# A modified Thompson sampling-based learning algorithm for unknown linear systems

Mukul Gagrani[a], Sagar Sudhakara[b], Aditya Mahajan[c], Ashutosh Nayyar[b], Ouyang Yi[d]

[a]Qualcomm, [b]USC, [c]McGill, [d]Preferred Networks

# Significant interest in RL for control

# Significant interest in RL for control


Robotics

# Significant interest in RL for control



Self driving cars

# Significant interest in RL for control



Smart Grids

# Significant interest in RL for control



Smart Grids

**Simplest setting:** Linear quadratic regulation

▷ Different classes of RL algorithms
▷ Provide different performance guarantees under different assumptions on the uncertainty

# Significant interest in RL for control



**Simplest setting:** Linear quadratic regulation

▷ Different classes of RL algorithms
▷ Provide different performance guarantees under
   different assumptions on the uncertainty

Relax the assumptions on uncertainty for a specific class of RL algorithms

1

# Learning in unknown linear systems

## Linear Quadratic Regulation

$$x_{t+1} = A_\theta x_t + B_\theta u_t + w_t, \quad w_t \sim N(0, \sigma_w^2 I)$$

$$c(x_t, u_t) = x_t^\top Q x_t + u_t^\top R u_t.$$

Given $\theta^\top = [A_\theta, B_\theta]$, choose a policy $\pi$ to minimize

$$J(\pi; \theta) = \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=1}^{T} c(x_t, w_t)\right].$$

# Learning in unknown linear systems

## Linear Quadratic Regulation

$$x_{t+1} = A_\theta x_t + B_\theta u_t + w_t, \quad w_t \sim N(0, \sigma_w^2 I)$$
$$c(x_t, u_t) = x_t^\top Q x_t + u_t^\top R u_t.$$

Given $\theta^\top = [A_\theta, B_\theta]$, choose a policy $\pi$ to minimize

$$J(\pi; \theta) = \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=1}^{T} c(x_t, w_t)\right].$$

## Optimal solution

When θ is known and $(A_\theta, B_\theta)$ is stabilizable, optimal policy $\pi^\star$ is given by

$$u_t = G(\theta) x_t$$

where

- $G(\theta) = -(R + B_\theta^\top S_\theta B_\theta)^{-1} B_\theta^\top S_\theta A_\theta$
- $S_\theta$ is the solution of the algebraic Riccati eqn

Moreover: $J(\pi_\theta^\star; \theta) = \sigma_w^2 \text{Tr}(S_\theta)$.

# Learning in unknown linear systems

## Linear Quadratic Regulation

$$x_{t+1} = A_\theta x_t + B_\theta u_t + w_t, \quad w_t \sim N(0, \sigma_w^2 I)$$
$$c(x_t, u_t) = x_t^\top Q x_t + u_t^\top R u_t.$$

Given $\theta^\top = [A_\theta, B_\theta]$, choose a policy $\pi$ to minimize
$$J(\pi; \theta) = \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=1}^{T} c(x_t, w_t)\right].$$

## Learning setup

- True parameter $\theta_\star$ is unknown
- Regret of any learning-based policy $\pi$:
$$R(T; \pi) = \mathbb{E}^\pi\left[\sum_{t=1}^{T} c(x_t, u_t) - TJ(\pi_{\theta_\star}^\star, \theta_\star)\right].$$

## Optimal solution

When $\theta$ is known and $(A_\theta, B_\theta)$ is stabilizable, optimal policy $\pi^\star$ is given by
$$u_t = G(\theta) x_t$$

where
- $G(\theta) = -(R + B_\theta^\top S_\theta B_\theta)^{-1} B_\theta^\top S_\theta A_\theta$
- $S_\theta$ is the solution of the algebraic Riccati eqn

Moreover: $J(\pi_\theta^\star; \theta) = \sigma_w^2 \text{Tr}(S_\theta)$.

Thompson sampling for LQ—(Gagrani et. al.)

