
Switching Latent Bandits

Alessio Russo
DEIB, Politecnico di Milano
alessio.russo@polimi.it

Alberto Maria Metelli
DEIB, Politecnico di Milano
albertomaria.metelli@polimi.it

Marcello Restelli
DEIB, Politecnico di Milano
marcello.restelli@polimi.it

Abstract

We consider a Latent Bandit problem where the latent state keeps changing in time according to an underlying Markov Chain and every state is represented by a specific Bandit instance. At each step, the agent chooses an arm and observes a random reward but is unaware of which MAB he is currently pulling. As typical in Latent Bandits, we assume to know the reward distribution of the arms of all the Bandit instances. Within this setting, our goal is to learn the transition matrix determined by the Markov process, so as to minimize the cumulative regret. We propose a technique to solve this estimation problem that exploits the properties of Markov Chains and results in solving a system of linear equations. We present an offline method that chooses the best subset of possible arms that can be used for matrix estimation, and we ultimately introduce the SL-EC learning algorithm based on an Explore Then Commit strategy that builds a belief representation of the current state and optimizes the instantaneous regret at each step. This algorithm achieves a regret of the order $\tilde{O}(T^{2/3})$ with T being the interaction horizon. Finally, we illustrate the effectiveness of the approach and compare it with state-of-the-art algorithms for non-stationary bandits.

1 Introduction

The Multi-Armed Bandit (MAB) framework is a well-known model used for sequential decision-making with little or no information. This framework has been successfully applied in a large number of fields, such as recommender systems, advertising, and networking. In the general MAB formulation, a learner sequentially selects an action among a finite set of different ones. The choice over the arm to select is made by properly balancing the exploration-exploitation trade-off with the goal of maximizing the expected total reward over a horizon T . Standard MAB literature requires the payoff of the available actions to be stationary (i.e., rewards come from a fixed distribution) in order to design efficient no-regret algorithms.

However, in many real-life applications, the stationarity assumption may not necessarily hold as data may be subjected to changes over time. In some applications, it is also possible to identify different data distributions each one corresponding to a specific working regime. In cases of large availability of historical data appearing in the form of past user interactions, it is possible to learn *offline* the observation models associated with the different arms for each working regime. Exploiting the knowledge on observation models leads to many advantages over the *fully online exploration* setting where no prior information is available at the beginning and a massive number of interactions is required to learn the observation models associated with each working regime. Even if the latent regime is not directly observable, by assuming to know the observation distributions, it can be inferred from the interaction process. Identifying the latent state accelerates the adaptation of the agent to the

environment leading to improved performances over time.

Past works focused on this state identification problem under the assumption of knowing the conditional observation models. Some works such as Maillard and Mannor [2014] and Zhou and Brunskill [2016] provided theoretically optimal UCB while others [Hong et al., 2020a] provided more practical algorithms based on Thompson Sampling.

The works cited above assume that the latent state does not change during the interaction process: once the real state is identified, the agent can act optimally. Differently, in this work, we embrace a more realistic scenario and assume that the latent state can change through time. In accordance with the latent bandits setting, we assume that the learning agent is aware of the observation models of the arms conditioned on each latent state. A setting similar to ours has been considered also in Hong et al. [2020b], the key difference is that they assume to have full or partial knowledge of both the observation and the transition models. We instead focus on the more challenging problem of learning the transition model given the knowledge of the observation models and maximizing the cumulative reward over T interaction steps. More specifically, our problem is modeled by assuming the existence of a set \mathcal{S} of different MABs all sharing the same set of finite arms \mathcal{I} , each generating rewards (our observations) in a finite set \mathcal{V} . Each state $s \in \mathcal{S} = \{s_1, \dots, s_S\}$ represents a different instance of a MAB. At each time step t , there is a transition from latent state s_{t-1} to the new latent state s_t according to the transition matrix governing the process. The action a_t selected in t will thus generate a reward conditioned on the latent state s_t .

Our Contribution We summarize here the main aspects and contributions related to this work:

- we design a procedure for the estimation of the transition matrix that converges to the true value under some mild assumptions. In order to obtain this result, we exploit the information derived from the conditional reward models, and we use some properties of Markov Chains;
- we provide high-probability confidence bounds for the proposed procedure using known results from statistical theory and novel estimation bounds of samples coming from Markov Chains;
- we propose the *Switching Latent Explore then Commit* (SL-EC) algorithm that uses the presented estimation method and then exploits the learned information achieving a $\tilde{O}(T^{2/3})$ regret bound on a finite horizon T ;
- we illustrate the effectiveness of the approach and compare it with state-of-the-art algorithms for the non-stationary bandits setting. Numerical simulations are reported in Appendix A.

2 Related Works

Non-stationary Bandits Non-stationary behaviors are closer to real-world scenarios, and this has induced a vast interest in the scientific community leading to the formulation of different methods that consider either abruptly changing environments [Garivier and Moulines, 2011], smoothly changing environments [Trovò et al., 2020], or settings with a bounded variation of the rewards [Besbes et al., 2014]. It is known that when rewards may arbitrarily change over time, the problem of Non-Stationary Bandits is intractable, meaning that only trivial bounds can be derived on the dynamic pseudo-regret. That is the main reason why in the literature there is a large focus on non-stationary settings enjoying some specific structure in order to design algorithms with better guarantees. Non-stationary MAB approaches typically include both passive methods in which arm selection is mainly driven by the most recent feedback [Auer et al., 2019, Besbes et al., 2014, Trovò et al., 2020] and active methods where a change detection layer is used to actively perceive a drift in the rewards and to discard old information [Liu et al., 2017, Cao et al., 2018]. A particular type of non-stationary Bandit problem related to our work includes the *restless Markov* setting [Ortner et al., 2014, Slivkins and Upfal, 2008] where each arm is associated with a different Markov process and the state of each arm evolves independently of the learner’s actions. Differently, Fiez et al. [2018] investigate MAB problems with rewards determined by an unobserved Markov Chain where the transition to the next state depends on the action selected at each time step, while Zhou et al. [2021] focus on MAB problems where the state transition dynamics evolves independently of the chosen action. The key distinction between their work and ours is that they do not assume prior knowledge of the conditional reward models and instead learn them concurrently with the transition matrix. They make use of spectral decomposition techniques [Anandkumar et al., 2014] and use this tool in a regret minimization algorithm achieving

a $\mathcal{O}(T^{2/3})$ regret bound. Their setting is more complex than ours but involves stronger assumptions, like the invertibility of the transition matrix. Furthermore, spectral methods need a vast amount of samples in order to provide reasonable estimation errors and can hardly be used in large problems.

Latent Bandits More similar lines of work are related to bandit studies where latent variables determine the distribution of rewards [Maillard and Mannor, 2014, Zhou and Brunskill, 2016]. In these works, the unobserved state is fixed across different rounds and the conditional rewards depend on the latent state. Maillard and Mannor [2014] developed UCB algorithms without context considering the two different cases in which the conditional rewards are either known or need to be estimated. This line of work has been extended to the contextual bandit case in Zhou and Brunskill [2016] where there is an offline procedure to learn the policies and a selection strategy to use them online. Hong et al. [2020a] proposed a TS procedure in the contextual case that updates a prior probability over the set of states in order to give a higher probability to the real latent state. A non-stationary variant of this setting is proposed in Hong et al. [2020b] where the latent states are assumed to change according to a Markov Chain. They develop TS algorithms under different cases when both the reward and transition models are completely known and when partial information about them is available. For the partial information case, they provide an algorithm based on particle filter which will be used for comparison in the experimental section in Appendix A. Differently from Hong et al. [2020b], we do not assume any prior information about the transition matrix and we learn it through interactions with the environment using the information about the reward models.

3 Switching Latent Bandits

3.1 Preliminaries

Markov Chains A Markov Chain (or Markov Process) [Feller, 1968] over the state space \mathcal{S} is a stochastic process $(S_t)_{t=1}^{\infty}$ satisfying the Markov property, meaning that for all $s_i, s_j \in \mathcal{S}$ and $t > 0$:

$$P(S_{t+1} = s_j | S_t = s_i, \dots, S_0 = s_0) = P(S_{t+1} = s_j | S_t = s_i).$$

More formally, a Markov chain is identified by a tuple $\langle \mathcal{S}, \mathbf{P}, \boldsymbol{\nu} \rangle$ with $\mathcal{S} = \{s_1, \dots, s_S\}$ being a (finite) set of states, \mathbf{P} is a state transition probability matrix with element $P_{ss'} = P(S_{t+1} = s' | S_t = s)$ and $\boldsymbol{\nu} \in \Delta^{\mathcal{S}-1}$ is the initial state distribution with $\nu_s = P(S_0 = s)$. Given the starting distribution $\boldsymbol{\nu}$ and the transition matrix \mathbf{P} , we can define the probability distribution over the state space after n steps as:

$$\boldsymbol{\nu}^{(n)} = \boldsymbol{\nu} \mathbf{P}^n.$$

We can classify Markov Chains according to the different properties they satisfy. In particular, a Markov Chain is *Regular* if some power n of the transition matrix \mathbf{P}^n has only positive elements [Puterman, 1994]. If a Markov Chain is Regular, it admits a unique stationary distribution, as can be seen in the following:

Proposition 3.1. *Let \mathbf{P} be the transition matrix of a Regular Markov Chain and \mathbf{v} an arbitrary probability vector. Then:*

$$\lim_{n \rightarrow \infty} \mathbf{v} \mathbf{P}^n = \boldsymbol{\pi},$$

where $\boldsymbol{\pi}$ is the unique stationary distribution of the chain, and the components of the vector $\boldsymbol{\pi}$ are all strictly positive.

