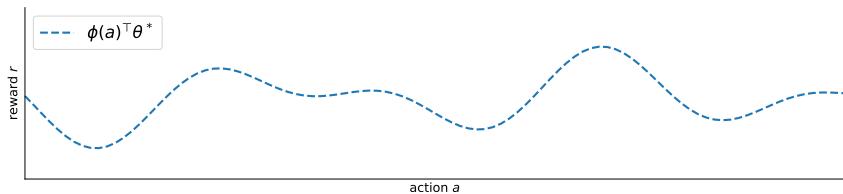


# Improved Algorithms for Stochastic Linear Bandits Using Tail Bounds for Martingale Mixtures

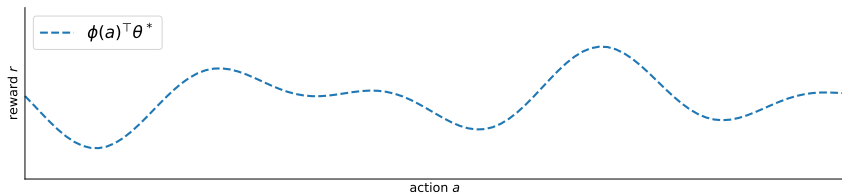
---

Hamish Flynn   David Reeb   Melih Kandemir   Jan Peters

# Stochastic Linear Bandits

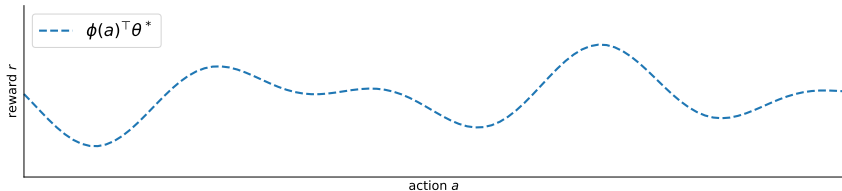


# Stochastic Linear Bandits



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

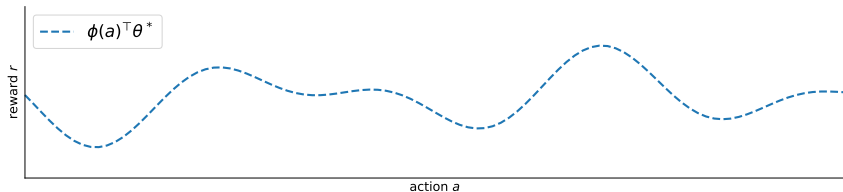
# Stochastic Linear Bandits



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

**Goal:** Minimise cumulative regret.

# Stochastic Linear Bandits



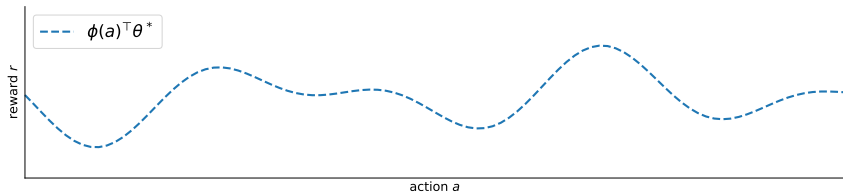
At round  $t$ , query any action  $a_t \in \mathcal{A}_t$ , receive a noisy reward  $r_t = \phi(a_t)^\top \theta^* + \epsilon_t$ .

**Goal:** Minimise cumulative regret.

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

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

# 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:** Minimise cumulative regret.

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

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

**Examples:** Black-box optimisation, recommendation systems, etc.

# Confidence Sets/Bounds for Stochastic Linear Bandits

**Confidence set:** A confidence set  $\Theta_t$  contains all  $\theta$ 's that could plausibly be  $\theta^*$  given data up to time  $t$ .

# Confidence Sets/Bounds for Stochastic Linear Bandits

**Confidence set:** A confidence set  $\Theta_t$  contains all  $\theta$ 's that could plausibly be  $\theta^*$  given data up to time  $t$ .

We want the smallest sequence of confidence sets  $\Theta_1, \Theta_2, \dots$  that satisfies the coverage condition

$$\mathbb{P}_{a_1, a_2, \dots, r_1, r_2, \dots} [\forall t \geq 1 : \theta^* \in \Theta_t] \geq 1 - \delta.$$



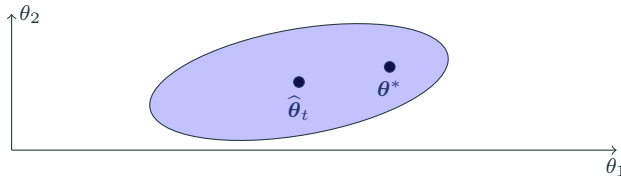
# Confidence Sets/Bounds for Stochastic Linear Bandits

**Confidence set:** A confidence set  $\Theta_t$  contains all  $\theta$ 's that could plausibly be  $\theta^*$  given data up to time  $t$ .

We want the smallest sequence of confidence sets  $\Theta_1, \Theta_2, \dots$  that satisfies the coverage condition

$$\mathbb{P}_{a_1, a_2, \dots}^{r_1, r_2, \dots} [\forall t \geq 1 : \theta^* \in \Theta_t] \geq 1 - \delta.$$

**Gold Standard (OFUL):**<sup>1</sup>  $\Theta_t$  is an ellipsoid centred at the regularised least squares/Ridge estimate  $\hat{\theta}_t$ , with a radius determined using self-normalised concentration and the method of mixtures.



---

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

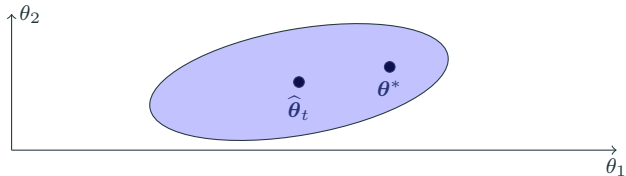
# Confidence Sets/Bounds for Stochastic Linear Bandits

**Confidence set:** A confidence set  $\Theta_t$  contains all  $\theta$ 's that could plausibly be  $\theta^*$  given data up to time  $t$ .

We want the smallest sequence of confidence sets  $\Theta_1, \Theta_2, \dots$  that satisfies the coverage condition

$$\mathbb{P}_{a_1, a_2, \dots}^{r_1, r_2, \dots} [\forall t \geq 1 : \theta^* \in \Theta_t] \geq 1 - \delta.$$

**Gold Standard (OFUL):**<sup>1</sup>  $\Theta_t$  is an ellipsoid centred at the regularised least squares/Ridge estimate  $\hat{\theta}_t$ , with a radius determined using self-normalised concentration and the method of mixtures.

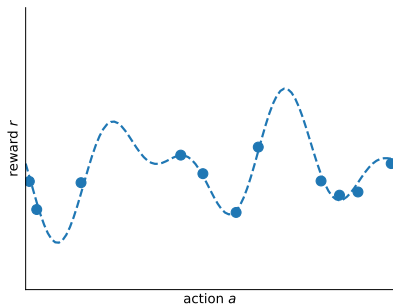


The corresponding upper confidence bound is  $\max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$ .

---

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

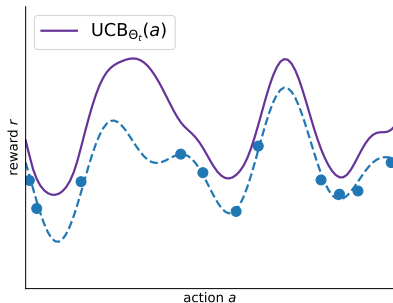
# UCB Algorithms for Stochastic Linear Bandits



**LinUCB:**

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

# UCB Algorithms for Stochastic Linear Bandits



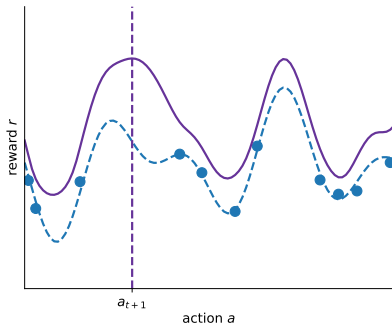
**LinUCB:**

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

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

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

# UCB Algorithms for Stochastic Linear Bandits



## LinUCB:

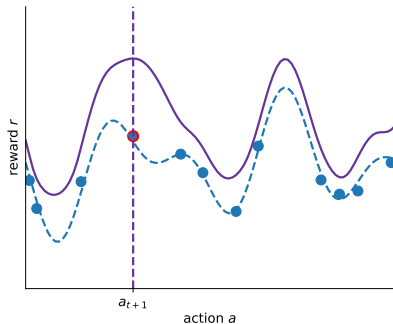
For  $t = 0, 1, 2, \dots$

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

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

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

# UCB Algorithms for Stochastic Linear Bandits



## LinUCB:

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

- Use  $\{(a_k, r_k)\}_{k=1}^t$  to construct a confidence set  $\Theta_t$  and the 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  $r_{t+1} = \phi(a_{t+1})^\top \theta^* + \epsilon_{t+1}$

# This Work

# This Work

**Question:** Is it possible to construct even tighter confidence sets/bounds for linear bandits?



# This Work

**Question:** Is it possible to construct even tighter confidence sets/bounds for linear bandits?

**Rest of the Talk:**

- Constructing confidence sets for linear bandits
- Computing confidence bounds/solving  $\max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$
- Regret bounds for LinUCB with our confidence sets
- In what sense is this better than OFUL (and why)?
- Some experimental results
- Open questions

