

Correlated Bandits



GROUP 11

ISHANK JUNEJA 16D070012

DHRUV VARSHNEY 16D070033

DEVANSHU SINGH GAHARWAR 16D070042

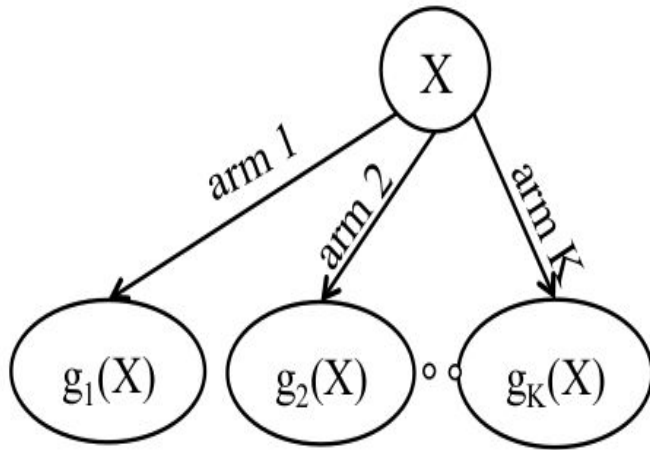




Introduction : What are correlated bandits ?

- Standard MAB Setup
- Independence assumption between arms relaxed
- Correlation between arms can be exploited if present
 - Skip pulling some arms based on correlation

Problem Formulation



- X is the hidden random variable
- $g_1(X), g_2(X), \dots, g_k(X)$ are the dependent reward functions
- g_1, g_2, \dots, g_k are known functions.
 - Assumption Valid ???



Motivating Example

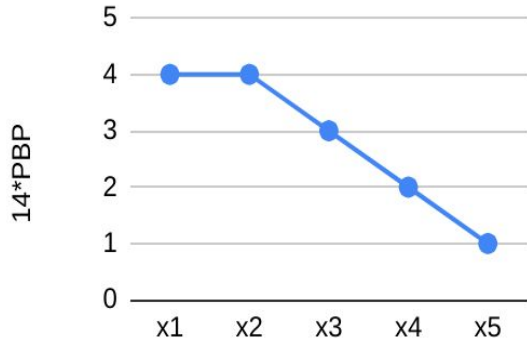
- Consider the case of Amazon expanding to a new country
- k arms $\equiv k$ mobile companies
 - $g_1, g_2 \dots g_k =$ product buying probability (PBP)
- Random variable \equiv Discrete Income levels
- $g_1, g_2 \dots g_k \rightarrow$ found using paid surveys.

Fact : Amazon to start operating in Bangladesh in 2020. Refer [this](#)



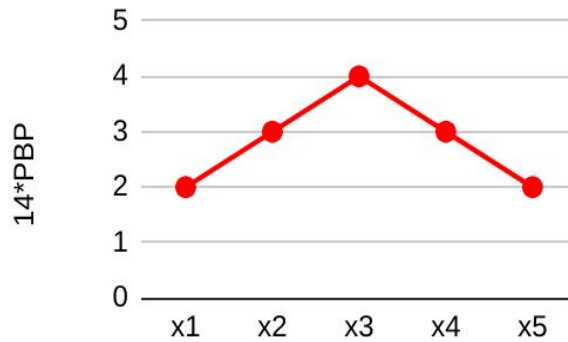
Motivating Example (Continued)

