# DYNAMICAL LINEAR BANDITS

## MARCO MUSSI, ALBERTO MARIA METELLI AND MARCELLO RESTELLI

{marco.mussi, albertomaria.metelli, marcello.restelli}@polimi.it

## MOTIVATION

- In real-world scenarios, actions leads both to **instantaneous** and **delayed** effects
- Delayed effects can be modeled by means of a **hidden** state
- The hidden state **evolves** depending on the previous hidden state and current **actions**

## CONTRIBUTIONS

- We define **Dynamical Linear Bandits** to represent sequential problems with a hidden state evolving with a **linear unknown dynamics**
  - We show how to represent the problem with a unique **Markov Parameter**
  - We show that the optimal policy is a **constant action**
  - We provide a **lower bound**, also proving that we cannot avoid knowing an upper bound on the **spectral radius** of the transition matrix
- We propose **DynLin-UCB**, an **optimistic regret minimization** algorithm and we provide:
  - a **regret analysis** resulting in $\widetilde{\mathcal{O}}(\sqrt{T})$ expected regret
  - an extensive **numerical validation**

## SETTING - OVERVIEW

### DYNAMICAL LINEAR BANDITS (DLB)

$$
\begin{aligned}
\underbrace{y_t}_{\text{Reward}} &= \underbrace{\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle}_{\substack{\text{Delayed} \\ \text{reward}}} + \underbrace{\langle \boldsymbol{\theta}, \mathbf{u}_t \rangle}_{\substack{\text{Immediate} \\ \text{reward}}} + \underbrace{\eta_t}_{\substack{\text{Reward} \\ \text{noise}}} \\
\underbrace{\mathbf{x}_{t+1}}_{\substack{\text{New} \\ \text{state}}} &= \underbrace{\mathbf{A}\mathbf{x}_t}_{\substack{\text{State} \\ \text{contribution}}} + \underbrace{\mathbf{B}\mathbf{u}_t}_{\substack{\text{Action} \\ \text{contribution}}} + \underbrace{\boldsymbol{\epsilon}_t}_{\substack{\text{State} \\ \text{noise}}}
\end{aligned}
$$

- The state $\mathbf{x}_t \in \mathbb{R}^n$ is **not observable**
- The action $\mathbf{u}_t$ can be chosen in action space $\mathcal{U} \subseteq \mathbb{R}^d$
- $\boldsymbol{\omega}, \boldsymbol{\theta}, \mathbf{A}$, and $\mathbf{B}$ are **unknown**

### ASSUMPTIONS

- Spectral radius: $\rho(\mathbf{A}) < 1$
- $\Phi(\mathbf{A}) = \sup\limits_{\tau \geqslant 0} \|\mathbf{A}^\tau\|_2 / \rho(\mathbf{A})^\tau < \infty$ } **STABILITY**

- $\| \cdot \|_2$ of $\boldsymbol{\theta}, \boldsymbol{\omega}, \mathbf{B}, \mathbf{u}, \mathbf{x}$ bounded
- $\sup\limits_{\mathbf{u}, \mathbf{u}' \in \mathcal{U}} \langle \boldsymbol{\theta}, \mathbf{u} - \mathbf{u}' \rangle \leqslant 1$ } **BOUNDEDNESS**

- $\eta_t$ and $\boldsymbol{\epsilon}_t$ are $\sigma^2$-subgaussian } **SUBGAUSSIANITY**

## SETTING - OPTIMAL POLICY AND LOWER BOUND

### CUMULATIVE MARKOV PARAMETER

$$
\mathbf{h} = \boldsymbol{\theta} + \mathbf{B}^\mathsf{T}(\mathbf{I} - \mathbf{A})^{-\mathsf{T}}\boldsymbol{\omega}
$$

### OPTIMAL POLICY

$$
\mathbf{u}^* \in \arg\max_{\mathbf{u} \in \mathcal{U}} \langle \mathbf{h}, \mathbf{u} \rangle
$$

### EXPECTED AVERAGE REWARD

$$
J := \liminf_{H \to +\infty} \mathbb{E}\left[ \frac{1}{H} \sum_{t=1}^{H} y_t \right]
$$