Having established the concept of stationary distribution, we give now another core definition, the one of *spectral gap*, that will be useful for what will follow. Before that, we define the set $(\lambda_i)_{i \in [S]}$ of ordered eigenvalues of \mathbf{P} , with $1 \geq |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_S|$. Assuming to consider a Regular Markov Chain, the system has a unique stationary distribution, and an eigenvalue $\lambda_1 = 1$.

Definition 3.1. *The spectral gap β of a Markov Process defined by transition matrix \mathbf{P} is $1 - |\lambda_2|$.*

The spectral gap provides valuable information about the process. For Regular Markov Chains, the spectral gap controls the rate of exponential decay to the stationary distribution [Saloff-Coste, 1997].

3.2 Problem Formulation

Consider a set \mathcal{S} of $S = |\mathcal{S}|$ different MAB problems. Each MAB has a finite set of discrete arms $\mathcal{I} := \{a_1, \dots, a_I\}$ with cardinality $I = |\mathcal{I}|$ and, by pulling an arm a , it is possible to get a reward r

taken from the set $\mathcal{V} = \{r_1, \dots, r_V\}$ of possible rewards. In our setting, we assume to have a finite set of rewards $V = |\mathcal{V}|$ with each reward $r \in \mathcal{V}$ bounded for simplicity in the range $[0, 1]$. All the considered MABs share the same sets of arms \mathcal{I} and rewards \mathcal{V} . At each step, the MABs alternate according to an underlying Markov Chain having transition probability \mathbf{P} with size $S \times S$. The interaction process is as follows: at each time instant t , the agent chooses an arm $I_t = a$ and observes a reward $R_t = r$ that is determined by the underlying state $S_t = s$ of the process. More formally, the probability associated to this event is

$$Q(r|s, a) := P(R_t = r | S_t = s, I_t = a). \quad (1)$$

For the moment, we will stick with the assumption that the distribution $Q(\cdot|s, i)$ is categorical. In Section 5.1, we will see how continuous distributions can also be handled in this setting. Given all the MABs, the actions and possible observations, we can define the three-dimensional observation tensor \mathbf{O} with size $S \times I \times V$ where the element $O_{s,a,r}$ represents the probability of observing the reward r being in state s and pulling arm a .

In particular, by fixing a state s and an action a , the vector $\mathbf{O}_{s,a,\cdot}$ contains the parameters of the categorical distribution associated with state s and action a . Motivated by the realistic scenario of massive availability of past interaction data in domains such as recommender systems that allows learning the reward models during an offline phase, we make the assumption of knowing the observation tensor \mathbf{O} while our objective is to learn the transition matrix \mathbf{P} that governs the Chain.

3.3 Reference Matrix Definition

We will introduce here some elements whose utility will be clarified in Section 4.

Let's consider the set $\mathcal{C}_S := \{(s_i, s_j) | s_i, s_j \in \mathcal{S}\}$ with $|\mathcal{C}_S| = S^2$ of all the ordered combinations of pairs of states. These combinations identify all the possible state transitions that can be seen from a generic time step t to the successive one $t + 1$. Analogously, we can define the sets $\mathcal{C}_I := \{(a_i, a_j) | a_i, a_j \in \mathcal{I}\}$ with $|\mathcal{C}_I| = I^2$ and $\mathcal{C}_V := \{(r_i, r_j) | r_i, r_j \in \mathcal{V}\}$ with $|\mathcal{C}_V| = V^2$ which are respectively the ordered combinations of pairs of all consecutive arms and of consecutive rewards that can be seen in two contiguous time intervals. From the knowledge of the observation tensor \mathbf{O} and for each $(s_i, s_j) \in \mathcal{C}_S, (a_i, a_j) \in \mathcal{C}_I, (r_i, r_j) \in \mathcal{C}_V$, we are able to compute the following probabilities:

$$P(R_t = r_i, R_{t+1} = r_j | S_t = s_i, S_{t+1} = s_j, I_t = a_i, I_{t+1} = a_j) = O_{s_i, a_i, r_i} O_{s_j, a_j, r_j}. \quad (2)$$

Equation 2 basically allows us to define the probability associated to each possible couple of rewards, actions and states that can occur in consecutive time steps. Hence, by fixing a specific combination of arms (a_h, a_k) from \mathcal{C}_I and by leveraging Equation 2, we can build matrix $\mathbf{H}^{a_h, a_k} \in \mathbb{R}^{V^2 \times S^2}$ where the elements along the rows are associated to combinations in \mathcal{C}_V and the elements along the columns are associated to combinations in \mathcal{C}_S . The element $H_{d,e}^{a_h, a_k}$ contains the value computed in Equation 2 associated to the d -th combination of rewards in \mathcal{C}_V and the e -th combination of states in \mathcal{C}_S assuming to have pulled actions (a_h, a_k) . Having established this procedure to build matrix \mathbf{H}^{a_h, a_k} for the couple of actions (a_h, a_k) , we can now build similar matrices associated with each of the other combinations of arms. By stacking all these matrices together, we get the matrix $\mathbf{A} \in \mathbb{R}^{I^2 V^2 \times S^2}$.

This matrix is a reformulation of the observation tensor \mathbf{O} that expresses the relation between pairs of different elements. The definition of matrix \mathbf{A} will be relevant for the proposed estimation method. In the following, we will refer to the matrix \mathbf{A} also with the name reference matrix.

3.4 Belief Update

As previously said, at each time step t , we only observe the reward realization, but we are unaware of the Bandit instance from which the arm has been pulled. However, it is possible to define a belief representation over the current state by using the information derived from the observation tensor \mathbf{O} and the transition matrix \mathbf{P} defining the Chain.

We introduce a belief vector $\mathbf{b}_t \in \Delta^{S-1}$ representing the probability distribution over the current state at time t . The belief update formulation includes a correction step that adjusts the current belief \mathbf{b}_t using the reward r_t obtained by pulling arm a_t and a prediction step that computes the new belief \mathbf{b}_{t+1} simulating a transition step. The overall update is as follows:

$$\mathbf{b}_{s,t+1} = \frac{\sum_{s'} \mathbf{b}_{s',t} Q(R_t = r_t | S_t = s', I_t = a_t) \mathbf{P}(s|s')}{\sum_{s''} Q(R_t = r_t | S_t = s'', I_t = a_t) \mathbf{b}_{s'',t}}. \quad (3)$$

The choice of the arm to pull is driven, at each step t , by

$$I_t = \arg \max_{a \in \mathcal{I}} \sum_{s \in \mathcal{S}} \sum_{r \in \mathcal{V}} r Q(r|s, a) \mathbf{b}_{s,t}. \quad (4)$$

In this case, the goal is to pull the arm that provides the highest instantaneous expected reward, given the belief representation \mathbf{b}_t of the states.

3.5 Assumptions

We need now to introduce some assumptions that should hold in our setting:

Assumption 3.1. *The smallest element of the transition matrix $\epsilon := \min_{i,j \in \mathcal{S}} P_{i,j} > 0$.*

Assumption 3.2. *The reference matrix $\mathbf{A} \in \mathbb{R}^{I^2 V^2 \times S^2}$ is full column rank.*

Basically, the first assumption gives a non-null probability of transitioning from any state to any other. It is needed for two main reasons. The former is that this assumption implies the regularity of the Chain and, consequently, the presence of a unique stationary distribution, as shown in Proposition 3.1, the latter is mainly a theoretical reason as in our regret analysis we use a result from De Castro et al. [2017] that builds on this condition.

The second assumption, instead, guarantees that the joint distribution of pairs of rewards and pairs of actions given a specific state transition is not the result of a linear combination of the distributions over other state transitions. In the following, we will show that this is a sufficient condition to recover the matrix \mathbf{P} since it makes all state transitions distinguishable from the joint pairs of rewards and actions, and it also implies that $I^2 V^2 \geq S^2$.

4 Proposed Approach

4.1 Markov Chain Estimation

As previously stated, the objective is to learn the transition matrix \mathbf{P} using the observations we get from the different pulled arms assuming to know the tensor $\mathbf{O} \in \mathbb{R}^{S \times I \times V}$. First of all, we start with a consideration about the transition matrix that defines the chain. Building on Assumption 3.1 and following Proposition 3.1, we can say that exists a unique stationary distribution. This distribution can be easily found by solving the equation below:

$$\pi \mathbf{P} = \pi.$$

From the stationary distribution π , we can define the diagonal matrix $\mathbf{\Pi} = \text{diag}(\pi)$ having the values of the stationary distribution along its diagonal, and we can define the matrix $\mathbf{W} = \mathbf{\Pi} \mathbf{P}$ satisfying $\sum_{i,j \in \mathcal{S}} W_{i,j} = 1$. We can see matrix \mathbf{W} as the transition matrix \mathbf{P} where the transition probabilities from each state (reported along the rows of the transition matrix) are scaled by the probability of the state, given by the stationary distribution. Having defined the matrix \mathbf{W} , we can interpret the element $W_{i,j}$ as the probability of seeing the transition from state s_i to state s_j when the two consecutive pairs of states are sampled from the mixed Chain. We will also refer to \mathbf{W} as the stationary transition distribution matrix. Our objective will be to build an estimate $\widehat{\mathbf{W}}$ of the \mathbf{W} matrix from which we will derive $\widehat{\mathbf{P}}$.