Mi



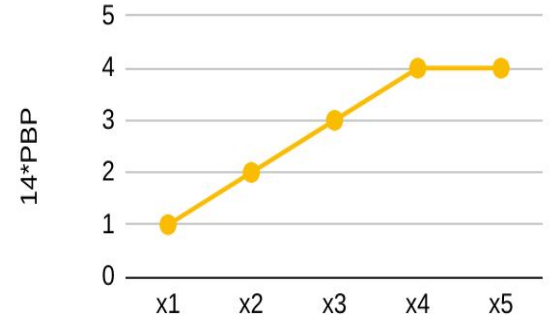
Income level -->

Samsung



Income level -->

One Plus



Income level -->



Approach 1.0

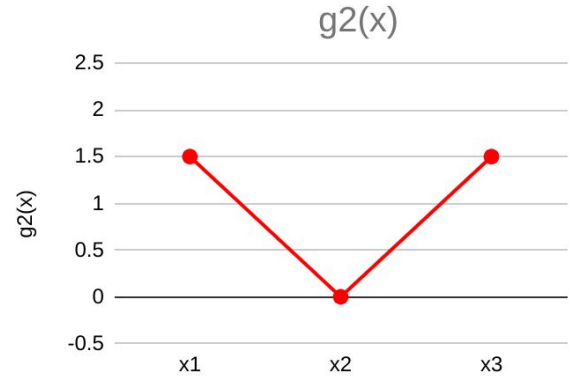
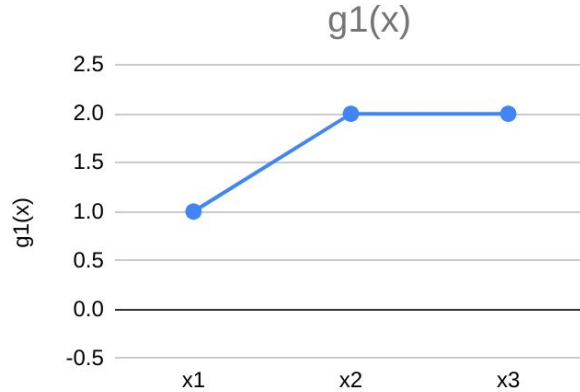
- C-UCB algorithm proposed
- Need to classify arms as Competitive and Non-Competitive
- How to decide whether an arm is competitive ?
 - Pseudo Reward of arm ℓ wrt k : $s_{\ell,k}(r) \triangleq \max_{x:g_k(x)=r} g_\ell(x)$
 - Expected Pseudo Gap of arm ℓ wrt k : $\Delta_{\ell,k} \triangleq \mu_k - E[s_{\ell,k}(g_k(X))]$
 - Arm ℓ is **non-competitive** wrt k if pseudo gap is positive.



Approach 1.0 (Continued)

- Compute the expected quantities empirically (Law of large numbers)
 - Empirical reward of arm k : $\mu_k = \frac{\sum_{\tau} 1_{k(\tau)=k} g_{k(\tau)}(X(\tau))}{n_k(t)}$
 - Empirical Pseudo Reward : $\Phi_{l,k}(t) \triangleq \frac{\sum_{\tau} 1_{k(\tau)=k} s_{l,k}(r_{\tau})}{n_k(t)}$
 - $\mu_k > \Phi_{l,k}(t) \Rightarrow$ Arm l is non-competitive wrt arm k

Example :



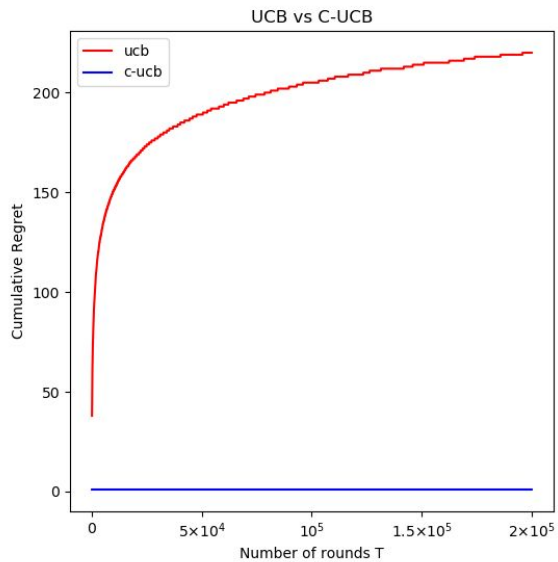
- Arm 1 pulled 10 times out of which reward 1 (3 times) and 2 (7 times)
 - $\mu_1 = (1 \cdot 3 + 2 \cdot 7) / 10 = 1.7$
 - $s_{21}(1) = s_{21}(2) = 1.5 \Rightarrow \Phi_{21}(t) = 1.5$
 - Clearly $\Delta_{21} = \mu_1 - \Phi_{21}(t) > 0 \Rightarrow$ Arm 2 is non - competitive



