

Comments on "Semi-Parametric Sampling for Stochastic Bandits with Many Arms"

Martyna Ksyta
ksyta.martyna@gmail.com

Abstract

In "Semi-parametric sampling for stochastic bandits with many arms" (Ou et al. 2019) the authors introduced a new contextual multi-armed bandit algorithm (called LSPS). It was claimed to be superior to the existing methods in situations with large number of candidate arms, when the reward depends on a context and arm-specific factor. The results were supported by better Bayesian regret bound and experimental performance.

In this comment, we prove that the formulation of Bayesian regret bound presented in the work mentioned above is incorrect - it is a decreasing function of number of steps and depends on a random variable associated with the estimated parameters. We also point and fix (with a proof) a minor mistake in a provided lemma. Additionally we reproduce experiments on synthetic data originally performed by the authors and obtain different results. Based on that we conclude that claims about superiority of LSPS might be exaggerated.

Introduction

Multi-armed bandit

Multi-armed bandit is an optimization problem. There is an agent and a set of possible actions (also called arms) she can choose from at each time step from 1 to T . An action brings a reward, which is revealed only after its selection. Agent's goal is to maximize the sum of rewards for all the time steps. In the contextual version of the problem before deciding on an action the agent also knows the context - feature vectors associated with the arms. Is it assumed that the reward is a function or a parametrized random variable. As time passes the agent knows more about the reward. At each moment she can decide to what extent use already gathered information (i.e. selecting arms with the context that has previously yielded high rewards) and to what extent explore (i.e. choosing a new context or arms) in the hope of finding arms giving even better rewards. This is a classical reinforcement learning problem known as the exploration-exploitation trade-off. The detailed assumptions about arms, rewards and the context depend on the algorithm. The work commented by us focus on a situation with the context when a number of arms to choose from might be large - in the next subsection we will mention work and assumption behind the models related to this problem.

Related work

There is a range of algorithms that do not take context into account assuming the rewards are independent between arms. A model presented in (Agrawal and Goyal 2017) uses parametrized random variables as rewards of the arms. Some prior distributions on the parameters are assumed. At any time step an arm is selected according to its posterior probability of being the best one. This algorithm is known in general as Thompson sampling (Thompson 1933) and it was empirically proved to be effective by multiple studies ((Granmo 2010), (Scott 2010), (Graepel et al. 2010), (Chapelle and Li 2011), (May and Leslie 2011), (Kaufmann, Korda, and Munos 2012)), as stated in (Agrawal and Goyal 2017). Other solutions not incorporating context comprise of Upper Confidence Bound algorithms studied by (Lai and Robbins 1985), (Auer, Cesa-Bianchi, and Fischer 2002), (Cappé et al. 2013), (Kaufmann, Cappé, and Garivier 2012) and (Agrawal 1995).

The disadvantage of such models is that they do not incorporate information about the context, which might be shared between arms. There are algorithms taking this into account. LinUCB (Li et al. 2010), analyzed later in (Chu et al. 2011), is based on the assumption that expected rewards are linear functions of the context, separate for each arm. With a prior distribution of their coefficient vectors we can derive formula for confidence bound of the expected reward for each arm and select the one maximizing it. (Agrawal and Goyal 2013) assume also a linear model but with a single coefficient vector shared between arms and distributed normally. The difference between arms is incorporated using feature vectors. The selected arm maximizes the reward calculated with the coefficient vector sampled from its posterior distribution. Other models based on the assumption that regret depends linearly on the context, were presented in (Auer 2002), (Dani, Hayes, and Kakade 2008), (Rusmevichientong and Tsitsiklis 2010), (Abbasi-Yadkori, Pál, and Szepesvári 2011) and (Kim and Paik 2019), which deals with the high-dimensional context.

There are also contextual multi-armed bandit algorithms based on general linear models studied by (Filippi et al. 2010), (Li, Lu, and Zhou 2017), (Jun et al. 2017) and (Kveton et al. 2020).

Assumption about linearity of the regret might not always be feasible. Recently non-linear bandit algorithms were pro-

posed. There are models based on neural networks: (Alessiardo, Féraud, and Bouneffouf 2014), (Collier and Llorens 2018), random forests: (Féraud et al. 2016), decision trees (Elmachetoub et al. 2017), piecewise constant estimators of the reward function: (Rigollet and Zeevi 2010), (Slivkins 2011), (Perchet, Rigollet et al. 2013) and Gaussian processes: (Krause and Ong 2011), (Srinivas et al. 2012). There is also a kernelized non-linear version of UCB: (Valko et al. 2013). The algorithm presented in (Foster et al. 2018) enables use of many predictor classes e.g. regularized linear functions or gradient-boosted regression trees. There are also attempts to combine multi-armed bandit with clustering (Gentile et al. 2017), (Li, Karatzoglou, and Gentile 2016) and (Wang et al. 2018). Authors of (Gupta, Joshi, and Yağan 2020) consider situation when the rewards are known functions of a common latent random variable. There is an algorithm testing linearity of regret and then selecting appropriate model proposed by (Ghosh, Chowdhury, and Gopalan 2017).

To overcome the discrepancy between reality and model assumptions the semi-parametric models were also introduced by (Krishnamurthy, Wu, and Syrgkanis 2018) and (Ou et al. 2019). These models assume that the reward for an arm come from a distribution whose mean is determined by two factors: one dependent on the features of an arm (parametric part) and the other dependent only on an arm (non-parametric part). In (Peng et al. 2019) the authors present a semi-parametric model in which the features depend also on the time step.

Linear Semi-Parametric Sampling model

In order to establish the notation for further use, we present below the model described in (Ou et al. 2019), called LSPS.