Let's now define an exploration policy θ that selects pairs of arms to be played in successive rounds. We use this policy for T_0 episodes on MABs that switch according to the underlying Markov Chain, and we obtain a sequence $\mathcal{D} = \{(a_1, r_1), (a_2, r_2), \dots, (a_{T_0}, r_{T_0})\}$. This sequence can also be represented by combining non-overlapping pairs of consecutive elements, thus obtaining $\text{Pairs}(\mathcal{D}) = \{(a_1, a_2, r_1, r_2), \dots, (a_{T_0-1}, a_{T_0}, r_{T_0-1}, r_{T_0})\}$.

We introduce now the vector $\mathbf{n}_{T_0} \in \mathbb{N}^{I^2 V^2}$ that counts the number of occurrences of elements in $\text{Pairs}(\mathcal{D})$. More formally, for each cell of the vector \mathbf{n}_{T_0} , we have:

$$\mathbf{n}_{T_0}(a_i, a_j, r_i, r_j) = \sum_{t=0}^{T_0/2} \mathbb{1}\{I_{2t} = a_i, I_{2t+1} = a_j, R_{2t} = r_i, R_{2t+1} = r_j\}.$$

Given the previous considerations, we are now ready to state a core result that links the stationary transition distribution matrix \mathbf{W} and the count vector \mathbf{n}_{T_0} as follows:

$$\begin{aligned} \mathbb{E}[\mathbf{n}_{T_0}(a_i, a_j, r_i, r_j)] &= \\ &= \sum_{s_i, s_j} W_{s_i, s_j} \sum_{t=0}^{T_0/2} \theta(I_{2t} = a_i, I_{2t+1} = a_j) P((R_{2t} = r_i, R_{2t+1} = r_j) | (a_i, a_j), (s_i, s_j)). \end{aligned} \quad (5)$$

This equation basically states that a specific couple of rewards will be observed after having pulled a specific couple of arms a number of times which depends on the conditional probabilities of rewards given the couple of arms and each couple of states, weighted by the probability W_{s_i, s_j} that each state transition occurs. We can write the previous formulation in matrix form as follows:

$$\mathbb{E}[\mathbf{n}_{T_0}] = \frac{T_0}{2} \mathbf{D} \mathbf{A} \mathbf{w}, \quad (6)$$

where the matrix \mathbf{A} is the reference matrix already defined in Section 3.3, vector $\mathbf{w} = \text{Vec}(\mathbf{W})$ is the vectorization of the matrix \mathbf{W} , while $\mathbf{D} \in \mathbb{R}^{I^2 V^2}$ is a diagonal matrix containing the probabilities (defined by policy θ) associated to each combination of arms, each appearing with multiplicity V^2 . Having defined Equation 6, we are able to compute an estimate of the vector $\hat{\mathbf{w}}$ based on the obtained vector count \mathbf{n}_{T_0} :

$$\hat{\mathbf{w}} = \mathbf{A}^\dagger \hat{\mathbf{D}}_{T_0}^{-1} \mathbf{n}_{T_0}, \quad (7)$$

where \mathbf{A}^\dagger is the Moore–Penrose inverse of reference matrix \mathbf{A} and matrix $\hat{\mathbf{D}}_{T_0}$ is the diagonal matrix that counts with multiplicity V^2 the number of occurrences of each combination of arms (we assume that each combination of arms has been pulled at least once, so $\hat{\mathbf{D}}_{T_0}$ is invertible). In the limit of infinite samples, Equation 7 has a fixed exact solution that is $\hat{\mathbf{w}} = \mathbf{w}$. After the computation of $\hat{\mathbf{w}}$, we obtain an estimate of $\hat{\mathbf{P}}$. The derivation implies two main steps: the first is to write back the vector $\hat{\mathbf{w}}$ in matrix form, reversing the vectorization operation and obtaining matrix $\hat{\mathbf{W}}$; the second step consists in normalizing each obtained row so that the values on each row sum to 1, thus deriving $\hat{\mathbf{P}}$.

4.2 SL-EC Algorithm

Having established an estimation procedure for the transition matrix $\hat{\mathbf{P}}$, we will now provide an algorithm that makes use of this approach in a regret minimization framework.

We consider a finite horizon T for our problem. We propose an algorithm called *Switching Latent Explore then Commit* (SL-EC) that proceeds using an EC approach where the exploration phase is devoted to finding the best estimation of the transition matrix $\hat{\mathbf{P}}$, while during the exploitation phase, we maximize the instantaneous expected reward using the information contained in the belief state \mathbf{b} with the formulation provided in Equation 4. The Exploration phase lasts for T_0 episodes, where T_0 is optimized w.r.t. the total horizon T , as will be seen in Equation 10. The presented approach is explained in the pseudocode of Algorithm 1.

Basically, a set of all the ordered combinations of pairs of arms is generated at the beginning of the exploration phase, and the pairs of arms are sequentially pulled in a round-robin fashion until the exploration phase is over. The choice of a round-robin approach allows the highlighting of some interesting properties in the theoretical analysis, as will be shown later in Section 5. When the exploration phase is over, an estimation of the transition matrix $\hat{\mathbf{P}}$ is computed using the procedure described in Section 4.1. After that, a belief vector \mathbf{b} is initialized, assigning a uniform probability to all the states, and it is updated using the estimated $\hat{\mathbf{P}}$, considering the history of samples collected during the exploration phase up to T_0 . Finally, the exploitation phase starts, as described in the pseudocode of the algorithm.

4.3 Arm selection policy

In Algorithm 1, we propose a simple approach for choosing the arms to pull. Each ordered combination of pairs of arms is indeed pulled the same number of times during the exploration phase by using a deterministic approach. However, the estimation framework proposed in Section 4.1 allows for a

more flexible arm selection policy. We may randomize the arm choice by assigning non-uniform probabilities to each combination. This aspect allows exploiting the knowledge of the known reward distribution of each arm, for example, giving a higher probability to the combinations of arms that are more rewarding (assuming an initial uniform distribution over state transitions).

Offline arm selection In problems with a large number of available arms, a round-robin approach among all possible combinations of pairs may be detrimental as it needs a longer exploration horizon to properly fill the vector count \mathbf{n} and to have better estimation results.

A more convenient approach, in this case, would be to select a subset of different arms, thus leading to a limited number of combinations of pairs of arms to use during the exploration phase. Clearly, in the general case, the removal of some arms may lead to a loss of the total information available. This is not the case when for example we remove a redundant arm, that is an arm that induces the same reward distribution as another arm, given all the latent states. Intuitively, the arm selection procedure tends to promote diversity among arms and remove redundant or similar ones. It turns out we are able to get an understanding of the information loss we suffer by selecting specific arms, given the knowledge of the reference matrix \mathbf{A} , that we are indeed able to compute beforehand. In particular, in Section 5 devoted to the theoretical analysis, we will see that the expression $\frac{1}{\sigma_{\min}(\mathbf{A})}$,

with $\sigma_{\min}(\mathbf{A})$ representing the minimum singular value of the reference matrix \mathbf{A} , is an index of the complexity of the problem and we can use this value to drive the choice of the best subset of arms to use. In particular, by fixing a number $J < I$ of arms to use among those available, the choice over the best subset of size J can be done as follows. For each possible subset of arms of size J , we can derive a new reference matrix \mathbf{G} from \mathbf{A} , by extracting from the reference matrix the rows associated with arms' combinations that are feasible using the new subset of arms. At this point, a good candidate subset of arms will be the one with the lowest $\frac{1}{\sigma_{\min}(\mathbf{G})}$. Understandably, this approach implies that the new reference matrix \mathbf{G} derived from the subset of selected arms should be full-column rank, thus satisfying Assumption 3.2.

5 Theoretical Analysis

We will now provide theoretical guarantees on the matrix estimation procedure presented in Section 4.1 and we will prove a regret bound for the SL-EC Algorithm.

We start with a concentration bound on the transition matrix $\hat{\mathbf{P}}$ estimated using samples coming from a round-robin collection policy.