**Notation, Etc.**

# Notation

## Notation:

- $\theta^*$  = true parameter vector,  $\theta$  = a candidate parameter vector
- $f_1, f_2, \dots$  = sequence of predictions
- $\mathcal{H}_t = (a_1, r_1, \dots, a_t, r_t, a_{t+1})$  = history of the bandit problem

# Notation

## Notation:

- $\theta^*$  = true parameter vector,  $\theta$  = a candidate parameter vector
- $f_1, f_2, \dots$  = sequence of predictions
- $\mathcal{H}_t = (a_1, r_1, \dots, a_t, r_t, a_{t+1})$  = history of the bandit problem

## Predictable sequences:

- I call a sequence of random variables  $x_1, x_2, \dots$  predictable if, given  $\mathcal{H}_{t-1}$ ,  $x_t$  is no longer random
- e.g.  $a_1, a_2, \dots$  are predictable,  $f_1, f_2, \dots$  are predictable,  $r_1, r_2, \dots$  are not predictable

# Notation

## Notation:

- $\theta^*$  = true parameter vector,  $\theta$  = a candidate parameter vector
- $f_1, f_2, \dots$  = sequence of predictions
- $\mathcal{H}_t = (a_1, r_1, \dots, a_t, r_t, a_{t+1})$  = history of the bandit problem

## Predictable sequences:

- I call a sequence of random variables  $x_1, x_2, \dots$  predictable if, given  $\mathcal{H}_{t-1}$ ,  $x_t$  is no longer random
- e.g.  $a_1, a_2, \dots$  are predictable,  $f_1, f_2, \dots$  are predictable,  $r_1, r_2, \dots$  are not predictable

## Matrix/vector notation:

- $\Phi_t = [\phi(a_1), \dots, \phi(a_t)]^\top \in \mathbb{R}^{t \times d}$  = matrix of first  $t$  feature vectors
- $\mathbf{r}_t = [r_1, \dots, r_t]^\top$  = vector of first  $t$  rewards
- $\boldsymbol{\epsilon}_t = [\epsilon_1, \dots, \epsilon_t]^\top$  = vector of first  $t$  noise variables
- $\mathbf{f}_t = [f_1, \dots, f_t]^\top$  = vector of first  $t$  predictions

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm



# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm
- E.g.,  $\theta_1, \theta_2, \dots$  could be a sequence of parameter estimates, and we could set  $f_t = \phi(a_t)^\top \theta_t$

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm
- E.g.,  $\theta_1, \theta_2, \dots$  could be a sequence of parameter estimates, and we could set  $f_t = \phi(a_t)^\top \theta_t$
- We could choose something boring like  $f_t \equiv 0$

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm
- E.g.,  $\theta_1, \theta_2, \dots$  could be a sequence of parameter estimates, and we could set  $f_t = \phi(a_t)^\top \theta_t$
- We could choose something boring like  $f_t \equiv 0$

**Randomised predictions:** Later on in the talk, we will consider distributions over sequences of predictions.

- $\mathbf{f}_t$  will be a random draw from  $P_t$ , which is distribution on  $\mathbb{R}^t$

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm
- E.g.,  $\theta_1, \theta_2, \dots$  could be a sequence of parameter estimates, and we could set  $f_t = \phi(a_t)^\top \theta_t$
- We could choose something boring like  $f_t \equiv 0$

**Randomised predictions:** Later on in the talk, we will consider distributions over sequences of predictions.

- $\mathbf{f}_t$  will be a random draw from  $P_t$ , which is distribution on  $\mathbb{R}^t$
- Each  $P_t$  can depend on the history  $\mathcal{H}_{t-1}$

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm
- E.g.,  $\theta_1, \theta_2, \dots$  could be a sequence of parameter estimates, and we could set  $f_t = \phi(a_t)^\top \theta_t$
- We could choose something boring like  $f_t \equiv 0$

**Randomised predictions:** Later on in the talk, we will consider distributions over sequences of predictions.

- $\mathbf{f}_t$  will be a random draw from  $P_t$ , which is distribution on  $\mathbb{R}^t$
- Each  $P_t$  can depend on the history  $\mathcal{H}_{t-1}$
- If  $P_t = \mathcal{N}(\mu_t, \mathbf{T}_t)$ ,  $\mu_t$  can still be thought of the first  $t$  predictions, and  $\mathbf{T}_t$  can be thought of as the uncertainty associated with the first  $t$  predictions

# Sequences of Predictions

Throughout the talk,  $f_1, f_2, \dots$  is a sequence of predictions for the rewards  $r_1, r_2, \dots$ , where each  $f_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

## Examples:

- For a fixed  $\theta$ , we could set  $f_t = \phi(a_t)^\top \theta$  for each  $t \geq 1$  (or  $\mathbf{f}_t = \Phi_t \theta$  in matrix notation)
- $f_1, f_2, \dots$  could be a sequence of predictions generated by running an online learning algorithm
- E.g.,  $\theta_1, \theta_2, \dots$  could be a sequence of parameter estimates, and we could set  $f_t = \phi(a_t)^\top \theta_t$
- We could choose something boring like  $f_t \equiv 0$

**Randomised predictions:** Later on in the talk, we will consider distributions over sequences of predictions.

- $\mathbf{f}_t$  will be a random draw from  $P_t$ , which is distribution on  $\mathbb{R}^t$
- Each  $P_t$  can depend on the history  $\mathcal{H}_{t-1}$
- If  $P_t = \mathcal{N}(\mu_t, \mathbf{T}_t)$ ,  $\mu_t$  can still be thought of the first  $t$  predictions, and  $\mathbf{T}_t$  can be thought of as the uncertainty associated with the first  $t$  predictions
- E.g. suppose  $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and define  $P_t$  to be the induced distribution on  $\Phi_t \theta$ , so  $P_t = \mathcal{N}(\mathbf{0}, \Phi_t \Phi_t^\top)$

## **Constructing Confidence Sets**

# General Plan (Part 1)

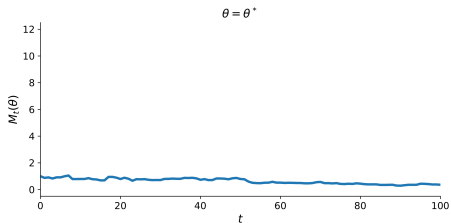
**Step 1:** Construct a collection of non-negative random processes  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  such that:



# General Plan (Part 1)

**Step 1:** Construct a collection of non-negative random processes  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  such that:

$$\mathbb{E}[M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) | \mathcal{H}_{t-1}] \leq M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)$$

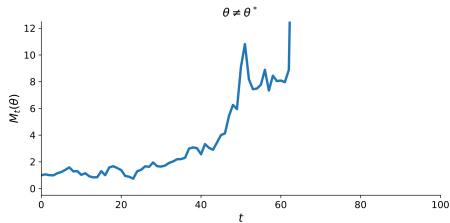
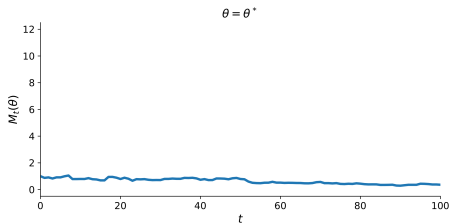


# General Plan (Part 1)

**Step 1:** Construct a collection of non-negative random processes  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  such that:

$$\mathbb{E}[M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) | \mathcal{H}_{t-1}] \leq M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)$$

If  $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$ ,  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  blows up

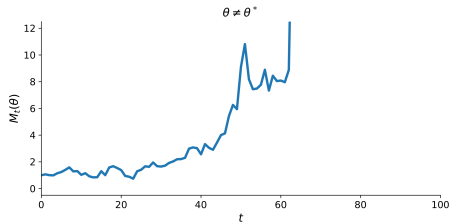
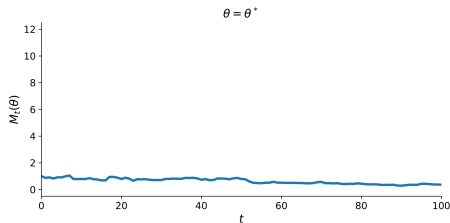


# General Plan (Part 1)

**Step 1:** Construct a collection of non-negative random processes  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  such that:

$$\mathbb{E}[M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) | \mathcal{H}_{t-1}] \leq M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)$$

If  $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$ ,  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  blows up



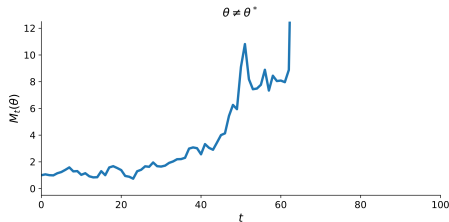
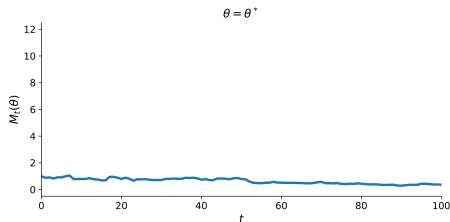
We want to maximise  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  w.r.t.  $\mathbf{f}_t$ , but the maximiser is not a predictable sequence.

# General Plan (Part 1)

**Step 1:** Construct a collection of non-negative random processes  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  such that:

$$\mathbb{E}[M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) | \mathcal{H}_{t-1}] \leq M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)$$

If  $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$ ,  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  blows up



We want to maximise  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  w.r.t.  $\mathbf{f}_t$ , but the maximiser is not a predictable sequence.