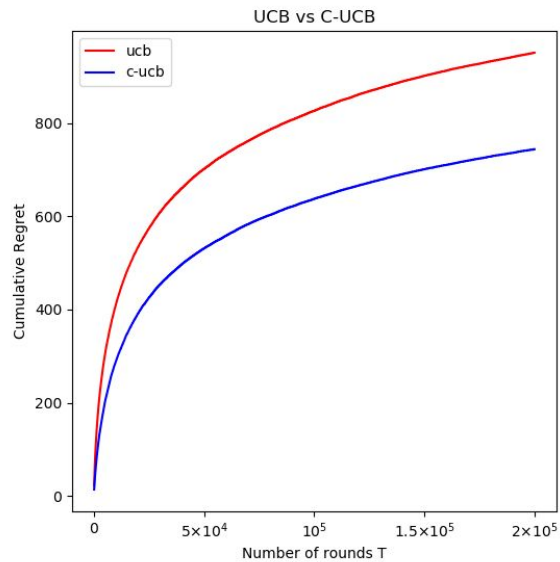
C-UCB ALGORITHM

- Initialise using standard UCB method ($n_k = 0, I_k = \infty$)
- For every iteration t do :
 1. Choose reference arm.
 2. Find the empirically competitive set wrt reference arm.
 3. Apply UCB over the set of competitive arms to get optimal arm k_t
 4. For all arms $k \neq k_t$ update the empirical pseudo rewards.
 5. Update the standard UCB parameters (n_k, I_k)
 6. Update the empirical reward for arm k_t

