



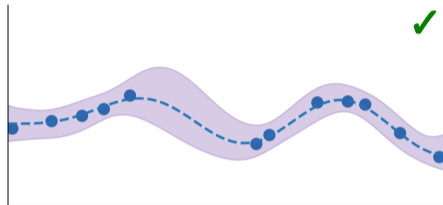
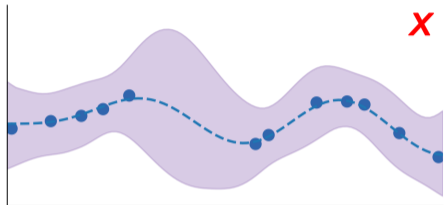
BOSCH



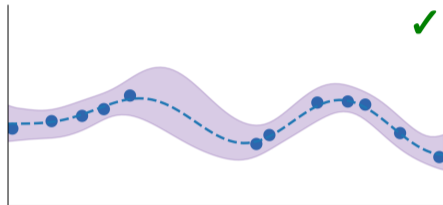
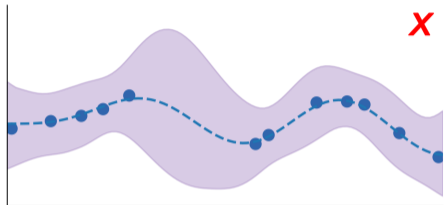
Improved Algorithms for Stochastic Linear Bandits Using Tail Bounds for Martingale Mixtures

Hamish Flynn David Reeb Melih Kandemir Jan Peters

Tighter Confidence Bounds Are Better Confidence Bounds

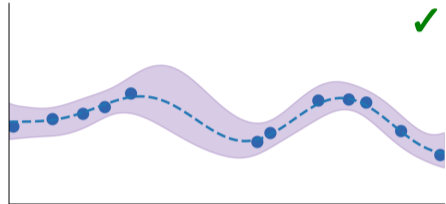
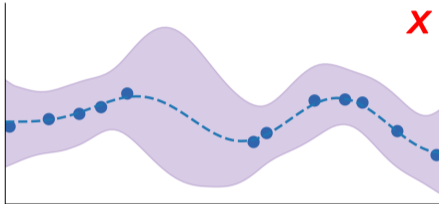


Tighter Confidence Bounds Are Better Confidence Bounds



Confidence bounds \rightarrow Upper Confidence Bound (UCB) algorithms for bandits.

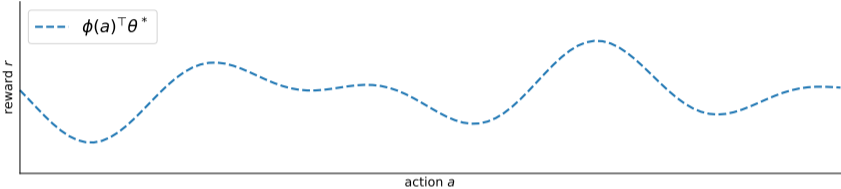
Tighter Confidence Bounds Are Better Confidence Bounds



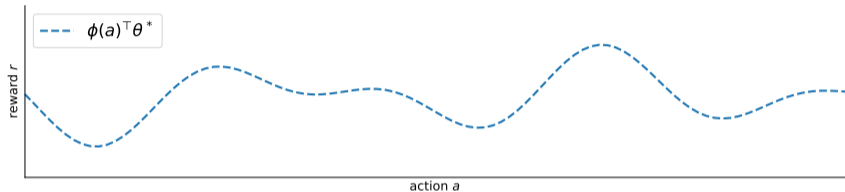
Confidence bounds \rightarrow Upper Confidence Bound (UCB) algorithms for bandits.

Tighter confidence bounds \rightarrow better UCB algorithms.