# Learning in unknown linear systems

## Linear Quadratic Regulation

$$x_{t+1} = A_\theta x_t + B_\theta u_t + w_t, \quad w_t \sim N(0, \sigma_w^2 I)$$

$$c(x_t, u_t) = x_t^\top Q x_t + u_t^\top R u_t.$$

Given $\theta^\top = [A_\theta, B_\theta]$, choose a policy $\pi$ to minimize

$$J(\pi; \theta) = \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=1}^{T} c(x_t, w_t)\right].$$

## Learning setup

- True parameter $\theta_\star$ is unknown
- Regret of any learning-based policy $\pi$:

$$R(T; \pi) = \mathbb{E}^\pi\left[\sum_{t=1}^{T} c(x_t, u_t) - T J(\pi_{\theta_\star}^\star, \theta_\star)\right].$$

## Optimal solution

When θ is known and $(A_\theta, B_\theta)$ is stabilizable, optimal policy $\pi^\star$ is given by

$$u_t = G(\theta) x_t$$

where

- $G(\theta) = -(R + B_\theta^\top S_\theta B_\theta)^{-1} B_\theta^\top S_\theta A_\theta$
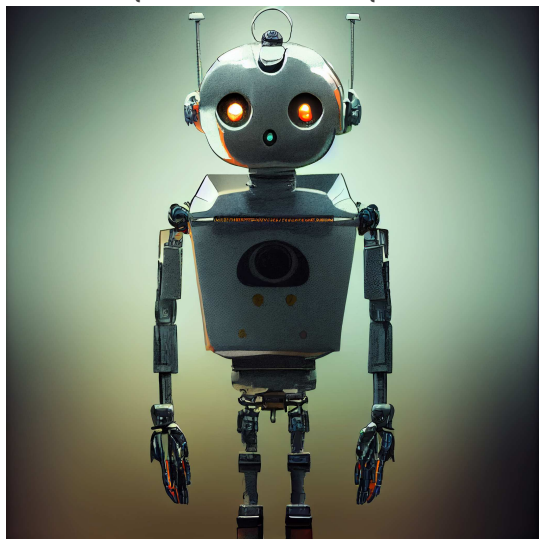- $S_\theta$ is the solution of the algebraic Riccati eqn

Moreover: $J(\pi_\theta^\star; \theta) = \sigma_w^2 \text{Tr}(S_\theta).$

## Key research question

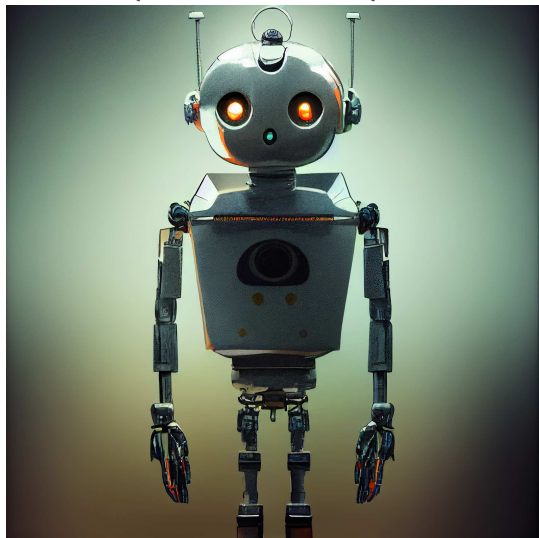▷ How does regret scale with horizon T?

# Different learning frameworks


Explore vs Exploit
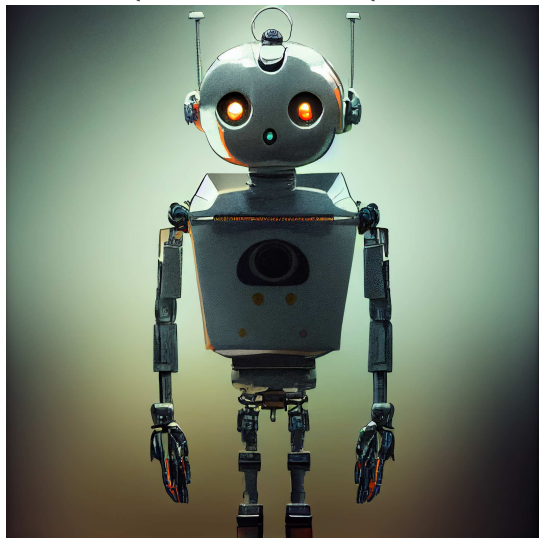
# Different learning frameworks


Explore vs Exploit

## Certainty equivalence

▷ Generate estimate $\hat{\theta}_t$ based on past observations

▷ Use controller: $u_t = G(\hat{\theta}_t)x_t + \varepsilon_t$ (exploration noise)