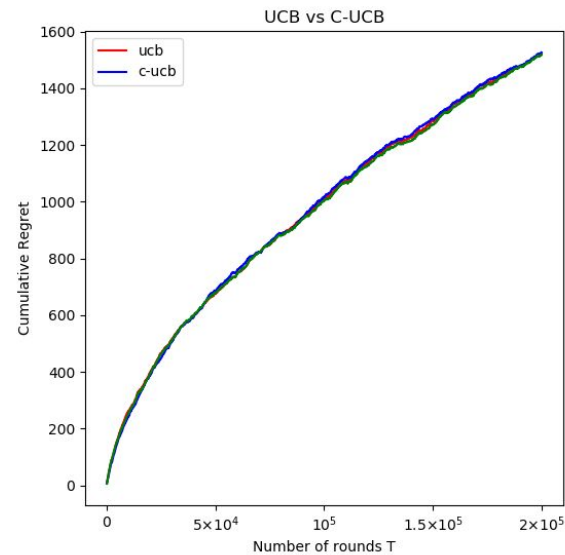
Simulation Results



No Competitive arm



1 Competitive arm

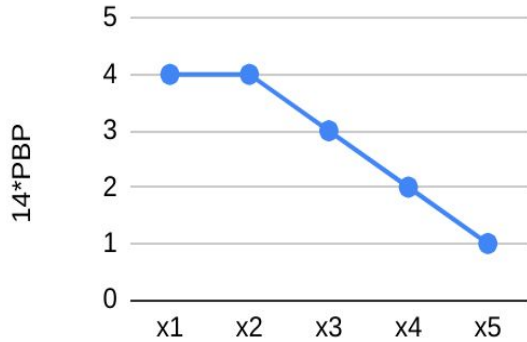


All arms are competitive



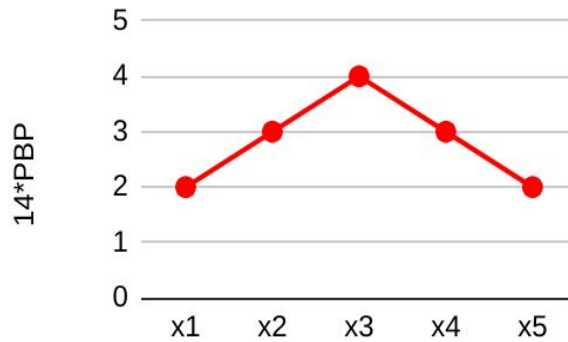
Reward Functions for Simulation Results

Mi



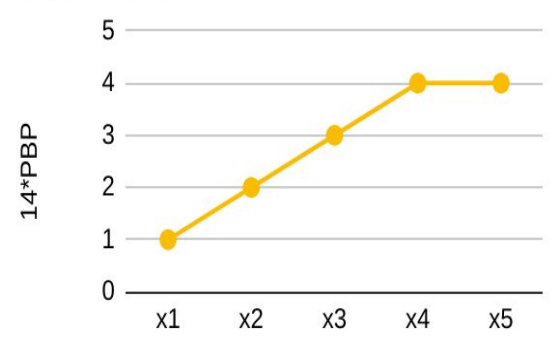
Income level -->

Samsung



Income level -->

One Plus



Income level -->



Our Contribution

A New Algorithm : Uni-C-UCB



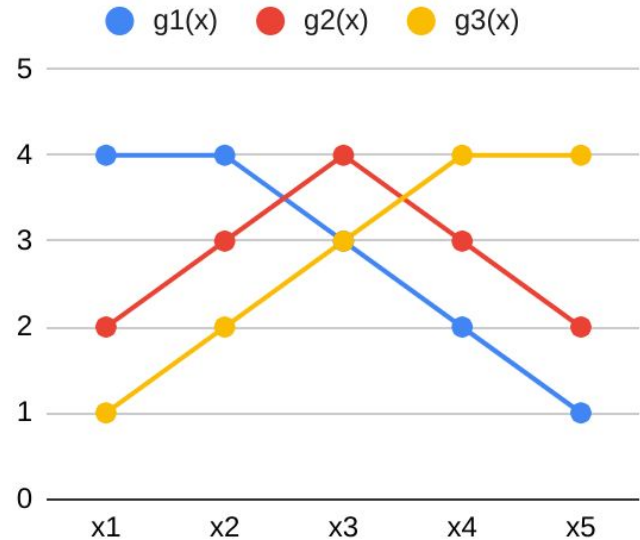
Approach 2.0

- Underlying distribution of X which can be learnt is used.
- How to decide whether an arm is competitive ?
 - **Pseudo distribution** - Empirically estimate the pmf(s)
 - **Confidence Set (C^*)** :
 - First sort the pmf in decreasing order of probability.
 - $C^* = \{1, 2 \dots j\}$ where j is the minimum k s.t. $\sum_{i=1}^k p(x_i) > 1 - \epsilon$
 - Arm k is **Non-Competitive** if $g_k(x) < g_j(x) \forall x \in C^*$ and some arm j .

Example :

Consider the following 3 cases :

- If $\{x_1, x_2\} \in C^* \Rightarrow$
Arms 1 is competitive
- If $\{x_2, x_3, x_4\} \in C^* \Rightarrow$
All arms are competitive
- If $\{x_1, x_2, x_3\} \in C^* \Rightarrow$
Arms 1,2 are competitive