Stochastic Linear Bandits

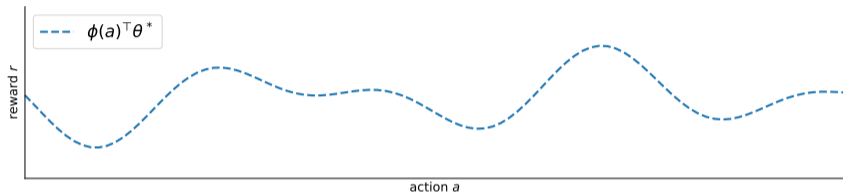


Stochastic Linear Bandits



At round t , query any action $a_t \in \mathcal{A}_t$, receive a noisy reward $r_t = \phi(a_t)^\top \boldsymbol{\theta}^* + \epsilon_t$.

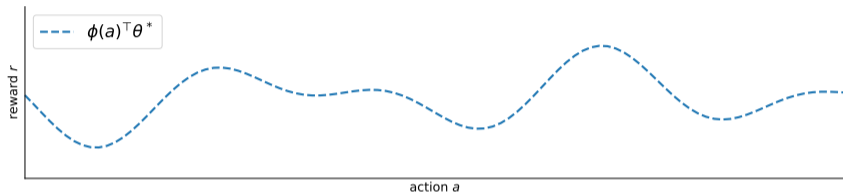
Stochastic Linear Bandits



At round t , query any action $a_t \in \mathcal{A}_t$, receive a noisy reward $r_t = \phi(a_t)^T \theta^* + \epsilon_t$.

Goal: Maximise total reward/minimise cumulative regret.

Stochastic Linear Bandits



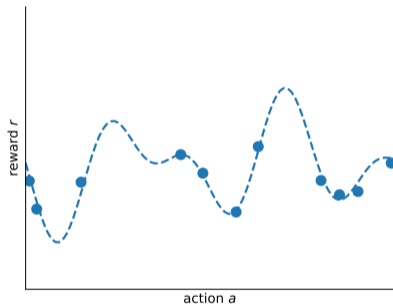
At round t , query any action $a_t \in \mathcal{A}_t$, receive a noisy reward $r_t = \phi(a_t)^\top \theta^* + \epsilon_t$.

Goal: Maximise total reward/minimise cumulative regret.

Assumptions: $\epsilon_1, \epsilon_2, \dots$ are (conditionally) σ -sub-Gaussian and $\|\theta^*\|_2 \leq B$.

$\theta^* \in \mathbb{R}^d$ is unknown, ϕ is known and upper bounds on σ and B are known.

UCB Algorithms for Stochastic Linear Bandits (e.g. OFUL¹)

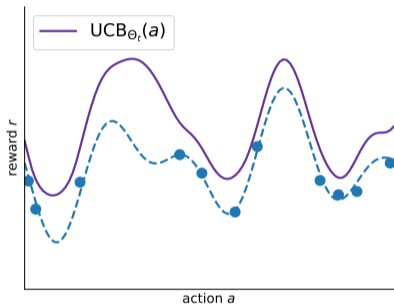


LinUCB:

For $t = 0, 1, 2, \dots$

¹Y. Abbasi-Yadkori et al. (2011) Improved algorithms for linear stochastic bandits. NeurIPS

UCB Algorithms for Stochastic Linear Bandits (e.g. OFUL¹)



LinUCB:

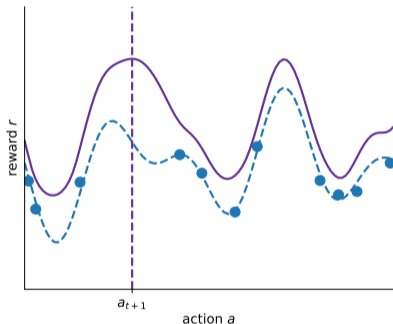
For $t = 0, 1, 2, \dots$

- Use $\{(a_k, r_k)\}_{k=1}^t$ to construct a confidence set Θ_t and the corresponding upper confidence bound

$$\text{UCB}_{\Theta_t}(a) := \max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$$