Assumptions

The agent has to select an action for each time step $t \in 1, 2, \dots, T$. The possible actions are fixed over time and indexed by $i = 1, 2, \dots, N$. There is a feature vector $x_i \in \mathbb{R}^d$, $\|x_i\| \leq 1$, associated with each action. The reward for an action, denoted as r_i , is a random variable. Each time an agent selects i -th action the observed reward is sampled from r_i . r_i may be represented as its expected value - γ_i - disturbed by a random noise η_t . Values of γ_i are initially and only once sampled for each arm and do not change over time. The means of their distributions are linear functions of x_i and coefficient vector θ , which must be found by the algorithm. Conditionally on θ variables γ_i are independent. There is a prior distribution of θ assumed. Values $\sigma_1, \sigma_2, \sigma_3 \in \mathbb{R}$ are hyperparameters of the model. All these assumptions, together with the distributions used by the authors, are summarized below:

$$r_i | \gamma_i \sim \mathcal{N}(\gamma_i, \sigma_1^2)$$

$$\gamma_i | (\theta, x_i) \sim \mathcal{N}(\theta^T x_i, \sigma_2^2)$$

$$\theta \sim \mathcal{N}(0, \sigma_3^2 I)$$

Thompson sampling

In order to select an arm at time t the LSPS algorithm uses Thompson sampling. It is described in the frame below. We follow a similar notation to (Russo et al. 2018).

Algorithm 1 Thompson sampling in LSPS

```

 $H_t \leftarrow \emptyset$ 
for  $t = 1, 2, \dots, T$  do
  sample  $\theta_t \sim \mathbb{P}(\theta | H_t)$ 
  for  $i = 1, 2, \dots, N$  do
    sample  $\gamma_{i,t} \sim \mathbb{P}(\gamma_{i,t} | \theta_t, H_t)$ 
  end for
   $i_t \leftarrow \operatorname{argmax}_{i=1,2,\dots,N} \gamma_{i,t}$ 
  observe reward  $r_{i_t,t}$  for  $i_t$ 
   $H_t \leftarrow H_t \cup \{(i_t, r_{i_t,t})\}$ 
end for

```

Posterior distributions

In order to use Thompson sampling one must know $\mathbb{P}(\theta | H_t)$ and $\mathbb{P}(\gamma_{i,t} | \theta_t, H_t)$. As calculated in (Ou et al. 2019) these distributions depends on:

$n_{i,t}$ - number of times the i -th arm was chosen up to time t ,

$\bar{r}_{i,t}$ - average observed reward of the i -th arm up to time t

The posterior distribution of θ can be represented as follows:

$$\theta | H_t \sim \mathcal{N}(\hat{\theta}_t, A_t^{-1}) \quad (1)$$

where:

$$\hat{\theta}_t = A_t^{-1} b_t$$

$$A_t = \frac{1}{\sigma_3^2} I + \sum_{i=1}^N \frac{n_{i,t}}{\sigma_1^2 + n_{i,t} \sigma_2^2} x_i x_i^T$$

$$b_t = \sum_{i=1}^N \frac{n_{i,t} \bar{r}_{i,t}}{\sigma_1^2 + n_{i,t} \sigma_2^2} x_i$$

The formula for posterior distribution of γ_i is:

$$\gamma_i | (\theta_t, H_t) \sim \mathcal{N}(\hat{\gamma}_{i,t}, \sigma_{i,t}^2) \quad (2)$$

where:

$$\hat{\gamma}_{i,t} = \frac{\sigma_2^2 n_{i,t} \bar{r}_{i,t} + \sigma_1^2 \theta_t^T x_i}{\sigma_1^2 + n_{i,t} \sigma_2^2}$$

$$\sigma_{i,t}^2 = \frac{\sigma_1^2 \sigma_2^2}{\sigma_1^2 + n_{i,t} \sigma_2^2}$$

Useful definitions

Before presenting our remarks to (Ou et al. 2019) we would like to start with useful definitions with some explanatory comments.

Bayesian regret

Bayesian regret is a measure of performance of posterior sampling algorithms. It was first proposed in (Russo and Van Roy 2014). It is based on a regret - the difference between the expected reward of the optimal arm and the expected reward of selected arms. For LSPS it can be written in the following way:

$$\text{Regret}(X, \theta, \gamma_1, \dots, \gamma_N, T, N, \sigma_1, \sigma_2, \sigma_3) = \sum_{t=1}^T \mathbb{E}[\max_{i=1, \dots, N} \gamma_i - r_{i_t, t} | \theta, \gamma_1, \dots, \gamma_N]$$

Bayesian regret is defined as the expected value of the regret calculated over the prior distribution of the parameters of the model. For LSPS we get:

$$\text{BayesRegret}(X, T, N, \sigma_1, \sigma_2, \sigma_3) = \mathbb{E}[\text{Regret}]$$

It is a function of X (which can be bound), time T , N , d and hyperparameters of the model: σ_1 , σ_2 and σ_3 . In (Ou et al. 2019) there is an additional parameter introduced, which does not fit into this formula. This could be explained by the fact, that the presented model is semi-parametric and the latter expectation is calculated conditionally on θ . However the definition of the Bayesian regret presented by the authors in equation (5) is similar to the one above and it does not assume some fixed parameters. If the definition of Bayesian regret was different for semi-parametric models it would not be possible to compare them to other Thompson sampling models.

R-sub-Gaussian random variable

A random variable η is called R-sub-Gaussian, if $\forall \lambda \geq 0$:

$$\mathbb{E}[e^{\lambda \eta}] \leq e^{\lambda^2 R^2 / 2}$$

Notation \tilde{O}

f, g - real valued functions, $g(x) > 0$.

If $\exists k > 0$ so that:

$$f(x) = O(g(x) \log^k(g(x))),$$

we can write

$$f(x) = \tilde{O}(g(x))$$

Main remarks

Comments on the Bayesian regret bound

In the Theorem 1 the authors of (Ou et al. 2019) provide the formulation of the Bayesian regret bound of the LSPS algorithm. We will first present the theorem and then raise our concerns about it.

