

Linear Bandits

Wen Sun

CS 6789: Foundations of Reinforcement Learning

Recap on MAB

Setting:

We have K many arms: a_1, \dots, a_K

Recap on MAB

Setting:

We have K many arms: a_1, \dots, a_K

Each arm has a unknown reward distribution, i.e., $\nu_i \in \Delta([0,1])$,
w/ mean $\mu_i = \mathbb{E}_{r \sim \nu_i}[r]$

Regret

More formally, we have the following learning objective:

$$\text{Regret}_T = T\mu^\star - \sum_{t=0}^{T-1} \mu_{I_t} \quad \mu^\star = \max_{i \in [K]} \mu_i$$

Regret

More formally, we have the following learning objective:

$$\text{Regret}_T = T\mu^\star - \sum_{t=0}^{T-1} \mu_{I_t} \quad \mu^\star = \max_{i \in [K]} \mu_i$$

Total expected reward if we pulled best arm over T rounds

Regret

More formally, we have the following learning objective:

$$\text{Regret}_T = T\mu^\star - \sum_{t=0}^{T-1} \mu_{I_t} \quad \mu^\star = \max_{i \in [K]} \mu_i$$

Total expected reward if we pulled best arm over T rounds

Total expected reward of the arms we pulled over T rounds

Regret

More formally, we have the following learning objective:

$$\text{Regret}_T = T\mu^\star - \sum_{t=0}^{T-1} \mu_{I_t} \quad \mu^\star = \max_{i \in [K]} \mu_i$$

Total expected reward if we pulled best arm over T rounds

Total expected reward of the arms we pulled over T rounds

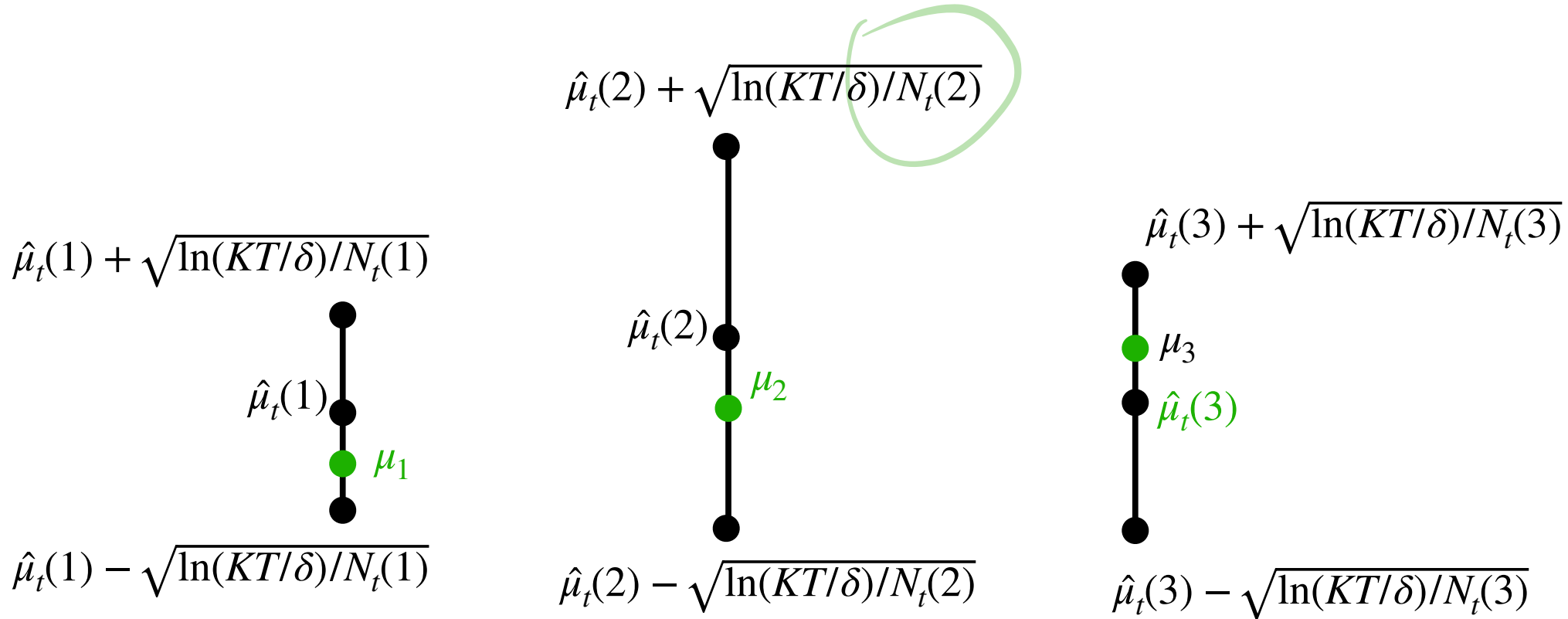
Goal: no-regret, i.e., $\text{Regret}_T/T \rightarrow 0$, as $T \rightarrow \infty$