# Different learning frameworks


Explore vs Exploit

## Certainty equivalence
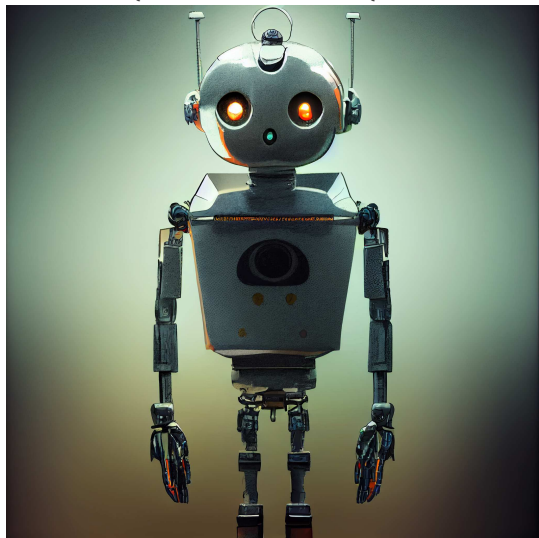
▷ Generate estimate $\hat{\theta}_t$ based on past observations

▷ Use controller: $u_t = G(\hat{\theta}_t)x_t + \varepsilon_t$ (exploration noise)

## Upper Confidence Bound (UCB)

▷ Generate UCB estimate $\bar{\theta}_t$ based on past observations.

▷ Use controller: $u_t = G(\bar{\theta}_t)x_t$

# Different learning frameworks



Explore vs Exploit

## Certainty equivalence

▷ Generate estimate $\hat{\theta}_t$ based on past observations

▷ Use controller: $u_t = G(\hat{\theta}_t)x_t + \varepsilon_t$ (exploration noise)
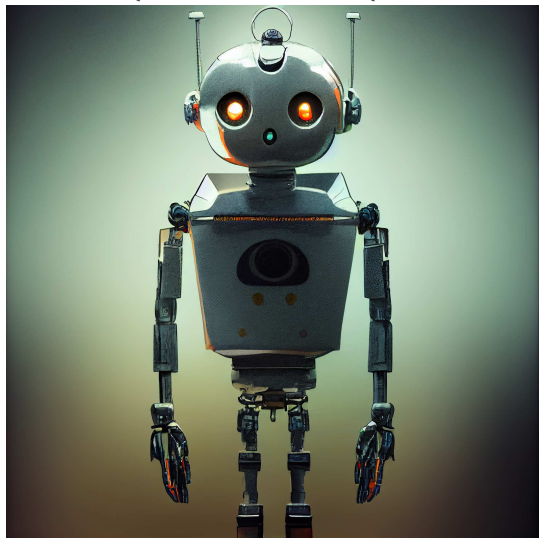
## Upper Confidence Bound (UCB)

▷ Generate UCB estimate $\bar{\theta}_t$ based on past observations.

▷ Use controller: $u_t = G(\bar{\theta}_t)x_t$

## Posterior/Thompson sampling

▷ Maintain posterior $\mu_t$ on $\theta_\star$

▷ Sample $\tilde{\theta}_t \sim \mu_t$

▷ Use controller: $u_t = G(\tilde{\theta}_t)x_t$

# Different learning frameworks


Explore vs Exploit

## Certainty equivalence

▷ Generate estimate $\hat{\theta}_t$ based on past observations

▷ Use controller: $u_t = G(\hat{\theta}_t)x_t + \varepsilon_t$ (exploration noise)

## Upper Confidence Bound (UCB)

▷ Generate UCB estimate $\bar{\theta}_t$ based on past observations.