Theorem 1. *If $\forall t \leq T$, η_t is R-sub-Gaussian, $d \leq \sqrt{N}$ and $\epsilon_{max} \leq \sqrt{\frac{N}{dT}} \left(\frac{d}{\sqrt{N}}\right)^{2\alpha}$, where $\alpha \in [0, 1]$, then the algorithm LSPS can achieve upper Bayesian regret bound:*

$$\tilde{O}(\sqrt{N}^{1-\alpha} d^\alpha \sqrt{T}) \quad (3)$$

with

$$\frac{\sigma_2^2}{\sigma_1^2} = \frac{T}{N} \left(\frac{d}{\sqrt{N}}\right)^\alpha \quad (4)$$

Formula for the bound of the Bayesian regret We can use (4) to substitute $\left(\frac{d}{\sqrt{N}}\right)^\alpha$ into the formula for the Bayesian regret bound, getting rid of α . This would give us the formula below, which is a decreasing function of T .

$$\tilde{O}\left(\frac{\sigma_2^2}{\sigma_1^2} N^{\frac{3}{2}} T^{-\frac{1}{2}}\right) \quad (5)$$

This can not be true. The detailed proof is provided in the Appendix. Intuitively the differences between expected rewards of an optimal arm and the selected arms are non-negative for each time step, so their sum can not decrease in time to zero, as the bounding function.

Impact of ϵ_{max} We would also like to comment the presence of ϵ_{max} . Firstly it is a random variable and its presence in the formulation of the Bayesian regret bound seems inappropriate. According to the definition presented by (Ou et al. 2019) in the section "Problem Formulation" in equation (1) we can interpret it in the following way: For i -th arm the variable ϵ_i is the difference between parametric part of the expected reward and the expected reward itself. ϵ_{max} is the maximal absolute difference of these values over all arms:

$$\epsilon_i = f(\theta, x_i) - \gamma_i$$

$$\epsilon_{max} = \max\{|\epsilon_1|, \dots, |\epsilon_N|\}$$

The Bayesian regret bound should be the expected value calculated over the prior distribution of θ . The detailed derivation of the results could cast some light on the reason for presence of ϵ_{max} . It is supposedly present in the supplementary materials. However these materials were not published anywhere and we were not able to obtain them from the Authors.

Secondly the constraint on ϵ_{max} depends on T . Using similar substitution as for the regret bound we get:

$$\epsilon_{max} \leq \frac{\sigma_2^4}{\sigma_1^4} d^{-1} N^{\frac{5}{2}} T^{-\frac{5}{2}}$$

so we see that the bound for ϵ_{max} narrows with increasing T . This limits the applicability of the Theorem 1. To consider how significant it is effect we can use the attributes of the e-commerce dataset presented by (Ou et al. 2019), as they should reflect the real-world scenario. In this case $d = 5$, $N = 1000$ and $T = 200000$. The bound for ϵ_{max} equals $3.5 * 10^{-7} \frac{\sigma_2^4}{\sigma_1^4}$, which seem to be low, but in order to validate its sensibility one must know the order of magnitude of σ_1 , σ_2 and X .

Claims about α coefficient We can notice that the claim that $\alpha \in [0, 1]$, also stated in the abstract of the publication, is not always true. From the equation (4) we see that:

$$\alpha = \frac{\ln T - \ln N + 2(\ln \sigma_1 - \ln \sigma_2)}{\frac{1}{2} \ln(N) - \ln d}$$

From the fact, that the denominator is greater than 0 (because $d \leq \sqrt{N}$), we get that α can be negative if $\frac{T}{N} \leq \frac{\sigma_2^2}{\sigma_1^2}$. This may happen in a situation when N is comparable to

T . This situation is described by the authors as one where "LSPS can achieve significant improvement even when bias is relatively large". In such case α is smaller than 0 for σ_2 accordingly greater than σ_1 . If $\alpha < 0$, the $(\sqrt{N}/d)^\alpha$ improvement declared by the authors is actually a deterioration. To prevent this we need additional assumption that $\frac{T}{N} \geq \frac{\sigma_2^2}{\sigma_1^2}$.

Value of α can be also > 1 if:

$$T > \frac{N\sqrt{N}}{d} \left(\frac{\sigma_2}{\sigma_1} \right)^2$$

This however does not have negative impact on the Theorem 1.

Comments on a lemma

The following claim is stated in (Ou et al. 2019) as Lemma 1:

Lemma 1. For any $t \leq T$, with probability $1 - \frac{\delta}{NT^2}$

$$|\gamma_{i,t} - \widehat{\gamma}_{i,t}| \leq \sqrt{2 \ln \frac{NT^2}{2\delta}} \sigma_{i,t} \quad (6)$$

This lemma is not completely correct. Taking $\frac{\delta}{NT^2} = \frac{1}{2}$ we get, that with probability $\frac{1}{2}$:

$$|\gamma_{i,t} - \widehat{\gamma}_{i,t}| \leq 0$$

which is equivalent to:

$$\mathbb{P}(\gamma_{i,t} = \widehat{\gamma}_{i,t}) = \frac{1}{2}$$

From the model formulation in (2) we know that:

$$\gamma_i(\theta_t, H_t) \sim \mathcal{N}(\widehat{\gamma}_{i,t}, \sigma_{i,t}^2)$$

so we can calculate:

$$\begin{aligned} \mathbb{P}(\gamma_{i,t} = \widehat{\gamma}_{i,t}) &= \mathbb{E}[\mathbb{P}(\gamma_{i,t} = \widehat{\gamma}_{i,t} | (\theta_t, H_t))] = \\ \mathbb{E}[0 | (\theta_t, H_t)] &= 0 \neq \frac{1}{2} \end{aligned}$$

The lemma is true if we use:

$$\sqrt{2 \ln \frac{NT^2}{\delta}} \sigma_{i,t}$$

in the right-hand side of the inequality (6) and require the probability to be at least $1 - \frac{\delta}{NT^2}$.