UCB: Optimism in the face of Uncertainty

Given the confidence interval, we pick arm that has the **highest Upper-Conf-Bound:**

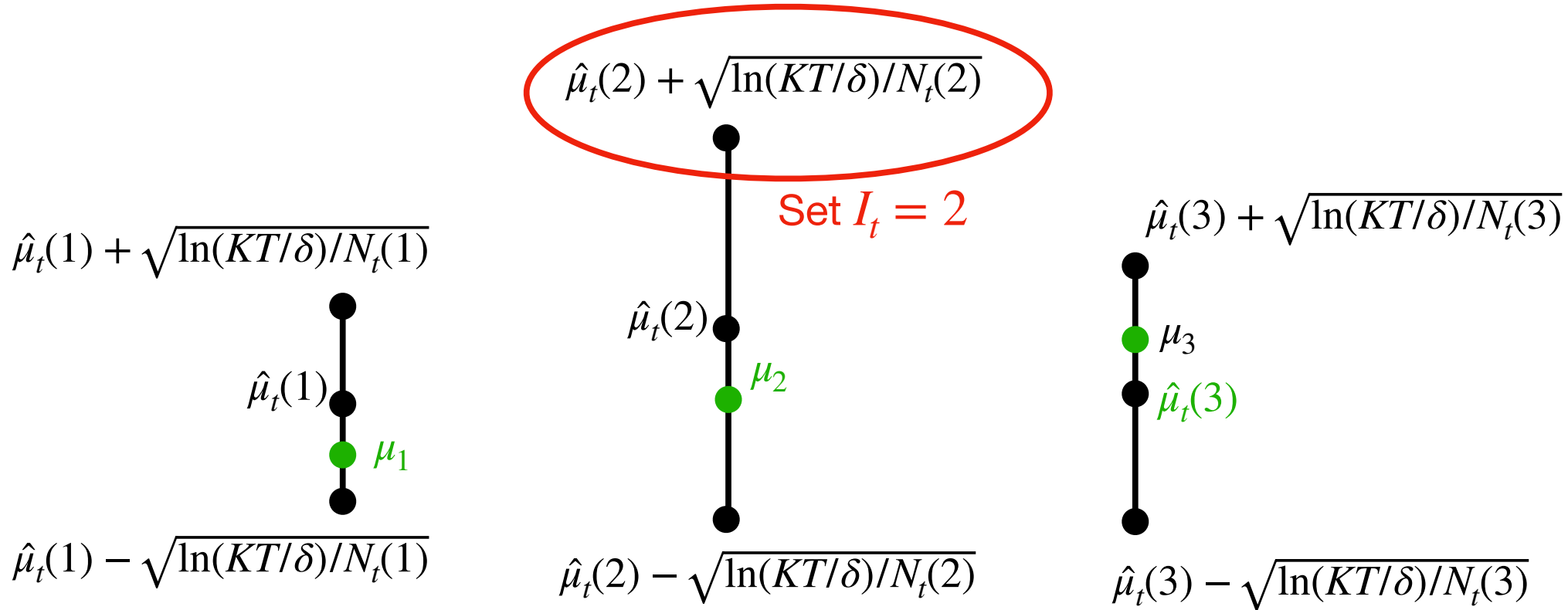
UCB: Optimism in the face of Uncertainty

Given the confidence interval, we pick arm that has the **highest Upper-Conf-Bound:**