▷ Use controller: $u_t = G(\bar{\theta}_t)x_t$

## Posterior/Thompson sampling

▷ Maintain posterior $\mu_t$ on $\theta_\star$

▷ Sample $\tilde{\theta}_t \sim \mu_t$

▷ Use controller: $u_t = G(\tilde{\theta}_t)x_t$

# Thompson sampling with dynamic episodes (TSDE)

[Ouyang, Gagrani, Jain 2020]

▷ Bayesian RL algorithm
▷ Generalization on Thompson sampling (or posterior sampling) for bandits

▷ Very simple algorithm which requires no hyper-parameter tuning and works well in practice

# Thompson sampling with dynamic episodes (TSDE)

▷ Bayesian RL algorithm
▷ Generalization on Thompson sampling (or posterior sampling) for bandits

▷ Very simple algorithm which requires no hyper-parameter tuning and works well in practice

## Assumptions on the true parameter

▷ $\theta_\star$ lies in a compact set.

$$\theta_\star^\top = \begin{bmatrix} A_\star & \bigg| & B_\star \end{bmatrix} \in \Omega$$

Compact set

# Thompson sampling with dynamic episodes (TSDE)

▷ Bayesian RL algorithm
▷ Generalization on Thompson sampling (or posterior sampling) for bandits

▷ Very simple algorithm which requires no hyper-parameter tuning and works well in practice

## Assumptions on the true parameter

▷ $\theta_\star$ lies in a compact set.
▷ Independent **truncated Gaussian prior** on each row of $\theta_\star^\top$:

$$\bar{\mu}_1(\theta) = \left[ \prod_{i=1}^{n} N(\hat{\theta}_1(i), \Sigma_1) \right]\Bigg|_\Omega$$

$$\theta_\star^\top = \begin{bmatrix} \rule{3cm}{0pt} \end{bmatrix} \in \Omega$$

Compact set

4

# Thompson sampling with dynamic episodes (TSDE)

## Properties of the posterior

▷ Posterior $\mu_t$ is also truncated Gaussian with $\mu_t(\theta) = \left[\prod_{i=1}^{n} N(\hat{\theta}_t(i), \Sigma_t)\right]\Big|_{\Omega}$ where
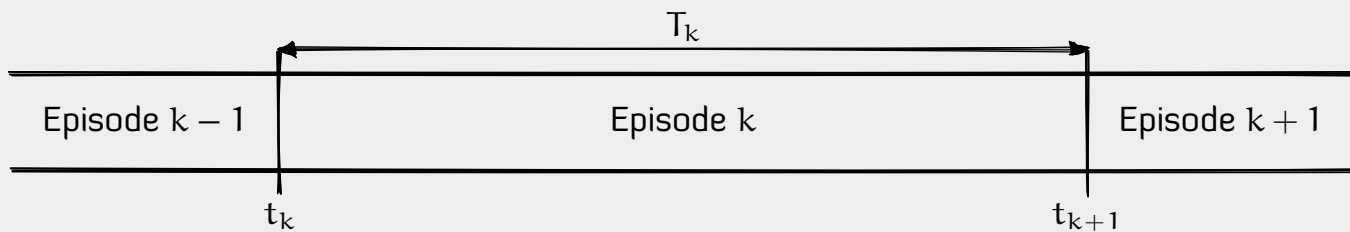
$$\hat{\theta}_{t+1}(i) = \hat{\theta}_t(i) + \frac{\Sigma_t z_t(x_{t+1}(i) - \hat{\theta}_t(i)^\top z_t)}{\sigma_w^2 + z_t^\top \Sigma_t z_t}$$

$$\Sigma_{t+1}^{-1} = \Sigma_t^{-1} + \frac{1}{\sigma_w^2} z_t z_t^\top.$$
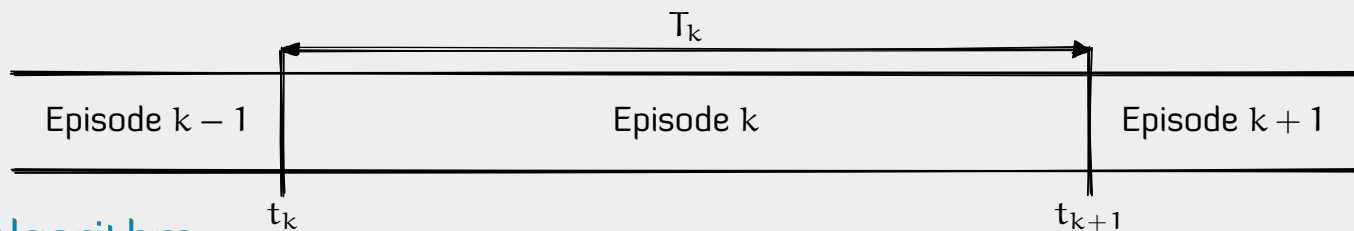
where $z_t = \text{vec}(x_t, u_t)$.

