)	66.14 (10.59)	0 (0)	0 (0)	5.58 (0.42)	34.14 (6.08)	59.60 (6.25)	0.33 (0.25)	0.33 (0.01)
SEEDA-Plateau	0.99 (0.31)	32.66 (10.12)	66.00 (10.36)	0 (0)	0 (0)	5.09 (2.05)	35.02 (7.78)	59.21 (6.64)	0.33 (0.02)	0.33 (0)
Independent TS	3.39 (2.56)	51.28 (9.92)	44.47 (10.34)	0.38 (0.33)	0.46 (0.37)	0.80 (0.60)	3.50 (2.70)	7.59 (5.21)	23.37 (10.19)	64.68 (12.51)
KL-UCB	0.07 (0.06)	55.74 (12.38)	28.67 (12.76)	5.43 (2.10)	0.07 (0.05)	98.0 (2.62)	0.42 (0.23)	0.47 (0.28)	0.52 (0.49)	0.55 (0.04)
UCB	0.81 (0.74)	41.68 (16.07)	57.24 (16.07)	0.23 (0.17)	0.01 (0.01)	6.88 (0.29)	15.09 (1.60)	20.10 (2.10)	25.75 (2.55)	32.16 (2.87)
3+3	0 (0)	2.40 (1.02)	12.00 (2.35)	17.60 (3.44)	45.20 (10.34)	16.04 (5.60)	17.82 (9.48)	20.19 (1.84)	18.12 (1.60)	16.81 (4.20)
CRM	0 (0)	0 (0)	0 (0)	100 (0)	0 (0)	0 (0)	0 (0)	0 (0)	99.66 (0.01)	0.33 (0.01)
MCRM	4.33 (0.25)	26.47 (1.80)	69.18 (1.86)	0 (0)	0 (0)	4.67 (0.25)	26.40 (1.80)	68.92 (2.10)	0 (0)	0 (0)
Multi-obj	0.24 (0.13)	15.33 (9.65)	17.59 (9.71)	0.12 (0.05)	0.03 (0.03)	24.33 (4.28)	26.12 (3.51)	18.95 (6.11)	16.05 (2.93)	14.52 (2.61)

package implementation provided by (Yoshida, 2019). The resulting models are as follows.

$$\text{neurodeg: } \text{resp} = 169.94 + \frac{12.95 \text{dose}}{1.85 + \text{dose}},$$

$$\text{IBScovars: } \text{resp} = 0.26 + \frac{0.68 \text{dose}}{4.01 + \text{dose}}.$$

As for the toxicity event, since it is not reported in the dataset, we resort to simulations with model (1).

The allocation and recommendation percentages of each dose for all the algorithms are shown in Table 3 and Table 4 for both datasets. We have similar observations as in the synthetic experiment that SEEDA and SEEDA-Plateau recommend the correct doses majority of the times, while the suboptimal recommendation is mostly safe in that the doses immediately below MTD are recommended second most. The same is true for allocation.

6. Related works

This work is concerned with adaptive phase I clinical trials, whose uptake in practice is starting to increase considerably. See (Bretz et al., 2017; Pallmann et al., 2018) for recent comprehensive surveys. The main motivation to use these adaptive designs is to learn as the trial progresses and use this learning to deliver more efficient or more ethically appealing trials. Adaptive clinical trial with sequential patient recruitment is considered in (Atan et al., 2019), but it does

not address the subsequent dose allocation. The 3+3 and the CRM designs or their variations remain the de facto adaptive designs in practice for dose-finding studies (Petroni et al., 2017; Pallmann et al., 2018), although new methodologies that aim at better safety protection are also proposed (Lee et al., 2017). In recent years, there is a growing interest in adaptive trial designs for MTA because of its different dose-response relationships (Zang et al., 2014; Riviere et al., 2018), but these studies do not explicitly enforce the safety constraints during the trial; neither do they provide theoretical guarantees on the trial performance.

Multi-armed bandit has long been considered as an important tool for learning in clinical trials, dating back to the earliest papers of (Thompson, 1933; Robbins, 1952). Developing bandit models and algorithms that better suit the specific requirements of adaptive clinical trials has attracted some attention in recent years. Villar et. al (Villar et al., 2015b; Villar & Rosenberger, 2018) adopted the (modified) forward-looking Gittins index rule for multi-arm clinical trials. The authors of (Wang et al., 2018) propose a regional bandit model that can be applied to learning the drug dosage and patient response relationship. The sample complexity of thresholding bandit is analyzed in (Garivier et al., 2017), which matches MTD identification. Furthermore, dose-finding clinical trials with heterogeneous groups are investigated in (Lee et al., 2020) from a MAB perspective. Probably the closest work to ours is (Aziz et al., 2019), which also considers both toxicity and efficacy. However,

Table 4: Recommendation & allocation percentages of the IBScovars datasets. In each cell the first row reports the mean value over 1000 repetitions, and the second row reports the (standard deviation).

	Recommended					Allocated				
	0.01	0.10	0.30	0.70	0.95	0.01	0.10	0.30	0.70	0.95
Toxicity probabilities	0.01	0.10	0.30	0.70	0.95	0.01	0.10	0.30	0.70	0.95
Efficacy probabilities	0.01	0.20	0.27	0.33	0.43	0.01	0.20	0.27	0.33	0.43
SEEDA	1.14 (1.31)	35.04 (7.58)	63.47 (7.60)	0 (0)	0 (0)	10.11 (0.82)	34.35 (5.42)	54.86 (5.57)	0.33 (0.17)	0.33 (0.01)
SEEDA-Plateau	2.08 (2.76)	36.51 (10.31)	61.06 (10.42)	0 (0)	0 (0)	8.97 (4.00)	34.60 (4.34)	55.75 (7.83)	0.33 (0.02)	0.33 (0.01)
Independent TS	7.52 (7.15)	48.47 (9.79)	43.30 (9.71)	0.31 (0.60)	0.39 (0.17)	1.82 (0.85)	3.89 (3.61)	26.66 (10.74)	20.65 (13.54)	23.17 (10.05)
KL-UCB	28.55 (9.21)	44.90 (9.95)	23.24 (10.06)	3.28 (2.60)	0 (0)	98.33 (0.35)	0.37 (0.44)	0.40 (0.87)	0.42 (0.06)	0.45 (0.60)
UCB	1.73 (1.21)	45.26 (9.21)	52.81 (9.25)	0.17 (0.08)	0.01 (0.01)	9.41 (0.38)	15.37 (1.45)	18.94 (1.87)	23.22 (2.26)	33.04 (2.61)
3+3	2.40 (0.88)	12.00 (7.65)	17.60 (6.87)	45.20 (6.86)	22.80 (8.87)	16.04 (2.85)	17.82 (5.29)	20.19 (8.29)	18.12 (5.52)	22.81 (5.45)
CRM	0 (0)	0 (0)	0 (0)	1.35 (0.10)	98.65 (0.04)	0 (0)	0 (0)	0 (0)	99.66 (0.86)	0.33 (0.03)
MCRM	4.34 (0.26)	26.91 (2.15)	68.74 (2.20)	0 (0)	0 (0)	4.67 (0.04)	26.83 (0.05)	68.49 (0.91)	0 (0)	0 (0)
Multi-obj	0.45 (0.25)	16.18 (5.49)	16.56 (10.53)	0.10 (0.03)	0.02 (0)	1.23 (1.14)	3.27 (3.12)	5.79 (5.12)	13.11 (6.43)	76.56 (6.04)

the safety constraint, which is an essential constraint of real-world phase I trials, has not been explicitly considered in these papers.

