# Stochastic Bandits
# The UCB algorithm

Yevgeny Seldin

# Quick recap of the last lecture



*environmental resistance*

adversarial

i.i.d.

full    bandit    *feedback*

Started last time, finishing today

HA1: $\bar{R}_T$ not growing with $T$

$K$ actions

time

- Regret: $R_T = \sum_{t=1}^T \ell_{t,A_t} - \min_a \sum_{t=1}^T \ell_{t,a}$

- Expected regret: $\mathbb{E}[R_T] = \mathbb{E}\left[\sum_{t=1}^T \ell_{t,A_t}\right] - \mathbb{E}\left[\min_a \sum_{t=1}^T \ell_{t,a}\right]$

- Pseudo-regret: $\bar{R}_T = \mathbb{E}\left[\sum_{t=1}^T \ell_{t,A_t}\right] - \min_a \mathbb{E}\left[\sum_{t=1}^T \ell_{t,a}\right] = \mathbb{E}\left[\sum_{t=1}^T \ell_{t,A_t}\right] - T\mu^*$

$\qquad = \sum_{a=1}^K \Delta(a)\mathbb{E}[N_T(a)]$

# Lower Confidence Bound (LCB) algorithm for losses (Originally Upper Confidence Bound (UCB) for rewards) ("Optimism in the face of uncertainty" approach)

- Define $L_t^{CB}(a) = \hat{\mu}_{t-1}(a) - \sqrt{\dfrac{3\ln t}{2N_{t-1}(a)}}$ lower confidence bound
  - (We will show that with high probability $L_t^{CB}(a) \leq \mu(a)$ for all $t$)

- LCB Algorithm:
  - Play each arm once
  - For $t = K + 1, K + 2, \ldots$:
    - Play $A_t = \arg\min_a L_t^{CB}(a)$

- No knowledge of $T$
- No knowledge of $\Delta$
- Works for any $K$

Rewards $\leftrightarrow$ Losses

$\ell_{t,a} = 1 - r_{t,a}$

$r_{t,a} = 1 - \ell_{t,a}$

- Theorem:
$$\bar{R}_T \leq 6 \sum_{a:\Delta(a)>0} \frac{\ln T}{\Delta(a)} + \left(1 + \frac{\pi^2}{3}\right) \sum_a \Delta(a)$$

Lower Confidence Bound (LCB) algorithm for losses
(Originally Upper Confidence Bound (UCB) for rewards)
("Optimism in the face of uncertainty" approach)

- Define $L_t^{CB}(a) = \hat{\mu}_{t-1}(a) - \sqrt{\dfrac{3 \ln t}{2N_{t-1}(a)}}$ lower confidence bound
  - (We will show that with high probability $L_t^{CB}(a) \leq \mu(a)$ for all $t$)

- LCB Algorithm:
  - Play each arm once
  - For $t = K + 1, K + 2, \dots$:
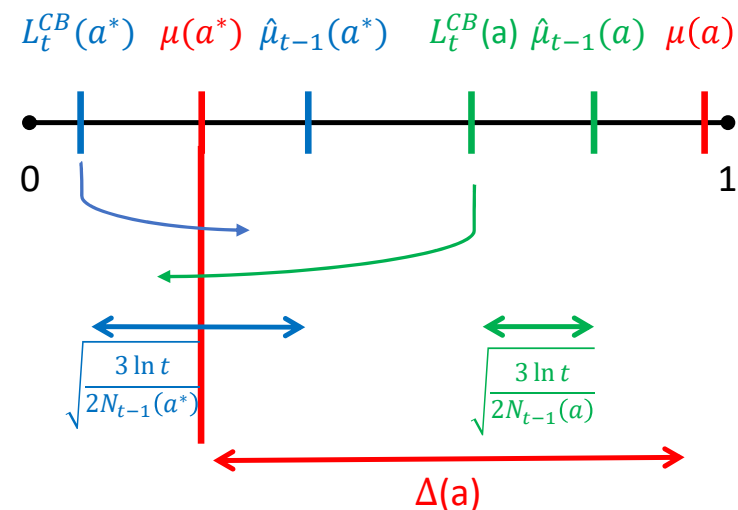    - Play $A_t = \arg\min_a L_t^{CB}(a)$