Proof. Let's denote $a = \sqrt{2 \ln \frac{NT^2}{\delta}}$ and $X = \frac{\gamma_{i,t} - \widehat{\gamma}_{i,t}}{\sigma_{i,t}}$. We know that:

$$X | (\theta_t, H_t) \sim \mathcal{N}(0, 1)$$

so we can apply Proposition 1 to get:

$$\mathbb{P}(|X| \leq a | (\theta_t, H_t)) \geq 1 - e^{-\frac{a^2}{2}}$$

After taking the expectation of both sides we get the formula to be proven. \square

Proposition 1. For any $a > 0$ and $X \sim \mathcal{N}(0, 1)$

$$\mathbb{P}(|X| \leq a) \geq 1 - e^{-\frac{a^2}{2}}$$

Proof. The density of normal distribution is an even function and $a > 0$, so we can write:

$$\begin{aligned} \mathbb{P}(|X| > a) &= 2\mathbb{P}(X > a) = \\ 2 \int_a^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx &= \frac{2}{\sqrt{\pi}} \int_{\frac{a}{\sqrt{2}}}^{+\infty} e^{-t^2} dt \end{aligned}$$

The integral can be bounded using Propositions 2 and 3:

$$\begin{aligned} \frac{2}{\sqrt{\pi}} e^{-\frac{a^2}{2}} e^{\frac{a^2}{2}} \int_{\frac{a}{\sqrt{2}}}^{+\infty} e^{-t^2} dt &\leq \\ \frac{2}{\sqrt{\pi}} e^{-\frac{a^2}{2}} \frac{1}{\frac{a}{\sqrt{2}} + \sqrt{\frac{a^2}{2} + \frac{4}{\pi}}} &< e^{-\frac{a^2}{2}} \end{aligned}$$

So we finally got that:

$$\mathbb{P}(|X| > a) \leq e^{-\frac{a^2}{2}}$$

and

$$\mathbb{P}(|X| \leq a) = 1 - \mathbb{P}(|X| > a) \geq 1 - e^{-\frac{a^2}{2}}$$

\square

Proposition 2. Formula 7.1.13 from (Abramowitz and Stegun 1964)

For any $z \geq 0$:

$$e^{z^2} \int_z^{+\infty} e^{-t^2} dt \leq \frac{1}{z + \sqrt{z^2 + \frac{4}{\pi}}}$$

Proposition 3. For any $a > 0$:

$$\frac{2}{\sqrt{\pi} \left(\frac{a}{\sqrt{2}} + \sqrt{\frac{a^2}{2} + \frac{4}{\pi}} \right)} < 1$$

Proof. The expression is decreasing with respect to a , so for any $a > 0$ is smaller than for $a = 0$:

$$\frac{2}{\sqrt{\pi} \sqrt{\frac{4}{\pi}}} = 1$$

\square

Not related theorem

In the section 5 "Regret analysis" (Ou et al. 2019) refer to the Proposition 3 from (Russo and Van Roy 2016). However this proposition is not connected to authors' statements. We suspect that authors had in mind Proposition 2 from (Russo and Van Roy 2014), which gives the Bayesian regret bound of $\tilde{O}(\sqrt{NT})$ if rewards are within interval $[0, 1]$. We can also find in (Russo and Van Roy 2014) claim that this proposition can be extended to cases where reward is not bounded but where instead its distribution is "light-tailed". In this case it might fit the statement of (Ou et al. 2019).

Experiment

We have reproduced the experiment performed by (Ou et al. 2019) on the synthetic data. The code can be found in [this](#) Github repository, files parametrizing the experiments and their results, with saved states of the models and rewards are stored on Google Drive [here](#). All the models were implemented from scratch, since the code was not provided with the publication.

Hyperparameters

The LSPS and one of the models used for the comparison utilise hyperparameters - their values used in the original experiment were not provided by the authors. They only mentioned to use the same corresponding hyperparameters of the LSPS and TS-Lin models. Because of this we decided to perform grid search over values from $[0.1, 1, 10]$ for reduced number of timesteps ($T = 25000$) and later run models once again with the best hyperparameters for all the timesteps.

We were not able to perform the experiment on e-commerce dataset, because it is not publicly available. We feel that even contradicting results on sythetic dataset are worth sharing with the machine learning community.

Compared models

LSPS was compared to the models listed below. All of them are based on Thompson sampling.

- TS-Gau (Agrawal and Goyal 2017). This model does not use the context. It assumes that the reward of each arm $\sim \mathcal{N}(\mu_i, 1)$. The prior distribution of μ_i is $\mathcal{N}(0, 1)$. The posterior distribution is calculated based on the rewards observed for each arm.
- TS-Beta (Agrawal and Goyal 2017). This is similar model to TS-Gau. The difference is in the used distributions: reward for each arm is assumed to follow the Bernoulli distribution. Its mean has $\mathcal{B}(1, 1)$ prior distribution.
- TS-Lin (Agrawal and Goyal 2013) This model assumes that the expected reward for each arm $\sim \mathcal{N}(x_i^T \mu, v^2)$. The prior of μ is $\mathcal{N}(0, v^2 I)$ and v is a hyperparameter.