On the other hand, the problem of safe exploration has attracted a lot of attention recently, albeit often in control (Koller et al., 2018) and general reinforcement learning (Berkenkamp et al., 2017). The authors in (Sui et al., 2015) propose the SAFEOPT algorithm for safe exploration in Gaussian processes, and (Kazerouni et al., 2017) presents a variant of linear UCB method for the contextual linear bandit problem. A different line of works (Maillard, 2013; Galichet et al., 2013) consider minimizing risk in MAB, but they are mostly casted in the mean-variance framework with respect to the reward distribution.

7. Conclusions

Learning in adaptive clinical trials faces several unique challenges that have not been well addressed, which may have contributed to their lack of adoption in actual clinical trials. In particular, the safety constraints resulting from ethical and societal considerations have been insufficiently researched, which has motivated us to develop the SEEDA algorithm that explicitly imposes safety constraints (in terms of toxicity) while also aiming for maximum patient response in a dose-finding study. Theoretical analysis of SEEDA is carried out and the proposed algorithm is further extended to the increase-then-plateau efficacy model and shown to

have smaller regret thanks to the unimodal structure. The performance advantages over state-of-the-art adaptive clinical trial designs are illustrated with experiments on both synthetic and real-world datasets.

8. Acknowledgements

CS acknowledges the funding support from Kneron, Inc. SSV thanks the funding received from the National Institute for Health Research Cambridge Biomedical Research Centre at the Cambridge University Hospitals NHS Foundation Trust and the UK Medical Research Council (grant number: MC_UU_00002/3). The research of MV has been supported by ONR and NSF 1524417 and 1722516.

References

- Atan, O., Zame, W. R., and van der Schaar, M. Sequential patient recruitment and allocation for adaptive clinical trials. In *Proceedings of The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1891–1900, Apr 2019.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2):235–256, May 2002.
- Aziz, M., Kaufmann, E., and Riviere, M.-K. On multi-armed bandit designs for phase I clinical trials. *arXiv e-prints*, art. arXiv:1903.07082, March 2019.

- Berkenkamp, F., Turchetta, M., Schoellig, A. P., and Krause, A. Safe model-based reinforcement learning with stability guarantees. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 908–919, Long Beach, California, USA, December 2017.
- Biesheuvel, E. and Hothorn, L. Many-to-one comparisons in stratified designs. *Biometrical Journal*, 44:101–116, 2002.
- Bornkamp, B., Bretz, F., Dette, H., and Pinheiro, J. C. Response-adaptive dose-finding under model uncertainty. *Annals of Applied Statistics*, 5:1611–1631, 2011.
- Bretz, F., Gallo, P., and Maurer, W. Adaptive designs: The swiss army knife among clinical trial designs? *Clinical Trials*, 14(5):417–424, 2017.
- Combes, R. and Proutière, A. Unimodal bandits: Regret lower bounds and optimal algorithms. In *Proceedings of the 31th International Conference on Machine Learning*, pp. 521–529, Beijing, China, June 2014.
- Galichet, N., Sebag, M., and Teytaud, O. Exploration vs exploitation vs safety: risk-aware multi-armed bandits. In *Proceedings of the 5th Asian Conference on Machine Learning*, pp. 245–260, November 2013.
- Garivier, A. and Cappè, O. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proceedings of Conference On Learning Theory (COLT)*, 2011.
- Garivier, A., Ménard, P., and Rossi, L. Thresholding bandit for dose-ranging: The impact of monotonicity. *arXiv e-prints*, art. arXiv:1711.04454, November 2017.
- Kazerouni, A., Ghavamzadeh, M., Abbasi, Y., and Van Roy, B. Conservative contextual linear bandits. In *Proceedings of Advances in Neural Information Processing Systems*, pp. 3910–3919, 2017.
- Koller, T., Berkenkamp, F., Turchetta, M., and Krause, A. Learning-based model predictive control for safe exploration. In *IEEE Conference on Decision and Control (CDC)*, pp. 6059–6066, December 2018.
- Lee, H.-S., Shen, C., Jordon, J., and van der Schaar, M. Contextual constrained learning for dose-finding clinical trials. In *Proceedings of The 23rd International Conference on Artificial Intelligence and Statistics*, Aug. 2020.
- Lee, S. M., Ursino, M., Cheung, Y. K., and Zohar, S. Dose-finding designs for cumulative toxicities using multiple constraints. *Biostatistics*, 20(1):17–29, Nov. 2017.
- Maillard, O.-A. Robust risk-averse stochastic multi-armed bandits. In *Proceedings of the 24th International Conference on Algorithmic Learning Theory*, pp. 218–233, Singapore, 2013.
- Montori, V. M. et al. Randomized trials stopped early for benefit: A systematic review. *JAMA*, 294(17):2203–2209, Nov. 2005.
- Neuenschwander, B., Branson, M., and Gsponer, T. Critical aspects of the Bayesian approach to phase I cancer trials. *Statistics in Medicine*, 27(13):2420–2439, 2008.
- O’Quigley, J., Pepe, M., and Fisher, L. Continual reassessment method: a practical design for phase I clinical trials in cancer. *Biometrics*, 43(1):33–48, 1990.
- Pallmann, P. et al. Adaptive designs in clinical trials: why use them, and how to run and report them. *BMC Medicine*, 16(1):29, Feb 2018.
- Paoletti, X. and Postel-Vinay, S. Phase I–II trial designs: how early should efficacy guide the dose recommendation process? *Annals of Oncology*, 29(3):540–541, Feb. 2018.
- Petroni, G. R., Wages, N. A., Paux, G., and Dubois, F. Implementation of adaptive methods in early-phase clinical trials. *Statistics in Medicine*, 36(2):215–224, 2017.
- Postel-Vinay, S. et al. Clinical benefit in phase-I trials of novel molecularly targeted agents: does dose matter? *British Journal of Cancer*, 100(9):1373–1378, May 2009.
- Riviere, M.-K., Yuan, Y., Jourdan, J.-H., Dubois, F., and Zohar, S. Phase I/II dose-finding design for molecularly targeted agent: Plateau determination using adaptive randomization. *Statistical Methods in Medical Research*, 27(2):466–479, 2018.
- Robbins, H. Some aspects of the sequential design of experiments. *Bull. Amer. Math. Soc.*, 58:527–535, 1952.
- Roberts, T. G. et al. Trends in the risks and benefits to patients with cancer participating in phase I clinical trials. *JAMA*, 292(17):2130–2140, Nov. 2004.
- Storer, B. E. Design and analysis of phase I clinical trials. *Biometrics*, 45:925–37, 1989.
- Sui, Y., Gotovos, A., Burdick, J. W., and Krause, A. Safe exploration for optimization with Gaussian processes. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 997–1005, 2015.
- Thiessen, B. et al. A phase I/II trial of GW572016 (lapatinib) in recurrent glioblastoma multiforme: clinical outcomes, pharmacokinetics and molecular correlation. *Cancer Chemotherapy and Pharmacology*, 65(2):353–361, Jan 2010.
- Thompson, W. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, December 1933.