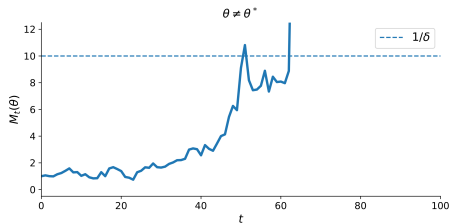
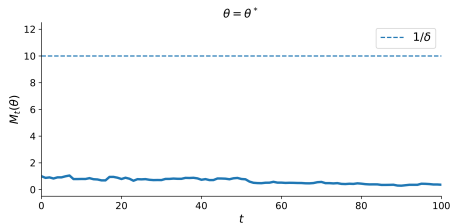
**Step 2:** Laplace's method/pseudo-maximisation/method of mixtures

$$\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta})] \approx \max_{\mathbf{f}_t} M_t(\mathbf{f}_t, \boldsymbol{\theta}).$$

# General Plan (Part 2)

**Step 3:** Use Ville's inequality to determine a threshold level

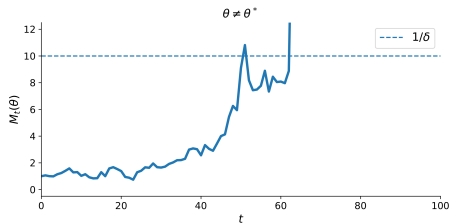
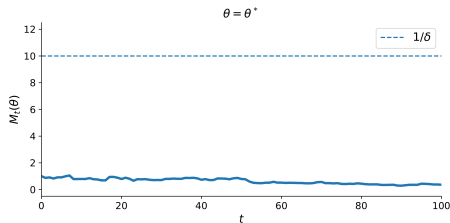
$$\mathbb{P}(\forall t \geq 1 : \mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta) \geq 1 - \delta.$$



# General Plan (Part 2)

**Step 3:** Use Ville's inequality to determine a threshold level

$$\mathbb{P}(\forall t \geq 1 : \mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta) \geq 1 - \delta.$$



**Step 4:** We define our confidence sets as

$$\Theta_t := \left\{ \boldsymbol{\theta} \in \mathbb{R}^d : \mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta})] \leq 1/\delta \right\} \cap \left\{ \boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_2 \leq B \right\}.$$

# What Do We Want From $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$ ?

We want to construct a collection of supermartingales  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  such that:

- $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is always non-negative
- $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)]$  has a closed-form expression whenever  $P_t$  is Gaussian
- $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  is a convex constraint for  $\boldsymbol{\theta}^*$

# What Do We Want From $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$ ?

We want to construct a collection of supermartingales  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  such that:

- $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is always non-negative
- $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)]$  has a closed-form expression whenever  $P_t$  is Gaussian
- $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  is a convex constraint for  $\boldsymbol{\theta}^*$

Look for  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  in the form

$$M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) = \exp \left( \sum_{k=1}^t \text{quad}(f_k, \phi(a_k)^\top \boldsymbol{\theta}^*) \right) = \prod_{k=1}^t \exp \left( \text{quad}(f_k, \phi(a_k)^\top \boldsymbol{\theta}^*) \right)$$



# What Do We Want From $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$ ?

We want to construct a collection of supermartingales  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  such that:

- $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is always non-negative
- $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)]$  has a closed-form expression whenever  $P_t$  is Gaussian
- $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  is a convex constraint for  $\boldsymbol{\theta}^*$

Look for  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  in the form

$$M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) = \exp \left( \sum_{k=1}^t \text{quad}(f_k, \phi(a_k)^\top \boldsymbol{\theta}^*) \right) = \prod_{k=1}^t \exp \left( \text{quad}(f_k, \phi(a_k)^\top \boldsymbol{\theta}^*) \right)$$

This is a supermartingale if, for all  $t \geq 1$ ,

$$\mathbb{E} \left[ \exp \left( \text{quad}(f_t, \phi(a_t)^\top \boldsymbol{\theta}^*) \right) \mid \mathcal{H}_{t-1} \right] \leq 1.$$

# Choosing The Quadratic Bits

Since  $\epsilon_1, \epsilon_2, \dots$  are (conditionally)  $\sigma$ -sub-Gaussian, we know that for any predictable  $\lambda_1, \lambda_2, \dots$ ,

$$\mathbb{E} \left[ \exp(\lambda_t (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*) \epsilon_t) | \mathcal{H}_{t-1} \right] \leq \exp \left( \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right).$$

# Choosing The Quadratic Bits

Since  $\epsilon_1, \epsilon_2, \dots$  are (conditionally)  $\sigma$ -sub-Gaussian, we know that for any predictable  $\lambda_1, \lambda_2, \dots$ ,

$$\mathbb{E} \left[ \exp(\lambda_t(f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)\epsilon_t) | \mathcal{H}_{t-1} \right] \leq \exp \left( \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right).$$

Therefore, we know that

$$\mathbb{E} \left[ \exp \left( \lambda_t(f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)\epsilon_t - \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right) | \mathcal{H}_{t-1} \right] \leq 1. \quad (1)$$

# Choosing The Quadratic Bits

Since  $\epsilon_1, \epsilon_2, \dots$  are (conditionally)  $\sigma$ -sub-Gaussian, we know that for any predictable  $\lambda_1, \lambda_2, \dots$ ,

$$\mathbb{E} \left[ \exp(\lambda_t (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*) \epsilon_t) | \mathcal{H}_{t-1} \right] \leq \exp \left( \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right).$$

Therefore, we know that

$$\mathbb{E} \left[ \exp \left( \lambda_t (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*) \epsilon_t - \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right) | \mathcal{H}_{t-1} \right] \leq 1. \quad (1)$$

Using  $\epsilon_t = r_t - \phi(a_t)^\top \boldsymbol{\theta}^*$ , the  $\exp(\dots)$  term in (1) can be re-written as

$$\exp \left( \frac{\lambda_t}{2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{\lambda_t}{2} (f_t - r_t)^2 + \frac{1}{2} (\lambda_t - \sigma^2 \lambda_t^2) (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2 \right).$$

# Choosing The Quadratic Bits

Since  $\epsilon_1, \epsilon_2, \dots$  are (conditionally)  $\sigma$ -sub-Gaussian, we know that for any predictable  $\lambda_1, \lambda_2, \dots$ ,

$$\mathbb{E} \left[ \exp(\lambda_t(f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)\epsilon_t) | \mathcal{H}_{t-1} \right] \leq \exp \left( \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right).$$

Therefore, we know that

$$\mathbb{E} \left[ \exp \left( \lambda_t(f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)\epsilon_t - \frac{\sigma^2 \lambda_t^2 (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2}{2} \right) | \mathcal{H}_{t-1} \right] \leq 1. \quad (1)$$

Using  $\epsilon_t = r_t - \phi(a_t)^\top \boldsymbol{\theta}^*$ , the  $\exp(\dots)$  term in (1) can be re-written as

$$\exp \left( \frac{\lambda_t}{2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{\lambda_t}{2} (f_t - r_t)^2 + \frac{1}{2} (\lambda_t - \sigma^2 \lambda_t^2) (f_t - \phi(a_t)^\top \boldsymbol{\theta}^*)^2 \right).$$

Setting  $\lambda_t \equiv 1/\sigma^2$ , this becomes

$$\exp \left( \frac{1}{2\sigma^2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{1}{2\sigma^2} (f_t - r_t)^2 \right).$$

# Putting Everything Together

Multiplying the  $\exp(\text{quad}(\dots))$  terms together, we obtain

$$\begin{aligned} M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) &= \prod_{k=1}^t \exp \left( \frac{1}{2\sigma^2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{1}{2\sigma^2} (f_t - r_t)^2 \right) \\ &= \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right). \end{aligned}$$

# Putting Everything Together

Multiplying the  $\exp(\text{quad}(\dots))$  terms together, we obtain

$$\begin{aligned} M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) &= \prod_{k=1}^t \exp \left( \frac{1}{2\sigma^2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{1}{2\sigma^2} (f_t - r_t)^2 \right) \\ &= \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right). \end{aligned}$$

**Closed-form integration.**  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is an unnormalised Gaussian density function (with mean  $\mathbf{r}_t$  and covariance  $\sigma^2 \mathbf{I}$ ), so we can use known tricks for integrating products of Gaussian densities.

# Putting Everything Together

Multiplying the  $\exp(\text{quad}(\cdots))$  terms together, we obtain

$$\begin{aligned} M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) &= \prod_{k=1}^t \exp \left( \frac{1}{2\sigma^2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{1}{2\sigma^2} (f_t - r_t)^2 \right) \\ &= \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right). \end{aligned}$$

**Closed-form integration.**  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is an unnormalised Gaussian density function (with mean  $\mathbf{r}_t$  and covariance  $\sigma^2 \mathbf{I}$ ), so we can use known tricks for integrating products of Gaussian densities.

**Convex constraint.**  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is the composition of  $\exp(\cdot)$  and a convex function of  $\boldsymbol{\theta}^*$ , which means  $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  is a convex constraint for  $\boldsymbol{\theta}^*$ .