¹Y. Abbasi-Yadkori et al. (2011) Improved algorithms for linear stochastic bandits. NeurIPS

UCB Algorithms for Stochastic Linear Bandits (e.g. OFUL¹)



LinUCB:

For $t = 0, 1, 2, \dots$

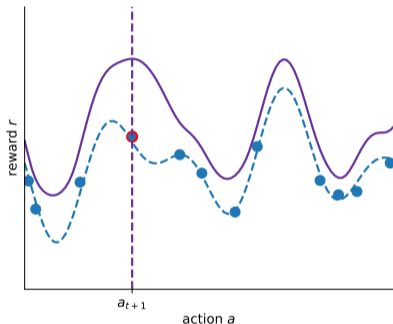
- Use $\{(a_k, r_k)\}_{k=1}^t$ to construct a confidence set Θ_t and the corresponding upper confidence bound

$$\text{UCB}_{\Theta_t}(a) := \max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$$

- Play $a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}_{t+1}} \{\text{UCB}_{\Theta_t}(a)\}$

¹Y. Abbasi-Yadkori et al. (2011) Improved algorithms for linear stochastic bandits. NeurIPS

UCB Algorithms for Stochastic Linear Bandits (e.g. OFUL¹)



LinUCB:

For $t = 0, 1, 2, \dots$

- Use $\{(a_k, r_k)\}_{k=1}^t$ to construct a confidence set Θ_t and the corresponding upper confidence bound

$$\text{UCB}_{\Theta_t}(a) := \max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$$

- Play $a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}_{t+1}} \{\text{UCB}_{\Theta_t}(a)\}$
- Observe reward $r_{t+1} = \phi(a_{t+1})^\top \theta^* + \epsilon_{t+1}$

¹Y. Abbasi-Yadkori et al. (2011) Improved algorithms for linear stochastic bandits. NeurIPS

In This Work

- New confidence sets Θ_t for stochastic linear bandits using a new tail bound for martingale mixtures

In This Work

- New confidence sets Θ_t for stochastic linear bandits using a new tail bound for martingale mixtures
- Provably tighter upper/lower confidence bounds than previous state-of-the-art (OFUL)

In This Work

- New confidence sets Θ_t for stochastic linear bandits using a new tail bound for martingale mixtures
- Provably tighter upper/lower confidence bounds than previous state-of-the-art (OFUL)
- LinUCB with our tighter confidence bounds leads to improved performance in hyperparameter tuning problems

In This Work

- New confidence sets Θ_t for stochastic linear bandits using a new tail bound for martingale mixtures
- Provably tighter upper/lower confidence bounds than previous state-of-the-art (OFUL)
- LinUCB with our tighter confidence bounds leads to improved performance in hyperparameter tuning problems
- LinUCB with our confidence sets has an $O(d\sqrt{T}\ln(T))$ worst-case cumulative regret bound (like OFUL)

Baseline Tail Bound

General Plan. Derive a data-dependent constraint for θ^* using tail bounds for non-i.i.d. data.

Baseline Tail Bound

General Plan. Derive a data-dependent constraint for $\boldsymbol{\theta}^*$ using tail bounds for non-i.i.d. data.

We want a bound on the sum of squared errors

$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 = \sum_{k=1}^t (\phi(a_k)^\top \boldsymbol{\theta}^* - r_k)^2 = \sum_{k=1}^t \epsilon_k^2.$$

Baseline Tail Bound

General Plan. Derive a data-dependent constraint for $\boldsymbol{\theta}^*$ using tail bounds for non-i.i.d. data.