UCB: Optimism in the face of Uncertainty

Given the confidence interval, we pick arm that has the **highest Upper-Conf-Bound:**



Optimism in the face of uncertainty

Denote the optimal arm $I^* = \arg \max_{i \in [K]} \mu_i$; recall $I_t = \arg \max_{i \in [K]} \hat{\mu}_t(i) + \sqrt{\frac{\ln(KT/\delta)}{N_t(i)}}$

Optimism in the face of uncertainty

Denote the optimal arm $I^\star = \arg \max_{i \in [K]} \mu_i$; recall $I_t = \arg \max_{i \in [K]} \hat{\mu}_t(i) + \sqrt{\frac{\ln(KT/\delta)}{N_t(i)}}$

$$\text{Regret-at-t} = \mu^\star - \mu_{I_t}$$

Optimism in the face of uncertainty

Denote the optimal arm $I^* = \arg \max_{i \in [K]} \mu_i$; recall $I_t = \arg \max_{i \in [K]} \hat{\mu}_t(i) + \sqrt{\frac{\ln(KT/\delta)}{N_t(i)}}$

$$\text{Regret-at-t} = \mu^* - \mu_{I_t}$$

$$\leq \underbrace{\hat{\mu}_t(I_t)}_{\Delta} + \sqrt{\frac{\ln(TK/\delta)}{N_t(I_t)}} - \underbrace{\mu_{I_t}}_{\Delta}$$

$$\leq 2 \sqrt{\frac{\ln(\dots)}{N_t(I_t)}}$$

Optimism in the face of uncertainty

Denote the optimal arm $I^* = \arg \max_{i \in [K]} \mu_i$; recall $I_t = \arg \max_{i \in [K]} \hat{\mu}_t(i) + \sqrt{\frac{\ln(KT/\delta)}{N_t(i)}}$

$$\text{Regret-at-t} = \mu^* - \mu_{I_t}$$

$$\leq \hat{\mu}_t(I_t) + \sqrt{\frac{\ln(TK/\delta)}{N_t(I_t)}} - \mu_{I_t} \leq 2\sqrt{\frac{\ln(TK/\delta)}{N_t(I_t)}}$$

optimism

$$\begin{aligned} \sum_{t=0}^T \text{Regret-at-} t &\leq \sum_{t=0}^T \sqrt{\frac{\ln(\dots)}{N_t(I_t)}} \\ &\leq \sqrt{KT} \end{aligned}$$

Today:

MAB w/ K arms has regret $O(\sqrt{KT})$

What if there are infinitely many actions?

Introducing structures in the reward function

Outline for Today:

1. Linear Bandit Setting

2. Algorithm: LinUCB

3. Regret analysis of LinUCB

Linear Bandit Setting

We have an action set $D \subset \mathbb{R}^d$

Linear Bandit Setting

We have an action set $D \subset \mathbb{R}^d$

Expected reward of each action $x \in D$ is linear:

$$\mathbb{E}[r|x] = (\mu^\star)^\top x$$

$$\mu^\star \in \mathbb{R}^d$$

Linear Bandit Setting

We have an action set $D \subset \mathbb{R}^d$

Expected reward of each action $x \in D$ is linear:

$$\mathbb{E}[r | x] = (\mu^\star)^\top x$$

Every time we pick an action $x \in D$, we observe a noisy reward

$$r = \mu^\star \cdot x + \eta$$

$$a \cdot b = a^T b$$

Linear Bandit Setting

We have an action set $D \subset \mathbb{R}^d$

Expected reward of each action $x \in D$ is linear:

$$\mathbb{E}[r|x] = (\mu^\star)^T x$$

$$\mu^\star \in \mathbb{R}^d$$

Every time we pick an action $x \in D$, we observe a noisy reward

$$r = \mu^\star \cdot x + \eta$$

Zero mean i.i.d noise

Learning protocol and goal:

For $t = 1$ to T :

Learning protocol and goal:

For $t = 1$ to T :

Learner selects $x_t \in D$ (based on history)

Learning protocol and goal:

For $t = 1$ to T :

Learner selects $x_t \in D$ (based on history)

Learner observes a noisy reward, i.e., $r_t = \mu^\star \cdot x_t + \eta_t$

Learning protocol and goal:

For $t = 1$ to T :

Learner selects $x_t \in D$ (based on history)

Learner observes a noisy reward, i.e., $r_t = \mu^\star \cdot x_t + \eta_t$

Goal: minimize regret

$$\text{Regret} := T\mu^\star \cdot x^\star - \sum_{t=0}^{T-1} \mu^\star \cdot x_t$$

Learning protocol and goal:

For $t = 1$ to T :

Learner selects $x_t \in D$ (based on history)

Learner observes a noisy reward, i.e., $r_t = \mu^\star \cdot x_t + \eta_t$

Goal: minimize regret

$$\text{Regret} := T\mu^\star \cdot x^\star - \sum_{t=0}^{T-1} \mu^\star \cdot x_t$$

Expected
reward
of x_t

$$x^\star = \arg \max_{x \in D} \mu^\star \cdot x$$

Outline for Today:

1. Linear Bandit Setting

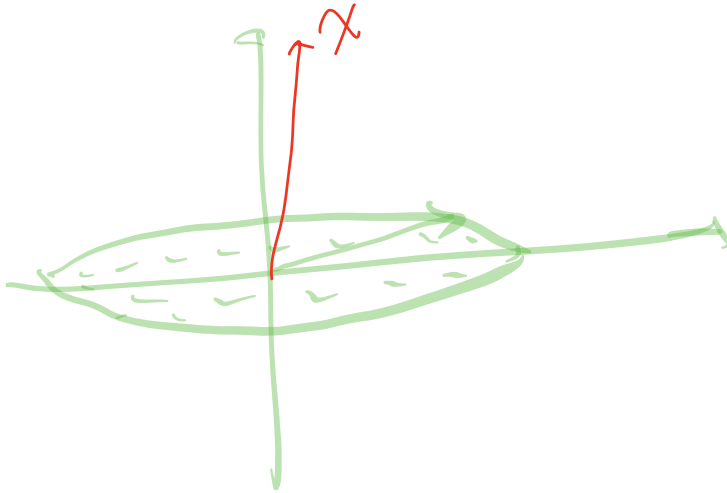
2. Algorithm: LinUCB

3. Regret analysis of LinUCB

LinUCB algorithm

Overall idea:

Ridge linear regression for learning μ^* + design exploration bonus



MAR:

$$\hat{\mu} + \sqrt{\frac{1}{N_t(F_t)}}$$

LinUCB algorithm

In iteration t :

1. Perform Ridge LR on data $\{x_i, r_i\}_{i=0}^{t-1}$:

$$\text{Set } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

LinUCB algorithm

In iteration t :

1. Perform Ridge LR on data $\{x_i, r_i\}_{i=0}^{t-1}$:

$$\text{Set } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

2: Set exploration bonus: $b_t(x) = \beta \sqrt{x^\top \Sigma_t^{-1} x}$
 \triangle

$$\Sigma_t = \sum_{i=0}^{t-1} x_i x_i^\top + \lambda I$$

LinUCB algorithm

In iteration t :

1. Perform Ridge LR on data $\{x_i, r_i\}_{i=0}^{t-1}$:

$$\text{Set } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

2: Set exploration bonus: $b_t(x) = \beta \sqrt{x^\top \Sigma_t^{-1} x}$

3: Play optimistically, i.e., $x_t = \arg \max_{x \in D} \hat{\mu}_t^\top x_t + b_t(x)$

Outline for Today:

1. Linear Bandit Setting

2. Algorithm: LinUCB

$$\tilde{O} \sqrt{d \cdot T}$$

↑
Discretization

3. Regret analysis of LinUCB

$$\tilde{O} \sqrt{d^2 T}$$

Analysis of Ridge Linear Regression

$$\text{Recall } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

Analysis of Ridge Linear Regression

$$\hat{\mu}_x - \mu^*$$

$$\text{Recall } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

$$\hat{\mu}_t = \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i r_i$$

$$\Sigma_t = \sum_{i=0}^{t-1} x_i x_i^\top + \lambda I$$

Regularization

$$r_i = \mu^{*\top} x_i + \eta_i$$

Analysis of Ridge Linear Regression

$$\text{Recall } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

$$\begin{aligned} \hat{\mu}_t &= \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i r_i \\ &= \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i (x_i^\top \mu^* + \eta_i) \end{aligned}$$