# Putting Everything Together

Multiplying the  $\exp(\text{quad}(\cdots))$  terms together, we obtain

$$\begin{aligned} M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) &= \prod_{k=1}^t \exp \left( \frac{1}{2\sigma^2} (\phi(a_t)^\top \boldsymbol{\theta}^* - r_t)^2 - \frac{1}{2\sigma^2} (f_t - r_t)^2 \right) \\ &= \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right). \end{aligned}$$

**Closed-form integration.**  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is an unnormalised Gaussian density function (with mean  $\mathbf{r}_t$  and covariance  $\sigma^2 \mathbf{I}$ ), so we can use known tricks for integrating products of Gaussian densities.

**Convex constraint.**  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  is the composition of  $\exp(\cdot)$  and a convex function of  $\boldsymbol{\theta}^*$ , which means  $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  is a convex constraint for  $\boldsymbol{\theta}^*$ .

**Blowing up when  $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$ .** If  $f_1, f_2, \dots$  predicts the rewards better than  $\phi(a_1)^\top \boldsymbol{\theta}, \phi(a_2)^\top \boldsymbol{\theta}, \dots$ , then  $M_t(\mathbf{f}_t, \boldsymbol{\theta})$  will grow exponentially with  $t$  (in expectation).

# Motivation for Mixing

We would like to maximise  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  w.r.t.  $\mathbf{f}_t$ , but the maximiser is not a predictable sequence.

$$\operatorname{argmax}_{\mathbf{f}_t \in \mathbb{R}^t} \left\{ \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right) \right\} = \mathbf{r}_t.$$

# Motivation for Mixing

We would like to maximise  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  w.r.t.  $\mathbf{f}_t$ , but the maximiser is not a predictable sequence.

$$\operatorname{argmax}_{\mathbf{f}_t \in \mathbb{R}^t} \left\{ \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right) \right\} = \mathbf{r}_t.$$

For a function  $g(x)$  with a minimiser  $x^*$ , Laplace's asymptotic formula tells us that,

$$\begin{aligned} \int_{-\infty}^{\infty} \exp(-\lambda g(x)) dx &\approx \int_{-\infty}^{\infty} \exp \left( -\lambda g(x^*) - \frac{\lambda}{2} g''(x^*) (x - x^*)^2 \right) dx \\ &= \exp(-\lambda g(x^*)) \sqrt{\frac{2\pi}{g''(x^*)}} = \max_x \{ \exp(-\lambda g(x)) \} \sqrt{\frac{2\pi}{g''(x^*)}} \end{aligned}$$

# Motivation for Mixing

We would like to maximise  $M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)$  w.r.t.  $\mathbf{f}_t$ , but the maximiser is not a predictable sequence.

$$\operatorname{argmax}_{\mathbf{f}_t \in \mathbb{R}^t} \left\{ \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right) \right\} = \mathbf{r}_t.$$

For a function  $g(x)$  with a minimiser  $x^*$ , Laplace's asymptotic formula tells us that,

$$\begin{aligned} \int_{-\infty}^{\infty} \exp(-\lambda g(x)) dx &\approx \int_{-\infty}^{\infty} \exp \left( -\lambda g(x^*) - \frac{\lambda}{2} g''(x^*) (x - x^*)^2 \right) dx \\ &= \exp(-\lambda g(x^*)) \sqrt{\frac{2\pi}{g''(x^*)}} = \max_x \{ \exp(-\lambda g(x)) \} \sqrt{\frac{2\pi}{g''(x^*)}} \end{aligned}$$

This suggests that we can perform “pseudo-maximisation” w.r.t.  $\mathbf{f}_t$  via integration w.r.t. a (probability) measure, i.e.

$$\mathbb{E}_{\mathbf{f}_t \sim P_t} \left[ \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right) \right] \approx \max_{\mathbf{f}_t} \left\{ \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\mathbf{f}_t - \mathbf{r}_t\|_2^2 \right) \right\}$$

# Which Mixture Distributions Are Allowed?

We can choose any sequence of mixture distributions as long as  $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)]$  is a supermartingale, i.e.

$$\mathbb{E} [\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] | \mathcal{H}_{t-1}] \leq \mathbb{E}_{\mathbf{f}_{t-1} \sim P_{t-1}} [M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)] .$$

# Which Mixture Distributions Are Allowed?

We can choose any sequence of mixture distributions as long as  $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)]$  is a supermartingale, i.e.

$$\mathbb{E} [\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] | \mathcal{H}_{t-1}] \leq \mathbb{E}_{\mathbf{f}_{t-1} \sim P_{t-1}} [M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)] .$$

Suppose that  $P_1, P_2, \dots$  satisfies

1.  $P_t$  depends on only the history  $\mathcal{H}_{t-1}$  (and not the future  $r_t, a_{t+1}, r_{t+1}, \dots$ )
2.  $P_t(\mathbf{f}_t) = p_t(f_t | \mathbf{f}_{t-1}) P_{t-1}(\mathbf{f}_{t-1})$

# Which Mixture Distributions Are Allowed?

We can choose any sequence of mixture distributions as long as  $\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)]$  is a supermartingale, i.e.

$$\mathbb{E} [\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] | \mathcal{H}_{t-1}] \leq \mathbb{E}_{\mathbf{f}_{t-1} \sim P_{t-1}} [M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)] .$$

Suppose that  $P_1, P_2, \dots$  satisfies

1.  $P_t$  depends on only the history  $\mathcal{H}_{t-1}$  (and not the future  $r_t, a_{t+1}, r_{t+1}, \dots$ )
2.  $P_t(\mathbf{f}_t) = p_t(f_t | \mathbf{f}_{t-1}) P_{t-1}(\mathbf{f}_{t-1})$

In this case, we have

$$\begin{aligned} \mathbb{E} [\mathbb{E}_{\mathbf{f}_t \sim P_t} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] | \mathcal{H}_{t-1}] &= \mathbb{E}_{\mathbf{f}_t \sim P_t} [\mathbb{E} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*) | \mathcal{H}_{t-1}]] & (1.) \\ &\leq \mathbb{E}_{\mathbf{f}_t \sim P_t} [M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)] & (M_t \text{ is a supermartingale}) \\ &= \mathbb{E}_{\mathbf{f}_{t-1} \sim P_{t-1}} [M_{t-1}(\mathbf{f}_{t-1}, \boldsymbol{\theta}^*)] & (2.) \end{aligned}$$

# Gaussian Mixture Distributions

A sequence  $P_1 = \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{T}_1), P_2 = \mathcal{N}(\boldsymbol{\mu}_2, \mathbf{T}_2), \dots$  satisfies 1. and 2. if

$$\boldsymbol{\mu}_t = \left[ \begin{array}{c} | \\ \boldsymbol{\mu}_{t-1} \\ | \\ \hline \boldsymbol{\mu}_t \end{array} \right], \quad \mathbf{T}_t = \left[ \begin{array}{ccc|c} & & & T_1 \\ & & & \vdots \\ & & & T_{t-1} \\ \hline T_1 & \cdots & T_{t-1} & T_t \end{array} \right],$$

Each new element  $\boldsymbol{\mu}_t$  and row/column  $T_1, \dots, T_t$  can depend on the history  $\mathcal{H}_{t-1}$ .



# Gaussian Mixture Distributions

A sequence  $P_1 = \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{T}_1), P_2 = \mathcal{N}(\boldsymbol{\mu}_2, \mathbf{T}_2), \dots$  satisfies 1. and 2. if

$$\boldsymbol{\mu}_t = \left[ \begin{array}{c|c} & \\ \hline \boldsymbol{\mu}_{t-1} & \\ \hline & \boldsymbol{\mu}_t \end{array} \right], \quad \mathbf{T}_t = \left[ \begin{array}{ccc|c} & & & T_1 \\ & & & \vdots \\ & & & T_{t-1} \\ \hline T_1 & \cdots & T_{t-1} & T_t \end{array} \right],$$

Each new element  $\boldsymbol{\mu}_t$  and row/column  $T_1, \dots, T_t$  can depend on the history  $\mathcal{H}_{t-1}$ .

With  $P_t = \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{T}_t)$ , the martingale mixture is

$$\mathbb{E}_{\mathbf{f}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{T}_t)} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] = \frac{1}{\sqrt{\det(\mathbf{I} + \mathbf{T}_t/\sigma^2)}} \exp \left( \frac{1}{2\sigma^2} \|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 - \frac{1}{2\sigma^2} \|\boldsymbol{\mu}_t - \mathbf{r}_t\|_{(\mathbf{I} + \mathbf{T}_t/\sigma^2)^{-1}}^2 \right).$$

# Martingale Mixture Tail Bound

The constraint  $\mathbb{E}_{\mathbf{f}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{T}_t)} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  can be rearranged into

$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 \leq (\boldsymbol{\mu}_t - \mathbf{r}_t)^\top \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right)^{-1} (\boldsymbol{\mu}_t - \mathbf{r}_t) + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta).$$

# Martingale Mixture Tail Bound

The constraint  $\mathbb{E}_{\mathbf{f}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{T}_t)} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  can be rearranged into

$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 \leq (\boldsymbol{\mu}_t - \mathbf{r}_t)^\top \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right)^{-1} (\boldsymbol{\mu}_t - \mathbf{r}_t) + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta).$$

Standard mixture distributions:  $P_t = \mathcal{N}(\mathbf{0}, c\Phi_t\Phi_t^\top)$

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

# Martingale Mixture Tail Bound

The constraint  $\mathbb{E}_{\mathbf{f}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{T}_t)} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  can be rearranged into

$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 \leq (\boldsymbol{\mu}_t - \mathbf{r}_t)^\top \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right)^{-1} (\boldsymbol{\mu}_t - \mathbf{r}_t) + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta).$$

Standard mixture distributions:  $P_t = \mathcal{N}(\mathbf{0}, c\Phi_t\Phi_t^\top)$

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

On the one hand...