We want a bound on the sum of squared errors

$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 = \sum_{k=1}^t (\phi(a_k)^\top \boldsymbol{\theta}^* - r_k)^2 = \sum_{k=1}^t \epsilon_k^2.$$

Baseline tail bound. Use the σ -sub-Gaussian property of ϵ_k : with probability $\geq 1 - \delta$

$$\forall t \geq 1: \quad \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 \leq \sigma^2 t + 2\sigma^2 \sqrt{t \ln \left(\frac{t^2 \pi^2}{6\delta} \right)} + 2\sigma^2 \ln \left(\frac{t^2 \pi^2}{6\delta} \right).$$

Martingale Mixture Tail Bound for Linear Bandits

Choose a sequence of incrementally updated mean vectors μ_1, μ_2, \dots and covariance matrices T_1, T_2, \dots

$$\mu_t = \left[\begin{array}{c|c} & \mu_t \end{array} \right]^\top, \quad T_t = \left[\begin{array}{ccc|c} & & & T_1 \\ & & & \vdots \\ & T_{t-1} & & T_{t-1} \\ \hline T_1 & \dots & T_{t-1} & T_t \end{array} \right].$$

Martingale Mixture Tail Bound for Linear Bandits

Choose a sequence of incrementally updated mean vectors μ_1, μ_2, \dots and covariance matrices T_1, T_2, \dots

$$\mu_t = \left[\begin{array}{c|c} \mu_{t-1} & \mu_t \end{array} \right]^\top, \quad T_t = \left[\begin{array}{ccc|c} & & & T_1 \\ & & & \vdots \\ & T_{t-1} & & T_{t-1} \\ \hline T_1 & \dots & T_{t-1} & T_t \end{array} \right].$$

μ_t and T_1, \dots, T_t can depend on the previous data $a_1, r_1, \dots, a_{t-1}, r_{t-1}, a_t$.

Martingale Mixture Tail Bound for Linear Bandits

Choose a sequence of incrementally updated mean vectors μ_1, μ_2, \dots and covariance matrices T_1, T_2, \dots

$$\mu_t = \left[\begin{array}{c|c} & \mu_{t-1} \\ \hline & \mu_t \end{array} \right]^\top, \quad T_t = \left[\begin{array}{ccc|c} & & & T_1 \\ & & & \vdots \\ & & & T_{t-1} \\ \hline T_1 & \dots & T_{t-1} & T_t \end{array} \right].$$

μ_t and T_1, \dots, T_t can depend on the previous data $a_1, r_1, \dots, a_{t-1}, r_{t-1}, a_t$.

Standard choice. $\mu_t = \mathbf{0}$, $T_t = \Phi_t \Phi_t^\top = (\phi(a_i)^\top \phi(a_j))_{1 \leq i, j \leq t}$.

Martingale Mixture Tail Bound for Linear Bandits

Choose a sequence of incrementally updated mean vectors μ_1, μ_2, \dots and covariance matrices $\mathbf{T}_1, \mathbf{T}_2, \dots$

$$\mu_t = \left[\begin{array}{c|c} & \mu_t \end{array} \right]^\top, \quad \mathbf{T}_t = \left[\begin{array}{c|c} \mathbf{T}_{t-1} & \begin{array}{c} T_1 \\ \vdots \\ T_{t-1} \end{array} \\ \hline T_1 \quad \dots \quad T_{t-1} & T_t \end{array} \right].$$

μ_t and T_1, \dots, T_t can depend on the previous data $a_1, r_1, \dots, a_{t-1}, r_{t-1}, a_t$.

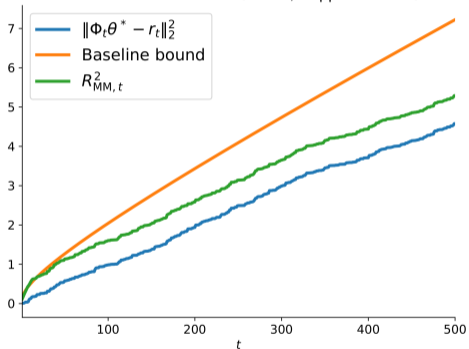
Standard choice. $\mu_t = \mathbf{0}$, $\mathbf{T}_t = \Phi_t \Phi_t^\top = (\phi(a_i)^\top \phi(a_j))_{1 \leq i, j \leq t}$.

