

# An Analysis of Ensemble Sampling

Presenter: Hao Liang

217019008@link.cuhk.edu.cn

The Chinese University of Hong Kong, Shenzhen, China

March 1, 2022

Original Paper: Qin, C., Wen, Z., Lu, X., & Van Roy, B. (2022). An Analysis of Ensemble Sampling. arXiv preprint arXiv:2203.01303.

# Outline

Introduction

Preliminaries

Regret Bound

Approximation Error under ES

Appendix

# Thompson Sampling in Online Decision Making

- ▶ Thompson sampling (TS) is an effective heuristic for trading off between exploration and exploitation in online decision making problems
  - bandit: MAB, linear bandit, contextual bandit
  - RL
- ▶ Mechanism
  - Maintain a **posterior** distribution of models (initialized with a prior)
  - Sample models from the posterior distribution as **statistically plausible** models
  - Choose **greedy/optimal** action w.r.t. the statistically plausible models
- ▶ Limitations
  - Requires conjugacy properties, e.g. Beta/Bernoulli, Gaussian/Gaussian
  - Difficult to apply to complex models like neural net

# From Thompson Sampling to Ensemble Sampling

- ▶ Approximate TS algorithms
  - Laplace approximation: limited to unimodal distribution
  - MCMC: computationally expensive for complex models
  - Ensemble sampling (ES)
  - Hypermodel: generalization/variation of ES
- ▶ ES as a practical approximation to TS
  - Fast incremental update/low computational cost
  - Applicable to neural net
- ▶ Applications of variations of ES
  - DRL [Osband et al., 2016, 2018, 2019]
  - Online recommendation [Lu et al., 2018, Hao et al., 2020, Zhu and Van Roy, 2021]
  - MARL [Dimakopoulou and Van Roy, 2018]

## Main contributions

- ▶ No rigorous theory of ES
- ▶ Lu and Van Roy [2017] provided the first regret bound of ES applied to the linear bandit, with a flaw in the analysis
- ▶ Contributions
  - The first rigorous regret analysis of ES for the linear bandit
  - A general Bayesian regret bound for any algorithm for the linear bandit

# Outline

Introduction

**Preliminaries**

Regret Bound

Approximation Error under ES

Appendix

# Linear Gaussian Bandit

- ▶ Reward at time  $t$  and is **linear** in action  $a$  with **Gaussian** noise

$$R_{t,a} = a^\top \theta + W_{t,a}$$

- $a \in \mathcal{A} \subset \mathbb{R}^d$  with  $K = |\mathcal{A}|$
- $\theta \in \mathbb{R}^d$  is the model parameter
- $W_t$  is an i.i.d. sequence with  $W_t \sim N(0, \sigma^2 I_K)$

- ▶ Bayesian framework

$$\theta \sim \mathbb{P}_1(\theta \in \cdot) = N(\mu_0, \Sigma_0)$$

- ▶ At each time  $t$ , the agent chooses  $A_t$  and only observes  $R_{t,A_t}$
- ▶ Bandit/partial feedback:  $R_{t,a}$  for  $a \neq A_t$  not revealed

## Bayesian Regret

- ▶ History at time  $t$

$$\mathcal{H}_t = (A_1, R_{1,A_1}, \dots, A_{t-1}, R_{t-1,A_{t-1}}).$$

- ▶ Given a model  $\theta$ , the optimal action

$$A_* := \arg \max_{a \in \mathcal{A}} \mathbb{E}[R_{t,a} | \theta] = \arg \max_{a \in \mathcal{A}} a^\top \theta$$

- ▶ Frequentist regret

$$\begin{aligned} \text{Regret}(T, \theta) &:= \sum_{t=1}^T \mathbb{E}[R_{t,A_*} - R_{t,A_t} | \theta] \\ &= \sum_{t=1}^T \mathbb{E}[A_*^\top \theta - A_t^\top \theta | \theta]. \end{aligned}$$

- ▶ Bayesian regret

$$\text{Regret}(T) := \mathbb{E}_{\theta \sim \mathbb{P}_1}[\text{Regret}(T, \theta)].$$



# Ensemble Sampling

- ▶ Without conjugacy properties, exact TS becomes computationally infeasible
- ▶ ES serves as a practical approximation to TS

---

## Algorithm 1 Thompson Sampling

---

- 1: **for**  $t \in [T]$  **do**
  - 2:   Sample  $\tilde{\theta}_t \sim \mathbb{P}(\theta \in \cdot | \mathcal{H}_t)$
  - 3:   Execute  $A_t \sim \arg \max_{a \in \mathcal{A}} a^\top \tilde{\theta}_t$
  - 4:   Observe  $R_{t,A_t}$
  - 5:   Update  $\mathbb{P}(\theta \in \cdot | \mathcal{H}_t) \rightarrow \mathbb{P}(\theta \in \cdot | \mathcal{H}_{t+1})$
- 