- Theorem:
$$\bar{R}_T \leq 6 \sum_{a:\Delta(a)>0} \frac{\ln T}{\Delta(a)} + \left(1 + \frac{\pi^2}{3}\right) \sum_a \Delta(a)$$

- Proof:

  - $\bar{R}_T = \sum_{a=1}^K \Delta(a)\mathbb{E}[N_T(a)]$

  - When can we play $a \neq a^*$?

  - Bound the expected number of times
    $$L_t^{CB}(a) \leq L_t^{CB}(a^*)$$

# Proof



$$\bar{R}_t(a) = \sum_a \Delta(a)\mathbb{E}[N_T(a)]$$

$$L_t^{CB}(a) = \hat{\mu}_{t-1}(a) - \sqrt{\frac{3\ln t}{2N_{t-1}(a)}}$$

- Bound the expected number of times $L_t^{CB}(a) \leq L_t^{CB}(a^*)$

- The expected number of times $L_t^{CB}(a) \leq L_t^{CB}(a^*)$ is bounded by
    1. The expected number of times $L_t^{CB}(a^*) \geq \mu(a^*)$
    2. Plus expected the number of times $L_t^{CB}(a) \leq \mu(a^*)$

# Proof continued

$L_t^{CB}(a^*)$ $\mu(a^*)$ $\hat{\mu}_{t-1}(a^*)$    $L_t^{CB}(a)$ $\hat{\mu}_{t-1}(a)$ $\mu(a)$

$$< \sqrt{\frac{3\ln t}{2N_{t-1}(a^*)}} \qquad < \sqrt{\frac{3\ln t}{2N_{t-1}(a)}}$$

$$\sqrt{\frac{3\ln t}{2N_{t-1}(a^*)}} \qquad \sqrt{\frac{3\ln t}{2N_{t-1}(a)}}$$

$\Delta(a)$

1. The expected number of times $L_t^{CB}(a^*) \geq \mu(a^*)$ is bounded by

   The expected number of times $\hat{\mu}_{t-1}(a^*) \geq \mu(a^*) + \sqrt{\frac{3\ln t}{2N_{t-1}(a^*)}}$

   ( ⟷ "breaks" – the confidence bound fails)

2. The expected the number of times $L_t^{CB}(a) \leq \mu(a^*)$ is bounded by

   2.1 The expected number of times $\hat{\mu}_{t-1}(a) \leq \mu(a) - \sqrt{\frac{3\ln t}{2N_{t-1}(a)}}$
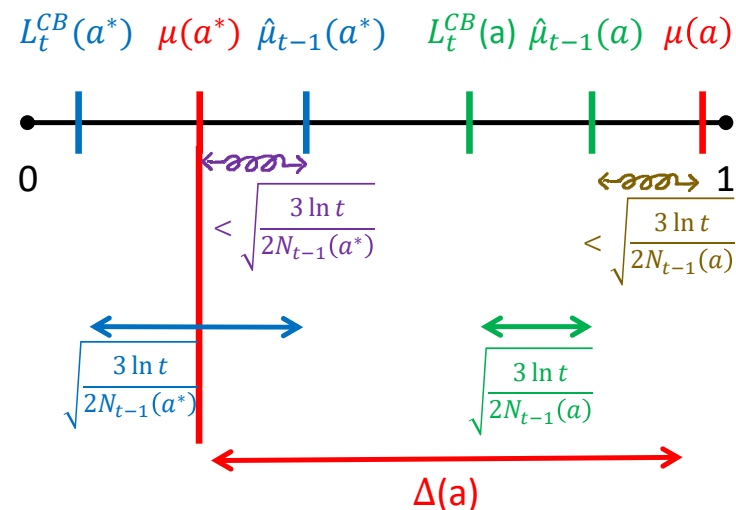
   ( ⟷ "breaks" – the confidence bound fails)