### ONLINE EXPECTED REGRET LOWER BOUND

$$
\mathbb{E}\, R(\boldsymbol{\pi}, T) = \mathbb{E}\left[ \sum_{t=1}^{T} J^* - y_t \right] \geqslant \Omega\left( \frac{d\sqrt{T}}{(1-\overline{\rho})^{1/2}} \right)
$$

## ALGORITHM - DYNAMICAL LINEAR UPPER CONFIDENCE BOUND (DynLin-UCB)

**DynLin-UCB** is an **optimistic regret minimization** algorithm that operates in **epochs**

- Played action is retrieved using an **optimistic** index
- **Epochs** are of increasing length $H_m$
  - Knowledge of an **upper bound on the spectral radius** $1 > \overline{\rho} \geqslant \rho(\mathbf{A})$
- The selected action is **persisted** for $H_m$ times
- The **regression** estimate Markov parameters $\widehat{\mathbf{h}}_t$ is performed only using the **last sample**
  - i.e., when the hidden state is approximately **steady**

---

**Algorithm 1** DynLin-UCB

Initialize $\mathbf{V}_0 = \lambda \mathbf{I}_d, \mathbf{b}_0 = \mathbf{0}_d, \widehat{\mathbf{h}}_0 = \mathbf{0}_d,$
$\qquad m \leftarrow 1, t \leftarrow 1$
**while** $t < T$ **do**
$\quad$ Compute $\mathbf{u}_t^*$ maximizing $\text{UCB}_t(\mathbf{u})$
$\quad$ Define $H_m = \lceil \log m / \log(1/\overline{\rho}) \rceil$
$\quad$ **for** $j \in \{1, \ldots, H_m\}$ **do**
$\quad\quad$ Play $\mathbf{u}_t^* = \mathbf{u}_{t-1}^*$ and Observe $y_t$
$\quad\quad t \leftarrow t + 1$
$\quad$ **end**
$\quad$ Update $\mathbf{V}_t = \mathbf{V}_{t-1} + \mathbf{u}_t \mathbf{u}_t^\mathsf{T}$
$\qquad\qquad\quad \mathbf{b}_t = \mathbf{b}_{t-1} + \mathbf{u}_t y_t$
$\quad$ Compute $\widehat{\mathbf{h}}_t = \mathbf{V}_t^{-1} \mathbf{b}_t$
$\quad m \leftarrow m + 1$
**end**

---

### OPTIMISTIC ACTIONS CHOICE

$$
\mathbf{u}_t^* = \arg\max_{\mathbf{u} \in \mathcal{U}} \text{UCB}_t(\mathbf{u}) := \langle \widehat{\mathbf{h}}_{t-1}, \mathbf{u} \rangle + \beta_{t-1} \|\mathbf{u}\|_{\mathbf{V}_{t-1}^{-1}}
$$

### BOUND FOR DynLin-UCB

$$
\forall t \in [1, T]: \qquad \beta_t = \frac{\overline{c}_1}{\sqrt{\lambda}} \log(e(t+1)) + \overline{c}_2\sqrt{\lambda} + \sqrt{2\overline{\sigma}^2 \left( \log\left(\frac{1}{\delta}\right) + \frac{d}{2}\log\left(1 + \frac{tU^2}{d\lambda}\right) \right)}
$$

where $\overline{c}_1, \overline{c}_2$, and $\overline{\sigma}^2$ are constants, and $\lambda > 0$ is a regularization parameter

### ONLINE EXPECTED REGRET UPPER BOUND

$$
\mathbb{E}\, R(\texttt{DynLin-UCB}, T) \leqslant \mathcal{O}\left( \frac{d\sigma\sqrt{T}(\log T)^{\frac{3}{2}}}{1-\overline{\rho}} + \frac{\sqrt{dT}(\log T)^2}{(1-\overline{\rho})^{\frac{3}{2}}} + \frac{1}{(1-\rho(\mathbf{A}))^2} \right)
$$