- $\mathcal{N}(\mathbf{0}, c\Phi_t\Phi_t^\top)$  is good enough to give us tighter confidence sets/bounds
- $c\Phi_t\Phi_t^\top$  is rank  $d$ , so  $R_{\text{MM},t}^2$  can be computed relatively cheaply

# Martingale Mixture Tail Bound

The constraint  $\mathbb{E}_{\mathbf{f}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{T}_t)} [M_t(\mathbf{f}_t, \boldsymbol{\theta}^*)] \leq 1/\delta$  can be rearranged into

$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 \leq (\boldsymbol{\mu}_t - \mathbf{r}_t)^\top \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right)^{-1} (\boldsymbol{\mu}_t - \mathbf{r}_t) + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{\mathbf{T}_t}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta).$$

Standard mixture distributions:  $P_t = \mathcal{N}(\mathbf{0}, c\Phi_t\Phi_t^\top)$

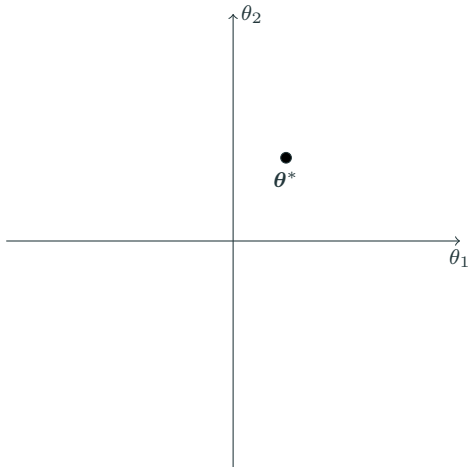
$$\|\Phi_t \boldsymbol{\theta}^* - \mathbf{r}_t\|_2^2 \leq \mathbf{r}_t^\top \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right)^{-1} \mathbf{r}_t + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta) =: R_{\text{MM},t}^2.$$

On the one hand...

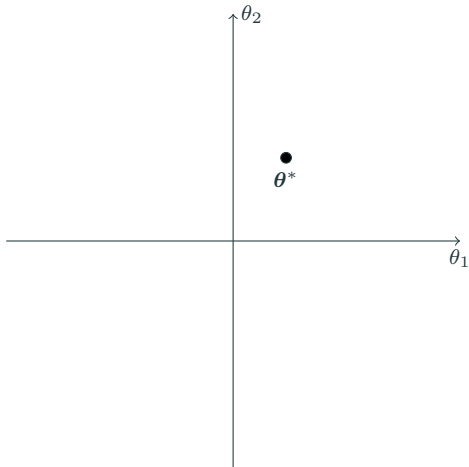
- $\mathcal{N}(\mathbf{0}, c\Phi_t\Phi_t^\top)$  is good enough to give us tighter confidence sets/bounds
- $c\Phi_t\Phi_t^\top$  is rank  $d$ , so  $R_{\text{MM},t}^2$  can be computed relatively cheaply

On the other hand,  $\boldsymbol{\mu}_t = \mathbf{0}$  seems bit a silly.

# Confidence Sets For Linear Bandits



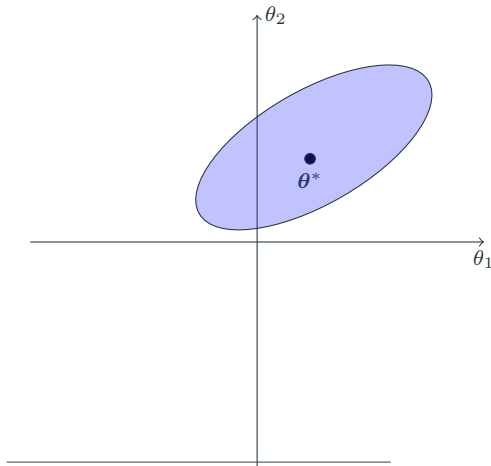
# Confidence Sets For Linear Bandits



Using our martingale tail bound, we have

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

# Confidence Sets For Linear Bandits



Using our martingale tail bound, we have

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

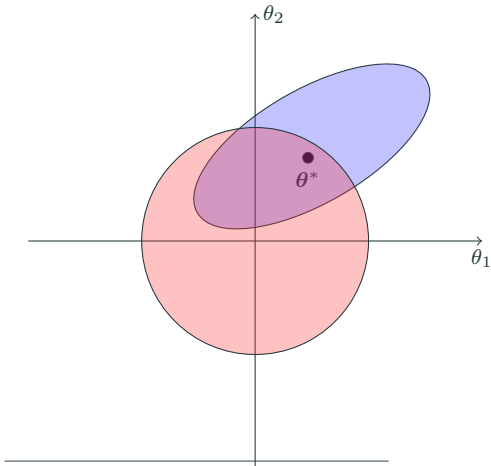
This means that  $\theta^*$  lies within the set\*

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

\* this set can be re-written as  $\{\theta \in \mathbb{R}^d : \|\theta - \hat{\theta}_t\|_{\Phi_t^\top \Phi_t} \leq \tilde{R}_t\}$ , where  $\hat{\theta}_t = \Phi_t^\dagger \mathbf{r}_t$  and  $\tilde{R}_t$  is some other radius quantity



# Confidence Sets For Linear Bandits



Using our martingale tail bound, we have

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

This means that  $\theta^*$  lies within the set\*

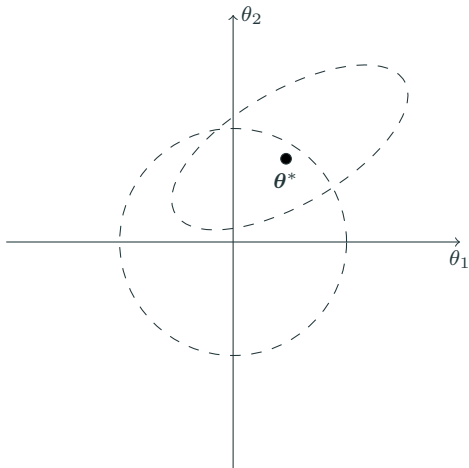
$$\{\theta \in \mathbb{R}^d : \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{\text{MM},t}\}.$$

Incorporating the smoothness assumption, we obtain

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

\* this set can be re-written as  $\{\theta \in \mathbb{R}^d : \|\theta - \hat{\theta}_t\|_{\Phi_t^\top \Phi_t} \leq \tilde{R}_t\}$ , where  $\hat{\theta}_t = \Phi_t^\dagger \mathbf{r}_t$  and  $\tilde{R}_t$  is some other radius quantity

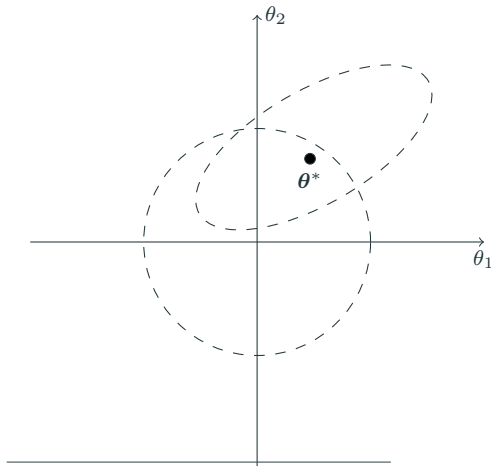
# Confidence Sets For Linear Bandits (Single Ellipsoid)



By taking a weighted (by  $\alpha > 0$ ) sum, we obtain a single quadratic constraint for  $\theta^*$

$$\|\Phi_t \theta^* - r_t\|_2^2 + \alpha \|\theta^*\|_2^2 \leq R_{\text{MM},t}^2 + \alpha B^2,$$

# Confidence Sets For Linear Bandits (Single Ellipsoid)



By taking a weighted (by  $\alpha > 0$ ) sum, we obtain a single quadratic constraint for  $\theta^*$

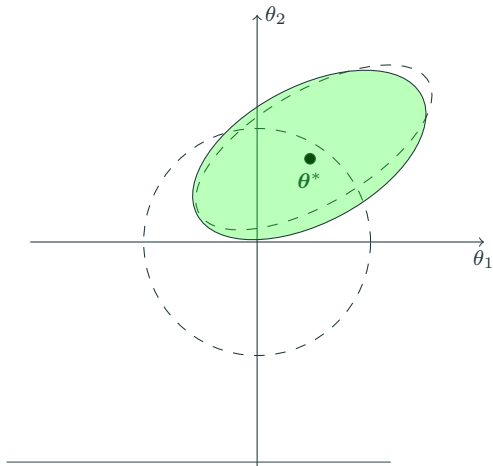
$$\|\Phi_t \theta^* - \mathbf{r}_t\|_2^2 + \alpha \|\theta^*\|_2^2 \leq R_{\text{MM},t}^2 + \alpha B^2,$$

By completing the square on the LHS, this constraint can be re-written as\*

$$\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})} \leq R_{\text{AMM},t}$$

\* where  $\hat{\theta}_{\alpha,t} = (\Phi_t^\top \Phi_t + \alpha \mathbf{I})^{-1} \Phi_t^\top \mathbf{r}_t$ ,  $R_{\text{AMM},t}^2 = R_{\text{MM},t}^2 + \alpha B^2 - \mathbf{r}_t^\top \mathbf{r}_t + \mathbf{r}_t^\top \Phi_t (\Phi_t^\top \Phi_t + \alpha \mathbf{I})^{-1} \Phi_t^\top \mathbf{r}_t$