# Thompson sampling with dynamic episodes (TSDE)

# Thompson sampling with dynamic episodes (TSDE)

## TSDE Algorithm

▷  At start of episode:  Sample $\tilde{\theta}_k \sim \mu_{t_k}$

▷  During the episode:  Use $u_t = G(\tilde{\theta}_k)\, x_t$

▷  Terminate episode if:  $(t - t_k > T_{k-1})$ or $(\det \Sigma_t < \frac{1}{2}\det \Sigma_{t_k})$
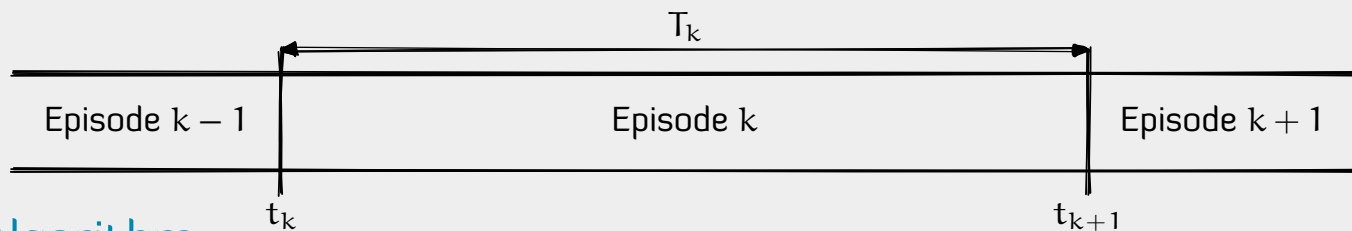
# Thompson sampling with dynamic episodes (TSDE)

## TSDE Algorithm

▷ At start of episode: Sample $\tilde{\theta}_k \sim \mu_{t_k}$

▷ During the episode: Use $u_t = G(\tilde{\theta}_k)\, x_t$

▷ Terminate episode if: $(t - t_k > T_{k-1})$ or $(\det \Sigma_t < \frac{1}{2} \det \Sigma_{t_k})$

**Intuition:** $\det \Sigma_t < \frac{1}{2} \det \Sigma_{t_k}$ implies that current posterior is much better than the posterior at the start of the episode. **Resample to exploit this knowledge**

# Thompson sampling with dynamic episodes (TSDE)

| Assumption A1 | There exists an $\delta \in (0,1)$ such that for any $\theta, \phi \in \Omega$, $\|A_\theta + B_\theta G(\phi)\| \leq \delta$. |
|---|---|

# Thompson sampling with dynamic episodes (TSDE)

| Assumption A1 | There exists an $\delta \in (0,1)$ such that for any $\theta, \phi \in \Omega$, $\|A_\theta + B_\theta G(\phi)\| \leq \delta$. |
|---|---|

## Discussion on assumptions

▷ A1 is a **strong assumption**.

▷ Requires that close loop system dynamics under any mismatched controller should have **spectral norm** less than one.

# Thompson sampling with dynamic episodes (TSDE)

| Assumption A1 | There exists an $\delta \in (0, 1)$ such that for any $\theta, \phi \in \Omega$, $\|A_\theta + B_\theta G(\phi)\| \leq \delta$. |
|---|---|

| Theorem | Under A1, $\quad R(T; \text{TSDE}) \leq C\sqrt{T}\,(\log T)^q$ |
|---|---|

## Discussion on assumptions

▷ A1 is a **strong assumption**.

▷ Requires that close loop system dynamics
under any mismatched controller should have
**spectral norm** less than one.

# Thompson sampling with dynamic episodes (TSDE)

| Assumption A1 | There exists an $\delta \in (0, 1)$ such that for any $\theta, \phi \in \Omega$, $\|A_\theta + B_\theta G(\phi)\| \leq \delta$. |
|---|---|

| Theorem | Under A1, $\quad R(T; \text{TSDE}) \leq C\sqrt{T}(\log T)^q$ |
|---|---|

## Discussion on assumptions

▷ A1 is a **strong assumption**.

▷ Requires that close loop system dynamics under any mismatched controller should have **spectral norm** less than one.

## Discussion on the results

▷ The regret is **Bayesian regret**, i.e., includes an expectation over the prior.

▷ Different from frequentist regret, which provides a high-probability bound on regret for the true parameter.