Lemma 5.1. *Suppose Assumptions 3.1 and 3.2 hold. By fixing an exploration parameter T_0 and by pulling each combination of pairs of arms in a round-robin fashion, with probability $1 - \delta$ the*

Algorithm 1: SL-EC Algorithm

Input: Reference Matrix \mathbf{A} , Exploration horizon T_0 ,
Total horizon T

Initialize vector of counts $\mathbf{n} \in \mathcal{N}^{I^2V^2}$ with zeroes

$t \leftarrow 0$

$\mathcal{D} \leftarrow \{\}$

while $t \leq T_0$ **do**

foreach $(a_i, a_j) \in I^2$ **do**

 Pull arm $I_t = a_i$

 Observe reward r_t

 Pull arm $I_{t+1} = a_j$

 Observe reward r_{t+1}

 Update \mathbf{n} with $(I_t, I_{t+1}, r_t, r_{t+1})$

$\mathcal{D}.add((I_t, r_t), (I_{t+1}, r_{t+1}))$

$t \leftarrow t + 2$

$\hat{\mathbf{w}} \leftarrow$ Use Equation 7

$\hat{\mathbf{P}} \leftarrow$ Compute Transition Matrix($\hat{\mathbf{w}}$)

$t \leftarrow 0$

$\mathbf{b}_0 \leftarrow Uniform()$

while $t \leq T$ **do**

if $t \leq T_0$ **then**

$I_t = \mathcal{D}.getAction(t)$

else

$I_t = \arg \max_{a \in \mathcal{I}} \sum_{s \in \mathcal{S}} \sum_{r \in \mathcal{V}} rQ(r|s, a)\mathbf{b}_{s,t}$

 Observe reward r_t

$\mathbf{b}_{t+1} \leftarrow UpdateBelief(\mathbf{b}_t, I_t, r_t)$

$t \leftarrow t + 1$

estimation error of the transition matrix \mathbf{P} will be:

$$\|\mathbf{P} - \widehat{\mathbf{P}}\|_F \leq \frac{2I^2V}{\sigma_{\min}(\mathbf{A})\pi_{\min}} \sqrt{\frac{2S \log \frac{2I^2V^2}{\delta}}{(1 - \lambda^{2I^2})T_0}}, \quad (8)$$

where $\|\cdot\|_F$ represents the Frobenius norm [Golub and Van Loan, 1996], σ_{\min} represents the minimum singular value of the reference matrix \mathbf{A} , π_{\min} is the minimum component in the probability vector representing the stationary distribution of the Chain, and λ represents the second highest eigenvalue of matrix \mathbf{P} . We will provide here a sketch of the proof of the presented Lemma. A more detailed version of this proof is reported in Appendix B.

Sketch of the proof The proof of Lemma 5.1 builds on two principal results. The former comprises a relation that links the estimation error of the matrix \mathbf{P} with the estimation error of the stationary transition distribution matrix \mathbf{W} , while the latter is a concentration bound on the estimated $\widehat{\mathbf{W}}$ from the true one \mathbf{W} . Concerning the first result, we can say that:

$$\|\mathbf{P} - \widehat{\mathbf{P}}\|_F \leq \frac{2\sqrt{S}\|\mathbf{W} - \widehat{\mathbf{W}}\|_F}{\pi_{\min}}.$$

This result follows from a sequence of algebraic manipulations, also involving a derivation from [Ramponi et al., 2020].

We now need to define a bound on $\|\mathbf{W} - \widehat{\mathbf{W}}\|_F$. In order to bound this quantity, we apply the vectorization operator $Vec(\cdot)$ to the two matrices obtaining respectively \mathbf{w} and $\widehat{\mathbf{w}}$ and use the fact that $\|\mathbf{W} - \widehat{\mathbf{W}}\|_F = \|\mathbf{w} - \widehat{\mathbf{w}}\|_2$. We proceed as follows:

$$\begin{aligned} \|\mathbf{w} - \widehat{\mathbf{w}}_{T_0}\|_2 &= \left\| \frac{2}{T_0} \mathbf{A}^\dagger \mathbf{D}^{-1} (\mathbb{E}[\mathbf{n}_{T_0}] - \mathbf{n}_{T_0}) \right\|_2 = \|\mathbf{A}^\dagger (\mathbf{z} - \widehat{\mathbf{z}})\|_2 \\ &\leq \|\mathbf{A}^\dagger\|_2 \|\mathbf{z} - \widehat{\mathbf{z}}\|_2 = \frac{1}{\sigma_{\min}(\mathbf{A})} \|\mathbf{z} - \widehat{\mathbf{z}}\|_2, \end{aligned}$$

where in the second equality we replaced the term $(2/T_0)\mathbf{D}^{-1}\mathbb{E}[\mathbf{n}_{T_0}]$ with the vector $\mathbf{z} \in \mathbb{R}^{I^2V^2}$ and similarly for $\widehat{\mathbf{z}}$ using in the expression the observed vector count \mathbf{n}_{T_0} instead of its expectation $\mathbb{E}[\mathbf{n}_{T_0}]$. In the inequality instead, we used the consistency property for the spectral norm of matrix \mathbf{A}^\dagger . Finally, we bound the remaining part as follows:

$$\begin{aligned} \|\mathbf{z} - \widehat{\mathbf{z}}\|_2 &= \sqrt{\sum_{i=1}^{I^2V^2} |z_i - \widehat{z}_i|^2} \leq \sqrt{\sum_{i=1}^{I^2V^2} \frac{(1 + \lambda^{2I^2}) \log \frac{2I^2V^2}{\delta}}{2(1 - \lambda^{2I^2}) \frac{T_0}{2I^2}}} \\ &\leq \sqrt{\frac{I^2V^2(1 + \lambda^{2I^2}) \log \frac{2I^2V^2}{\delta}}{2(1 - \lambda^{2I^2}) \frac{T_0}{2I^2}}} \leq \frac{I^2V}{\sigma_{\min}(\mathbf{A})} \sqrt{\frac{2 \log \frac{2I^2V^2}{\delta}}{(1 - \lambda^{2I^2})T_0}}, \end{aligned}$$

where, on the first inequality, we used Hoeffding's inequality with probability $1 - \frac{\delta}{I^2V^2}$ for each component of the vector $\widehat{\mathbf{z}}$ and a union bound in the second inequality. In our case, in which samples are generated from a Markov Process, we employed a variant of Hoeffding's inequality that accounts for non-independent samples. We utilized the formulation presented in Fan et al. [2021] which incorporates an additional term $\frac{1+\lambda}{1-\lambda}$ in the bound. More details on this can be found in Proposition C.2 in Appendix C. It is important to note that this proposition holds when the starting distribution of the chain corresponds to the stationary distribution $\boldsymbol{\mu}_0 = \boldsymbol{\pi}$, an assumption we can make in our problem. However, if this is not the case, we would suffer a further logarithmic term in the regret (See Theorem 12 in Fan et al. [2021]).

We were able to improve this result by introducing an exponential term $2I^2$ to the second highest eigenvalue λ . This is possible thanks to the adoption of a round-robin procedure for the choice of combinations of arms. Notably, each combination is pulled every $2I^2$ steps of the Markov Process, resulting in a faster mixing of the chain. A more formal result of this aspect can be found in Corollary C.1 in Appendix C.

Having established the results on the estimation matrix \mathbf{P} , we can now provide regret guarantees for Algorithm 1. The oracle we use is aware of both the observation tensor \mathbf{O} and the transition matrix \mathbf{P}

but does not observe the hidden state. As well as our algorithm, it builds a belief over the states, using the formulation defined in Equation 3 and selects the arm maximizing the expected instantaneous reward. The derived regret upper bound is provided in the following:

Theorem 5.1. *Suppose Assumptions 3.1 and 3.2 hold. By considering a finite horizon T , there exists a constant T_0 , with $T > T_0$, such that with probability $1 - \delta$, the regret of the SL-EC Algorithm satisfies:*

$$\mathfrak{R}(T) \leq 2 \left(\frac{LI^2V}{\pi_{\min}\sigma_{\min}(\mathbf{A})} \sqrt{\frac{2S \log \frac{2I^2V^2}{\delta}}{1 - \lambda^{2I^2}}} \cdot T \right)^{2/3}, \quad (9)$$

where L is a constant that depends on the ϵ value appearing in Assumption 3.1 (More details in Appendix C). The presented regret has an order of $\mathcal{O}(T^{2/3})$ w.r.t the horizon T , as common when using an Explore-Then-Commit algorithm. A detailed proof of this theorem can be found in Appendix B. The presented bound on the regret can be achieved by appropriately choosing the exploration horizon T_0 . More specifically, we set it as follows:

$$T_0 = \left(\frac{LTI^2V}{\sigma_{\min}(\mathbf{A})\pi_{\min}} \sqrt{\frac{2S \log \frac{2I^2V^2}{\delta}}{(1 - \lambda^{2I^2})}} \right)^{2/3}. \quad (10)$$

5.1 Continuous Reward Distributions and Dependency on the Number of Observations

By analyzing the results on the bound of the regret, we can observe that it scales with I^2V , which can be concerning for problems with many arms or a large number of observations. To address the high number of arms, we proposed the offline arm selection procedure highlighted in Section 4.3. It is indeed likely that when $I \gg S$, some arms contain redundant information and can be easily discarded for the estimation procedure.

Regarding the number of observations, continuous reward models pose a challenge as the number of observations would be infinite, making the construction of the reference matrix unfeasible. However, we can discretize the distribution into U distinct segments and obtain a matrix with dimension $I^2U^2 \times S^2$. Hence, we assign to each segment a probability value that represents the likelihood of a particular sample originating from the continuous distribution and belonging to that segment.

The discretization of a continuous distribution paves the way for important considerations because the number of different segments U determines the size of the reference matrix \mathbf{A} . In principle, we can choose U such that $U^2 \geq S^2$ and this allows us to estimate the transition matrix even by using a unique combination of arms (as long as Assumption 3.2 is satisfied). It is an interesting problem to determine in this setting the number of suitable splits and the location of the split points that lead to a faster estimation of the transition matrix.