---

## Algorithm 2 Ensemble Sampling

---

- 1: Sample:  $\tilde{\theta}_{1,1}, \dots, \tilde{\theta}_{1,M} \sim \mathbb{P}_1(\theta \in \cdot)$
  - 2: **for**  $t \in [T]$  **do**
  - 3:   Sample  $m \sim \text{unif}\{1, \dots, M\}$
  - 4:   Execute  $A_t \sim \arg \max_{a \in \mathcal{A}} a^\top \tilde{\theta}_{t,m}$
  - 5:   Observe  $R_{t,A_t}$
  - 6:   Update  $\tilde{\theta}_{t,1:M} \rightarrow \tilde{\theta}_{t+1,1:M}$
-

## Update Details

- ▶ The posterior distribution at time  $t + 1$  is still Gaussian

$$\Sigma_{t+1} = \left( \Sigma_t^{-1} + \frac{1}{\sigma^2} A_t A_t^\top \right)^{-1} \quad \text{and} \quad \mu_{t+1} = \Sigma_{t+1} \left( \Sigma_t^{-1} \mu_t + \frac{R_{t,A_t}}{\sigma^2} A_t \right)$$

- ▶ Satisfies conjugacy properties though
- ▶ ES updates each  $m$ -th model according to

$$\tilde{\theta}_{t+1,m} = \Sigma_{t+1} \left( \Sigma_t^{-1} \tilde{\theta}_{t,m} + \frac{R_{t,A_t} + \tilde{W}_{t,m}}{\sigma^2} A_t \right),$$

where each  $\tilde{W}_t = (\tilde{W}_{t,1}, \dots, \tilde{W}_{t,M}) \sim N(0, \sigma^2 I_M)$  is an **independent** random perturbation.

# Outline

Introduction

Preliminaries

**Regret Bound**

Approximation Error under ES

Appendix

# Regret Bound

- ▶ Regret bound for ES
- ▶ General regret bound for any algorithms
- ▶ From general regret to regret bound for ES

## Regret Bound for ES

### Theorem 1.

Algorithm 2 for linear bandit with prior  $N(\mu_0, \Sigma_0)$  and  $M$  models satisfies

$$\begin{aligned} \text{Regret}(T) &\leq \underbrace{\iota \sqrt{dT \mathbb{H}(A_*)}}_{(a)} + \underbrace{\kappa T \sqrt{\frac{K \log(6TM)}{M}}}_{(b)} \\ &= \tilde{O}\left(\sqrt{dT \mathbb{H}(A_*)} + T \sqrt{K/M}\right), \end{aligned}$$

where  $\mathbb{H}(A_*)$  is the entropy of the optimal action  $A_*$  under the prior, and

$$\iota := \sqrt{2 \left( \max_{a \in \mathcal{A}} a^\top \Sigma_0 a + \sigma^2 \right)}$$

and

$$\kappa := 2 \sqrt{(4 \log K + 5) \max_{a \in \mathcal{A}} a^\top \Sigma_0 a + \max_{a \in \mathcal{A}} (a^\top \mu_0)^2 + \sigma^2} = \tilde{O}(\sqrt{\log K}).$$

## Comparison with the regret bound on TS

The regret bound for ES

$$\text{Regret}(T) \leq \tilde{O} \left( \underbrace{\sqrt{dT\mathbb{H}(A_*)}}_{(a)} + \underbrace{T\sqrt{K/M}}_{(b)} \right)$$

- ▶ Term (a) is exactly the regret bound achieved by TS [Russo and Van Roy, 2016]
- ▶ Term (b) accounts for posterior distribution mismatch
- ▶ As  $M \rightarrow \infty$ , term (b) converges to 0, and the regret bound reduces to that of TS
- ▶ When  $M$  is finite and satisfies  $M = \Omega(KT/d)$ , the regret bound matches TS

## Notations

- ▶ For two discrete distributions  $P = (p_1, \dots, p_n)$  and  $Q = (q_1, \dots, q_n)$ , the KL divergence and Hellinger distance between  $P$  and  $Q$  are defined as

$$\mathbf{d}_{\text{KL}}(P\|Q) := \sum_{i \in [n]} p_i \log(p_i/q_i) \quad \text{and} \quad \mathbf{d}_{\text{H}}(P\|Q) := \sqrt{\sum_{i \in [n]} (\sqrt{p_i} - \sqrt{q_i})^2}$$