Setup

We have used the same setup as one described in (Ou et al. 2019) whenever possible. All x_i and θ were sampled from $\mathcal{N}(0, I)$, then we took the absolute values over all of their elements. Later x_i were divided all by the largest of their norms, so $\forall i \|x_i\| \leq 1$. Parameter θ was normalized and then multiplied by a predefined constant a . The expected rewards for each arm - γ_i - were set to $x_i^T \theta + \epsilon_i$. Values of ϵ_i were sampled from $\mathcal{U}[0, 1 - a]$. (Ou et al. 2019) claim to use $\mathcal{U}[1 - a, 1]$ instead, but it is not feasible. In this case we could get values of expected reward higher than 1. The authors claims to model reward using Gaussian and binomial distributions. For the binomial distribution (which has two parameters: p and n) n was not specified, so we assumed that the authors meant $n = 1$, which is equivalent to Bernoulli distribution and in line with the assumption of TS-Beta model. In this case the expected rewards can not be

greater than 1. Also if we use parameter a to control the impact of bias ϵ_i on the expected reward, the norm of θ should increase and the possible values of ϵ_i should decrease with increasing a . This is not the case if we use $\mathcal{U}[1 - a, 1]$.

The experiment was performed for 1000 arms, 5-dimensional feature vectors and a equal 0.5 (reflecting semi-parametric case) and 1 (purely linear case). We repeated calculations for 3 values of random seeds: 1, 2, 3 for hyperparameters tuning and 3 different values for final calculations: 4, 5, 6. This was not originally performed by the authors.

Results

We share results from all four cases tested in the simulation. In Figure 1 there is a cumulative regret shown for Gaussian distribution and $a \in \{0.5, 1\}$. In Figure 2 we present binomial distribution. Bold lines represent values averaged over all seeds. The shaded areas reflect maximal and minimal values. The cumulative regret of LSPS in some cases is very close to TS-Lin. In order to compare these models we include numerical results in Table 1. (Ou et al. 2019) present some results on graphs. It is not mentioned which distribution they refer to, so we can not compare all the results directly. However in all cases presented by the authors LSPS is the best model. For $a = 0.5$ the difference is clear, for $a = 1$ is much smaller. We obtained similar results only in case of Gaussian distribution. For $a = 1$ TS-Lin was better than LSPS but its cumulative regret was within bounds formed by minimum and maximum of LSPS results. Results for binomial distribution show much more variability - for $a = 0.5$ the best model is TS-Lin and the range of cumulative regret for LSPS is very wide. For $a = 1$ LSPS gives the best performance. In order to get certainty which of these models: LSPS or TS-Lin is better more experiments are needed. It can not be decided based only on one run of the simulation, as done in (Ou et al. 2019). We also obtained differences in performance of TS-Lin in the semi-parametric case. (Ou et al. 2019) shows that it is significantly worse comparing to other models. We did not get such results. It could be explained by the fact that TS-Lin is able to quickly estimate the rewards of the best arms even in semi-parametric environment ($a = 0.5$). Then it selects them, obtaining more observations and improving the fit for these best arms, at the expense of fit to other ones.