- Tighiouart, M., Liu, Y., and Rogatko, A. Escalation with overdose control using time to toxicity for cancer phase I clinical trials. *PLOS ONE*, 9:1–13, 03 2014.
- Villar, S. S. and Rosenberger, W. F. Covariate-adjusted response-adaptive randomization for multi-arm clinical trials using a modified forward looking gittins index rule. *Biometrics*, 74(1):49–57, 2018.
- Villar, S. S., Bowden, J., and Wason, J. Multi-armed bandit models for the optimal design of clinical trials: Benefits and challenges. *Statistical Science*, 30(2):199–215, May 2015a.
- Villar, S. S., Wason, J., and Bowden, J. Response-adaptive randomization for multi-arm clinical trials using the forward looking Gittins index rule. *Biometrics*, 71(4):969–978, 2015b.
- Wang, Z., Zhou, R., and Shen, C. Regional multi-armed bandits. In *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 510–518, Playa Blanca, Lanzarote, Canary Islands, Apr. 2018.
- Whitehead, J. et al. A novel phase I/IIa design for early phase oncology studies and its application in the evaluation of MK-0752 in pancreatic cancer. *Statistics in Medicine*, 31(18):1931–1943, 2012.
- Yahyaa, S. and Manderick, B. Thompson sampling for multi-objective multi-armed bandits problem. In *Proceedings of European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pp. 47–52, Bruges, Belgium, April 2015.
- Yan, F., Thall, P. F., Lu, K. H., Gilbert, M. R., and Yuan, Y. Phase I–II clinical trial design: a state-of-the-art paradigm for dose finding. *Annals of Oncology*, 29(3):694–699, Dec. 2017.
- Yap, C. et al. Implementation of adaptive dose-finding designs in two early phase haematological trials: clinical, operational, and methodological challenges. *Trials*, 14(1):O75, Nov 2013.
- Yoshida, K. *Emax Model Analysis with 'Stan'*. Columbia University, New York, USA, 2019. URL <https://cran.r-project.org/web/packages/rstanemax>.
- Zang, Y., Lee, J. J., and Yuan, Y. Adaptive designs for identifying optimal biological dose for molecularly targeted agents. *Clinical Trials*, 11(3):319–327, 2014.

Supplementary Material: Learning for Dose Allocation in Adaptive Clinical Trials with Safety Constraints

Cong Shen, Zhiyang Wang, Sofía S. Villar, Mihaela van der Schaar

A. Preliminaries

Before presenting the technical proofs, we introduce some notations and regularity assumptions on the dose-toxicity model, which can be verified to hold for Eqn. (1). For a general toxicity function $p_k(a)$ of an unknown parameter $a \in \mathcal{A}$, the following regularities are imposed:

Assumption 2 1) *Monotonicity:* For each $k \in \mathcal{K}$ and $a, a' \in \mathcal{A}$ there exists $C_{1,k} > 0$ and $1 < \gamma_{1,k}$, such that $|p_k(a) - p_k(a')| \geq C_{1,k}|a - a'|^{\gamma_{1,k}}$.

2) *Hölder continuity:* For each $k \in \mathcal{K}$ and $a, a' \in \mathcal{A}$ there exists $C_{2,k} > 0$ and $0 < \gamma_{2,k} \leq 1$, such that $|p_k(a) - p_k(a')| \leq C_{2,k}|a - a'|^{\gamma_{2,k}}$.

We note that both monotonicity and continuity assumptions are mild and standard in the literature; see (Wang et al., 2018). Proposition 1 immediately follows with Assumption 2.

Proposition 1 For functions $p_k(a), \forall k \in \mathcal{K}$ that satisfy Assumption 2, we have:

1) $p_k(a)$ is invertible;

2) For each $k \in \mathcal{K}$ and $d, d' \in \mathcal{P}$, we have $|p_k^{-1}(d) - p_k^{-1}(d')| \leq \bar{C}_{1,k}|d - d'|^{\bar{\gamma}_{1,k}}$, where $\bar{\gamma}_{1,k} = \frac{1}{\gamma_{1,k}}$, $\bar{C}_{1,k} = (\frac{1}{C_{1,k}})^{\frac{1}{\gamma_{1,k}}}$.

For ease of exposition, we denote $C_1 = \min C_{1,k}$, $C_2 = \max C_{2,k}$, $\gamma_1 = \max \gamma_{1,k}$, $\gamma_2 = \min \gamma_{2,k}$, $\bar{\gamma}_1 = 1/\gamma_1$, and $\bar{C}_1 = C_1^{-\bar{\gamma}_1}$.

B. Select Design Parameters

The parameters appeared in Assumption 2 collectively determine the confidence interval in Eqn. (3). We take function (1) as an example to show how to select these parameters. We have

$$\begin{aligned} |p_k(a) - p_k(a')| &\geq C_{1,k}|a - a'|^{\gamma_{1,k}}, \\ \frac{|p_k(a) - p_k(a')|}{|a - a'|} &\geq C_{1,k}|a - a'|^{\gamma_{1,k}-1}, \\ \min_{a \in \mathcal{A}} p'_k(a) &\geq C_{1,k}|\mathcal{A}|^{\gamma_{1,k}-1}, \\ \log\left(\frac{\tanh(d_k) + 1}{2}\right) &\geq C_{1,k}|\mathcal{A}|^{\gamma_{1,k}-1}. \end{aligned}$$

Therefore, we can first set $\gamma_{1,k}$ as $\frac{3}{2}$ and find the corresponding $C_{1,k}$. Then, with the known function $p_k(a)$, parameters can be approximately calculated.

C. Proof of Lemma 1

$$\begin{aligned}
 \mathbb{P}[\hat{a}(t) + \alpha(t) < p_i^{-1}(\theta)] &\leq \mathbb{P}[\hat{a}(t) + \alpha(t) < a^*] \\
 &\leq \mathbb{P}[|a^* - \hat{a}(t)| > \alpha(t)] \\
 &\leq \mathbb{P}\left[\sum_{k=1}^K w_k(t-1)\bar{C}_1|\hat{p}_k(t) - p_k(a^*)|^{\gamma_1} > \alpha(t)\right] \\
 &\leq \sum_{k=1}^K \mathbb{P}\left[|\hat{p}_k(t) - p_k(a^*)| > \left(\frac{\alpha(t)}{w_k(t-1)\bar{C}_1 K}\right)^{\gamma_1}\right] \\
 &\leq \sum_{k=1}^K 2 \exp\left(-2N_k(t) \left(\frac{\alpha(t)}{w_k(t)\bar{C}_1 K}\right)^{2\gamma_1}\right) \tag{11} \\
 &\leq 2K \exp\left(-2 \left(\frac{\alpha(t)}{\bar{C}_1 K}\right)^{2\gamma_1} t\right) = \delta. \tag{12}
 \end{aligned}$$

Inequality (11) is from the Hoeffding's inequality and (12) is derived from the definition of $N_k(t) = tw_k(t)$ and Assumption 2 with $\gamma_1 > 1$.

D. Proof of Lemma 2

From the Hoeffding's Inequality and Eqn. (6), we have:

$$\alpha(t) \leq p_k^{-1}(\theta) - a^* - \epsilon = \Delta_k - \epsilon,$$

where $\Delta_k = |a^* - p_k^{-1}(\theta)|$ denotes the gap between the true value of parameter a and the parameter corresponding to when the toxicity of dose level d_k is exactly at the MTD threshold θ . When $t > t_1$ and with the definition of $\alpha(t)$ in Eqn. (3), the lemma can be immediately derived.

E. Proof of Theorem 1

Depending on whether the optimal dose level is included in the admissible set or not, we can decompose the regret into two parts:

$$\begin{aligned}
 R(n) &= \sum_{t=1}^n \mathbb{P}[k^* \notin \mathcal{D}_1(t)]Q + \mathbb{P}[k^* \in \mathcal{D}_1(t)]R_2(n) \\
 &\leq n\delta Q + R_2(n).
 \end{aligned}$$

The probability of the first error event $\{k^* \notin \mathcal{D}_1(t)\}$ can be bounded by Lemma 1, which indicates that at each step t the probability of a safe dose level being excluded from the admissible set is bounded by δ . For the second part, $R_2(n)$ represents the regret when the optimal dose is included in the admissible set. In this case, the error event is due to the inaccuracy of parameter estimation at the beginning as well as the limited efficacy information provided by each sample. Using Lemma 2, we have:

$$\begin{aligned}
 R_2(n) &\leq t_1 + (K - M) \sum_{t=1}^n \exp(-2t\epsilon^2) + \sum_{t=t_1+1}^n \sum_{d_k: p_k \leq \theta} \mathbb{1}\{I(t) = k\} \\
 &\leq t_1 + \frac{K - M}{2\epsilon^2} + \sum_{d_k: p_k \leq \theta} \frac{c \log(n)}{q^* - q_k}.
 \end{aligned}$$

Putting the regret from both error events together leads to (7), which completes the proof.