   2.2 Plus the number of times $\sqrt{\frac{3\ln t}{2N_{t-1}(a)}} + \sqrt{\frac{3\ln t}{2N_{t-1}(a)}} = \sqrt{\frac{6\ln t}{N_{t-1}(a)}} \geq \Delta(a)$

   (The confidence bound holds, but it is not tight enough)

   $\Rightarrow \quad N_{t-1}(a) \leq \frac{6\ln t}{\Delta(a)^2} \leq \frac{6\ln T}{\Delta(a)^2}$

- Mid-summary: $\mathbb{E}[N_T(a)] \leq \frac{6\ln T}{\Delta(a)^2} + 1 + \mathbb{E}[1.] + \mathbb{E}[2.1]$

# Proof continued

$L_t^{CB}(a^*)$   $\mu(a^*)$   $\hat{\mu}_{t-1}(a^*)$    $L_t^{CB}(a)$   $\hat{\mu}_{t-1}(a)$   $\mu(a)$

0       1

$< \sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}}$      $< \sqrt{\dfrac{3\ln t}{2N_{t-1}(a)}}$

$\sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}}$     $\sqrt{\dfrac{3\ln t}{2N_{t-1}(a)}}$

$\Delta(a)$

- $\mathbb{E}[N_T(a)] \leq \dfrac{6\ln T}{\Delta(a)^2} + 1 + \mathbb{E}[\;\;] + \mathbb{E}[\;\;]$

- Let $F(a^*)$ be the expected number of times $\hat{\mu}_{t-1}(a^*) \geq \mu(a^*) + \sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}}$

- Bound $\mathbb{P}\left( \hat{\mu}_{t-1}(a^*) - \mu(a^*) \geq \sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}} \right)$   $N_{t-1}(a^*)$ is a random variable dependent on $\hat{\mu}_t(a^*)$!
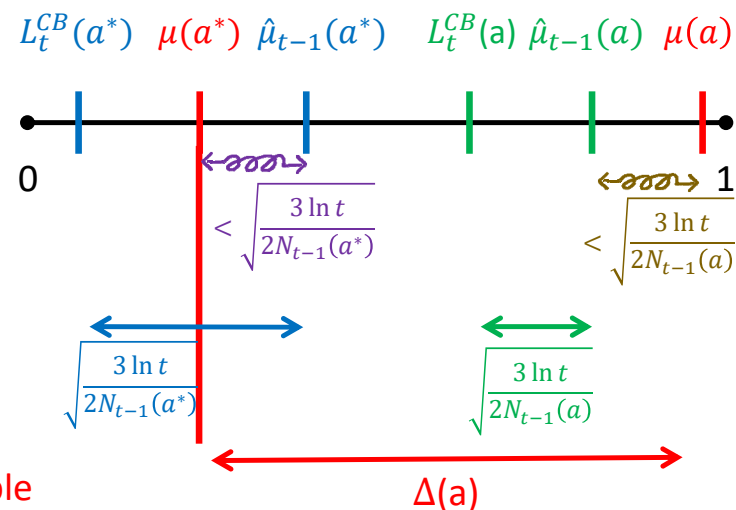
- Idea: break dependent events into independent events and take a union bound
- Introduce $X_1, \dots, X_T$ r.v. with the same distribution as $\ell_{t,a^*}$
- Let $\bar{\mu}_s = \dfrac{1}{s}\sum_{i=1}^{s} X_i$
- $\mathbb{P}\left( \hat{\mu}_{t-1}(a^*) - \mu(a^*) \geq \sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}} \right) \leq \mathbb{P}\left( \exists s \in \{1, \dots, t\}: \bar{\mu}_s - \mu(a^*) \geq \sqrt{\dfrac{\ln t^3}{2s}} \right)$
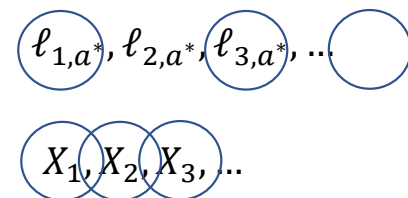
$$\underset{\text{union}}{\leq} \sum_{s=1}^{t} \mathbb{P}\left( \bar{\mu}_s - \mu(a^*) \geq \sqrt{\dfrac{\ln t^3}{2s}} \right)$$

$$\underset{\text{Hoeffding}}{\leq} \sum_{s=1}^{t} \dfrac{1}{t^3} = \dfrac{1}{t^2}$$