## SIMILAR SETTINGS

Dynamical Linear Bandits can be seen as:

- **Partially Observable Markov Decision Processes** (Littman et al., 1995), in which the state $\mathbf{x}_t$ is non-observable, and the learner has access to the noisy observation $y_t$
- Multiple Input Single Output discrete–time **Linear Time–Invariant System** (Kalman, 1963)
- Non–contextual **Linear Bandit** (Abe and Long, 1999) when the hidden state does not affect the reward, i.e., when $\boldsymbol{\omega} = \mathbf{0}$
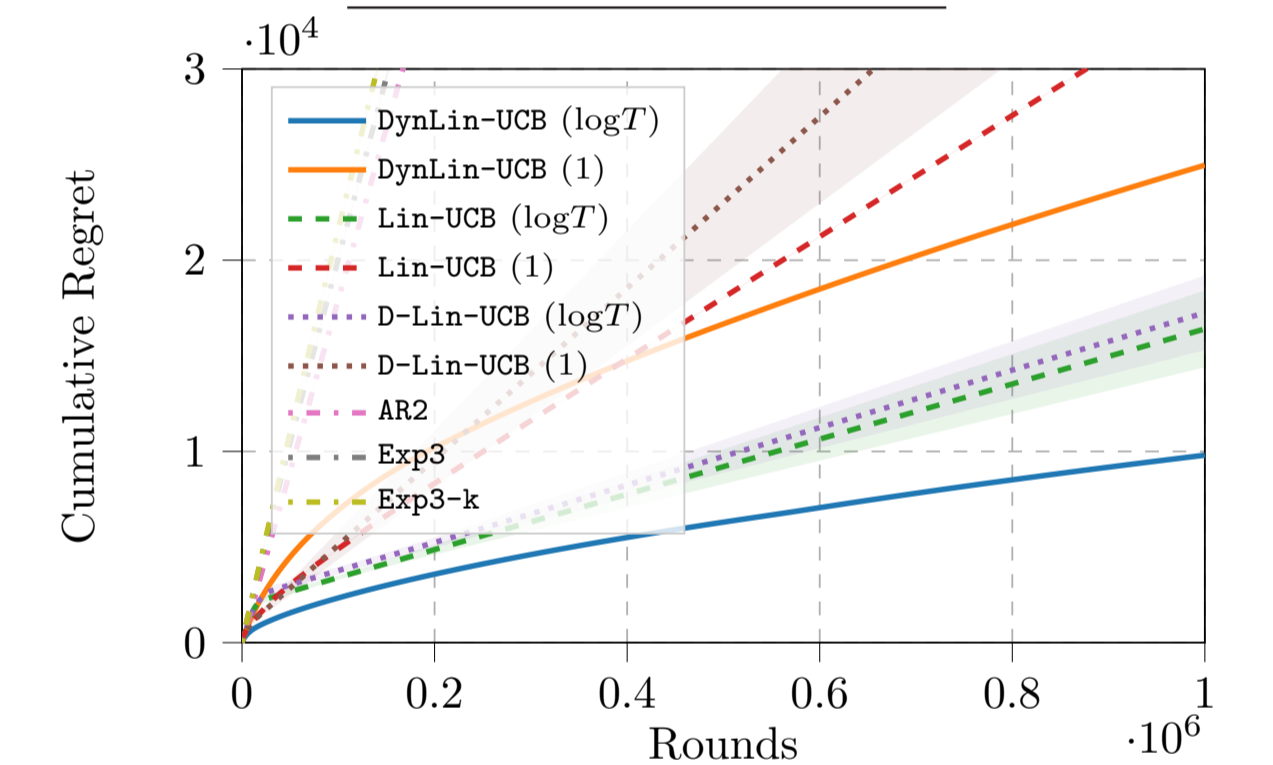
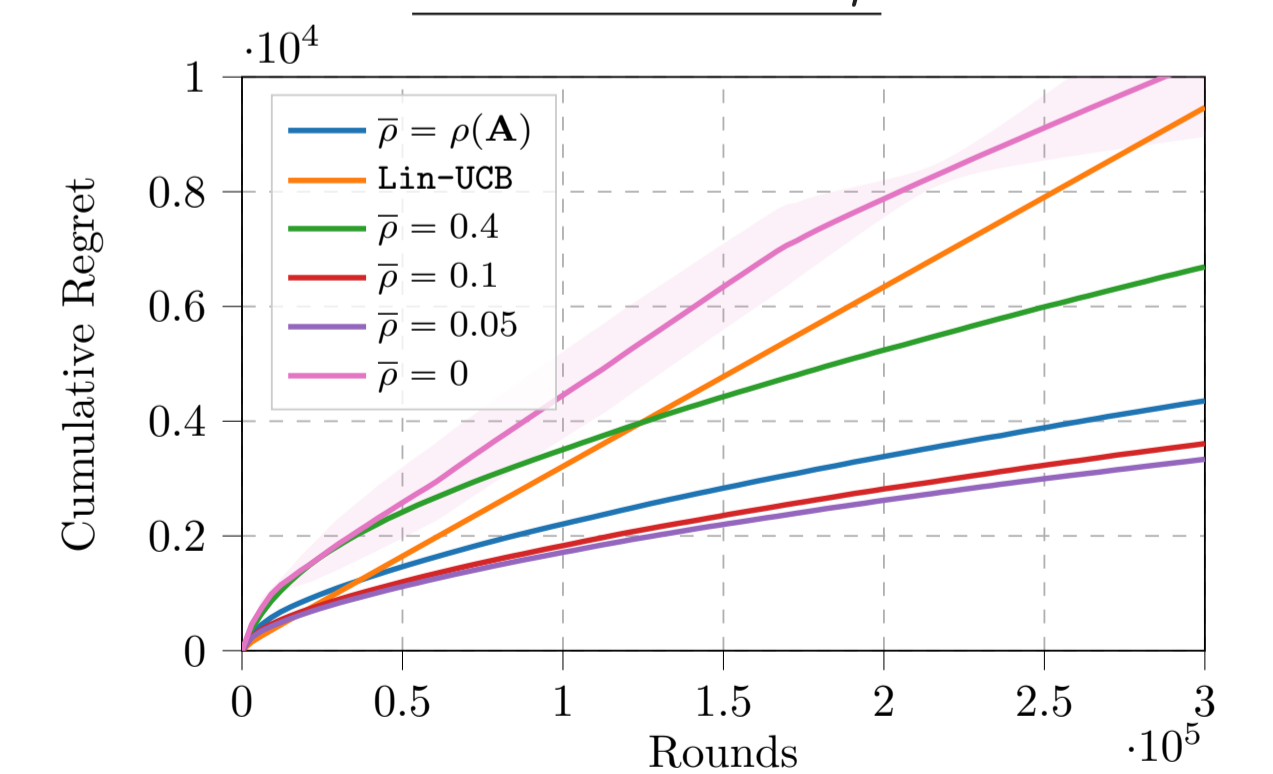## EXPERIMENTAL VALIDATION

### EXPERIMENTAL SETTINGS

$$
\mathbf{A} = \text{diag}(0.2, 0, 0.1) \qquad \mathbf{B} = \text{diag}(0.25, 0, 0.1)
$$
$$
\boldsymbol{\theta} = (0, 0.5, 0.1)^\mathsf{T} \qquad \boldsymbol{\omega} = (1, 0, 0.1)^\mathsf{T} \qquad \sigma^2 = 10^{-2}
$$

### CUMULATIVE REGRET



### SENSITIVITY TO $\overline{\rho}$



## REFERENCES

Naoki Abe and Philip M. Long. Associative reinforcement learning using linear probabilistic concepts. In *ICML*, pages 3–11, 1999.

Rudolf Emil Kalman. Mathematical description of linear dynamical systems. *Journal of the Society for Industrial and Applied Mathematics*, 1(2):152–192, 1963.

Michael L. Littman, Anthony R. Cassandra, and Leslie Pack Kaelbling. Learning policies for partially observable environments: Scaling up. In *ICML*, pages 362–370, 1995.