Analysis of Ridge Linear Regression

$$\text{Recall } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

$$\begin{aligned} \hat{\mu}_t &= \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i r_i \\ &= \Sigma_t^{-1} \left(\sum_{i=0}^{t-1} x_i (x_i^\top \mu^* + \eta_i) \right) = \Sigma_t^{-1} (\Sigma_t - \lambda I) \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i \end{aligned}$$

$\hookrightarrow \sum_{i=0}^{t-1} x_i x_i^\top$

Analysis of Ridge Linear Regression

$$\text{Recall } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

$$\hat{\mu}_t = \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i r_i$$

$$= \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i (x_i^\top \mu^\star + \eta_i) = \cancel{\Sigma_t^{-1} (\Sigma_t - \lambda I) \mu^\star} + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

$$= \mu^\star - \lambda \Sigma_t^{-1} \mu^\star + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Analysis of Ridge Linear Regression

$$\text{Recall } \hat{\mu}_t := \arg \min_{\mu} \sum_{i=0}^{t-1} (\mu^\top x_i - r_i)^2 + \lambda \|\mu\|_2^2$$

$$\hat{\mu}_t = \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i r_i$$

$$= \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i (x_i^\top \mu^\star + \eta_i) = \Sigma_t^{-1} (\Sigma_t - \lambda I) \mu^\star + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

$$= \mu^\star - \lambda \Sigma_t^{-1} \mu^\star + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

$$\hat{\mu}_t - \mu^\star = -\lambda \Sigma_t^{-1} \mu^\star + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Analysis of Ridge Linear Regression

$$\Sigma_t = \sum_{i=0}^{t-1} x_i x_i^T + \lambda I$$
$$\hat{\mu}_t - \mu^* = -\lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Let us look at the training error:

$$\sqrt{(\hat{\mu}_t - \mu^*)^T \Sigma_t (\hat{\mu}_t - \mu^*)}$$

$$\Rightarrow \sum_{i=1}^{t-1} \left((\hat{\mu}_x - \mu^*)^T x_i \right)^2 + |(\hat{\mu}_x - \mu^*)|^2 \cdot \lambda$$

Analysis of Ridge Linear Regression

$$\hat{\mu}_t - \mu^* = -\lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Let us look at the training error:

$$\sqrt{(\hat{\mu}_t - \mu^*)^T \Sigma_t (\hat{\mu}_t - \mu^*)} \leq \left\| \lambda \Sigma_t^{-1/2} \mu^* \right\| + \left\| \Sigma_t^{-1/2} \sum_{i=0}^{t-1} \eta_i x_i \right\|$$

Triangle Ineq

$$\left\| \Sigma_t^{-1/2} (\hat{\mu}_t - \mu^*) \right\| = \left\| -\lambda \Sigma_t^{-1/2} \mu^* + \Sigma_t^{-1/2} \sum_{i=0}^{t-1} x_i \eta_i \right\|$$

Analysis of Ridge Linear Regression

$$\hat{\mu}_t - \mu^* = -\lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Let us look at the training error:

$$\sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} \leq \left\| \lambda \Sigma_t^{-1/2} \mu^* \right\| + \left\| \Sigma_t^{-1/2} \sum_{i=0}^{t-1} \eta_i x_i \right\|$$
$$\leq \sqrt{\lambda} \|\mu^*\| + ???$$

$$\sigma_{\max} \left(\Sigma_t^{-1/2} \right) \leq \frac{1}{\sqrt{\lambda}}$$

Analysis of Ridge Linear Regression

$$\hat{\mu}_t - \mu^* = -\lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Let us look at the training error:

$$\begin{aligned} \sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} &\leq \left\| \lambda \Sigma_t^{-1/2} \mu^* \right\| + \left\| \Sigma_t^{-1/2} \sum_{i=0}^{t-1} \eta_i x_i \right\| \\ &\leq \sqrt{\lambda} \|\mu^*\| + ??? \end{aligned}$$