# Confidence Sets For Linear Bandits (Single Ellipsoid)



By taking a weighted (by  $\alpha > 0$ ) sum, we obtain a single quadratic constraint for  $\theta^*$

$$\|\Phi_t \theta^* - r_t\|_2^2 + \alpha \|\theta^*\|_2^2 \leq R_{\text{MM},t}^2 + \alpha B^2,$$

By completing the square on the LHS, this constraint can be re-written as\*

$$\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)} \leq R_{\text{AMM},t}$$

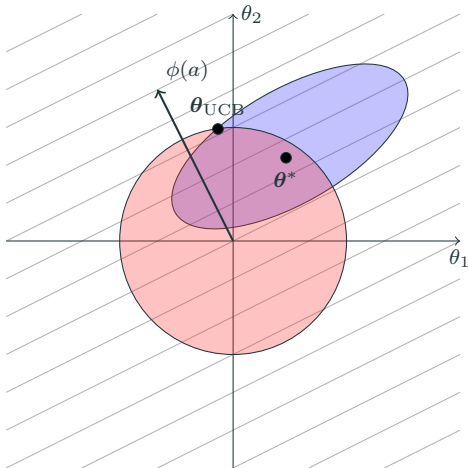
This means that  $\theta^*$  lies within the ellipsoid

$$\Theta_t^\alpha = \{\theta \in \mathbb{R}^d : \|\theta - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)} \leq R_{\text{AMM},t}\}.$$

\* where  $\hat{\theta}_{\alpha,t} = (\Phi_t^\top \Phi_t + \alpha I)^{-1} \Phi_t^\top r_t$ ,  $R_{\text{AMM},t}^2 = R_{\text{MM},t}^2 + \alpha B^2 - r_t^\top r_t + r_t^\top \Phi_t (\Phi_t^\top \Phi_t + \alpha I)^{-1} \Phi_t^\top r_t$

## **Computing and Maximising Confidence Bounds**

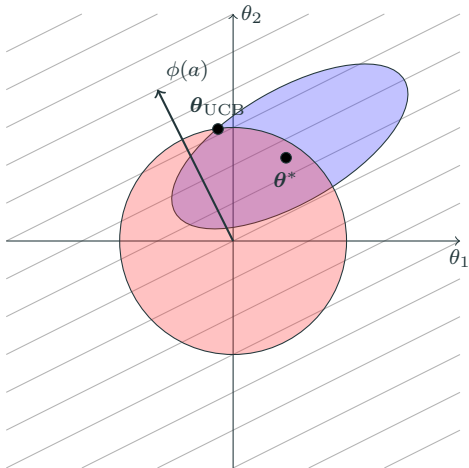
# Convex Martingale Mixture UCB



The UCB for our double ellipsoid confidence set is

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

# Convex Martingale Mixture UCB

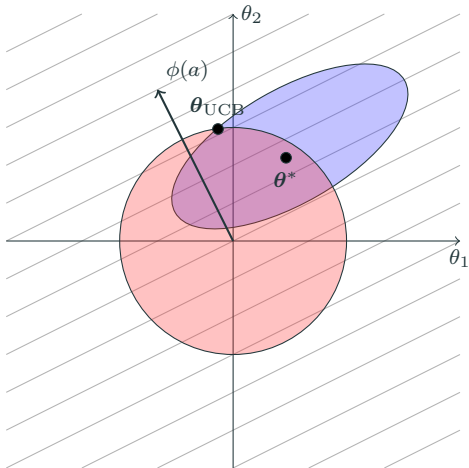


The UCB for our double ellipsoid confidence set is

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

$\text{UCB}_{\Theta_t}(a)$  can be computed in  $\mathcal{O}(d^3)$  time complexity via interior point methods.

# Convex Martingale Mixture UCB



The UCB for our double ellipsoid confidence set is

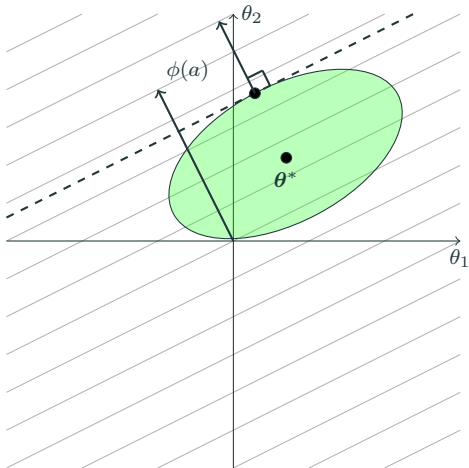
$$\begin{aligned}\text{UCB}_{\Theta_t}(a) &= \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta \\ &\text{s.t. } \|\Phi_t \theta - \mathbf{r}_t\|_2 \leq R_{\text{MM},t} \\ &\text{and } \|\theta\|_2 \leq B \\ &= \phi(a)^\top \theta_{\text{UCB}}.\end{aligned}$$

$\text{UCB}_{\Theta_t}(a)$  can be computed in  $\mathcal{O}(d^3)$  time complexity via interior point methods.

We call LinUCB with these confidence sets/bounds Convex Martingale Mixture (CMM-)UCB.



# Analytic Martingale Mixture UCB

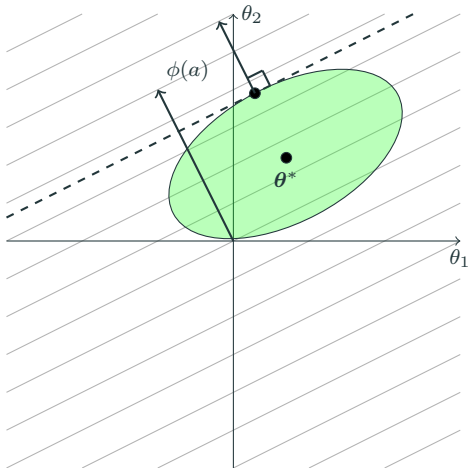


The UCB for our single ellipsoid confidence set is

$$\text{UCB}_{\Theta_t^\alpha}(a) = \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta$$

$$\text{s.t. } \|\theta - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)}^2 \leq R_{\text{AMM},t}^2.$$

# Analytic Martingale Mixture UCB



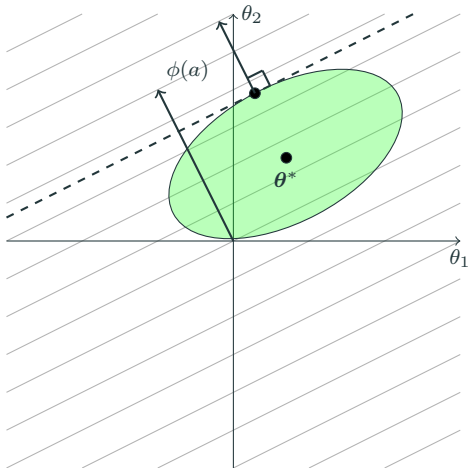
The UCB for our single ellipsoid confidence set is

$$\begin{aligned} \text{UCB}_{\Theta_t^\alpha}(a) &= \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta \\ \text{s.t. } &\|\theta - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)}^2 \leq R_{\text{AMM},t}^2. \end{aligned}$$

This time, there is a closed-form solution.

$$\text{UCB}_{\Theta_t^\alpha}(a) = \phi(a)^\top \hat{\theta}_{\alpha,t} + R_{\text{AMM},t} \|\phi(a)\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}}.$$

# Analytic Martingale Mixture UCB



The UCB for our single ellipsoid confidence set is

$$\begin{aligned} \text{UCB}_{\Theta_t^\alpha}(a) &= \max_{\theta \in \mathbb{R}^d} \phi(a)^\top \theta \\ \text{s.t. } &\|\theta - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)}^2 \leq R_{\text{AMM},t}^2. \end{aligned}$$

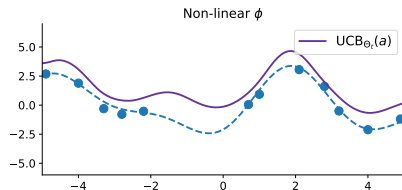
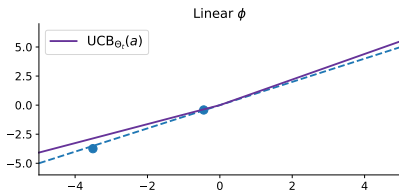
This time, there is a closed-form solution.

$$\text{UCB}_{\Theta_t^\alpha}(a) = \phi(a)^\top \hat{\theta}_{\alpha,t} + R_{\text{AMM},t} \|\phi(a)\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}}.$$

We call LinUCB with these confidence sets/bounds Analytic Martingale Mixture (AMM-)UCB.