Martingale mixture tail bound. With probability $\geq 1 - \delta$, for all $t \geq 1$

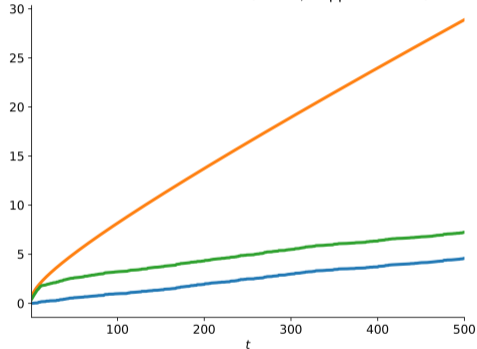
$$\|\Phi_t \theta^* - \mathbf{r}_t\|_2^2 \leq (\mu_t - \mathbf{r}_t)^\top \left(\mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right)^{-1} (\mu_t - \mathbf{r}_t) + \sigma^2 \ln \left(\det \left(\mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta) =: R_{\text{MM},t}^2.$$

Our Tail Bound Against The Baseline Tail Bound

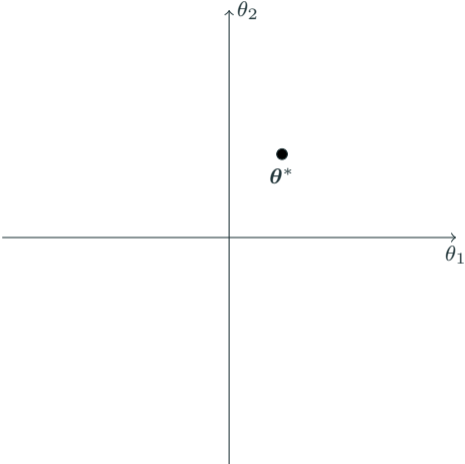
True Sub-Gaussian Parameter: $\sigma = 0.1$, Upper Bound: $\sigma = 0.1$



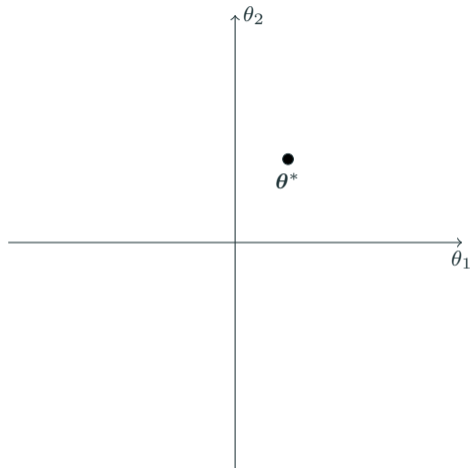
True Sub-Gaussian Parameter: $\sigma = 0.1$, Upper Bound: $\sigma = 0.2$