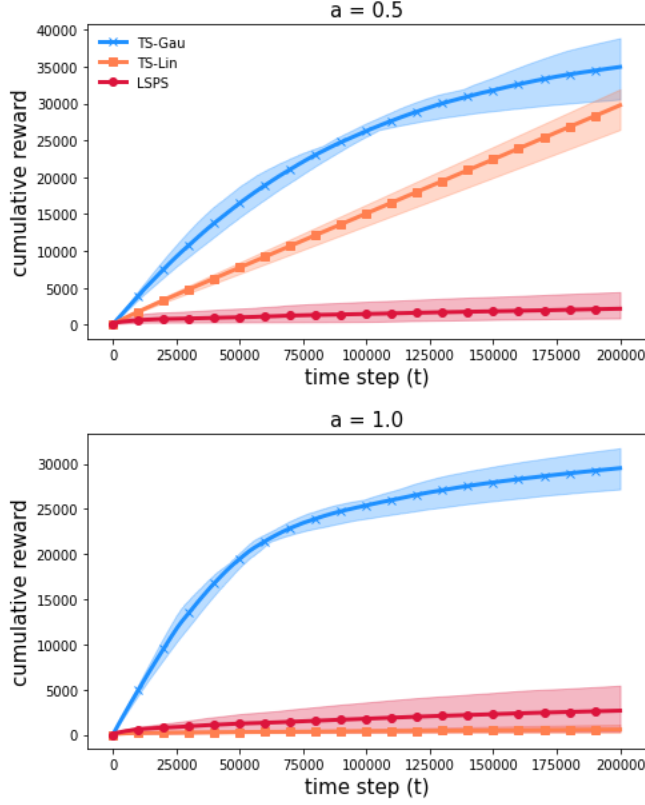


Figure 1: Results of the experiment for Gaussian distribution

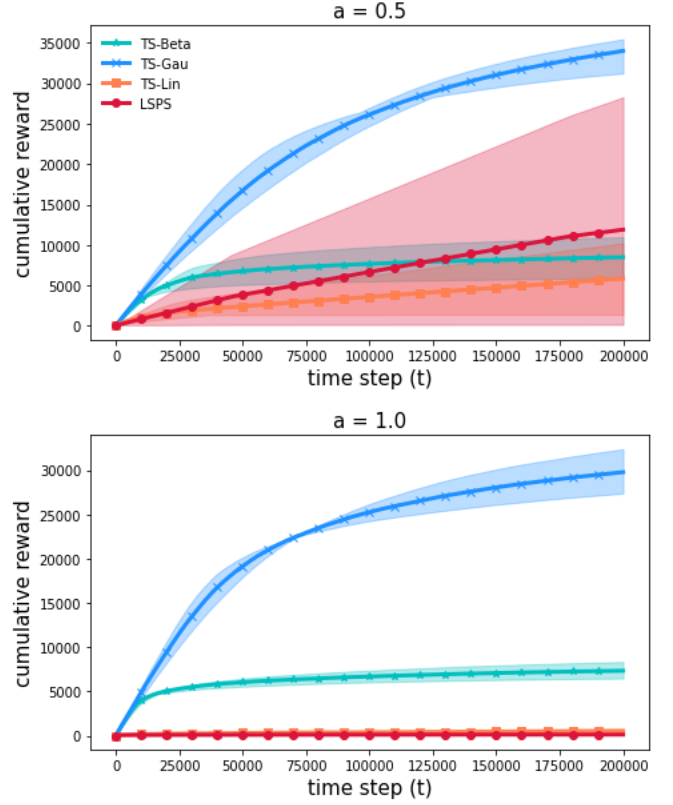


Figure 2: Results of the experiment for binomial distribution

Table 1: Comparison of the models on the synthetic data

Model	Cumulative regret		
	Mean	Min	Max
reward Gaussian, $a = 0.5$			
TS-Gau	34950.22	30551.22	38854.90
TS-Lin	29757.08	26442.41	31937.36
LSPS	2180.79	883.14	4445.42
reward Gaussian, $a = 1.0$			
TS-Gau	29532.66	27199.17	31751.64
TS-Lin	587.06	249.70	1121.74
LSPS	2675.45	370.63	5452.13
reward binomial, $a = 0.5$			
TS-Beta	8485.57	5918.24	10979.35
TS-Gau	33996.48	31228.52	35459.35
TS-Lin	5850.11	1392.79	10219.05
LSPS	11887.84	146.80	28273.94
reward binomial, $a = 1.0$			
TS-Beta	7355.23	6493.95	8359.40
TS-Gau	29788.95	27389.01	32396.54
TS-Lin	537.31	241.48	818.74
LSPS	143.27	80.40	224.21

Table 2: Hyperparameters used in the experiment

Model	Hyperparameters
reward Gaussian, $a = 0.5$	
TS-Gau	-
TS-Lin	$v = 1.0$
LSPS	$\sigma_1 = 1.0, \sigma_2 = 0.1, \sigma_3 = 0.1$
reward Gaussian, $a = 1.0$	
TS-Gau	-
TS-Lin	$v = 1.0$
LSPS	$\sigma_1 = 1.0, \sigma_2 = 0.1, \sigma_3 = 0.1$
reward binomial, $a = 0.5$	
TS-Beta	-
TS-Gau	-
TS-Lin	$v = 1.0$
LSPS	$\sigma_1 = 0.1, \sigma_2 = 0.1, \sigma_3 = 0.1$
reward binomial, $a = 1.0$	
TS-Beta	-
TS-Gau	-
TS-Lin	$v = 1.0$
LSPS	$\sigma_1 = 0.1, \sigma_2 = 0.1, \sigma_3 = 1.0$

Appendix

Bound of the Bayesian regret

Here we will show that the bound for the Bayesian regret provided by (Ou et al. 2019) cannot be correct. First we will notice that the Bayesian regret is non-decreasing function of time horizon T . Later we will show that the proposed formula for the bound approaches 0 as T goes to infinity. These two facts contradict each other, except for the situation when Bayesian regret is always 0. We will show that it is not the case for LSPS algorithm. In order to simplify the notation we will omit all arguments of BayesRegret except time and write BayesRegret(T).

Proposition 4. For any $T \geq 1$:

BayesRegret(T) ≥ 0

BayesRegret(T) is non-decreasing function of T .

Proof.

$$\begin{aligned} \text{Regret}(T) &= \sum_{t=1}^T \mathbb{E}[\max_{i=1,\dots,N} \gamma_i - r_{i,t} | \theta, \gamma_1, \dots, \gamma_N] = \\ &= \sum_{t=1}^T \mathbb{E}[\max_{i=1,\dots,N} \gamma_i | \theta, \gamma_1, \dots, \gamma_N] - \mathbb{E}[r_{i,t} | \theta, \gamma_1, \dots, \gamma_N] = \\ &= \sum_{t=1}^T \max_{i=1,\dots,N} \gamma_i - \mathbb{E}[r_{i,t} | \theta, \gamma_1, \dots, \gamma_N] = \sum_{t=1}^T \max_{i=1,\dots,N} \gamma_i - \gamma_{i_t} \end{aligned}$$

Each element of the sum is greater or equal 0. Its expected value is also greater or equal 0. The Bayesian regret is sum of such T expected values, so it is non-decreasing function of T and BayesRegret(T) ≥ 0 for any T . \square

Proposition 5. For any $k > 0$ and any $\alpha > 0$:

$$\lim_{T \rightarrow +\infty} \alpha T^{-\frac{1}{2}} \log^k(\alpha T^{-\frac{1}{2}}) = 0$$

Proof. Let $u = \frac{T^{\frac{1}{2}}}{\alpha}$. Then the limit is equivalent to:

$$\begin{aligned} \lim_{u \rightarrow +\infty} \frac{\log^k(u^{-1})}{u} &= \lim_{u \rightarrow +\infty} \frac{(-1)^k \log^k(u)}{u} \\ &\stackrel{\text{H}}{=} \lim_{u \rightarrow +\infty} \frac{(-1)^k k \log^{k-1}(u)}{u} \end{aligned}$$

If $k - 1 \leq 0$ we can rewrite it as:

$$\lim_{u \rightarrow +\infty} \frac{(-1)^k k}{u \log^{1-k}(u)} = 0$$

If $k - 1 > 0$ we can apply L'hospital's rule $k-1$ more times until the exponent above $\log(u)$ is smaller or equal 0. Then we can rewrite the expression as in the first case and get 0 in the limit. \square

Proposition 6. It is not possible that:

$$\text{BayesRegret}(T) \in \tilde{O}\left(\frac{\sigma_2^2}{\sigma_1^2} N^{\frac{3}{2}} T^{-\frac{1}{2}}\right)$$

unless BayesRegret(T) = 0 for all T .

Proof. Assume, by contradiction, that there exists T_0 for which BayesRegret(T_0) > 0 (from Proposition 4 we know that it cannot be smaller than 0) and the bound is possible.