- ▶ They satisfy

$$\mathbf{d}_{\text{H}}^2(P\|Q) \leq \min \{ \mathbf{d}_{\text{KL}}(P\|Q), \mathbf{d}_{\text{KL}}(Q\|P) \}$$

## Notations

- ▶ The Shannon entropy of  $X$

$$\mathbb{H}(X) := - \sum_{x \in \mathcal{X}} \mathbb{P}(X = x) \log \mathbb{P}(X = x)$$

- ▶ The entropy of  $X$  conditional on  $Y = y$

$$\mathbb{H}(X | Y = y) := - \sum_{x \in \mathcal{X}} \mathbb{P}(X = x | Y = y) \log \mathbb{P}(X = x | Y = y)$$

- ▶ The conditional entropy of  $X$  given  $Y$

$$\mathbb{H}(X | Y) := \mathbb{E}_Y \left[ - \sum_{x \in \mathcal{X}} \mathbb{P}(X = x | Y) \log \mathbb{P}(X = x | Y) \right]$$

- ▶ The mutual information between  $X$  and  $Y$

$$\mathbb{I}(X; Y) := \mathbf{d}_{\text{KL}}(P(X, Y) \| P(X)P(Y))$$

- ▶ The conditional mutual information between  $X$  and  $Y$  given  $Z$

$$\mathbb{I}(X; Y | Z) := \mathbb{E}_Z [\mathbf{d}_{\text{KL}}(P(X, Y | Z) \| P(X | Z)P(Y | Z))]$$



## General Regret Bound

- ▶ Derive a general Bayesian regret bound for any learning algorithm
- ▶ It is of independent interest and might be used to analyze other bandit algorithms
- ▶ Use subscript  $t$  to denote conditioning on  $H_t$ ,

$$\mathbb{P}_t(\cdot) := \mathbb{P}(\cdot \mid H_t) \quad \text{and} \quad \mathbb{E}_t[\cdot] := \mathbb{E}[\cdot \mid H_t]$$

- ▶ Define

$$\bar{p}_t(\cdot) := \mathbb{P}_t(A_t = \cdot) \quad \text{and} \quad p_t(\cdot) := \mathbb{P}_t(A_* = \cdot)$$

- ▶ Both  $\bar{p}_t$  and  $p_t$  are specified by the algorithm
  - Under TS,  $\bar{p}_t = p_t$
  - For approximate TS,  $\bar{p}_t \approx p_t$

## General Regret Bound

### Theorem 2.

*Any learning algorithm for linear bandit satisfies*

$$\text{Regret}(T) \leq \iota \sqrt{dT \mathbb{H}(A_*)} + \eta \sum_{t=1}^T \sqrt{\mathbb{E}[\mathbf{d}_{\mathbb{H}}^2(\bar{p}_t \| p_t)]},$$

*where*

$$\eta := 2 \sqrt{\mathbb{E} \left[ \max_{a \in \mathcal{A}} (a^\top \theta)^2 \right] + \sigma^2}.$$

*Note that the expectation is taken w.r.t. the prior over  $\theta$ .*

## Discussion on the General Regret Bound

- ▶ The first term matches TS
- ▶ The second term quantifies the cumulative difference between  $\bar{p}_t$  and  $p_t$
- ▶  $\mathbf{d}_{\text{H}}^2(\bar{p}_t \| p_t)$ 
  - vanishes for TS
  - should be small for well approximated TS
- ▶  $\eta = 2\sqrt{\mathbb{E} \left[ \max_{a \in \mathcal{A}} (a^\top \theta)^2 \right] + \sigma^2}$  depends on the prior. For Gaussian prior,
  - $a^\top \theta$  is Gaussian r.v. for any  $a \in \mathcal{A}$
  - the expectation of the maximum of  $K$  squares of Gaussian r.v.s is  $\mathcal{O}(\log K)$
- ▶ Indeed,  $\eta \leq \kappa = 2\sqrt{(4 \log K + 5) \max_{a \in \mathcal{A}} a^\top \Sigma_0 a + \max_{a \in \mathcal{A}} (a^\top \mu_0)^2 + \sigma^2}$

## Regret Bound in terms of KL divergence

- ▶ Analyzing the KL divergence are sometimes easier
- ▶ Recall

$$\mathbf{d}_H^2(P\|Q) \leq \min \{ \mathbf{d}_{\text{KL}}(P\|Q), \mathbf{d}_{\text{KL}}(Q\|P) \}$$