Self-normalized Martingale bound

$$\left\| \Sigma_t^{-1/2} \sum_{i=0}^{t-1} \eta_i x_i \right\|^2 \approx d \log(t)$$

Self-normalized Bound for Vector-valued Martingales

Suppose $\{\eta_i\}_{i=0}^{\infty}$ are mean zero random variables, and $|\eta_i| \leq \sigma$;

Let $\{x_i\}_{i=0}^{\infty}$ be any sequence of random vectors with $\|x_i\| \leq 1$, then w/
prob $1 - \delta$, for all $t \geq 1$,

$$\left\| \Sigma_t^{-1/2} \sum_{i=0}^{t-1} x_i \eta_i \right\|^2 \leq \sigma^2 d \cdot \left(\ln \left(\frac{t}{\lambda} + 1 \right) + \ln(1/\delta) \right)$$

$$\left(\sum_{i=0}^{t-1} x_i \eta_i \right)^T \Sigma_t \left(\sum_{i=0}^{t-1} x_i \eta_i \right) \leq \sum_{i=0}^{t-1} x_i x_i^T + \lambda I$$

Analysis of Ridge Linear Regression (Continue)

$$\hat{\mu}_t - \mu^* = -\lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

Let us look at the training error:

$$\begin{aligned} \sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} &\leq \left\| \lambda \Sigma_t^{-1/2} \mu^* \right\| + \left\| \Sigma_t^{-1/2} \sum_{i=0}^{t-1} \eta_i x_i \right\| \approx \sqrt{d \ln(T)} \\ &\lesssim \underbrace{\sqrt{\lambda}} + \underbrace{\sigma \sqrt{d \cdot \ln(T/(\lambda \delta))}} \end{aligned}$$

Summary for Ridge Linear Regression

$$\hat{\mu}_t - \mu^* = -\lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{i=0}^{t-1} x_i \eta_i$$

$$\sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} \lesssim \sqrt{\lambda} + \sigma \sqrt{d \ln(T/(\lambda \delta))}$$

Optimism

$$\sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} \lesssim \sqrt{\lambda} + \sigma^2 d \ln(T/(\lambda\delta))$$

Let's construct uncertainty quantification for each action $x \in D$

$$|\hat{\mu}_t \cdot x - \mu^* \cdot x|$$

Optimism

$$\sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} \lesssim \sqrt{\lambda} + \sqrt{\sigma^2 d \ln(T/(\lambda\delta))}$$

Let's construct uncertainty quantification for each action $x \in D$

$$\begin{aligned} |\hat{\mu}_t \cdot x - \mu^* \cdot x| &\stackrel{CS}{\leq} \|\hat{\mu}_t - \mu^*\|_{\Sigma_t} \cdot \|x\|_{\Sigma_t^{-1}} \\ &\leq \sqrt{\lambda} + \sqrt{\sigma^2 d \ln(T)} \end{aligned}$$

Optimism

$$\sqrt{(\hat{\mu}_t - \mu^\star)^\top \Sigma_t (\hat{\mu}_t - \mu^\star)} \lesssim \sqrt{\lambda} + \sigma^2 d \ln(T/(\lambda\delta))$$

Let's construct uncertainty quantification for each action $x \in D$

$$\begin{aligned} |\hat{\mu}_t \cdot x - \mu^\star \cdot x| &\leq \|\hat{\mu}_t - \mu^\star\|_{\Sigma_t} \cdot \|x\|_{\Sigma_t^{-1}} \\ &\lesssim \left(\sqrt{\lambda} + \sigma \sqrt{d \ln(T/(\lambda\delta))} \right) \cdot \|x\|_{\Sigma_t^{-1}} \end{aligned}$$

Optimism

$$\sqrt{(\hat{\mu}_t - \mu^\star)^\top \Sigma_t (\hat{\mu}_t - \mu^\star)} \lesssim \sqrt{\lambda} + \sigma^2 d \ln(T/(\lambda\delta))$$