Confidence Sets For Linear Bandits



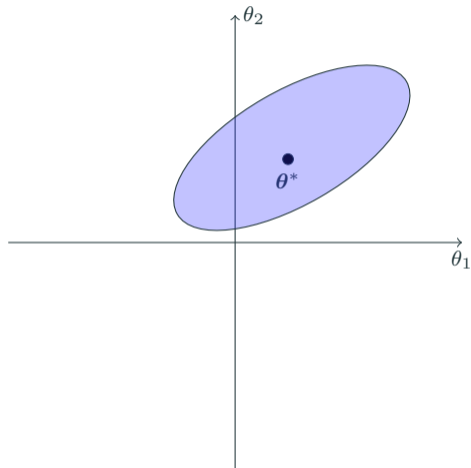
Confidence Sets For Linear Bandits



Using our martingale mixture tail bound, we have

$$\|\Phi_t \theta^* - \mathbf{r}_t\|_2 \leq R_{\text{MM},t},$$

Confidence Sets For Linear Bandits



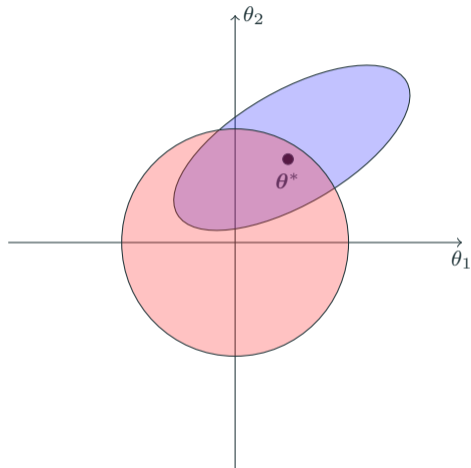
Using our martingale mixture tail bound, we have

$$\|\Phi_t \theta^* - \mathbf{r}_t\|_2 \leq R_{\text{MM},t},$$

This means that θ^* lies within the set

$$\{\theta \in \mathbb{R}^d : \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{\text{MM},t}\}.$$

Confidence Sets For Linear Bandits



Using our martingale mixture tail bound, we have

$$\|\Phi_t \theta^* - \mathbf{r}_t\|_2 \leq R_{\text{MM},t},$$

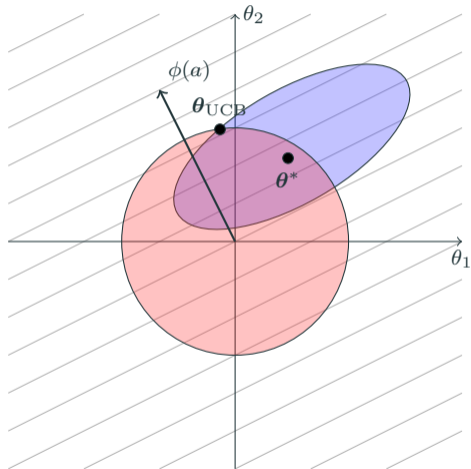
This means that θ^* lies within the set

$$\{\theta \in \mathbb{R}^d : \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{\text{MM},t}\}.$$

Incorporating the smoothness assumption, we obtain

$$\Theta_t = \{\theta \in \mathbb{R}^d : \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{\text{MM},t}, \|\theta\|_2 \leq B\}.$$

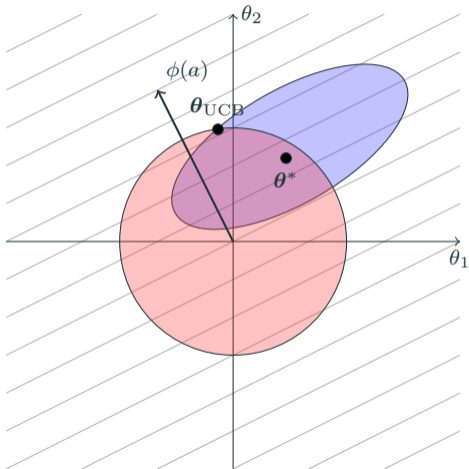
Convex Martingale Mixture UCB Algorithm



To run LinUCB with our confidence sets, we need to maximise $UCB_{\Theta_t}(a)$ w.r.t. a , where

$$\begin{aligned}UCB_{\Theta_t}(a) &= \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta \\ &\text{s.t. } \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{MM,t} \\ &\text{and } \|\theta\|_2 \leq B \\ &= \phi(a)^\top \theta_{UCB}.\end{aligned}$$

