

---

# Routine Bandits: Minimizing Regret on Recurring Problems

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study a variant of the multi-armed bandit problem in which a learner faces  
2 every day one of  $B$  many bandit instances, and call it a routine bandit. More  
3 specifically, at each period  $h \in [1, H]$ , the same bandit  $b^h$  is considered during  
4  $T > 1$  consecutive time steps, but the identity  $b^h$  is unknown to the learner. We  
5 assume all rewards distribution are Gaussian standard. Such a situation typically  
6 occurs in recommender systems when a learner may repeatedly serve the same  
7 user whose identity is unknown due to privacy issues. By combining bandit-  
8 identification tests with a KLUCB type strategy, we introduce the KLUCB for  
9 Routine Bandits (KLUCB-RB) algorithm. While independently running KLUCB  
10 algorithm at each period leads to a cumulative expected regret of  $\Omega(H \log T)$   
11 after  $H$  many periods, KLUCB-RB benefits from previous periods by aggregating  
12 observations from similar identified bandits, which yields a non-trivial scaling of  
13  $o(H \log T)$ . We provide numerical illustration that confirm the benefit of KLUCB-  
14 RB.

## 15 1 Introduction

## 16 2 ...

## 17 3 Regret analysis

## 18 4 Sketch of proof

19 **O-A:** bla

20 **A:** bla

21 **H:** bla

22 **L:** bla

23 bla

## 24 5 Numerical experiments

## 25 6 Discussion



27 **References**

28 **A Proof of ...**