### Corollary 3.

*Under then setting of Theorem 2,*

$$\text{Regret}(T) \leq \iota \sqrt{dT\mathbb{H}(A_*)} + \eta \sum_{t=1}^T \sqrt{\mathbb{E} [\min \{ \mathbf{d}_{\text{KL}}(\bar{p}_t\|p_t), \mathbf{d}_{\text{KL}}(p_t\|\bar{p}_t) \}]}$$

- ▶ Will show that for ES,  $\mathbf{d}_{\text{KL}}(\bar{p}_t\|p_t)$  can be bounded in terms of  $M$

## Proof Sketch of Theorem 2

- ▶ Step 1: rewrite cumulative regret

$$\text{Regret}(T) = \sum_{t=1}^T \mathbb{E} [\mathbb{E}_t [R_{t+1,A_*} - R_{t+1,A_t}]]$$

- ▶ Step 2: regret decomposition as sum of “main regret” and “approximation error”,

$$\mathbb{E}_t [R_{t,A_*} - R_{t+1,A_t}] = G_t + D_t$$

where

$$G_t \triangleq \sum_{a \in \mathcal{A}} \sqrt{\bar{p}_t(a)p_t(a)} (\mathbb{E}_t [R_{t,a} | A_* = a] - \mathbb{E}_t [R_{t+,a}])$$

and

$$D_t \triangleq \sum_{a \in \mathcal{A}} \left( \sqrt{p_t(a)} - \sqrt{\bar{p}_t(a)} \right) \left( \sqrt{p_t(a)} \mathbb{E}_t [R_{t,a} | A_* = a] + \sqrt{\bar{p}_t(a)} \mathbb{E}_t [R_{t,a}] \right)$$

## Proof Sketch of Theorem 2

- Step 3: Bound  $\sum_{t=1}^T \mathbb{E}[G_t]$

$$G_t \leq \iota \sqrt{d \cdot \mathbb{I}_t(A_*; (A_t, R_{t,A_t}))} \quad (\text{information-theoretic})$$

$$\implies \sum_{t=1}^T \mathbb{E}[G_t] \leq \iota \sqrt{dT\mathbb{H}(A_*)} \quad (\text{Cauchy-Schwartz inequality + chain rule})$$

- Step 4: Bound  $\sum_{t=1}^T \mathbb{E}[G_t]$

$$\mathbb{E}[D_t] \leq \eta \sqrt{\mathbb{E}[\mathbf{d}_H^2(\bar{p}_t \| p_t)]} \quad (\text{Cauchy-Schwartz inequality})$$

# Outline

Introduction

Preliminaries

Regret Bound

Approximation Error under ES

Appendix

$$\text{Bound } \sum_{t=1}^T \sqrt{\mathbb{E} [\mathbf{d}_H^2 (\bar{p}_t \| p_t)]}$$

**Lemma 4.**

*Under ES, for all  $t \in [T]$ ,*

$$\mathbb{E} [\mathbf{d}_{KL} (\bar{p}_t \| p_t)] \leq \frac{K \log(6(t+1)M)}{M}.$$

- ▶ Plugging Lemma 4 into the general regret bound in Theorem 2, we achieve the regret bound for ES in Theorem 1



# Outline

Introduction

Preliminaries

Regret Bound

Approximation Error under ES

Appendix

## Proof of Lemma 4

- ▶ ES first **uniformly** samples  $m \in [M]$ , and then samples the action  $A_t$  corresponding to  $\tilde{\theta}_{t,m}$  **uniformly** from the optimal action set

$$\tilde{\mathcal{A}}_{t,m} := \arg \max_{a \in \mathcal{A}} a^\top \tilde{\theta}_{t,m}$$

- ▶ Define the following approximation of  $p_t(a)$

$$\hat{p}_t(a) := \frac{1}{M} \sum_{m=1}^M \frac{1}{|\tilde{\mathcal{A}}_{t,m}|} \mathbb{I}\{a \in \tilde{\mathcal{A}}_{t,m}\}$$

- ▶ History  $H_t$  does not include  $\tilde{W}_t$ , and

$$\tilde{\theta}_{t,m} = \Sigma_t \left( \Sigma_{t-1}^{-1} \tilde{\theta}_{t-1,m} + \frac{R_{t,A_t} + \tilde{W}_{t,m}}{\sigma^2} A_t \right)$$

- ▶ Given  $H_t$ ,  $\hat{p}_t(a)$  is still random

$$\bar{p}_t(a) = \mathbb{E}_t [\hat{p}_t(a)]$$

## Proof of Lemma 4

- ▶ By convexity of KL divergence, the per-period approximation error

$$\mathbf{d}_{\text{KL}}(\bar{p}_t \| p_t) = \mathbf{d}_{\text{KL}}(\mathbb{E}_t[\hat{p}_t] \| p_t) \leq \mathbb{E}_t[\mathbf{d}_{\text{KL}}(\hat{p}_t \| p_t)]$$