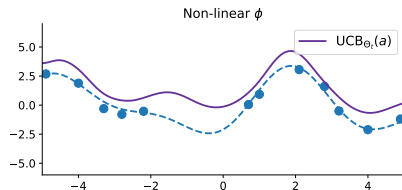
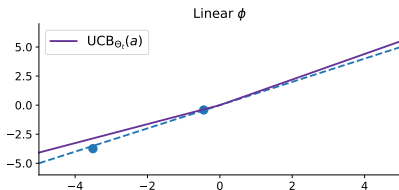
# Confidence Bound Maximisation

To run LinUCB with our confidence sets, we need to maximise  $UCB_{\Theta_t}(a) = \max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$  w.r.t.  $a$



# Confidence Bound Maximisation

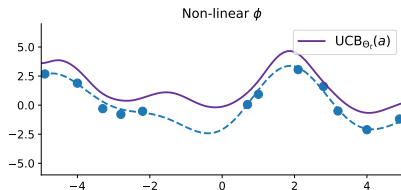
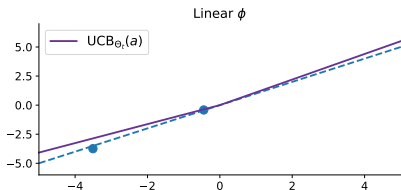
To run LinUCB with our confidence sets, we need to maximise  $UCB_{\Theta_t}(a) = \max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$  w.r.t.  $a$



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

# Confidence Bound Maximisation

To run LinUCB with our confidence sets, we need to maximise  $UCB_{\Theta_t}(a) = \max_{\theta \in \Theta_t} \{\phi(a)^\top \theta\}$  w.r.t.  $a$



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

For CMM-UCB,  $UCB_{\Theta_t}(a)$  and  $\nabla_a UCB_{\Theta_t}(a)$  can be computed numerically using differentiable convex optimisation (surprisingly easy with `cvxpylayers`<sup>2</sup>).

---

<sup>2</sup>A. Agrawal et al. (2019) Differentiable convex optimization layers. NeurIPS

## Regret Bounds

# Bounding the Radius

The full expression for the (squared) AMM radius is

$$\begin{aligned} R_{\text{AMM},t}^2 = & \mathbf{r}_t \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right)^{-1} \mathbf{r}_t + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta) \\ & + \alpha B^2 - \mathbf{r}_t^\top \mathbf{r}_t + \mathbf{r}_t^\top \Phi_t \left( \Phi_t^\top \Phi_t + \alpha \mathbf{I} \right)^{-1} \Phi_t^\top \mathbf{r}_t. \end{aligned}$$



# Bounding the Radius

The full expression for the (squared) AMM radius is

$$\begin{aligned} R_{\text{AMM},t}^2 = & \mathbf{r}_t \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right)^{-1} \mathbf{r}_t + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta) \\ & + \alpha B^2 - \mathbf{r}_t^\top \mathbf{r}_t + \mathbf{r}_t^\top \Phi_t \left( \Phi_t^\top \Phi_t + \alpha \mathbf{I} \right)^{-1} \Phi_t^\top \mathbf{r}_t. \end{aligned}$$

Using the Matrix Inversion Lemma, we have

$$\mathbf{r}_t \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right)^{-1} \mathbf{r}_t = \mathbf{r}_t^\top \mathbf{r}_t - \mathbf{r}_t^\top \Phi_t \left( \Phi_t^\top \Phi_t + \frac{\sigma^2}{c} \mathbf{I} \right)^{-1} \Phi_t^\top \mathbf{r}_t.$$

# Bounding the Radius

The full expression for the (squared) AMM radius is

$$R_{\text{AMM},t}^2 = \mathbf{r}_t \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right)^{-1} \mathbf{r}_t + \sigma^2 \ln \left( \det \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right) \right) + 2\sigma^2 \ln(1/\delta) \\ + \alpha B^2 - \mathbf{r}_t^\top \mathbf{r}_t + \mathbf{r}_t^\top \Phi_t \left( \Phi_t^\top \Phi_t + \alpha \mathbf{I} \right)^{-1} \Phi_t^\top \mathbf{r}_t.$$

Using the Matrix Inversion Lemma, we have

$$\mathbf{r}_t \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right)^{-1} \mathbf{r}_t = \mathbf{r}_t^\top \mathbf{r}_t - \mathbf{r}_t^\top \Phi_t \left( \Phi_t^\top \Phi_t + \frac{\sigma^2}{c} \mathbf{I} \right)^{-1} \Phi_t^\top \mathbf{r}_t.$$

If we set  $\alpha = \sigma^2/c$ , the quadratic terms cancel, and we can use the determinant-trace inequality<sup>\*</sup>

$$R_{\text{AMM},t}^2 = \sigma^2 \left( \ln \left( \det \left( \mathbf{I} + \frac{c\Phi_t\Phi_t^\top}{\sigma^2} \right) \right) + \frac{B^2}{c} + 2 \ln(1/\delta) \right) \leq \sigma^2 \left( d \ln \left( 1 + \frac{ctL^2}{\sigma^2 d} \right) + \frac{B^2}{c} + 2 \ln(1/\delta) \right)$$

---

<sup>\*</sup> Assuming  $\|\phi(a)\|_2 \leq L$ .

# OFUL Analysis

**Step 1.** Use optimism to bound the cumulative regret by the confidence bound widths.

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

# OFUL Analysis

**Step 1.** Use optimism to bound the cumulative regret by the confidence bound widths.

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

**Step 2.** For both CMM-UCB and AMM-UCB, we have

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

# OFUL Analysis

**Step 1.** Use optimism to bound the cumulative regret by the confidence bound widths.

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

**Step 2.** For both CMM-UCB and AMM-UCB, we have

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

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

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

## **Comparison With OFUL**

# How Does AMM-UCB Compare To OFUL?

We derive and use a bound on the norm of the noise vector

$$\|\epsilon_t\|_2 = \|\Phi_t \theta^* - r_t\|_2 \leq R_{\text{MM},t}.$$

# How Does AMM-UCB Compare To OFUL?

We derive and use a bound on the norm of the noise vector

$$\|\epsilon_t\|_2 = \|\Phi_t \theta^* - r_t\|_2 \leq R_{\text{MM},t}.$$

For  $c = \sigma^2/\alpha$  and any  $\alpha > 0$ , this leads to the inequality

$$\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})} \leq R_{\text{AMM},t} = \sqrt{\sigma^2 \ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + \mathbf{I} \right) \right) + \alpha B^2 + 2\sigma^2 \ln(1/\delta)}.$$



# How Does AMM-UCB Compare To OFUL?

We derive and use a bound on the norm of the noise vector

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

For  $c = \sigma^2/\alpha$  and any  $\alpha > 0$ , this leads to the inequality

$$\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})} \leq R_{\text{AMM},t} = \sqrt{\sigma^2 \ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + \mathbf{I} \right) \right) + \alpha B^2 + 2\sigma^2 \ln(1/\delta)}.$$

OFUL uses a bound on the (weighted) norm of the projection of the noise vector

$$\|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})^{-1}} \leq \sigma \sqrt{\ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + \mathbf{I} \right) \right) + 2 \ln(1/\delta)}.$$

# How Does AMM-UCB Compare To OFUL?

We derive and use a bound on the norm of the noise vector

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

For  $c = \sigma^2/\alpha$  and any  $\alpha > 0$ , this leads to the inequality

$$\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})} \leq R_{\text{AMM},t} = \sqrt{\sigma^2 \ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + \mathbf{I} \right) \right) + \alpha B^2 + 2\sigma^2 \ln(1/\delta)}.$$

OFUL uses a bound on the (weighted) norm of the projection of the noise vector

$$\|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})^{-1}} \leq \sigma \sqrt{\ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + \mathbf{I} \right) \right) + 2 \ln(1/\delta)}.$$

This leads to a similar, but looser (due to  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ ) inequality

$$\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha \mathbf{I})} \leq \sigma \sqrt{\ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + \mathbf{I} \right) \right) + 2 \ln(1/\delta) + \sqrt{\alpha} B}.$$

# Why Is AMM-UCB Better Than OFUL?

Bounds on  $\|\Phi_t \theta^* - r_t\|_2$  and  $\|\theta^*\|_2$  fit together better than bounds on  $\|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}}$  and  $\|\theta^*\|_2$ .

# Why Is AMM-UCB Better Than OFUL?

Bounds on  $\|\Phi_t \theta^* - r_t\|_2$  and  $\|\theta^*\|_2$  fit together better than bounds on  $\|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}}$  and  $\|\theta^*\|_2$ .

**OFUL:** Using the definition of  $\hat{\theta}_{\alpha,t}$ , and then the triangle inequality,

$$\begin{aligned}\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)} &= \|\Phi_t^\top \epsilon_t + \alpha \theta^*\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} \\ &\leq \|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} + \alpha \|\theta^*\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} \\ &\leq \sigma \sqrt{\ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + I \right) \right) + 2 \ln(1/\delta)} + \sqrt{\alpha} B\end{aligned}$$

# Why Is AMM-UCB Better Than OFUL?

Bounds on  $\|\Phi_t \theta^* - r_t\|_2$  and  $\|\theta^*\|_2$  fit together better than bounds on  $\|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}}$  and  $\|\theta^*\|_2$ .

**OFUL:** Using the definition of  $\hat{\theta}_{\alpha,t}$ , and then the triangle inequality,

$$\begin{aligned}\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)} &= \|\Phi_t^\top \epsilon_t + \alpha \theta^*\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} \\ &\leq \|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} + \alpha \|\theta^*\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} \\ &\leq \sigma \sqrt{\ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + I \right) \right) + 2 \ln(1/\delta)} + \sqrt{\alpha} B\end{aligned}$$

The triangle inequality step causes the  $\ln \det$  and  $\alpha B^2$  terms to appear under separate square roots.

# Why Is AMM-UCB Better Than OFUL?

Bounds on  $\|\Phi_t \theta^* - \mathbf{r}_t\|_2$  and  $\|\theta^*\|_2$  fit together better than bounds on  $\|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}}$  and  $\|\theta^*\|_2$ .