Another issue arises when the environment comprises numerous but finite observations. In such scenarios, we can employ the inverse approach by clustering some observations, thereby reducing the problem's scale. By selecting a number of clusters $C < V$, we can divide the observations into distinct groups. This allows us to utilize cluster-level probabilities (obtained by summing probabilities of the single observations) to construct a new reference matrix and consider counts at the cluster-level for the count vector \mathbf{n} . Of course, this approach may lead to a loss of information due to the clustering procedure but it may be beneficial in scenarios with limited availability of memory resources.

6 Numerical Simulations

In this section, we provide numerical simulations on synthetic data, demonstrating the effectiveness of the proposed Markov Chain estimation procedure. We compare the regret suffered by our SL-EC approach with other algorithms specifically designed for non-stationary environments. Different types of experiments are also reported in Appendix A.

Following the recent work of Zhou et al. [2021], we consider the subsequent baseline algorithms: the simple ϵ -greedy heuristics, a sliding-window algorithm such as *SW-UCB* [Garivier and Moulines, 2011] that is generally able to deal with non-stationary settings and the *Exp3.S* [Auer et al., 2002] algorithm. The parameters for all the baseline algorithms have been properly tuned according to the different considered settings. More details on this can be found in Appendix A. It is worth noting that unlike our Algorithm, the baseline algorithms do not have knowledge of the observation tensor or

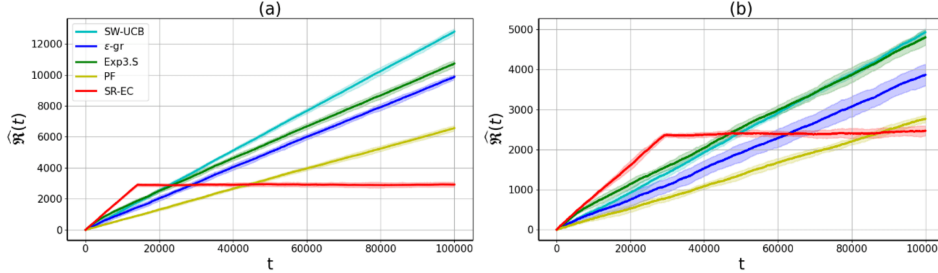


Figure 1: Plots of regret comparing the SL-EC Algorithm with some non-stationary bandit algorithms with a different number of problem parameters: (a) $S = 3, I = 4, V = 5$ (5 runs, 95% c.i.); (b) $S = 8, I = 5, V = 10$. (5 runs, 95% c.i.).

the underlying Markov Chain. In contrast, our approach utilizes the observation tensor to estimate the transition matrix and to update the belief over the current state. Additionally, we compare our approach with a particle filter algorithm proposed in Hong et al. [2020b] about non-stationary Latent Bandits. We provide their algorithm with full information about the observation model (as it is for our case) and an informative prior about the true transition model. The comparison is made in terms of the empirical cumulative regret $\hat{\mathfrak{R}}(t)$, which is the empirical counterpart of the expected cumulative regret $\mathfrak{R}(t)$ averaged over multiple independent runs. The regret results for some problem configurations are shown in Figure 1. Both plots exhibit similar patterns, with most of the baseline algorithms displaying a linear time dependence. This is expected since these algorithms do not take into account the underlying Markov Chain that governs the process. The particle filter algorithm, despite being given a good initial prior on the transition model, is unable to achieve the performance of SL-EC in the long run. Conversely, we can notice a quite different behavior for our algorithm that, in line with an Explore-Then-Commit approach, initially accumulates a large regret and then experiences a drastic slope change when the exploitation phase begins. The regret shown in each plot is the average over all the runs. As a remark, our algorithm outperforms the others when the spectral gap β of the chain is not close to zero. Indeed, if this is not the case, simple exploration heuristics such as ϵ -greedy would lead to comparable performance. A clear example is when the transition matrix \mathbf{P} defining the chain assigns equal probability to all transitions. In this scenario, all states can be considered independent and identically distributed, and we get no advantage from the knowledge of the matrix \mathbf{P} over the use of an algorithm such as ϵ -greedy.

7 Discussion and Conclusions

This paper studies a Latent Bandit problem with latent states changing in time according to an underlying unknown Markov Process. Each state is represented by a different Bandit instance that is unobserved by the agent. As common in the latent Bandit literature, we assumed to know the observation tensor relating each MAB to the reward distribution of its actions, and by using some mild assumptions, we presented a novel estimation technique using the information derived from consecutive pulls of pairs of arms. As far as we know, we are the first to present an estimation procedure of this type aiming at directly estimating the probabilities of the state transitions encoded in the matrix \mathbf{W} . We have shown that our approach is flexible as it allows choosing combinations of pairs of arms with non-uniform probability and easy as it does not require specific hyperparameters to be set. We also provided some offline techniques for the selection of the best subsets of arms to speed up the estimation process. We analyzed the dependence of the parameters on the complexity of the problem and we showed how our approach can be extended to handle models with continuous observation distributions. We used the presented estimation approach in our SL-EC algorithm that uses an Explore-Then-Commit approach and for which we proved a $\mathcal{O}(T^{2/3})$ regret bound. The experimental evaluation confirmed our theoretical findings showing advantages over some algorithms designed for non-stationary MABs and showing good estimation performances even in scenarios with larger problems (Appendix A). A natural future research direction consists of designing new algorithms that are able to exploit the flexibility in the exploration policy determined by the defined procedure, allegedly in an optimistic way.

References

- Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *J. Mach. Learn. Res.*, 15(1):2773–2832, jan 2014.
- Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- Peter Auer, Pratik Gajane, and Ronald Ortner. Adaptively tracking the best bandit arm with an unknown number of distribution changes. In *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 138–158, 2019.
- Kamyar Azizzadenesheli, Alessandro Lazaric, and Anima Anandkumar. Reinforcement learning of pomdps using spectral methods. In *Annual Conference Computational Learning Theory*, 2016.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Stochastic multi-armed-bandit problem with non-stationary rewards. In *Advances in Neural Information Processing Systems*, 2014.
- Yang Cao, Zheng Wen, Branislav Kveton, and Yao Xie. Nearly optimal adaptive procedure with change detection for piecewise-stationary bandit. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- Yohann De Castro, Élisabeth Gassiat, and Sylvain Le Corff. Consistent estimation of the filtering and marginal smoothing distributions in nonparametric hidden markov models. *IEEE Transactions on Information Theory*, 63(8):4758–4777, 2017. doi: 10.1109/TIT.2017.2696959.
- Jianqing Fan, Bai Jiang, and Qiang Sun. Hoeffding’s inequality for general markov chains and its applications to statistical learning. *Journal of Machine Learning Research*, 22(139):1–35, 2021.
- William Feller. *An Introduction to Probability Theory and its Applications Vol. I*. Wiley, 1968.
- Tanner Fiez, Shreyas Sekar, and Lillian J. Ratliff. Multi-armed bandits for correlated markovian environments with smoothed reward feedback. *arXiv: Learning*, 2018.
- Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for switching bandit problems. In *Proceedings of the 22nd International Conference on Algorithmic Learning Theory*, page 174–188, 2011.
- Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, third edition, 1996.
- Joey Hong, Branislav Kveton, Manzil Zaheer, Yinlam Chow, Amr Ahmed, and Craig Boutilier. Latent bandits revisited. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 13423–13433. Curran Associates, Inc., 2020a.
- Joey Hong, Branislav Kveton, Manzil Zaheer, Yinlam Chow, Amr Ahmed, Mohammad Ghavamzadeh, and Craig Boutilier. Non-stationary latent bandits. *CoRR*, abs/2012.00386, 2020b.
- Fang Liu, Joohyung Lee, and Ness B. Shroff. A change-detection based framework for piecewise-stationary multi-armed bandit problem. In *AAAI Conference on Artificial Intelligence*, 2017.
- Odalric-Ambrym Maillard and Shie Mannor. Latent bandits. *31st International Conference on Machine Learning, ICML 2014*, 1, 05 2014.
- Ronald Ortner, Daniil Ryabko, Peter Auer, and Rémi Munos. Regret bounds for restless markov bandits. *Theoretical Computer Science*, 558:62–76, 2014. ISSN 0304-3975. Algorithmic Learning Theory.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley; Sons, Inc., USA, 1st edition, 1994. ISBN 0471619779.

- Giorgia Ramponi, Amarildo Likmeta, Alberto Maria Metelli, Andrea Tirinzoni, and Marcello Restelli. Truly batch model-free inverse reinforcement learning about multiple intentions. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, 2020.
- Laurent Saloff-Coste. *Lectures on finite Markov chains*, pages 301–413. Springer Berlin Heidelberg, Berlin, Heidelberg, 1997. ISBN 978-3-540-69210-2.
- Aleksandrs Slivkins and Eli Upfal. Adapting to a changing environment: the brownian restless bandits. In *Annual Conference Computational Learning Theory*, 2008.
- Francesco Trovò, Stefano Paladino, Marcello Restelli, and Nicola Gatti. Sliding-window thompson sampling for non-stationary settings. *Journal of Artificial Intelligence Research*, 68:311–364, 05 2020. doi: 10.1613/jair.1.11407.
- Li Zhou and Emma Brunskill. Latent contextual bandits and their application to personalized recommendations for new users. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, IJCAI’16, page 3646–3653. AAAI Press, 2016. ISBN 9781577357704.
- Xiang Zhou, Yi Xiong, Ningyuan Chen, and Xuefeng Gao. Regime switching bandits. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

A Additional Numerical Simulations

In this section, we provide additional numerical simulations on synthetic data. Specifically, we show the efficiency of the offline arm selection procedure described in Section 4.3 and conduct a comparison between our approach and a modified technique based on Spectral Decomposition methods in order to highlight the performance difference.