- $\mathbb{E}[F(a^*)] = \sum_{t=1}^{\infty} \mathbb{P}\left( L_t^{CB}(a^*) \geq \mu(a^*) \right) \leq \sum_{t=1}^{\infty} \dfrac{1}{t^2} \leq \dfrac{\pi^2}{6}$

Hoeffding:

$$\mathbb{P}\left( \dfrac{1}{n}\sum_{i=1}^{n} Z_i - \mu \geq \sqrt{\dfrac{\ln \frac{1}{\delta}}{2n}} \right) \leq \delta$$

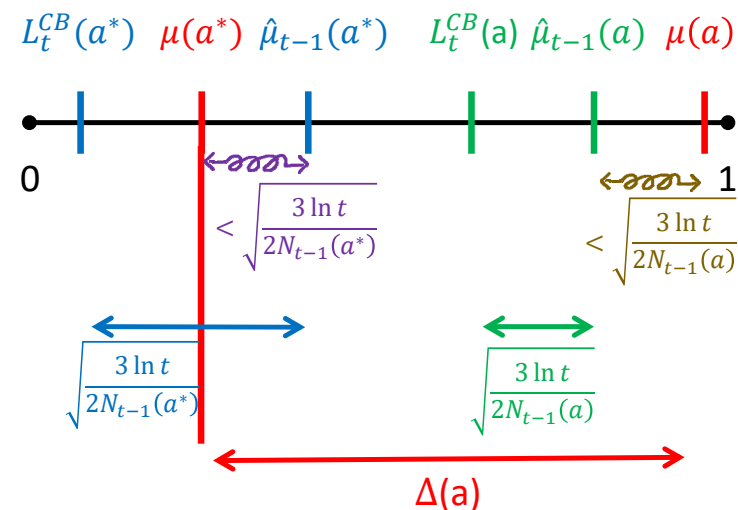$\ell_{1,a^*}, \ell_{2,a^*}, \ell_{3,a^*}, \dots$

$X_1, X_2, X_3, \dots$

# Proof summary

Top right diagram labels:

$L_t^{CB}(a^*)$  $\mu(a^*)$  $\hat{\mu}_{t-1}(a^*)$  $L_t^{CB}(\text{a})$  $\hat{\mu}_{t-1}(a)$  $\mu(a)$