**OFUL:** Using the definition of  $\hat{\theta}_{\alpha,t}$ , and then the triangle inequality,

$$\begin{aligned}\|\theta^* - \hat{\theta}_{\alpha,t}\|_{(\Phi_t^\top \Phi_t + \alpha I)} &= \|\Phi_t^\top \epsilon_t + \alpha \theta^*\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} \\ &\leq \|\Phi_t^\top \epsilon_t\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} + \alpha \|\theta^*\|_{(\Phi_t^\top \Phi_t + \alpha I)^{-1}} \\ &\leq \sigma \sqrt{\ln \left( \det \left( \frac{1}{\alpha} \Phi_t^\top \Phi_t + I \right) \right) + 2 \ln(1/\delta)} + \sqrt{\alpha} B\end{aligned}$$

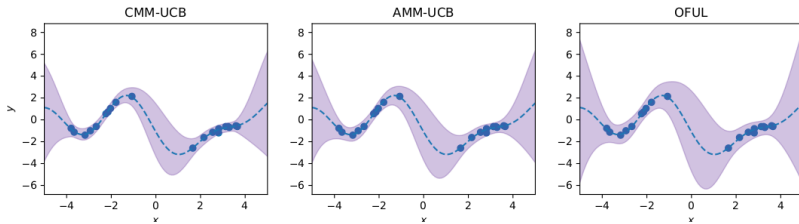
The triangle inequality step causes the  $\ln \det$  and  $\alpha B^2$  terms to appear under separate square roots.

**Ours:** We combine our constraints by completing the square on the LHS of

$$\|\Phi_t \theta^* - \mathbf{r}_t\|_2^2 + \alpha \|\theta^*\|_2^2 \leq R_{\text{MM},t}^2 + \alpha B^2$$

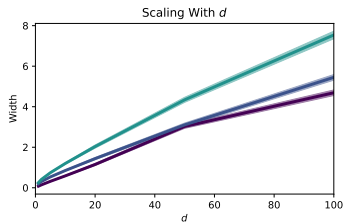
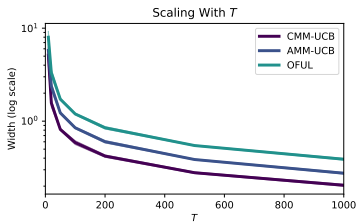
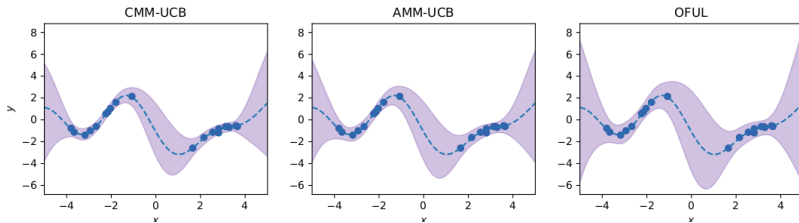
## **Some Experimental Results**

# Confidence Bound Comparison



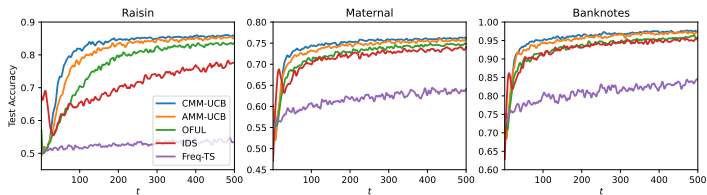


# Confidence Bound Comparison



# Hyperparameter Tuning

	Raisin		Maternal		Banknotes	
	Mean Acc	Max Acc	Mean Acc	Max Acc	Mean Acc	Max Acc
CMM-UCB (Ours)	<b>0.818</b> $\pm$ 0.018	<b>0.893</b> $\pm$ 0.019	<b>0.744</b> $\pm$ 0.020	<b>0.829</b> $\pm$ 0.023	<b>0.954</b> $\pm$ 0.005	<b>1.000</b> $\pm$ 0.000
AMM-UCB (Ours)	0.800 $\pm$ 0.017	0.892 $\pm$ 0.020	0.736 $\pm$ 0.020	<b>0.829</b> $\pm$ 0.023	0.948 $\pm$ 0.005	<b>1.000</b> $\pm$ 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	<b>1.000</b> $\pm$ 0.000
IDS <sup>3</sup>	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	<b>1.000</b> $\pm$ 0.000
Freq-TS <sup>4</sup>	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	<b>1.000</b> $\pm$ 0.000



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

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

## Open Questions

# When Do More Adaptive Mixture Distributions Help?

The means  $\boldsymbol{\mu}_t$  and covariances  $\boldsymbol{T}_t$  of the standard mixture distributions can be written in the form

$$\boldsymbol{\mu}_t = \begin{bmatrix} m(a_1) \\ m(a_2) \\ \vdots \\ m(a_t) \end{bmatrix}, \quad \boldsymbol{T}_t = \begin{bmatrix} k(a_1, a_1) & k(a_1, a_2) & \cdots & k(a_1, a_t) \\ k(a_2, a_1) & k(a_2, a_2) & \cdots & k(a_2, a_t) \\ \vdots & \vdots & \ddots & \vdots \\ k(a_t, a_1) & k(a_t, a_2) & \cdots & k(a_t, a_t) \end{bmatrix},$$

where  $m(a) = 0$  and  $k(a, a') = \phi(a)^\top \phi(a')$ .

# When Do More Adaptive Mixture Distributions Help?

The means  $\mu_t$  and covariances  $T_t$  of the standard mixture distributions can be written in the form

$$\mu_t = \begin{bmatrix} m(a_1) \\ m(a_2) \\ \vdots \\ m(a_t) \end{bmatrix}, \quad T_t = \begin{bmatrix} k(a_1, a_1) & k(a_1, a_2) & \cdots & k(a_1, a_t) \\ k(a_2, a_1) & k(a_2, a_2) & \cdots & k(a_2, a_t) \\ \vdots & \vdots & \ddots & \vdots \\ k(a_t, a_1) & k(a_t, a_2) & \cdots & k(a_t, a_t) \end{bmatrix},$$

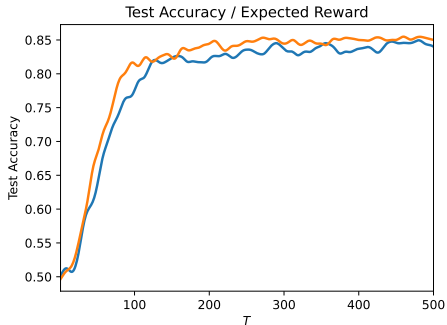
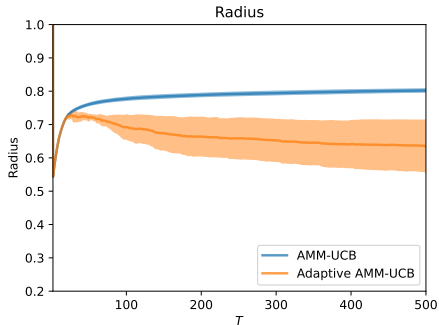
where  $m(a) = 0$  and  $k(a, a') = \phi(a)^\top \phi(a')$ .

We also tried out

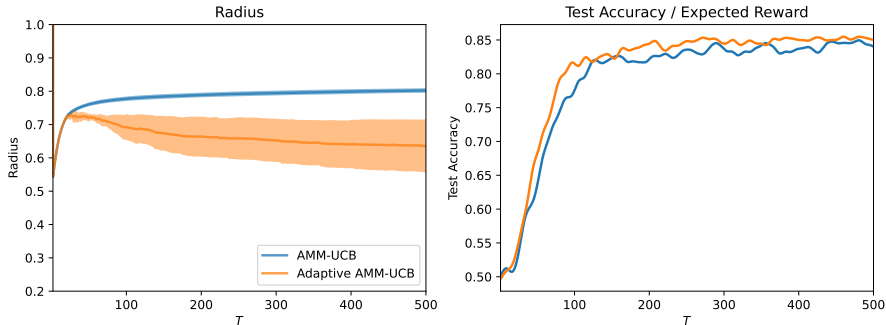
$$\mu_t = \begin{bmatrix} m_0(a_1) \\ m_1(a_2) \\ \vdots \\ m_{t-1}(a_t) \end{bmatrix}, \quad T_t = \begin{bmatrix} k_0(a_1, a_1) & k_1(a_1, a_2) & \cdots & k_{t-1}(a_1, a_t) \\ k_1(a_2, a_1) & k_1(a_2, a_2) & \cdots & k_{t-1}(a_2, a_t) \\ \vdots & \vdots & \ddots & \vdots \\ k_{t-1}(a_t, a_1) & k_{t-1}(a_t, a_2) & \cdots & k_{t-1}(a_t, a_t) \end{bmatrix},$$

where  $m_t(a) = \mathbf{k}_t(a)^\top (\mathbf{K}_t + \beta \mathbf{I})^{-1} \mathbf{r}_t$  and  $k_t(a, a') = k(a, a') - \mathbf{k}_t(a)^\top (\mathbf{K}_t + \beta \mathbf{I})^{-1} \mathbf{k}_t(a')$ .

# When Do More Adaptive Mixture Distributions Help?



# When Do More Adaptive Mixture Distributions Help?



Adaptive mixture distributions don't always help this much though.

Thank you for listening!