Convex Martingale Mixture UCB Algorithm

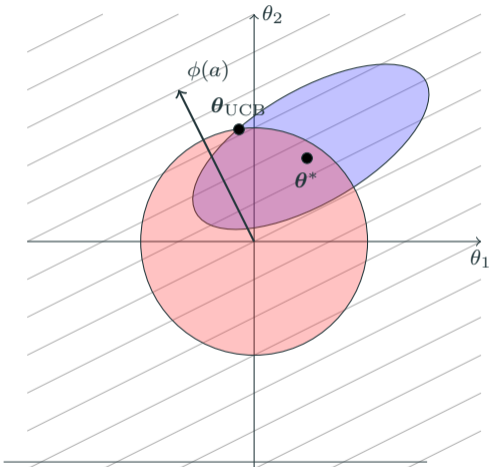


To run LinUCB with our confidence sets, we need to maximise $UCB_{\Theta_t}(a)$ w.r.t. a , where

$$\begin{aligned}UCB_{\Theta_t}(a) &= \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta \\ &\text{s.t. } \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{MM,t} \\ &\text{and } \|\theta\|_2 \leq B \\ &= \phi(a)^\top \theta_{UCB}.\end{aligned}$$

For continuous action sets, we approximately maximise $UCB_{\Theta_t}(a)$ w.r.t a using gradient-based methods.

Convex Martingale Mixture UCB Algorithm



To run LinUCB with our confidence sets, we need to maximise $UCB_{\Theta_t}(a)$ w.r.t. a , where

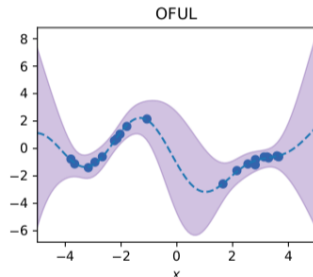
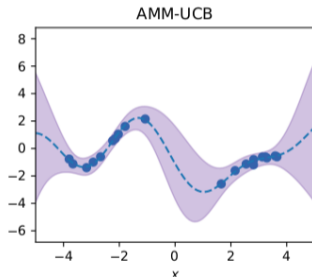
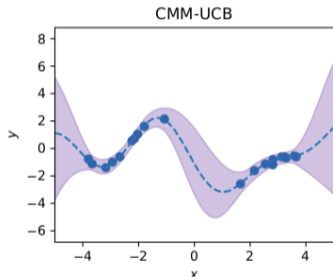
$$\begin{aligned}UCB_{\Theta_t}(a) &= \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta \\ &\text{s.t. } \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{MM,t} \\ &\text{and } \|\theta\|_2 \leq B \\ &= \phi(a)^\top \theta_{UCB}.\end{aligned}$$

For continuous action sets, we approximately maximise $UCB_{\Theta_t}(a)$ w.r.t. a using gradient-based methods.

We calculate $UCB_{\Theta_t}(a)$ and $\nabla_a UCB_{\Theta_t}(a)$ numerically using differentiable convex optimisation (cvxpylayers)².

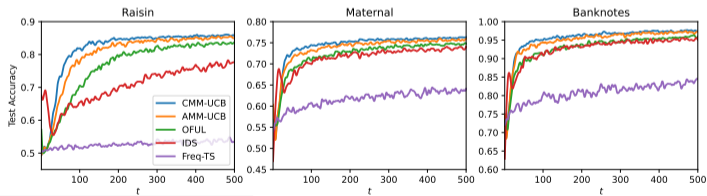
²A. Agrawal et al. (2019) Differentiable convex optimization layers. NeurIPS

Tighter Confidence Bounds



Hyperparameter Tuning

	Raisin		Maternal		Banknotes	
	Mean Acc	Max Acc	Mean Acc	Max Acc	Mean Acc	Max Acc
CMM-UCB (Ours)	0.818 \pm 0.018	0.893 \pm 0.019	0.744 \pm 0.020	0.829 \pm 0.023	0.954 \pm 0.005	1.000 \pm 0.000
AMM-UCB (Ours)	0.800 \pm 0.017	0.892 \pm 0.020	0.736 \pm 0.020	0.829 \pm 0.023	0.948 \pm 0.005	1.000 \pm 0.000
OFUL	0.764 \pm 0.019	0.891 \pm 0.019	0.722 \pm 0.019	0.827 \pm 0.022	0.929 \pm 0.006	1.000 \pm 0.000
IDS ³	0.706 \pm 0.048	0.891 \pm 0.020	0.714 \pm 0.019	0.827 \pm 0.024	0.926 \pm 0.007	1.000 \pm 0.000
Freq-TS ⁴	0.527 \pm 0.022	0.884 \pm 0.019	0.616 \pm 0.018	0.823 \pm 0.022	0.808 \pm 0.012	1.000 \pm 0.000