A.1 Estimation Error of Transition Matrix

The first set of experiments is devoted to showing the error incurred by the estimation procedure of the transition matrix in relation to the number of samples considered and the set of actions used for estimation. The left side of Figure 2 illustrates the estimation error of the transition matrix given different instances of Switching Bandits with increasing number of states. In particular, we fix the number of total actions $I = 10$ and number of observations $V = 10$ and consider three instances with $S = 5$, $S = 10$ and $S = 15$ number of states. As it is expected, we can see that as the number of states increases the problem becomes more complex, and more samples are needed in order to improve the estimation. Figure 2 reports the $\|\cdot\|_1$ of the error between the true and the estimated transition matrix, scaled by the number of states. We can see that the estimation procedure is particularly efficient leading to low error values even with a limited number of samples, as can be seen from the steep error drop experienced in the first part of the plot.

The right plot in Figure 2, instead, shows the estimation error obtained by using a different subset of arms. As mentioned in previous sections, it is not always beneficial to use all the available actions during the estimation procedure, but selecting a subset of actions may be preferable. Furthermore, we show that by selecting specific subsets of arms we can improve the estimation w.r.t using other subsets. For this experiment, we consider $J = 3$ arms among the $I = 8$ available for a Switching MAB instance with $S = 5$ states. We then identify the optimal subset of arms of size J and initiate the estimation process using the selected subset. In order to find the best one, we generate all matrices of type \mathbf{G} , as described in Section 4.3 and choose the matrix with lowest $\frac{1}{\sigma_{\min}(\mathbf{G})}$. The subset of arms generating that matrix will be used for estimation. The estimation error of the best subset of arms is represented in the plot with the red line, while we represent in green the estimation error of the subset having the lowest $\sigma_{\min}(\mathbf{G})$. The figure clearly exhibits the performance difference between the two choices, thereby validating our claims.

Experimental Details For the experiments related in Figure 2, we generated a set of transition and observation matrices with the following characteristics.

- for the plot on the left, we fixed the number $I = 10$ of possible actions and $V = 10$ of finite observations. We then consider the estimation procedure for problems of different sizes with respectively $S = 5$, $S = 10$ and $S = 15$ number of states;
- for the plot on the right, the considered estimated problem has $S = 5$ states, $I = 8$ possible actions, $V = 10$ finite observations.

Starting from the presented parameters, the transition and observation matrices have been generated as follows. An initial version of transition and observation matrices is generated with random elements and, subsequently:

- regarding the transition matrix, we add a tuned diagonal matrix to the initial random version and then normalize. In this way we give more probability on self transitions;
- regarding the observation tensor, for each pair of states and actions, we choose a specific reward that will be drawn with higher probability, in order to avoid having too much stochasticity in the reward distributions.

The scheme just presented is also used for the generation of matrices in the experiments showing the regret of the different algorithms.

For the experiments in the plot on the right, the expression referring to the complexity of the used subset of arms $c = \frac{1}{\sigma_{\min}(\mathbf{G})}$ has values respectively $c_g = 51.5$ for the green plot and $c_r = 13.03$ for the red plot, thus validating the intuition that a large c factor leads to a more difficult estimation procedure.

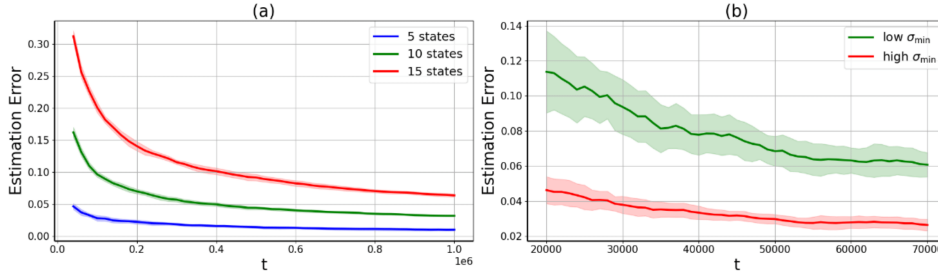


Figure 2: (a) Difference between the estimated and real transition matrix with an increasing number of samples. Metric used is $\|\cdot\|_1$ divided by the number of states (10 runs, 95% c.i.), (b) Difference between real and estimated transition matrix using two different subsets of arms of size $J = 3$ arms from the 8 available on a problem with 5 states. Metric used is $\|\cdot\|_1$ divided by the number of states (10 runs, 95% c.i.).

A.2 Comparisons with Modified Spectral Decomposition Techniques

The focus of this last set of experiments is to show the difference between a modified Spectral Decomposition (SD) technique and our approach. Among the various applications, SD techniques are typically used for learning with Hidden Markov Models (HMM) where no information about the observation and transition model is provided. In particular, Zhou et al. [2021] makes use of these techniques to get an estimation of both the observation and the transition model. It is important to highlight that SD methods are hardly used in practice because of their computational and sample complexity. Indeed, both the related works of Zhou et al. [2021] and Azzadenezsheli et al. [2016] include only proof-of-concept experiments with 2 hidden states and 2 possible actions. Given that our algorithm requires knowledge about the observation model, we consider a slightly different algorithm for performing SD estimation in order to help the estimation process and make the comparison fairer. The original SD technique to which we refer follows the procedures highlighted in Anandkumar et al. [2014] for HMM and makes use of the Robust Tensor Power (RTP) method for orthogonal tensor decomposition. In typical SD techniques, data is collected by sampling an action at each time step and adding the count of the observed reward to the computed statistics. With the presented modified technique, at each step, we do not simply provide the count of the observed reward but the whole reward distribution associated with every arm according to the current state. In this way, it is like pulling at each step all the arms and receiving full information about their associated reward distributions.

We perform various experiments by fixing the number of arms ($I = 20$) and the number of possible rewards ($V = 5$) for each arm and by changing the number of states. Each experiment is performed over 10 different runs, where for each run a transition matrix and observation tensor is generated. For our algorithm, we selected for each experiment 3 arms among the 20 available using our offline arms selection strategy. The transition and observation matrices are created in two different ways: we will see a first set of experiments (Table 1) where the two matrices are almost deterministic, hence having high probability on a specific observation/state and low probabilities for all the others. For transition matrices, the highest probability is assigned to the probability of staying in the same state. From the observation tensor point of view, this problem is easier as it makes the states more distinguishable. In the second set of experiments, the generated matrices have less peaked distributions and higher stochasticity, for both the transition and the observation models (Table 2). What we show in the table is the error in the estimation of the observation matrix with SD techniques (*SD O*), the error for the transition matrix (*SD T*), and the estimation error of our approach (*Our*). For *SD O*, we report the error in 1-norm between the real observation model and the estimated one averaged over all the distributions contained in the observation tensor, we provide this information just to show that indeed our modified procedure allows performing a good estimation of the observation model. This information is indeed separated with a dashed line from the errors in the estimation of the transition distribution, which is our focus in this set of experiments. For *SD T* and *Our* approach, we simply show the 1-norm error between real and estimated transition matrix. For each experiment, we report the mean over 10 runs and one standard deviation between parenthesis and we show in bold the experiments with lower estimation error. Commenting on the results, we can see that the error values of the *SD T* technique are comparable to our approach only in the case of 2 states when the

Table 1: Comparison with Nearly Deterministic Models

2 States	3K samples	6K samples	9K samples	15K samples
<i>SD O</i>	0.0052 (0.0008)	0.0040 (0.0008)	0.0035 (0.0007)	0.0027 (0.0007)
<i>SD T</i>	0.0583 (0.0280)	0.0313 (0.0162)	0.0263 (0.0057)	0.0183 (0.029)
<i>OUR</i>	0.0464 (0.0243)	0.0265 (0.0142)	0.0210 (0.0122)	0.0157 (0.0097)
3 States	150K samples	300K samples	600K samples	900K samples
<i>SD O</i>	0.0038 (0.0013)	0.0026 (0.0005)	0.0017 (0.0003)	0.0015 (0.0002)
<i>SD T</i>	0.3801 (0.1369)	0.3695 (0.1343)	0.3688 (0.1345)	0.3683 (0.1352)
<i>OUR</i>	0.0152 (0.0062)	0.0106 (0.0005)	0.0075 (0.0003)	0.0066 (0.0002)
5 States	150K samples	300K samples	600K samples	900K samples
<i>SD O</i>	0.0104 (0.0024)	0.0075 (0.0015)	0.0053 (0.0017)	0.0042 (0.0010)
<i>SD T</i>	0.9415 (0.2165)	0.9572 (0.2099)	0.9360 (0.2131)	0.9315 (0.2107)
<i>OUR</i>	0.1026 (0.0160)	0.0672 (0.0159)	0.0482 (0.0144)	0.0396 (0.0072)

Table 2: Comparison with Higher Model Stochasticity

2 States	150K samples	210K samples	270K samples
<i>SD O</i>	0.0405 (0.0734)	0.0341 (0.0605)	0.0319 (0.0598)
<i>SD T</i>	0.2900 (0.3061)	0.3273 (0.3364)	0.3077 (0.3326)
<i>OUR</i>	0.0319 (0.0174)	0.0228 (0.0184)	0.0198 (0.0118)
3 States	300K samples	600K samples	900K samples
<i>SD O</i>	0.0531 (0.0409)	0.0510 (0.0367)	0.0415 (0.0366)
<i>SD T</i>	0.9853 (0.7026)	1.1157 (0.6506)	1.0289 (0.6508)
<i>OUR</i>	0.0184 (0.0065)	0.0138 (0.0367)	0.0125 (0.0366)

transition matrix is almost deterministic. Besides the lower performances, the SD technique requires higher computational power and experiments with higher number of states were not able to reach convergence. In particular, experiments with more states and more stochastic models were not able to reach convergence with the number of samples of the order 10^5 , and by increasing their number, there were memory space problems with the used hardware (Intel i7-11th and 16G RAM).