F. Proof of Theorem 2

First we note:

$$\begin{aligned} p_{I(t)}(a^*) - \theta &\leq p_{I(t)}(a^*) - \theta + \theta - p_{I(t)}(a^* - \alpha(t)) \\ &\leq C_2 |a^* - \hat{a}(t) + \alpha(t)|^{\gamma_2}. \end{aligned}$$

Thus, the probability can be upper bounded as:

$$\mathbb{P}[\hat{a}(t) - a^* > \alpha(t) + \epsilon] \leq \exp(-2t(\alpha(t) + \epsilon)^2).$$

Reorganizing the terms, we finally have

$$\mathbb{P}\left[\frac{1}{n} \sum_{t=1}^n p_{I(t)}(a^*) - \theta < C_2 \epsilon^{\gamma_2}\right] \geq 1 - \exp(-2t(\alpha(t) + \epsilon)^2) \geq 1 - \delta.$$

G. Proof of Corollary 2

$$\begin{aligned} \mathbb{P}[|\hat{a}(n) - a^*| \geq \Delta_M] &\leq \sum_{k=1}^K \mathbb{P}\left[|\hat{p}_k(t) - p_k(a^*)| > \left(\frac{\Delta_M}{w_k(t)\bar{C}_1 K}\right)^{\gamma_1}\right] \\ &\leq \sum_{k=1}^K 2 \exp\left(-2N_k(n) \left(\frac{\Delta_M}{w_k(t)\bar{C}_1 K}\right)^{2\gamma_1}\right) \\ &\leq 2K \exp\left(-2 \left(\frac{\Delta_M}{\bar{C}_1 K}\right)^{2\gamma_1} n\right). \end{aligned}$$

H. Proof of Theorem 3

We first establish Lemma 3, whose proof directly follow Theorem C.1 in (Combes & Proutière, 2014).

Lemma 3 $\mathbb{E}[l_k(n)] = O(\log(\log(n))),$ for each $k \neq k^*$.

Then, following the similar proof steps in Theorem 1, we have the bound in (9).

I. Proof of Theorem 4

Since $k^* = \min\{M, N\}$ and $L_1(n)$ and $L_2(n)$ are the estimations for N and M respectively, $\{\hat{d}_r(n) \neq k^*\} \subseteq E_1 \cup E_2$, where $E_1 = \{L_1(n) \neq N\}$, $E_2 = \{L_2(n) \neq M\}$. The latter can be bounded by Corollary 2. With the notation $\beta_k(n) = \sqrt{\frac{c \log(n)}{N_k(n)}}$, the probability of E_1 can be bounded as follows:

$$\begin{aligned} \mathbb{P}[L_1(n) < M] &\leq \mathbb{P}[|\hat{q}_N(n) - \hat{q}_{N-1}(n)| \leq \beta_{N-1}(n) + \beta_N(n)] \\ &\leq \mathbb{P}[\hat{q}_{N-1}(n) - q_k + q_N - \hat{q}_N(n) \leq q_N - q_{N-1} - \beta_{N-1}(n) - \beta_N(n)] \\ &\leq 2 \exp\left(-2N_{N-1}(n) \left(\frac{q_N - q_{N-1} - \beta_{N-1}(n) - \beta_N(n)}{2}\right)^2\right) \\ &\leq 2 \exp\left(-2f(N-1) \log(n) \left(\frac{\Delta_{N-1,N} - \beta_{N-1}(n) - \beta_N(n)}{2}\right)^2\right) \\ &= o\left(n^{-\frac{5}{2}}\right). \end{aligned}$$

Furthermore,

$$\begin{aligned}
 \mathbb{P}[L_1(n) > M] &\leq \mathbb{P}[|\hat{q}_N(n) - \hat{q}_{N+1}(n)| > \beta_N(n) + \beta_{N+1}(n)] \\
 &\leq \mathbb{P}[|\hat{q}_N(n) - q_N| + |q_{N+1} - \hat{q}_{N+1}(n)| > \beta_N(n) + \beta_{N+1}(n)] \\
 &\leq \mathbb{P}[|\hat{q}_N(n) - q_N| > \beta_N(n)] + \mathbb{P}[|\hat{q}_{N+1}(n) - q_{N+1}| > \beta_{N+1}(n)] \\
 &\leq \frac{2}{n^c}.
 \end{aligned}$$

Lastly, $f(N-1)$ is the coefficient of the lower bound of $N_{N-1}(n)$, and can be written as (see Theorem 4.1 in (Combes & Proutière, 2014))

$$f(N-1) = \frac{1}{I(q_{N-1}, q_N)}.$$

This completes the proof.

J. Baseline designs in the experiments

The following baseline designs are used for comparison to SEEDA and SEEDA-Plateau in the experiments.

- **KL-UCB** (Garivier & Cappè, 2011): This approach ignores the safety constraint and focuses entirely on efficacy during allocation, as for each patient it allocates the dose level with the highest efficacy index. The efficacy performance for each dose level is characterized by the KL-UCB index. However, at the end of the experiment, a dose level is recommended according to $\hat{d}(n) = \arg \max_{k: \hat{p}_k(n) \leq \theta} \hat{q}_k(n)$, where $\hat{q}_k(n)$ and $\hat{p}_k(n)$ are the last empirical estimations of toxicity and efficacy for dose level d_k . This suggests that safety constraint is considered in recommendation. Accordingly, type I and type II errors are defined as:

$$\begin{aligned}
 e_1 &= \sum_{k \in \mathcal{K}} \mathbb{1}\{p_k \leq \theta\} \mathbb{1}\{\hat{p}_k(n) > \theta\}, \\
 e_2 &= \sum_{k \in \mathcal{K}} \mathbb{1}\{p_k > \theta\} \mathbb{1}\{\hat{p}_k(n) \leq \theta\}.
 \end{aligned}$$

- **UCB-1** (Auer et al., 2002): The allocation and recommendation rules are similar to KL-UCB above, with the only difference that the dose level with the highest UCB-1 index of efficacy is allocated to the patient.
- **Independent Thompson Sampling (TS)** (Thompson, 1933; Aziz et al., 2019): Toxicity and efficacy are estimated with Bayesian indices:

$$\tilde{p}_k(t) \sim \text{Beta}(S_k^p(t) + 1, N_k(t) - S_k^p(t) + 1),$$

and

$$\tilde{q}_k(t) \sim \text{Beta}(S_k^q(t) + 1, N_k(t) - S_k^q(t) + 1),$$

where $S_k^p(t)$ counts the number of toxic outcomes of dose level k among the first t patients and $S_k^q(t)$ counts the number of effective responses. The dose with maximum $\tilde{q}_k(t)$ is allocated to the t -th patient and $\hat{d}(n) = \arg \max_{k: \tilde{p}_k(n) \leq \theta} \tilde{q}_k(n)$ is recommended. Definitions of type I and type II errors are slightly modified to:

$$\begin{aligned}
 e_1 &= \sum_{k \in \mathcal{K}} \mathbb{1}\{p_k \leq \theta\} \mathbb{1}\{\tilde{p}_k(n) > \theta\}, \\
 e_2 &= \sum_{k \in \mathcal{K}} \mathbb{1}\{p_k > \theta\} \mathbb{1}\{\tilde{p}_k(n) \leq \theta\}.
 \end{aligned}$$