Let's construct uncertainty quantification for each action $x \in D$

$$|\hat{\mu}_t \cdot x - \mu^\star \cdot x| \leq \|\hat{\mu}_t - \mu^\star\|_{\Sigma_t} \cdot \|x\|_{\Sigma_t^{-1}}$$

$$\lesssim \left(\sqrt{\lambda} + \sigma \sqrt{d \ln(T/(\lambda\delta))} \right) \cdot \|x\|_{\Sigma_t^{-1}}$$

$$b_t(x) := \beta \cdot \|x\|_{\Sigma_t^{-1}}$$

Optimism

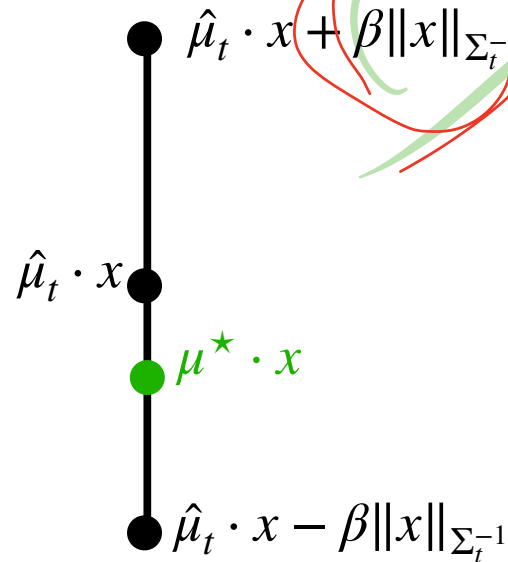
$$\sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} \lesssim \sqrt{\lambda} + \sigma^2 d \ln(T/(\lambda\delta))$$

$\forall x \in D$ Let's construct uncertainty quantification for each action $x \in D$

$$|\hat{\mu}_t \cdot x - \mu^* \cdot x| \leq \|\hat{\mu}_t - \mu^*\|_{\Sigma_t} \cdot \|x\|_{\Sigma_t^{-1}}$$

$$\lesssim \left(\sqrt{\lambda} + \sigma \sqrt{d \ln(T/(\lambda\delta))} \right) \cdot \|x\|_{\Sigma_t^{-1}}$$

$$b_t(x) := \beta \cdot \|x\|_{\Sigma_t^{-1}}$$



Optimism

$$\text{Optimism: } \mu^* \cdot x^* \leq \hat{\mu}_t \cdot x_t + \beta \|x_t\|_{\Sigma_t^{-1}}$$

\triangle

$\forall x \in D$

Proof:

$$x_t = \operatorname{argmax}_x \hat{\mu}_t \cdot x + \beta \|x\|_{\Sigma_t^{-1}}$$

$$\hat{\mu}_t \cdot x + \beta \|x\|_{\Sigma_t^{-1}}$$



$$\hat{\mu}_t \cdot x$$



$$\mu^* \cdot x$$



$$\hat{\mu}_t \cdot x - \beta \|x\|_{\Sigma_t^{-1}}$$



Regret

$$\text{Regret-at-t} = \mu^* \cdot x^* - \mu^* \cdot x_t$$

Regret

$$\text{Regret-at-t} = \mu^* \cdot x^* - \mu^* \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^* \cdot x_t$$

$$\left| \hat{\mu}_t^\top x_t - \mu^{*T} \cdot x_t \right| \leq \beta \|x_t\|_{\Sigma_t^{-1}}$$

Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

Intuitively this should be convincing already:

$$\mu^\star \cdot x^\star - \mu^\star \cdot x_t \geq \delta$$

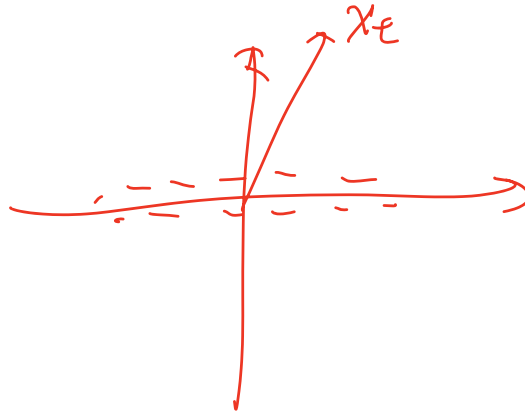
Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

Intuitively this should be convincing already:

Case 1: x_t is a bad arm, i.e., $2\beta \|x_t\|_{\Sigma_t^{-1}} \geq \mu^\star \cdot (x^\star - x_t) \geq \delta$



Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

Intuitively this should be convincing already:

Case 1: x_t is a bad arm, i.e., $2\beta \|x_t\|_{\Sigma_t^{-1}} \geq \mu^\star \cdot (x^\star - x_t) \geq \delta$

x_t falls in the subspace where “data is sparse”, i.e., we explored!