# Thompson sampling with dynamic episodes (TSDE)

[Ouyang, Gagrani, Jain 2020]

**Assumption** $\leq \delta$.

## Why bother with TSDE

▷ Works very well in practice. Requires no parameter tuning.

▷ Continues to work well when A1 is violated.

**Theorem** Under A1, $\quad R(T; TSDE) \leq C\sqrt{T}\,(\log T)^q$

## Discussion on assumptions

▷ A1 is a **strong assumption**.

▷ Requires that close loop system dynamics under any mismatched controller should have **spectral norm** less than one.

## Discussion on the results

▷ The regret is **Bayesian regret**, i.e., includes an expectation over the prior.

▷ Different from frequentist regret, which provides a high-probability bound on regret for the true parameter.

The strong assumption appears to be a limitation of the proof technique (and not the algorithm).

Can we relax it?

# How should the stability assumption be relaxed?

Ideally, should only require the true $\theta_\star$ to be stabilizable

▷ Bayesian equivalent:

$$\mathbb{P}(\theta \in \Omega : \theta \text{ is stabilizable}) = 1$$

# How should the stability assumption be relaxed?

Ideally, should only require the true $\theta_\star$ to be stabilizable

  ▷ Bayesian equivalent:

  $$\mathbb{P}(\theta \in \Omega : \theta \text{ is stabilizable}) = 1$$

. . .and be able to construct a stabilizing controller in finite time

  ▷ Don't know how to do that in Bayesian setting
  ▷ Guaranteeing stability with high probability is not sufficient

# How should the stability assumption be relaxed?

Ideally, should only require the true $\theta_\star$ to be stabilizable

> ▷ Bayesian equivalent:
>
> $$\mathbb{P}(\theta \in \Omega : \theta \text{ is stabilizable}) = 1$$

. . . and be able to construct a stabilizing controller in finite time

> ▷ Don't know how to do that in Bayesian setting
> ▷ Guaranteeing stability with high probability is not sufficient

First step in weakening the stability assumption

> ▷ Assumption A1 is defined in terms of **spectral norm**
> ▷ A natural relaxation is to replace spectral norm by **spectral radius**.
> ▷ . . . which is what we do in this paper

# This paper: Natural relaxation of Assumption A1

**Assumption A2**     There exists an $\delta \in (0, 1)$ such that for any $\theta, \phi \in \Omega$, $\rho(A_\theta + B_\theta G(\phi)) \leq \delta$.

▷ Controller for system $\phi$ stabilizes system $\theta$

▷ Still a strong assumption, but weaker (and more natural) than A1.

# This paper: Natural relaxation of Assumption A1

| Assumption A2 | There exists an $\delta \in (0, 1)$ such that for any $\theta, \phi \in \Omega$, $\rho(A_\theta + B_\theta G(\phi)) \leq \delta$. |
|---|---|

▷ Controller for system $\phi$ stabilizes system $\theta$

▷ Still a strong assumption, but weaker (and more natural) than A1.
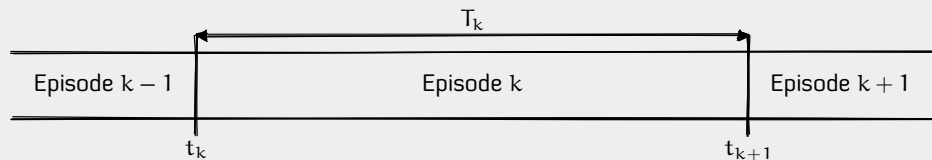
## Proof of regret bound of TSDE breaks down

▷ Proof relies on showing that there is some constant $\alpha_0$ such that

$$(\star) \qquad \mathbb{E}\Big[ \max_{1 \leq t \leq T} \|x_t\| \Big] \leq \sigma_w + \alpha_0 \mathbb{E}\Big[ \max_{1 \leq t \leq T} \|w_t\| \Big]$$

▷ Under (A1), $\quad \mathbb{E}[\|x_{t+1}\|] \leq \delta \mathbb{E}[\|x_t\|] + \mathbb{E}[\|w_t\|]$, which implies $\alpha_0 = 1/(1 - \delta)$.