- **CRM** (O'Quigley et al., 1990): We here employ the CRM algorithm with the same one-parameter toxicity model in our paper:

$$p_k(a) = \left(\frac{\tanh(d_k) + 1}{2} \right)^a.$$

We choose a typical prior distribution as $a \sim \exp(0.5)$. Therefore, d_k can be solved with $prior_{tox}$ and the prior mean of a . $\pi_t(a)$ denotes the posterior distribution of a after observing the outcomes of the first t patients. The allocation rule is a greedy one:

$$I_t^{CRM} = \arg \min_{k \in \mathcal{K}} |\theta - p_k(\hat{a}(t))|,$$

$$\hat{a}(t) = \int_0^\infty a d\pi_t(a),$$

where $\hat{a}(t)$ is the posterior mean value. With this estimation, the final recommendation rule can be written as:

$$\hat{d}(n) = \arg \min_{k \in \mathcal{K}} |\theta - p_k(\hat{a}(n))|.$$

- **3+3** (Storer, 1989): The lowest dose is first given to 3 patients. If none reports a toxic outcome, the next lowest dose level is given to the next 3 patients. If there are less than 2 among these 6 patients who report toxic outcome, the next lowest dose level is given to the next 3 patients; otherwise the experiment is stopped and the dose level used before stopping is recommended as MTD.
- **MCRM** (Neuenschwander et al., 2008): This algorithm classifies the probability of toxicity into four categories. For our simulated setting, the categories are set as:

Under-dosing:	$\pi_a(d) \in (0, 0.20]$
Targeted toxicity:	$\pi_a(d) \in (0.20, 0.35]$
Excessive toxicity:	$\pi_a(d) \in (0.35, 0.60]$
Unacceptable toxicity:	$\pi_a(d) \in (0.60, 1.00]$

The recommendation and the allocation rules are to maximize the probability of targeted toxicity while controlling the probability of excessive or unacceptable toxicity at $P^{thre} = 25\%$. Based on the posterior distribution of the toxicity, the probability that the toxicity falls in the above four categories can be calculated. The probability that it falls in Targeted category is denoted as P_i^t while falls in Excessive and Unacceptable categories as P_i^e . The selection rule is therefore $I_t = \arg \max_{P_i^e \leq P^{thre}} P_i^t$.

- **Multi-objective Bandits** (Yahyaa & Manderick, 2015): We implement the Pareto Thompson Sampling algorithm of (Yahyaa & Manderick, 2015) in our experiments. Specifically, after getting the estimations of toxicity and efficacy of each dose from running the Independent TS design, the algorithm computes the Pareto optimal dose level set \mathcal{I}^* , which means $\forall i \in \mathcal{I}^*, \forall j \notin \mathcal{I}^*, \tilde{p}_i(t) \leq \tilde{p}_j(t)$ or $\tilde{q}_i(t) \geq \tilde{q}_j(t)$.

Other policies designed for MTA, such as MTA-RA, depend on a different truncated two-parameter logistic efficacy model (Riviere et al., 2018). In our setting, the exact efficacy model is assumed to be unknown – we only make the increase-then-plateau assumption.

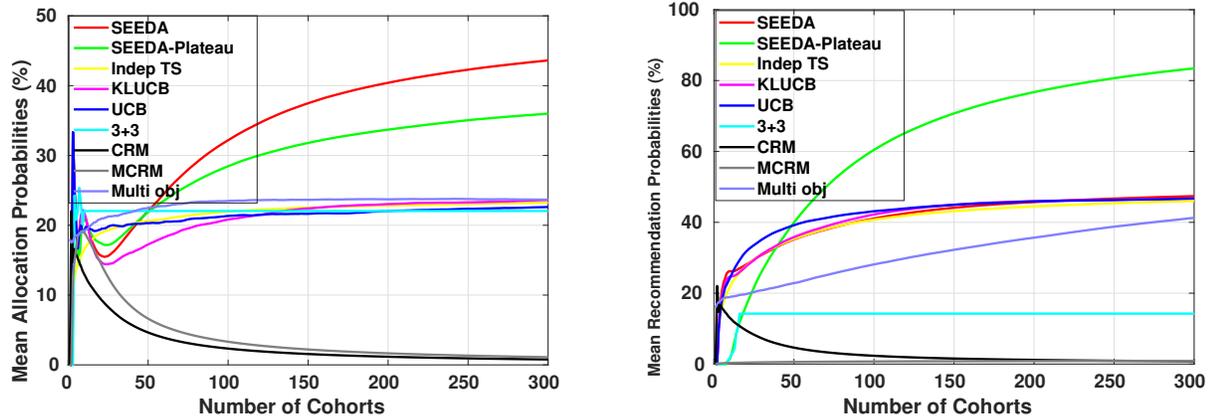
K. Additional experiment results under the same setting as in Section 5

Due to space limitations, we were not able to include all the experiment results of the setting in Section 5. These additional results are provided here.

In particular, Table 2 only reports the recommendation and allocation percentages for a given $n = 100$. It is of interest to see how these metrics change with n . We plot the mean allocation and recommendation probabilities as a function of n in Fig. 4. It can be seen that SEEDA-Plateau outperforms all other methods across a large range of n .

L. Experiment of a new setting and its comprehensive results

In the main paper, a setting that has the efficacy reaching the maximal value (the optimal dose) before toxicity hits MTD threshold is used. A different setting can be considered when maximum efficacy dose exceeds the MTD threshold. The


 Figure 4: Mean allocation (left) and recommendation (right) probabilities versus number of patients n .

experiment results for this setting (called “setting 2”) is reported in this section. Unless otherwise stated, the parameters are the same as in Section 5 of the main paper.

Table 5 presents the setting as well as the allocation and recommendation percentages of each dose for all considered algorithms. For this scenario, dose level 3 is the optimal one. We note that a large portion of the previous conclusions in the main paper still hold. However, the gain of SEEDA-Plateau is less significant over SEEDA, but still outperforms all the comparing designs. The corresponding Type I and Type II error rates are similarly plotted in Fig. 5.

Table 5: Recommendation & allocation percentages of different designs for setting 2.

	Recommended						Allocated					
Toxicity probabilities	0.1	0.2	0.25	0.4	0.5	0.6	0.1	0.2	0.25	0.4	0.5	0.6
Efficacy probabilities	0.3	0.4	0.5	0.7	0.7	0.7	0.3	0.4	0.5	0.7	0.7	0.7
SEEDA	9.54 (3.40)	19.34 (10.09)	52.66 (10.43)	16.00 (9.95)	2.12 (1.70)	0 (0)	6.82 (3.34)	17.61 (5.56)	48.99 (9.60)	21.77 (1.07)	3.47 (1.32)	1.33 (0.61)
SEEDA-Plateau	5.15 (3.72)	34.51 (5.96)	53.27 (6.80)	5.84 (2.64)	1.05 (0.50)	0.01 (0)	3.61 (2.28)	11.79 (1.79)	70.30 (7.51)	11.97 (5.12)	2.16 (0.42)	0.17 (0.12)
Independent TS	22.61 (5.61)	22.12 (7.43)	29.05 (8.24)	19.22 (5.96)	4.50 (2.41)	2.50 (2.01)	2.58 (1.90)	3.17 (2.23)	5.56 (3.72)	30.35 (4.73)	32.92 (4.62)	25.43 (3.82)
KL-UCB	19.72 (3.65)	21.03 (4.14)	29.19 (9.27)	24.02 (5.44)	5.46 (1.88)	0.59 (0.38)	2.13 (0.48)	2.50 (0.78)	3.37 (1.35)	32.80 (3.77)	30.63 (8.16)	28.58 (6.99)
UCB	8.95 (3.77)	22.45 (7.99)	41.04 (8.20)	21.61 (3.65)	4.83 (4.56)	1.11 (1.18)	8.12 (0.88)	10.31 (1.13)	13.20 (1.47)	22.90 (2.13)	22.58 (1.75)	22.89 (2.89)
3+3	6.80 (0.12)	20 (13.40)	23.80 (10.24)	29.80 (8.45)	16.40 (5.45)	3.20 (3.12)	26.99 (2.89)	27.50 (3.25)	19.59 (1.45)	13.14 (0.25)	5.01 (1.25)	0.76 (0.75)
CRM	0 (0)	0 (0)	0 (0)	0 (0)	99.10 (0.42)	0.90 (0.36)	0 (0)	0 (0)	0 (0)	0 (0)	99.11 (0.23)	0.89 (0.14)
MCRM	0 (0)	0.60 (0.93)	28.40 (13.29)	67.80 (13.95)	3.20 (3.06)	0 (0)	0.60 (0.09)	0.33 (0.12)	29.17 (9.47)	52.37 (13.95)	11.35 (4.34)	3.18 (1.92)
Multi-obj	6.57 (2.64)	13.38 (8.12)	50.95 (9.92)	22.71 (6.95)	4.44 (1.27)	1.95 (0.55)	20.17 (5.32)	14.78 (2.02)	19.05 (3.95)	20.29 (3.25)	14.57 (5.56)	11.17 (3.58)