$0$ ........................ $1$

$< \sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}}$

$< \sqrt{\dfrac{3\ln t}{2N_{t-1}(a)}}$

$\sqrt{\dfrac{3\ln t}{2N_{t-1}(a^*)}}$

$\sqrt{\dfrac{3\ln t}{2N_{t-1}(a)}}$

$\Delta(a)$

- $\bar{R}_t(a) = \sum_a \Delta(a)\mathbb{E}[N_T(a)]$

- $\mathbb{E}[N_T(a)] \leq \underbrace{\dfrac{6\ln T}{\Delta(a)^2}}_{\substack{\text{The time it takes} \\ \text{for confidence} \\ \text{intervals to} \\ \text{start working}}} + 1 + \underbrace{\dfrac{\pi^2}{6} + \dfrac{\pi^2}{6}}_{\substack{\text{The expected number} \\ \text{of times confidence} \\ \text{intervals fail}}}$
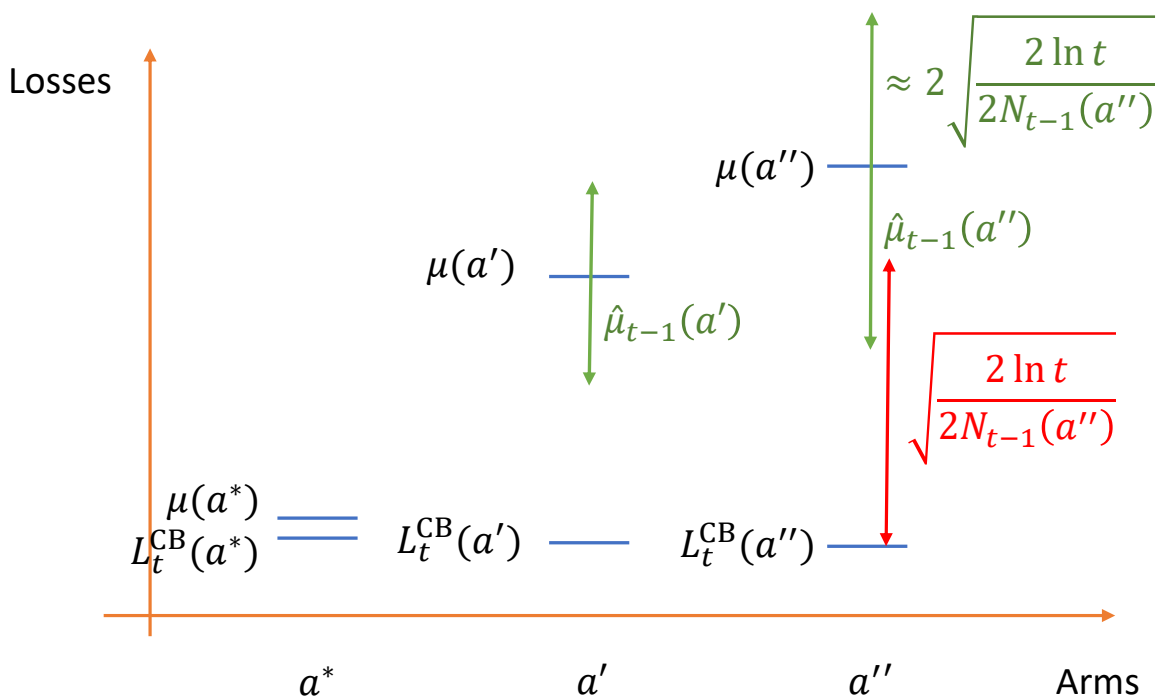
- $\bar{R}_T \leq 6\sum_{a:\Delta(a)>0} \dfrac{\ln T}{\Delta(a)} + \left(1 + \dfrac{\pi^2}{3}\right)\sum_a \Delta(a)$

- Home assignment:
  - Take $L_t^{CB}(a) = \hat{\mu}_{t-1}(a) - \sqrt{\dfrac{2\ln t}{2N_{t-1}(a)}}$ (instead of $L_t^{CB}(a) = \hat{\mu}_{t-1}(a) - \sqrt{\dfrac{3\ln t}{2N_{t-1}(a)}}$; i.e. confidence $\dfrac{1}{t^2}$ instead $\dfrac{1}{t^3}$)
  - Show $\bar{R}_T \leq 4\sum_{a:\Delta(a)>0} \dfrac{\ln T}{\Delta(a)} + (2\ln T + 3)\sum_a \Delta(a)$

# LCB algorithm dynamics (with $L_t^{CB}(a) = \hat{\mu}_{t-1}(a) - \sqrt{\frac{2\ln t}{2N_{t-1}(a)}}$)



- Confidence interval of the played arm shrinks ($N_{t-1}(a)$ grows)

- Confidence intervals of all other arms grow ($\ln t$ grows)

- $\Rightarrow$ all LCBs are roughly at the same level

- Most of the time $L_t^{CB}(a^*) \leq \mu(a^*)$

- $a^*$ is played a lot, so $L_t^{CB}(a^*)$ is very close to $\mu(a^*)$

- All other arms are played just enough to keep $\sqrt{\frac{2\ln t}{2N_{t-1}(a)}} = \theta(\Delta(a))$, i.e. $N_t(a) = \theta\left(\frac{\ln t}{\Delta(a)^2}\right)$