Denote $\alpha = \frac{\sigma_2^2}{\sigma_1^2} N^{\frac{3}{2}}$. We can use the definitions to get:

$$\exists T' \geq 1, c \geq 0, k \geq 0 :$$

$$\forall T \geq T' : \text{BayesRegret}(T) \leq c * \alpha T^{-\frac{1}{2}} \log^k(\alpha T^{-\frac{1}{2}}) \quad (7)$$

From Proposition 5 we know that right-hand side of the inequality approaches 0, so we know from the definition of limit that for $\frac{\text{BayesRegret}(T_0)}{2}$:

$$\exists T'' \geq 1 :$$

$$\forall T \geq T'' :$$

$$|c * \alpha T^{-\frac{1}{2}} \log^k(\alpha T^{-\frac{1}{2}}) - 0| \leq \frac{\text{BayesRegret}(T_0)}{2} \quad (8)$$

From the inequalities (7) and (8) we get that:

$$\begin{aligned} \forall T \geq \max(T', T'', T_0) : \\ \text{BayesRegret}(T) &\leq c * \alpha T^{-\frac{1}{2}} \log^k(\alpha T^{-\frac{1}{2}}) \\ &\leq |c * \alpha T^{-\frac{1}{2}} \log^k(\alpha T^{-\frac{1}{2}})| \leq \frac{\text{BayesRegret}(T_0)}{2} \\ &< \text{BayesRegret}(T_0) \end{aligned}$$

We obtained that for $T > T_0$:

$$\text{BayesRegret}(T) < \text{BayesRegret}(T_0)$$

which contradicts the Proposition 4. \square

Proposition 7. In case of the model formulation presented in section Assumptions:

$$\text{BayesRegret}(T_1) > 0$$

Proof.

$$\text{BayesRegret}(T_1) = \mathbb{E} \left[\mathbb{E} \left[\max_{i=1,\dots,N} \gamma_i - r_{i,1} | \theta, \gamma_1, \dots, \gamma_N \right] \middle| \theta \right]$$

Calculating the inner expected value we get:

$$\begin{aligned} \mathbb{E}[\max_{i=1,\dots,N} \gamma_i - r_{i,1} | \theta, \gamma_1, \dots, \gamma_N] &= \max_{i=1,\dots,N} \gamma_i \\ &- \sum_{i=1,\dots,N} \mathbb{E}[r_i | i_1 = i, \theta, \gamma_1, \dots, \gamma_N] * \mathbb{P}(i_1 = i | \theta, \gamma_1, \dots, \gamma_N) \end{aligned} \quad (9)$$

For $t = 1$ distributions used by the LSPS algorithm are the same for all arms, so: $\mathbb{P}(i_1 = i | \theta, \gamma_1, \dots, \gamma_N) = \frac{1}{N}$. Also $r_i | \gamma_i \sim \mathcal{N}(\gamma_i, \sigma_1^2)$. Using these facts we can write (9) as:

$$\max_{i=1,\dots,N} \gamma_i - \frac{1}{N} \sum_{i=1,\dots,N} \gamma_i$$

Value of this expression is always greater or equal zero. It is equal 0 if all γ_i have the same value. They are, conditionally on θ , independent and normally distributed, so probability of this event is 0 and we can infer that:

$$\mathbb{E} \left[\max_{i=1,\dots,N} \gamma_i - \frac{1}{N} \sum_{i=1,\dots,N} \gamma_i \middle| \theta \right] > 0$$

From this it follows that also the most outer \mathbb{E} is greater than 0. \square

References

- Abbasi-Yadkori, Y.; Pál, D.; and Szepesvári, C. 2011. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, 2312–2320.
- Abramowitz, M.; and Stegun, I. A. 1964. Handbook of mathematical functions, Natl. Bur. Stand. Appl. Math. Ser 55: 1046.
- Agrawal, R. 1995. Sample mean based index policies with $O(\log n)$ regret for the multi-armed bandit problem. *Advances in Applied Probability* 1054–1078.
- Agrawal, S.; and Goyal, N. 2013. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, 127–135.
- Agrawal, S.; and Goyal, N. 2017. Near-optimal regret bounds for thompson sampling. *Journal of the ACM (JACM)* 64(5): 1–24.
- Allesiardo, R.; Féraud, R.; and Bouneffouf, D. 2014. A neural networks committee for the contextual bandit problem. In *International Conference on Neural Information Processing*, 374–381. Springer.
- Auer, P. 2002. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research* 3(Nov): 397–422.
- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 47(2-3): 235–256.
- Cappé, O.; Garivier, A.; Maillard, O.-A.; Munos, R.; Stoltz, G.; et al. 2013. Kullback–leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics* 41(3): 1516–1541.
- Chapelle, O.; and Li, L. 2011. An empirical evaluation of thompson sampling. In *Advances in neural information processing systems*, 2249–2257.
- 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.
- Collier, M.; and Llorens, H. U. 2018. Deep Contextual Multi-armed Bandits. *arXiv preprint arXiv:1807.09809*.
- Dani, V.; Hayes, T. P.; and Kakade, S. M. 2008. Stochastic linear optimization under bandit feedback.
- Elmachetoub, A. N.; McNellis, R.; Oh, S.; and Petrik, M. 2017. A practical method for solving contextual bandit problems using decision trees. *arXiv preprint arXiv:1706.04687*.
- Féraud, R.; Allesiardo, R.; Urvoy, T.; and Clérot, F. 2016. Random forest for the contextual bandit problem. In *Artificial Intelligence and Statistics*, 93–101.
- Filippi, S.; Cappé, O.; Garivier, A.; and Szepesvári, C. 2010. Parametric bandits: The generalized linear case. In *Advances in Neural Information Processing Systems*, 586–594.
- Foster, D. J.; Agarwal, A.; Dudík, M.; Luo, H.; and Schapire, R. E. 2018. Practical contextual bandits with regression oracles. *arXiv preprint arXiv:1803.01088*.