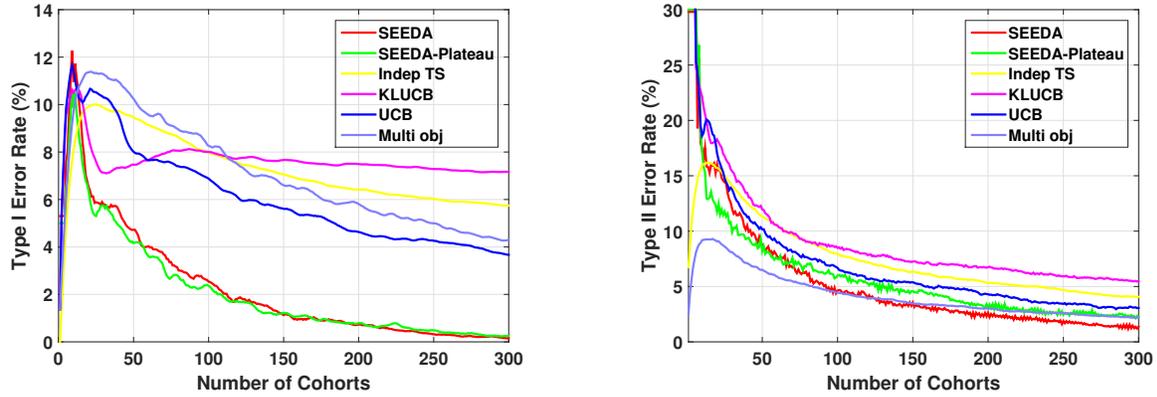


Figure 5: Type I and type II error rates in setting 2.

An in-depth look at the mean allocation and recommendation probabilities versus number of patients n for this new setting is given in Fig. 6. The same observation as in Section K holds.

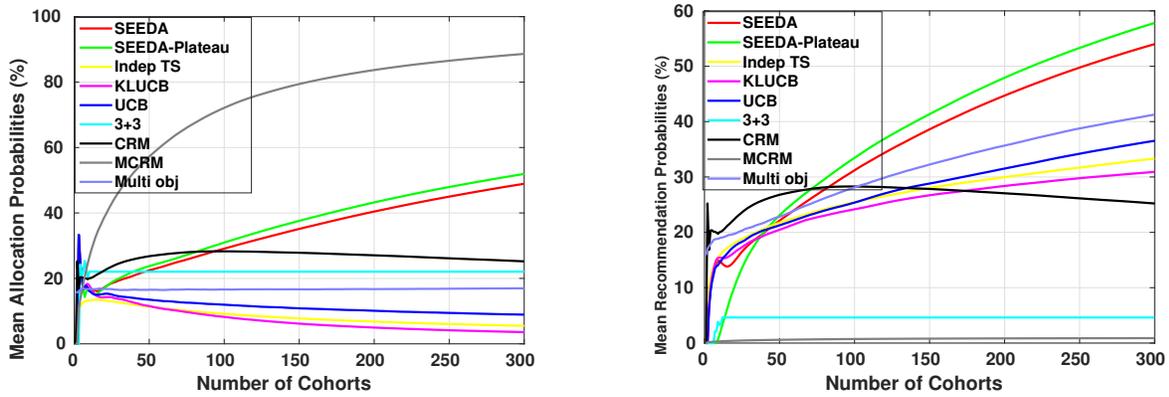


Figure 6: Mean allocation (left) and recommendation (right) probabilities versus number of patients n in setting 2.

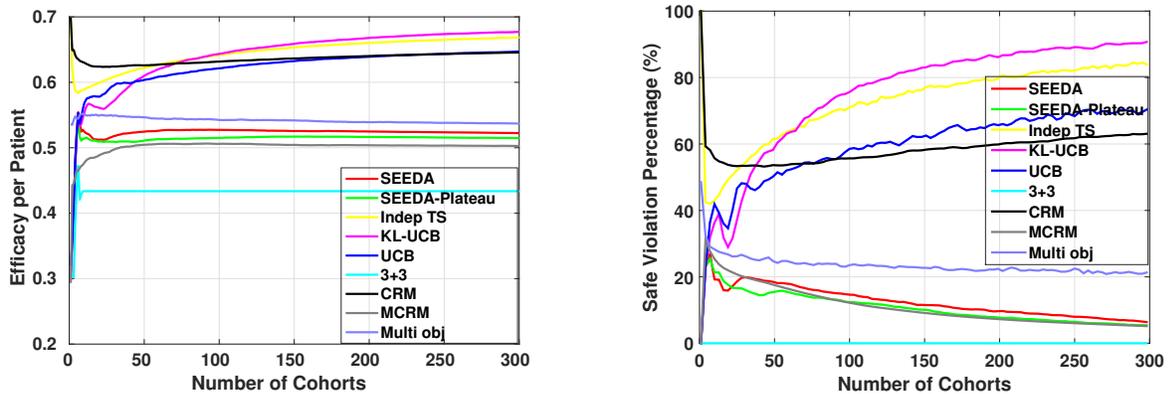


Figure 7: Comparison of efficacy per patient and the safety violation percentage in setting 2.

The convergence of efficacy and toxicity as t increases for setting 2 is plotted in Fig. 7. There is a notable difference to the previous result in Fig. 2, in that now SEEDA and SEEDA-Plateau converge to a different (but correct) dose than the other considered designs, which only emphasize maximum efficacy. It is clear that with such aggressive pursue of efficacy, they succeed in obtaining better treatment effect than SEEDA(-Plateau), but at the significant cost of frequent violation of the safety constraint: as opposed to safety violation percentage hovering between 40% and 50% in Fig. 2, now we face a violation in the range of 70% to 90% as shown in Fig. 7.

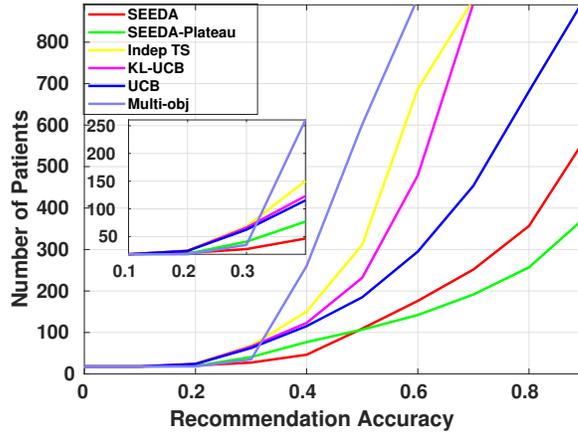


Figure 8: Sample size comparison in setting 2.