- ▶ Taking expectation on both sides

$$\mathbb{E}[\mathbf{d}_{\text{KL}}(\bar{p}_t \| p_t)] \leq \mathbb{E}[\mathbf{d}_{\text{KL}}(\hat{p}_t \| p_t)]$$



$$\mathbb{E}[\mathbf{d}_{\text{KL}}(\hat{p}_t \| p_t)] = \int_0^\infty \mathbb{P}(\mathbf{d}_{\text{KL}}(\hat{p}_t \| p_t) > \epsilon) d\epsilon$$

- ▶ Derive an upper bound on  $\mathbb{P}(\mathbf{d}_{\text{KL}}(\hat{p}_t \| p_t) > \epsilon)$  for any  $\epsilon > 0$

## Proof of Lemma 4

- ▶ For simplicity, consider **deterministic** action sequence  $a_{1:t} := (a_1, \dots, a_t)$ . Write

$$p_t^{a_{1:t-1}}(\cdot) := \mathbb{P}_t (A_* = \cdot \mid a_1, R_{1,a_1}, \dots, a_{t-1}, R_{t,a_{t-1}})$$

- ▶ Under **deterministic** action sequence  $a_{1:t}$

$$\tilde{\theta}_{t,1}^{a_{0:t-1}}, \dots, \tilde{\theta}_{t,M}^{a_{0:t-1}} \mid R_{1,a_1}, \dots, R_{t,a_{t-1}} \sim \mathbb{P}_t(\theta \in \cdot)$$

- ▶  $\hat{p}_t^{a_{0:t-1}}$  is an empirical distribution for  $p_t^{a_{0:t-1}}$

**Fact 5 (Sanov's theorem).**

## Relate to stochastic action sequence

- ▶ Using Sanov's theorem

$$\mathbb{P}(\mathbf{d}_{\text{KL}}(\hat{p}_t^{a_{0:t-1}} \| p_t^{a_{0:t-1}}) > \epsilon \mid \theta) \leq (M+1)^K e^{-M\epsilon}$$

- ▶ By applying the union bound over action counts, instead of action sequences

$$\mathbb{P}\left(\max_{a_{0:t-1} \in \mathcal{A}^t} \mathbf{d}_{\text{KL}}(\hat{p}_t^{a_{0:t-1}} \| p_t^{a_{0:t-1}}) > \epsilon \mid \theta\right) \leq (t+1)^K (M+1)^K e^{-M\epsilon}$$

- ▶ Relate to the KL divergence associated with **stochastic** action sequence

$$\mathbb{P}(\mathbf{d}_{\text{KL}}(\hat{p}_t \| p_t) > \epsilon \mid \theta) \leq \mathbb{P}\left(\max_{a_{0:t-1} \in \mathcal{A}^t} \mathbf{d}_{\text{KL}}(\hat{p}_t^{a_{0:t-1}} \| p_t^{a_{0:t-1}}) > \epsilon \mid \theta\right)$$

## Proof

- ▶ Fix  $t$ . For any  $\epsilon > 0$ ,

$$\mathbb{P}(\mathbf{d}_{KL}(\hat{p}_t \| p_t) > \epsilon) = \mathbb{E}[\mathbb{P}(\mathbf{d}_{KL}(\hat{p}_t \| p_t) > \epsilon | \theta) | \theta] \leq (t+1)^K (M+1)^K e^{-M\epsilon}.$$

- ▶ For any threshold  $\delta \geq 0$ ,

$$\begin{aligned}\mathbb{E}[\mathbf{d}_{KL}(\hat{p}_t \| p_t)] &= \int_0^\infty \mathbb{P}(\mathbf{d}_{KL}(\hat{p}_t \| p_t) > \epsilon) d\epsilon \\ &\leq \delta + (t+1)^K (M+1)^K \int_\delta^\infty e^{-M\epsilon} d\epsilon \\ &= \delta + \frac{(t+1)^K (M+1)^K e^{-M\delta}}{M}\end{aligned}$$

- ▶ Choosing the optimal  $\delta^* = \frac{K[\log(t+1) + \log(M+1)]}{M}$

$$\mathbb{E}[\mathbf{d}_{KL}(\hat{p}_t \| p_t)] \leq \frac{K[\log(t+1) + \log(M+1)] + 1}{M} \leq \frac{K \log(6(t+1)M)}{M}$$