Regret

$$\begin{aligned} \text{Regret-at-t} &= \mu^\star \cdot x^\star - \mu^\star \cdot x_t \\ &\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq \underbrace{2\beta \|x_t\|_{\Sigma_t^{-1}}}_{\text{Small}} \end{aligned}$$

Intuitively this should be convincing already:

Case 1: x_t is a bad arm, i.e., $2\beta \|x_t\|_{\Sigma_t^{-1}} \geq \mu^\star \cdot (x^\star - x_t) \geq \delta$

x_t falls in the subspace where “data is sparse”, i.e., we explored!

Case 2: confidence interval $\|x_t\|_{\Sigma_t^{-1}}$ is small

Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

Intuitively this should be convincing already:

Case 1: x_t is a bad arm, i.e., $2\beta \|x_t\|_{\Sigma_t^{-1}} \geq \mu^\star \cdot (x^\star - x_t) \geq \delta$

x_t falls in the subspace where “data is sparse”, i.e., we explored!

Case 2: confidence interval $\|x_t\|_{\Sigma_t^{-1}}$ is small

Then regret at this round is small too, i.e., we exploited!

Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

More formally, we can show:

$$\text{Regret} \leq \beta \sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}}$$

Regret

$$\text{Regret-at-t} = \mu^* \cdot x^* - \mu^* \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^* \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

More formally, we can show:

$$\text{Regret} \leq \beta \sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}} \leq \beta \sqrt{T} \cdot \sqrt{\sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}}^2}$$

$\beta = \sqrt{\lambda} + 6\sqrt{d \ln(T)}$

$\sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}}^2 \leq d \ln(T)$

$= d \sqrt{T}$

Regret

$$\text{Regret-at-t} = \mu^\star \cdot x^\star - \mu^\star \cdot x_t$$

$$\leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}} - \mu^\star \cdot x_t \leq 2\beta \|x_t\|_{\Sigma_t^{-1}}$$

More formally, we can show:

$$\begin{aligned} \text{Regret} &\leq \beta \sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}} \leq \beta \sqrt{T} \cdot \sqrt{\sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}}^2} \\ &\lesssim \beta \sqrt{T} \cdot \sqrt{d \ln(T/\lambda + 1)} \quad \forall \lambda \geq 1 \end{aligned}$$

$$\sum_{t=0}^{T-1} \|x_t\|_{\Sigma_t^{-1}}^2 \leq O(d \ln T)$$

Summary

1. To deal w/ infinitely many arms, we introduce linear structure in rewards

$$\begin{matrix} \lambda \\ \phi(x) \end{matrix}$$

$$\begin{aligned} K(x, x') \\ = \langle \phi(x), \phi(x') \rangle \end{aligned}$$

Summary

1. To deal w/ infinitely many arms, we introduce linear structure in rewards
2. Analysis of Ridge LR gives us bound on $|(\mu^\star - \hat{\mu}_t)^\top x|$

Summary

1. To deal w/ infinitely many arms, we introduce linear structure in rewards

2. Analysis of Ridge LR gives us bound on $|(\mu^\star - \hat{\mu}_t)^\top x|$

3. Optimism in the face of uncertainty: $\mu^\star \cdot x^\star \leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}}$

Summary

1. To deal w/ infinitely many arms, we introduce linear structure in rewards

2. Analysis of Ridge LR gives us bound on $|(\mu^\star - \hat{\mu}_t)^\top x|$

3. Optimism in the face of uncertainty: $\mu^\star \cdot x^\star \leq \hat{\mu}_t^\top x_t + \beta \|x_t\|_{\Sigma_t^{-1}}$

4. Regret is upper bounded by $\beta \sum_t \|x_t\|_{\Sigma_t} \leq \beta \sqrt{T} \sqrt{\sum_t \|x_t\|_{\Sigma_t^{-1}}^2}$