▷ Such a bound does not work under (A2).
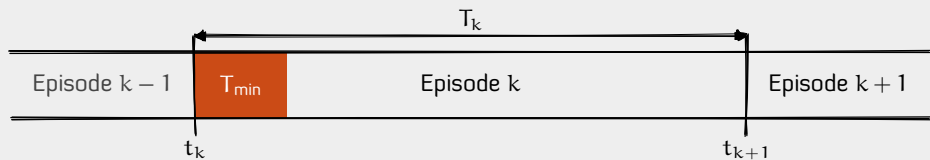
Need to modify the algorithm

# Modified TSDE



$T_k$ — Episode $k-1$ | Episode $k$ | Episode $k+1$, with markers $t_k$ and $t_{k+1}$

## Intuition

▷ Under (A2), in each episode the system is asymptotically stable.

▷ Asymptotic stability implies exponential stability.

▷ So, if the episode is sufficiently large, we can show that

$$\mathbb{E}[\|x_{t_{k+1}}\|] \leq \beta\,\mathbb{E}[\|x_{t_k}\|] + \bar{\alpha}\,\mathbb{E}\Big[\max_{t_k \leq t \leq t_{k+1}} \|w_t\|\Big]$$

which implies $(\star)$.

# Modified TSDE



## Intuition

▷ Under (A2), in each episode the system is asymptotically stable.

▷ Asymptotic stability implies exponential stability.

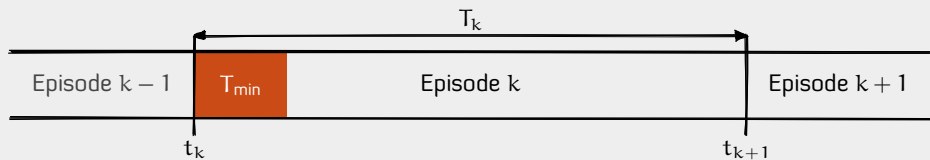▷ So, if the episode is sufficiently large, we can show that

$$\mathbb{E}[\|x_{t_{k+1}}\|] \leq \beta\,\mathbb{E}[\|x_{t_k}\|] + \bar{\alpha}\,\mathbb{E}\Big[\max_{t_k \leq t \leq t_{k+1}} \|w_t\|\Big]$$

which implies $(\star)$.

## Proposed modification

▷ To ensure that each episode is sufficiently large, do not stop in the first $T_{min}$ steps of an episode

▷ See paper for choice of $T_{min}$.

# Modified TSDE



## Intuition

▷ Under (A2), in each episode the system is asymptotically stable.

▷ Asymptotic stability implies exponential stability.

▷ So, if the episode is sufficiently large, we can show that

$$\mathbb{E}[\|x_{t_{k+1}}\|] \leq \beta \mathbb{E}[\|x_{t_k}\|] + \bar{\alpha} \mathbb{E}\Big[\max_{t_k \leq t \leq t_{k+1}} \|w_t\|\Big]$$

which implies $(\star)$.

## Proposed modification

▷ To ensure that each episode is sufficiently large, do not stop in the first $T_{min}$ steps of an episode

▷ See paper for choice of $T_{min}$.

## Implication

▷ The second stopping condition is not triggered for the $T_{min}$ steps of each episode.

▷ Requires other changes in the proof argument. See paper for details.

# Main results

| Assumption A2 | There exists an $\delta \in (0, 1)$ such that for any $\theta, \phi \in \Omega$, $\rho(A_\theta + B_\theta G(\phi)) \leq \delta$. |
|---|---|

| Theorem | Under A2,    $R(T; m\text{-TSDE}) \leq C\sqrt{T}\,(\log T)^q$ |
|---|---|

# Main results

| Assumption A2 | There exists an $\delta \in (0, 1)$ such that for any $\theta, \phi \in \Omega$, $\rho(A_\theta + B_\theta G(\phi)) \leq \delta$. |
|---|---|

| Theorem | Under A2, $\quad R(T; m\text{-TSDE}) \leq C\sqrt{T}\,(\log T)^q$ |
|---|---|

## Conclusion

▷ Relaxed a technical assumption for TSDE.

▷ Although A2 is weaker than A1, it still a **strong assumption**.

▷ Numerical experiments suggest that regret scales $\tilde{O}(\sqrt{T})$ even when A2 is not satisfied.

▷ Open question: How to further relax the stability assumption?

Thompson sampling for LQ—(Gagrani et. al.)

11

Thank you