- Gentile, C.; Li, S.; Kar, P.; Karatzoglou, A.; Zappella, G.; and Etrue, E. 2017. On context-dependent clustering of bandits. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 1253–1262. JMLR. org.
- Ghosh, A.; Chowdhury, S. R.; and Gopalan, A. 2017. Misspecified linear bandits. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Graepel, T.; Candela, J. Q.; Borchert, T.; and Herbrich, R. 2010. Web-scale bayesian click-through rate prediction for sponsored search advertising in microsoft’s bing search engine. Omnipress.
- Granmo, O.-C. 2010. Solving two-armed Bernoulli bandit problems using a Bayesian learning automaton. *International Journal of Intelligent Computing and Cybernetics*.
- Gupta, S.; Joshi, G.; and Yağın, O. 2020. Correlated multi-armed bandits with a latent random source. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 3572–3576. IEEE.
- Jun, K.-S.; Bhargava, A.; Nowak, R.; and Willett, R. 2017. Scalable generalized linear bandits: Online computation and hashing. In *Advances in Neural Information Processing Systems*, 99–109.
- Kaufmann, E.; Cappé, O.; and Garivier, A. 2012. On Bayesian upper confidence bounds for bandit problems. In *Artificial intelligence and statistics*, 592–600.
- 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, G.-S.; and Paik, M. C. 2019. Doubly-Robust Lasso Bandit. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems* 32, 5877–5887. Curran Associates, Inc. URL <http://papers.nips.cc/paper/8822-doubly-robust-lasso-bandit.pdf>.
- Krause, A.; and Ong, C. S. 2011. Contextual gaussian process bandit optimization. In *Advances in neural information processing systems*, 2447–2455.
- Krishnamurthy, A.; Wu, Z. S.; and Syrgkanis, V. 2018. Semiparametric contextual bandits. *arXiv preprint arXiv:1803.04204*.
- Kveton, B.; Zaheer, M.; Szepesvari, C.; Li, L.; Ghavamzadeh, M.; and Boutilier, C. 2020. Randomized exploration in generalized linear bandits. In *International Conference on Artificial Intelligence and Statistics*, 2066–2076.
- Lai, T. L.; and Robbins, H. 1985. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics* 6(1): 4–22.
- Li, L.; Chu, W.; Langford, J.; and Schapire, R. E. 2010. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, 661–670. ACM.
- Li, L.; Lu, Y.; and Zhou, D. 2017. Provably optimal algorithms for generalized linear contextual bandits. *arXiv preprint arXiv:1703.00048*.
- Li, S.; Karatzoglou, A.; and Gentile, C. 2016. Collaborative filtering bandits. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, 539–548.
- May, B. C.; and Leslie, D. S. 2011. Simulation studies in optimistic Bayesian sampling in contextual-bandit problems. *Statistics Group, Department of Mathematics, University of Bristol* 11(02).
- Ou, M.; Li, N.; Yang, C.; Zhu, S.; and Jin, R. 2019. Semi-parametric sampling for stochastic bandits with many arms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 7933–7940.
- Peng, Y.; Xie, M.; Liu, J.; Meng, X.; Li, N.; Yang, C.; Yao, T.; and Jin, R. 2019. A practical semi-parametric contextual bandit. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 3246–3252. AAAI Press.
- Perchet, V.; Rigollet, P.; et al. 2013. The multi-armed bandit problem with covariates. *The Annals of Statistics* 41(2): 693–721.
- Rigollet, P.; and Zeevi, A. 2010. Nonparametric bandits with covariates. *arXiv preprint arXiv:1003.1630*.
- Rusmevichientong, P.; and Tsitsiklis, J. N. 2010. Linearly parameterized bandits. *Mathematics of Operations Research* 35(2): 395–411.
- Russo, D.; and Van Roy, B. 2014. Learning to optimize via posterior sampling. *Mathematics of Operations Research* 39(4): 1221–1243.
- Russo, D.; and Van Roy, B. 2016. An information-theoretic analysis of Thompson sampling. *The Journal of Machine Learning Research* 17(1): 2442–2471.
- Russo, D. J.; Roy, B. V.; Kazerouni, A.; Osband, I.; and Wen, Z. 2018. A Tutorial on Thompson Sampling. *Foundations and Trends® in Machine Learning* 11(1): 1–96. ISSN 1935-8237. doi:10.1561/22000000070. URL <http://dx.doi.org/10.1561/22000000070>.
- Scott, S. L. 2010. A modern Bayesian look at the multi-armed bandit. *Applied Stochastic Models in Business and Industry* 26(6): 639–658.
- Slivkins, A. 2011. Contextual bandits with similarity information. In *Proceedings of the 24th annual Conference On Learning Theory*, 679–702. JMLR Workshop and Conference Proceedings.
- Srinivas, N.; Krause, A.; Kakade, S. M.; and Seeger, M. W. 2012. Information-theoretic regret bounds for gaussian process optimization in the bandit setting. *IEEE Transactions on Information Theory* 58(5): 3250–3265.
- Thompson, W. R. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25(3/4): 285–294.

Valko, M.; Korda, N.; Munos, R.; Flaounas, I.; and Cristianini, N. 2013. Finite-time analysis of kernelised contextual bandits. *arXiv preprint arXiv:1309.6869*.

Wang, Q.; Zeng, C.; Zhou, W.; Li, T.; Iyengar, S. S.; Shwartz, L.; and Grabarnik, G. Y. 2018. Online interactive collaborative filtering using multi-armed bandit with dependent arms. *IEEE Transactions on Knowledge and Data Engineering* 31(8): 1569–1580.