We would like to highlight that SD techniques are explicitly meant to work in a different setting, intrinsically more complex, where no information about either the transition or the observation model is provided. However, we want to show that if instead we have knowledge about the observation model, directly using this information in the SD techniques does not lead to performances comparable to our approach.

A.3 Details for Experiments on Algorithms Comparisons

For the set of experiments on Algorithm Comparison reported in the main body of the work, the parameters used for the generation of the transition and observation matrices are:

- for the plot on the left, the problem has $S = 3$ states, $I = 4$ possible actions, $V = 5$ finite observations;
- for the plot on the right, the problem has $S = 8$ states, $I = 5$ possible actions, $V = 10$ finite observations.

The generation of the matrices is not completely random and follows the same procedure explained in the previous paragraph for the experiment on the matrix estimation error. For the specific experiments considered, we adopted scaled values for the exploration horizon T_0 w.r.t. the result derived from the theory. However, despite a reduced number of samples, the estimation still presents good performances. For the plots shown in the main paper, the hyperparameters used are $\epsilon = 0.05$ for the ϵ -greedy approach, a value of $L_w = 1000$ for the sliding-window UCB algorithm, and the suggested value $1/T$ for the α parameter in the Exp3.S algorithm. For the particle filter algorithm, we used 100 different particles and a resampling threshold of 25 for the *Effective Sample Size*.

B Theoretical Results

In this Section, we will provide the proofs of the Theorems and Lemmas presented in the main paper. We will start by reporting here the main assumptions used for the presented problem.

Assumption 3.1. *The smallest element of the transition matrix $\epsilon := \min_{i,j \in S} P_{i,j} > 0$.*

Assumption 3.2. *The reference matrix $\mathbf{A} \in \mathbb{R}^{I^2 V^2 \times S^2}$ is full column rank.*

We will start by reporting Lemma 5.1 of the main paper and its proof.

Lemma 5.1. *Suppose Assumptions 3.1 and 3.2 hold. By fixing an exploration parameter T_0 and by pulling each combination of pairs of arms in a round-robin fashion, with probability $1 - \delta$ the estimation error of the transition matrix \mathbf{P} will be:*

$$\|\mathbf{P} - \widehat{\mathbf{P}}\|_F \leq \frac{2I^2 V}{\sigma_{\min}(\mathbf{A})\pi_{\min}} \sqrt{\frac{2S \log \frac{2I^2 V^2}{\delta}}{(1 - \lambda^{2I^2})T_0}}, \quad (8)$$

Proof. The proof of the presented bound can be decomposed into two main parts. On one side, we can define the bound of the estimation error of the vector $\widehat{\mathbf{w}}$ from the real value \mathbf{w} and, secondly, the error of the transition matrix \mathbf{P} that derives from the estimated vector $\widehat{\mathbf{w}}$. We first will tackle this last part:

$$\begin{aligned} \|\mathbf{P} - \widehat{\mathbf{P}}\|_F &= \sqrt{\sum_i \sum_j (P_{ij} - \widehat{P}_{ij})^2} = \sqrt{\sum_i \|P_i - \widehat{P}_i\|_2^2} \\ &= \sqrt{\sum_i \left\| \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|_1} - \frac{\widehat{\mathbf{w}}_i}{\|\widehat{\mathbf{w}}_i\|_1} \right\|_2^2} \\ &\leq \sqrt{\sum_i \left\| \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|_2} - \frac{\widehat{\mathbf{w}}_i}{\|\widehat{\mathbf{w}}_i\|_2} \right\|_2^2} \\ &\leq \sqrt{\sum_i \frac{4\|\mathbf{w}_i - \widehat{\mathbf{w}}_i\|_2^2}{\max\{\|\mathbf{w}_i\|_2, \|\widehat{\mathbf{w}}_i\|_2\}^2}} \\ &\leq \sqrt{\frac{4\|\mathbf{W} - \widehat{\mathbf{W}}\|_F^2}{\min_i \max\{\|\mathbf{w}_i\|_2, \|\widehat{\mathbf{w}}_i\|_2\}^2}} \\ &\leq \sqrt{\frac{4S\|\mathbf{W} - \widehat{\mathbf{W}}\|_F^2}{\pi_{\min}^2}} = \frac{2\sqrt{S}\|\mathbf{W} - \widehat{\mathbf{W}}\|_F}{\pi_{\min}}. \end{aligned} \quad (11)$$

In the presented derivation, we have rewritten each row of the matrix \mathbf{P} and $\widehat{\mathbf{P}}$ using, respectively, the vector w_i and \widehat{w}_i of the rows of the stationary transition distribution matrices \mathbf{W} and $\widehat{\mathbf{W}}$, while the first inequality derives from the fact that values in vector w_i are smaller than 1, thus implying $\|\mathbf{w}_i\|_1 \geq \|\mathbf{w}_i\|_2$ and the second inequality is obtained by using Lemma C.1.

The last line is instead derived from the following observations:

$$\|\mathbf{w}_i\|_2^2 = \sum_j w_{ij}^2 = \sum_j P(i,j)^2 \geq \frac{1}{S} \left(\sum_j P(i,j) \right)^2 = \frac{1}{S} \pi(i)^2 \geq \frac{\pi_{\min}^2}{S},$$

where the first inequality in the expression above follows from the fact that $\forall \mathbf{y} \in \mathbb{R}^Y$ in the unit simplex, $\|\mathbf{y}\|_2 \geq \frac{\|\mathbf{y}\|_1}{Y}$.

We will now derive the first part of the proof by defining a high probability bound on the estimated vector of stationary transition distribution $\widehat{\mathbf{w}}$ from its true value. In order to do that, we use the relation $\|\mathbf{W} - \widehat{\mathbf{W}}\|_F = \|\mathbf{w} - \widehat{\mathbf{w}}\|_2$. The bound is obtained assuming that the policy used for the

arms pulls all the combinations of pairs of arms in a round-robin fashion. The derivation is as follows:

$$\begin{aligned}
\|\mathbf{w} - \widehat{\mathbf{w}}_{T_0}\|_2 &= \left\| \frac{2}{T_0} \mathbf{A}^\dagger \mathbf{D}^{-1} (\mathbb{E}[\mathbf{n}_{T_0}] - \mathbf{n}_{T_0}) \right\|_2 \\
&= \|\mathbf{A}^\dagger (\mathbf{z} - \widehat{\mathbf{z}})\|_2 \\
&\leq \|\mathbf{A}^\dagger\|_2 \|\mathbf{z} - \widehat{\mathbf{z}}\|_2 \\
&= \frac{1}{\sigma_{\min}(\mathbf{A})} \sqrt{\sum_{i=1}^{I^2 V^2} |z_i - \widehat{z}_i|^2} \\
&\leq \frac{1}{\sigma_{\min}(\mathbf{A})} \sqrt{\sum_{i=1}^{I^2 V^2} \frac{(1 + \lambda^{2I^2}) \log \frac{2I^2 V^2}{\delta}}{2(1 - \lambda^{2I^2}) \frac{T_0}{2I^2}}} \\
&\leq \frac{1}{\sigma_{\min}(\mathbf{A})} \sqrt{\frac{I^2 V^2 (1 + \lambda^{2I^2}) \log \frac{2I^2 V^2}{\delta}}{2(1 - \lambda^{2I^2}) \frac{T_0}{2I^2}}} \\
&\leq \frac{I^2 V}{\sigma_{\min}(\mathbf{A})} \sqrt{\frac{2 \log \frac{2I^2 V^2}{\delta}}{(1 - \lambda^{2I^2}) T_0}}, \tag{12}
\end{aligned}$$

where in the second equality we have introduced the vector $\mathbf{z} \in \mathbb{R}^{I^2 V^2}$ which can be expressed as $\mathbf{z} = (2/T_0) \mathbf{D}^{-1} \mathbb{E}[\mathbf{n}_{T_0}]$ and can be seen as the expected vector count $\mathbb{E}[\mathbf{n}_{T_0}]$ where each component $\mathbb{E}[n_i]$ associated to a specific combination of pairs of arms, is normalized by the number of times that combination of arms has been pulled. The second inequality is obtained by the consistency property of matrices: the first norm represents the spectral norm of matrix \mathbf{A} , while the second is a $\|\cdot\|_2$ of a vector. The second inequality derives by applying Hoeffding's inequality on each component \widehat{z}_i of the vector $\widehat{\mathbf{z}}$ with probability $1 - \frac{\delta}{I^2 V^2}$, while the multiplicative term $(1 + \lambda^{2I^2})/(1 - \lambda^{2I^2})$ derives from Corollary C.1 and is the result of the dependence among samples originated from the Markov Chain.