Lastly, the sample efficiency is evaluated. Fig. 8 plots the minimum number of patients to achieve a given a recommendation accuracy for different algorithms.

M. Experiment setting 3 to 8 with evaluation of allocation and recommendation percentages

This section reports the allocation and recommendation percentages of each dose for all considered algorithms under different toxicity/efficacy probabilities. We reuse the same 6 scenarios as those in the experiments of (Zang et al., 2014). See Table 6 to 11 for the detailed results. They are in line with the conclusions of the main paper.

Table 6: Recommended & allocated percentages for Scenario 1 of (Zang et al., 2014).

	Recommended					Allocated				
Toxicity probability	0.08	0.12	0.2	0.3	0.4	0.08	0.12	0.2	0.3	0.4
Efficacy probability	0.2	0.4	0.6	0.8	0.55	0.2	0.4	0.6	0.8	0.55
SEEDA	2.72 (1.01)	4.88 (2.14)	21.72 (7.50)	69.52 (10.11)	1.16 (0.62)	2.84 (0.78)	4.67 (1.95)	18.55 (6.04)	71.20 (7.65)	2.74 (2.74)
Indep TS	2.34 (0.25)	4.38 (1.31)	12.91 (6.34)	76.83 (7.03)	3.54 (1.49)	1.67 (0.62)	2.99 (0.64)	7.93 (0.36)	81.18 (2.55)	6.23 (2.44)
KL-UCB	9.58 (1.57)	23.99 (3.53)	39.35 (8.10)	24.27 (9.13)	2.81 (2.28)	3.24 (0.34)	13.89 (0.51)	30.91 (1.64)	22.35 (2.14)	29.61 (1.12)
UCB	3.04 (0.91)	12.41 (3.11)	46.91 (8.68)	35.24 (7.68)	2.40 (1.99)	10.91 (0.72)	18.41 (1.31)	33.34 (2.10)	28.32 (2.67)	9.02 (1.85)
3+3	4 (2.65)	10.40 (4.73)	20 (5.94)	22.80 (2.73)	42.80 (6.95)	23.38 (5.79)	22.81 (1.22)	20.92 (4.63)	15.80 (2.14)	10.79 (1.26)
CRM	0.09 (0.02)	0.20 (0.02)	1.72 (0.02)	42.51 (2.38)	55.48 (2.38)	0.09 (0.02)	0.20 (0.02)	1.72 (0.02)	42.51 (2.38)	55.48 (2.38)
MCRM	1.09 (1.01)	2.26 (2.20)	26.69 (7.69)	65.68 (9.26)	4.28 (2.10)	2.09 (1.31)	2.26 (2.20)	26.50 (6.68)	64.88 (8.25)	4.28 (0.13)
Multi-obj	1.41 (1.13)	4.56 (3.97)	22.69 (8.44)	67.31 (9.93)	4.03 (3.29)	18.42 (1.31)	20.69 (2.40)	22.51 (6.67)	31.41 (8.25)	6.97 (1.23)

Table 7: Recommended & allocated percentages for Scenario 2 of (Zang et al., 2014).

	Recommended					Allocated				
Toxicity probability	0.01	0.05	0.10	0.15	0.3	0.01	0.05	0.10	0.15	0.3
Efficacy probability	0.6	0.8	0.5	0.4	0.2	0.6	0.8	0.5	0.4	0.2
SEEDA	6.3 (0.90)	91.23 (3.18)	1.45 (1.02)	0.53 (0.34)	0.08 (0.08)	5.56 (3.11)	87.26 (3.94)	2.95 (2.09)	2.14 (1.43)	2.09 (0.63)
Indep TS	5.31 (4.95)	92.09 (1.32)	1.47 (1.08)	0.64 (0.56)	0.48 (0.16)	7.99 (2.55)	83.18 (5.34)	4.27 (4.34)	2.91 (2.30)	1.65 (1.05)
KL-UCB	9.68 (2.73)	87.66 (2.98)	1.91 (1.20)	0.66 (0.44)	0.09 (0.04)	7.01 (1.57)	81.93 (1.94)	3.03 (0.82)	2.31 (0.51)	5.72 (0.31)
UCB	8.58 (3.98)	89.80 (4.18)	1.26 (1.24)	0.34 (0.24)	0.03 (0.03)	21.06 (2.20)	46.31 (2.69)	15.07 (1.68)	11.16 (1.28)	6.40 (0.73)
3+3	0.20 (0)	1.80 (0.32)	5.40 (0.78)	13.80 (2.37)	78.80 (8.34)	16.71 (3.35)	18.81 (3.65)	19.40 (3.78)	19.88 (3.14)	19.75 (4.65)
CRM	0 (0)	0 (0)	0 (0)	9.98 (0.42)	90.02 (0.42)	0 (0)	0 (0)	0 (0)	9.97 (1.25)	90.03 (1.43)
MCRM	0.08 (0)	0.17 (0.02)	1.15 (1.00)	13.47 (0.44)	85.13 (0.04)	1.08 (0.27)	0.17 (0.09)	1.15 (0.61)	13.44 (5.76)	84.16 (6.73)
Multi-obj	6.07 (1.74)	90.85 (1.86)	1.93 (0.54)	0.92 (0.30)	0.22 (0.11)	34.88 (7.26)	51.74 (6.81)	7.11 (2.41)	4.34 (1.20)	1.93 (0.50)

Table 8: Recommended & allocated percentages for Scenario 3 of (Zang et al., 2014).

	Recommended					Allocated				
Toxicity probability	0.06	0.08	0.14	0.2	0.3	0.06	0.08	0.14	0.2	0.3
Efficacy probability	0.2	0.4	0.6	0.8	0.55	0.2	0.4	0.6	0.8	0.55
SEEDA	1.84 (0.71)	1.97 (1.10)	6.15 (2.86)	88.12 (3.22)	1.58 (1.00)	2.27 (0.71)	2.54 (1.11)	6.27 (2.86)	85.46 (3.22)	3.46 (0.99)
Indep TS	0.76 (0.45)	1.55 (0.93)	5.49 (3.71)	89.85 (5.09)	2.35 (1.73)	1.67 (0.48)	2.98 (1.33)	8.17 (4.48)	81.28 (4.96)	5.89 (1.79)
KL-UCB	2.64 (0.54)	7.29 (1.15)	28.47 (3.22)	57.18 (3.58)	4.43 (1.41)	2.62 (0.54)	6.58 (1.15)	26.85 (3.22)	55.07 (3.58)	8.87 (1.41)
UCB	1.71 (0.48)	3.57 (1.33)	19.04 (4.48)	72.89 (4.96)	2.79 (1.79)	8.33 (0.48)	13.17 (1.33)	22.75 (4.48)	44.71 (4.96)	11.04 (1.79)
3+3	2.20 (1.93)	4.80 (2.10)	10.60 (3.22)	18.80 (3.92)	63.60 (9.33)	19.77 (3.54)	20.08 (5.93)	20.43 (5.12)	18.67 (3.95)	15.29 (3.45)
CRM	0 (0)	0 (0)	0 (0)	4.37 (0.69)	95.63 (0.69)	0 (0)	0 (0)	0 (0)	4.37 (0.66)	95.63 (0.66)
MCRM	0.60 (0.54)	0.87 (0.15)	3.57 (3.22)	31.89 (3.58)	63.07 (1.41)	1.60 (0.98)	0.87 (1.26)	3.57 (2.97)	31.68 (8.86)	62.28 (10.58)
Multi-obj	0.78 (0.20)	2.07 (0.45)	8.67 (1.97)	84.99 (2.40)	3.49 (1.02)	16.43 (0.20)	20.56 (0.45)	21.56 (1.97)	34.45 (2.41)	7.00 (1.02)

Table 9: Recommended & allocated percentages for Scenario 4 of (Zang et al., 2014).