³ J. Kirschner and A. Krause. (2018) Information directed sampling and bandits with heteroscedastic noise, COLT

⁴ S. Agrawal and N. Goyal. (2013) Thompson sampling for contextual bandits with linear payoffs, ICML

Regret Analysis

Step 1. Cumulative regret is bounded by the confidence bound widths (UCB minus LCB).

$$\sum_{t=1}^T \phi(a_t^*)^\top \boldsymbol{\theta}^* - \phi(a_t)^\top \boldsymbol{\theta}^* \leq \sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t^*) - \text{LCB}_{\Theta_{t-1}}(a_t) \leq \sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t) - \text{LCB}_{\Theta_{t-1}}(a_t).$$

Regret Analysis

Step 1. Cumulative regret is bounded by the confidence bound widths (UCB minus LCB).

$$\sum_{t=1}^T \phi(a_t^*)^\top \boldsymbol{\theta}^* - \phi(a_t)^\top \boldsymbol{\theta}^* \leq \sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t^*) - \text{LCB}_{\Theta_{t-1}}(a_t) \leq \sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t) - \text{LCB}_{\Theta_{t-1}}(a_t).$$

Step 2. Via weak duality for the convex program $\max_{\boldsymbol{\theta} \in \Theta_t} \{\phi(a)^\top \boldsymbol{\theta}\}$, we obtain

$$\sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t) - \text{LCB}_{\Theta_{t-1}}(a_t) \leq \sum_{t=1}^T 2R_{\text{AMM},t-1} \sqrt{\phi(a_t)^\top \left(\Phi_{t-1}^\top \Phi_{t-1} + \alpha \mathbf{I} \right)^{-1} \phi(a_t)}.$$

Regret Analysis

Step 1. Cumulative regret is bounded by the confidence bound widths (UCB minus LCB).

$$\sum_{t=1}^T \phi(a_t^*)^\top \boldsymbol{\theta}^* - \phi(a_t)^\top \boldsymbol{\theta}^* \leq \sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t^*) - \text{LCB}_{\Theta_{t-1}}(a_t) \leq \sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t) - \text{LCB}_{\Theta_{t-1}}(a_t).$$

Step 2. Via weak duality for the convex program $\max_{\boldsymbol{\theta} \in \Theta_t} \{\phi(a)^\top \boldsymbol{\theta}\}$, we obtain

$$\sum_{t=1}^T \text{UCB}_{\Theta_{t-1}}(a_t) - \text{LCB}_{\Theta_{t-1}}(a_t) \leq \sum_{t=1}^T 2R_{\text{AMM},t-1} \sqrt{\phi(a_t)^\top \left(\Phi_{t-1}^\top \Phi_{t-1} + \alpha \mathbf{I} \right)^{-1} \phi(a_t)}.$$

Step 3. Separately upper bound $R_{\text{AMM},T-1}$ and $\sum_{t=1}^T \sqrt{\phi(a_t)^\top \left(\Phi_{t-1}^\top \Phi_{t-1} + \alpha \mathbf{I} \right)^{-1} \phi(a_t)}$, to obtain

$$\sum_{t=1}^T \phi(a_t^*)^\top \boldsymbol{\theta}^* - \phi(a_t)^\top \boldsymbol{\theta}^* \leq \mathcal{O}(d\sqrt{T} \ln(T)).$$

Poster Session Information

Our paper title: “Improved Algorithms for Stochastic Linear Bandits Using Tail Bounds for Martingale Mixtures”

Poster #1801

Poster Session 6, Thu 14 December, 17:00 - 19:00 CST

We'll be happy to talk at the poster!