The third inequality results from the union bound and holds with probability $1 - \delta$.

By combining the bounds derived in 11 and in 12, we get to the final concentration result. \square

We are now ready to derive the main result related to the regret of the SL-EC Algorithm. We will report here Theorem 5.1 of the main paper.

Theorem 5.1. *Suppose Assumptions 3.1 and 3.2 hold. By considering a finite horizon T , there exists a constant T_0 , with $T > T_0$, such that with probability $1 - \delta$, the regret of the SL-EC Algorithm satisfies:*

$$\mathfrak{R}(T) \leq 2 \left(\frac{L I^2 V}{\pi_{\min} \sigma_{\min}(\mathbf{A})} \sqrt{\frac{2S \log \frac{2I^2 V^2}{\delta}}{1 - \lambda^{2I^2}}} \cdot T \right)^{2/3}, \tag{9}$$

Proof. The proof of the regret of the SL-EC Algorithm makes use of some of the results previously derived and it can be divided into the regret from the exploration and regret from the exploitation phase.

Considering an exploration phase of length T_0 , the regret initially suffered can be trivially bounded as:

$$\mathfrak{R}_{1:T_0} = \sum_{t=1}^{T_0} \max_a \langle \boldsymbol{\mu}_a, \mathbf{b}_t \rangle - r_t \leq \sum_{i=1}^{T_0} 1 = T_0. \tag{13}$$

For the exploitation phase, we use the estimate of the transition matrix $\widehat{\mathbf{P}}$ and use this matrix to define a belief vector that is initialized uniformly over the states and updated starting from the initial samples. Before proceeding with the analysis, we introduce the generic vector $\boldsymbol{\mu}_a$ of size S where the element μ_{a,s_i} referred to state $s_i \in \mathcal{S}$ contains the expected reward of pulling arm a while being

in state s_i . More formally, $\mu_{a,s_i} = \langle \mathbf{O}_{s_i,a,:}, \mathbf{r} \rangle$, with \mathbf{r} being the vector of size V of possible rewards. The analysis of the regret in this part is as follows:

$$\begin{aligned}
\mathfrak{R}_{T_0:T} &= \sum_{t=T_0+1}^T \max_a \langle \boldsymbol{\mu}_a, \mathbf{b}_t \rangle - \max_a \langle \boldsymbol{\mu}_a, \widehat{\mathbf{b}}_t \rangle \\
&\leq \sum_{t=T_0+1}^T \max_a |\langle \boldsymbol{\mu}_a, \mathbf{b}_t - \widehat{\mathbf{b}}_t \rangle| \\
&\leq \sum_{t=T_0+1}^T \|\boldsymbol{\mu}_a\|_\infty \|\mathbf{b}_t - \widehat{\mathbf{b}}_t\|_1 \\
&\leq \sum_{t=T_0+1}^T \|\mathbf{b}_t - \widehat{\mathbf{b}}_t\|_1 \\
&\leq \sum_{t=T_0+1}^T L \|\mathbf{P} - \widehat{\mathbf{P}}_{T_0}\|_F \\
&\leq \frac{2LTI^2V}{\sigma_{\min}(\mathbf{A})\pi_{\min}} \sqrt{\frac{2S \log \frac{2I^2V^2}{\delta}}{(1-\lambda^2I^2)T_0}}, \tag{14}
\end{aligned}$$

where in the second inequality we applied Hölder's inequality with norms ∞ and 1, while the third inequality is obtained from $\|\boldsymbol{\mu}_a\|_\infty \leq 1 \forall a$ and $\forall s \in \mathcal{S}$. The fourth inequality is obtained by applying Proposition C.2, while the last inequality uses the concentration derived in Lemma 5.1.

Combining together the regrets of the two phases derived in Equation equation 13 and in Equation equation 14 we have:

$$\mathcal{R}(T) \leq T_0 + \frac{2LTI^2V}{\sigma_{\min}(\mathbf{A})\pi_{\min}} \sqrt{\frac{2S \log \frac{2I^2V^2}{\delta}}{(1-\lambda^2I^2)T_0}}. \tag{15}$$

We can now optimize this bound w.r.t. the exploration length T_0 by vanishing the derivative of the right-hand side of Equation equation 15. What we get is the following term:

$$T_0 = \left(\frac{LTI^2V}{\sigma_{\min}(\mathbf{A})\pi_{\min}} \sqrt{\frac{2S \log \frac{2I^2V^2}{\delta}}{(1-\lambda^2I^2)}} \right)^{2/3}.$$

In order to be able to compute this T_0 , we need to have information about the minimum value of the stationary distribution π_{\min} and about the second highest eigenvalue λ

By substituting this value of T_0 into Equation equation 15, we get the result of the Theorem. \square

C Useful Lemmas and Deviation Inequalities

This section is devoted to the presentation of some results that are useful in understanding some proofs appearing in AppendixB.

Lemma C.1. (Lemma A.1 in Ramponi et al. [2020]) Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ any pair of vectors, then it holds that:

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \frac{\mathbf{y}}{\|\mathbf{y}\|_2} \right\|_2 \leq \frac{2\|\mathbf{x} - \mathbf{y}\|_2}{\max\{\|\mathbf{x}\|_2, \|\mathbf{y}\|_2\}}$$

Proof. The presented result follows from a sequence of algebraic manipulations:

$$\begin{aligned}
\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \frac{\mathbf{y}}{\|\mathbf{y}\|_2} \right\|_2 &= \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \frac{\mathbf{y}}{\|\mathbf{y}\|_2} \pm \frac{\mathbf{y}}{\|\mathbf{x}\|_2} \right\|_2 \\
&\leq \frac{\|\mathbf{x} - \mathbf{y}\|_2}{\|\mathbf{x}\|_2} + \frac{|\|\mathbf{x}\|_2 - \|\mathbf{y}\|_2|}{\|\mathbf{x}\|_2} \\
&\leq 2 \frac{\|\mathbf{x} - \mathbf{y}\|_2}{\|\mathbf{x}\|_2},
\end{aligned}$$

where the triangular inequality has been applied in the second line and the reverse triangular inequality in the last one, i.e. $\|\mathbf{x}\|_2 - \|\mathbf{y}\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2$. The result in the lemma can be derived by observing that, for symmetry reasons, the same derivation can be performed getting $\|\mathbf{y}\|_2$. \square

Proposition C.1. (Hoeffding's inequality with Markov Chains Fan et al. [2021]) Let $\{X_i\}_{i \geq 1}$ be a Markov Chain with stationary distribution π and absolute spectral gap $1 - \lambda > 0$. For any $t \in \mathbb{R}$, uniformly for all bounded functions $f_i : \mathcal{X} \rightarrow [a_i, b_i]$,

$$\mathbb{E}_\pi [e^{t(\sum_{i=1}^n f_i(X_i) - \sum_{i=1}^n \pi(f_i))}] \leq \exp \left(\frac{1 + \lambda}{1 - \lambda} \cdot \sum_{i=1}^n \frac{(b_i - a_i)^2}{4} \cdot \frac{t^2}{2} \right). \quad (16)$$

It follows that for any $\epsilon > 0$:

$$\mathbb{P}_\pi \left(\sum_{i=1}^n f_i(X_i) - \sum_{i=1}^n \pi(f_i) > \epsilon \right) \leq \exp \left(-\frac{1 - \lambda}{1 + \lambda} \cdot \frac{\epsilon^2}{2 \sum_{i=1}^n (b_i - a_i)^2 / 4} \right). \quad (17)$$

From the previous proposition, we can derive the following corollary:

Corollary C.1. Under the same considerations of Proposition C.1, if we consider n samples from the Markov Chain, each collected after a number k of transitions, we get:

$$\mathbb{P}_\pi \left(\sum_{i=1}^n f_i(X_i) - \sum_{i=1}^n \pi(f_i) > \epsilon \right) \leq \exp \left(-\frac{1 - \lambda^k}{1 + \lambda^k} \cdot \frac{\epsilon^2}{2 \sum_{i=1}^n (b_i - a_i)^2 / 4} \right). \quad (18)$$

Proposition C.2. (Controlling the belief error Zhou et al. [2021]) Assuming to know the smallest value ϵ in a transition matrix \mathbf{P} , given an estimator $\hat{\mathbf{P}}$ of the true transition matrix \mathbf{P} , for an arbitrary reward-action sequence $\{r_{1:t}, a_{i:t}\}_{t \geq 1}$, let $\hat{\mathbf{b}}_t$ and \mathbf{b}_t be the corresponding beliefs in period t under $\hat{\mathbf{P}}$ and \mathbf{P} respectively. Then there exists a constant L such that:

$$\|\hat{\mathbf{b}} - \mathbf{b}\|_1 \leq L \|\hat{\mathbf{P}} - \mathbf{P}\|_F,$$

where $L = \frac{4S(1-\epsilon)^2}{\epsilon^3} + \sqrt{S}$, and $\|\cdot\|_F$ is the Frobenius norm.