	Recommended					Allocated				
Toxicity probability	0.05	0.1	0.25	0.5	0.6	0.05	0.1	0.25	0.5	0.6
Efficacy probability	0.2	0.4	0.6	0.8	0.55	0.2	0.4	0.6	0.8	0.55
SEEDA	3.43 (1.26)	12.15 (3.69)	79.72 (4.25)	4.37 (1.90)	0 (0)	3.40 (1.24)	11.05 (3.48)	79.44 (4.28)	5.00 (1.75)	1.12 (0.45)
Indep TS	11.53 (9.17)	24.58 (10.80)	58.58 (12.42)	2.66 (1.53)	2.65 (3.42)	1.68 (0.99)	3.02 (2.39)	8.50 (5.40)	81.01 (16.00)	5.79 (6.50)
KL-UCB	24.60 (6.65)	37.78 (14.78)	28.34 (14.62)	6.91 (2.00)	2.37 (2.78)	1.91 (0.32)	2.43 (0.52)	3.41 (1.41)	51.61 (1.89)	40.64 (1.06)
UCB	4.87 (5.17)	32.53 (10.80)	60.34 (14.42)	1.84 (1.52)	0.42 (0.42)	14.29 (0.72)	26.93 (1.31)	40.69 (2.11)	9.15 (2.63)	8.94 (1.85)
3+3	3 (1.46)	6.20 (4.64)	34.20 (6.85)	40.40 (7.10)	16.20 (4.16)	22.57 (7.69)	22.82 (6.98)	26.70 (7.89)	17.10 (6.79)	4.29 (0.68)
CRM	0 (0)	0 (0)	0 (0)	95.56 (0.14)	4.44 (0.14)	0 (0)	0 (0)	0 (0)	95.23 (2.12)	4.77 (2.12)
MCRM	0.84 (0.83)	3.77 (1.73)	88.03 (3.92)	7.17 (3.58)	0.19 (0.01)	1.84 (0.83)	3.77 (1.73)	87.04 (3.91)	7.16 (3.57)	0.19 (0.01)
Multi-obj	3.64 (0.66)	19.66 (4.87)	70.79 (5.13)	4.80 (1.18)	1.11 (0.41)	19.93 (6.11)	23.96 (4.93)	23.97 (4.69)	26.16 (4.17)	5.98 (2.25)

Table 10: Recommended & allocated percentages for Scenario 5 of (Zang et al., 2014).

	Recommended					Allocated				
Toxicity probability	0.1	0.2	0.4	0.5	0.6	0.1	0.2	0.4	0.5	0.6
Efficacy probability	0.1	0.3	0.5	0.5	0.5	0.1	0.3	0.5	0.5	0.5
SEEDA	7.20 (1.10)	74.95 (4.42)	15.01 (4.84)	2.50 (1.46)	0 (0)	6.86 (0.96)	67.46 (3.49)	21.21 (4.11)	3.22 (1.58)	1.26 (0.61)
SEEDA-Plateau	12.60 (2.12)	82.20 (5.45)	4.60 (2.12)	0.60 (0.40)	0 (0)	19.50 (5.12)	56.46 (9.23)	15.49 (4.56)	7.56 (1.23)	1.00 (0.54)
Indep TS	21.59 (7.05)	50.75 (10.00)	21.15 (11.41)	4.73 (1.91)	1.78 (1.42)	2.67 (1.52)	6.34 (6.28)	29.19 (6.32)	30.59 (6.41)	31.22 (6.14)
KL-UCB	23.64 (4.52)	40.01 (10.18)	21.58 (10.82)	10.92 (2.16)	3.85 (0.81)	3.80 (0.75)	2.24 (1.38)	23.69 (10.60)	40.19 (9.94)	30.10 (10.93)
UCB	13.71 (1.63)	75.24 (9.14)	8.66 (5.79)	1.85 (0.79)	0.54 (0.96)	18.75 (0.64)	36.38 (4.19)	16.49 (2.57)	14.23 (2.63)	14.14 (2.55)
3+3	7.40 (1.42)	21.20 (12.30)	42.60 (6.42)	21.80 (3.06)	7.00 (4.12)	29.03 (0.79)	29.38 (3.32)	23.97 (2.14)	8.35 (1.15)	1.76 (0.42)
CRM	0 (0)	0 (0)	0 (0)	94.72 (0.04)	5.28 (0.05)	0 (0)	0 (0)	0 (0)	94.39 (0.02)	5.61 (0.02)
MCRM	2.86 (0.80)	62.72 (1.66)	33.03 (4.13)	1.25 (4.05)	0.14 (0)	3.86 (0.80)	62.02 (1.66)	32.73 (4.11)	1.25 (0.42)	0.14 (0.02)
Multi-obj	9.56 (0.58)	60.18 (3.92)	23.51 (4.17)	5.38 (1.00)	1.38 (0.39)	23.42 (6.89)	25.22 (5.30)	22.55 (5.28)	16.27 (6.61)	12.54 (5.79)

Table 11: Recommended & allocated percentages for Scenario 6 of (Zang et al., 2014).

	Recommended					Allocated				
Toxicity probability	0.01	0.03	0.05	0.1	0.2	0.01	0.03	0.05	0.1	0.2
Efficacy probability	0.1	0.3	0.45	0.6	0.6	0.1	0.3	0.45	0.6	0.6
SEEDA	1.47 (0.45)	1.79 (1.16)	5.12 (3.94)	48.97 (10.31)	42.32 (12.35)	3.59 (0.56)	2.93 (1.55)	5.89 (3.19)	45.65 (6.51)	41.94 (6.62)
SEEDA-Plateau	0 (0)	0.20 (0.05)	3 (1.38)	96 (5.72)	0.80 (0.56)	4.20 (3.75)	5.64 (2.45)	13.73 (5.42)	40.22 (9.85)	36.18 (4.75)
Indep TS	0.42 (0.31)	1.24 (0.86)	5.20 (3.13)	47.46 (12.35)	45.67 (12.22)	13.71 (1.06)	18.37 (3.55)	22.33 (5.87)	28.10 (8.80)	17.48 (8.57)
KL-UCB	1.96 (0.50)	2.55 (1.46)	9.57 (3.46)	54.30 (10.30)	31.62 (10.06)	3.78 (0.77)	3.32 (0.76)	9.42 (2.14)	52.03 (10.56)	31.45 (10.53)
UCB	1.31 (0.37)	2.06 (1.22)	9.47 (4.06)	56.47 (10.82)	30.69 (10.74)	8.18 (0.58)	12.58 (1.30)	19.54 (2.02)	32.85 (2.83)	26.84 (2.93)
3+3	0 (0)	1.40 (0.23)	2.20 (1.23)	8.20 (1.27)	88.20 (7.21)	17.14 (6.79)	18.15 (7.90)	18.32 (7.45)	20.07 (6.52)	20.74 (6.48)
CRM	0 (0)	0 (0)	0 (0)	65.39 (2.29)	34.61 (2.29)	0 (0)	0 (0)	0 (0)	65.15 (6.79)	34.85 (6.41)
MCRM	0.06 (0.02)	0.08 (0.04)	0.49 (0.50)	2.92 (1.21)	96.45 (3.00)	1.06 (0.25)	0.08 (0.04)	0.48 (0.29)	2.92 (2.01)	95.45 (3.00)
Multi-obj	0.63 (0.17)	1.60 (0.36)	6.78 (1.52)	49.01 (10.01)	41.98 (10.03)	13.71 (7.55)	18.37 (8.18)	22.33 (8.00)	28.10 (7.23)	17.48 (7.41)