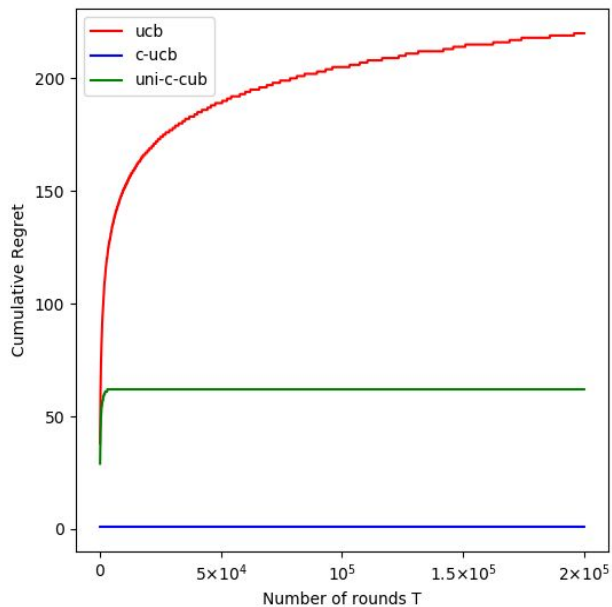
Uniform C-UCB ALGORITHM

- Initialise using standard UCB method ($n_k = 0, I_k = \infty$)
- Initialise $C^* = \text{support of r.v. } X$ and $\epsilon = 0.1$ (tuneable)
- For every iteration t do :
 1. Find the competitive set using C^*
 2. Apply UCB over the set of competitive arms to get optimal arm k_t
 3. Update the pseudo-distribution using Bayesian updates
 4. Update the Confidence Set (C^*)
 5. Update the standard UCB parameters (n_k, I_k)

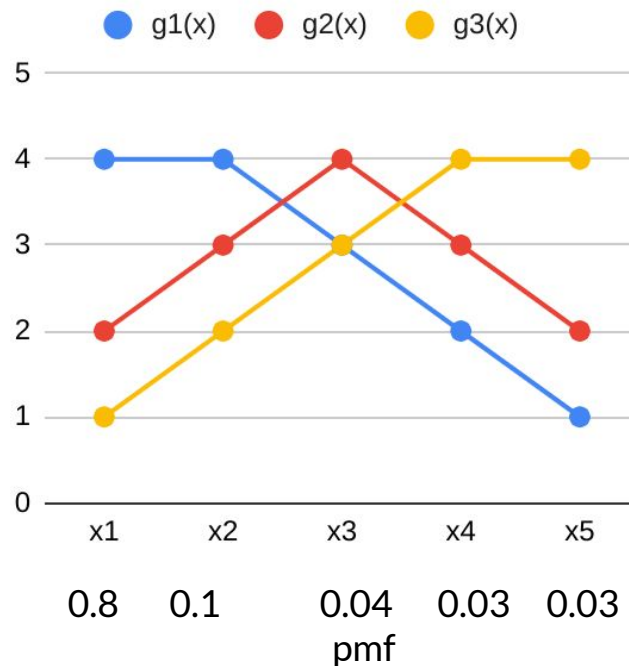
Simulation Results

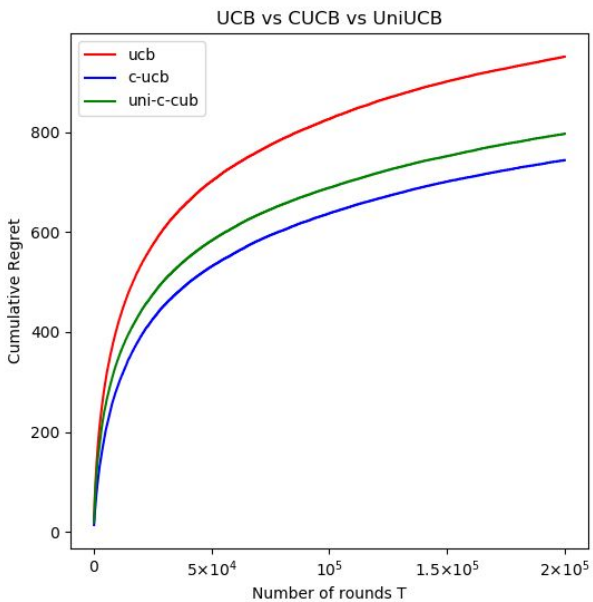


UCB vs CUCB vs UniUCB

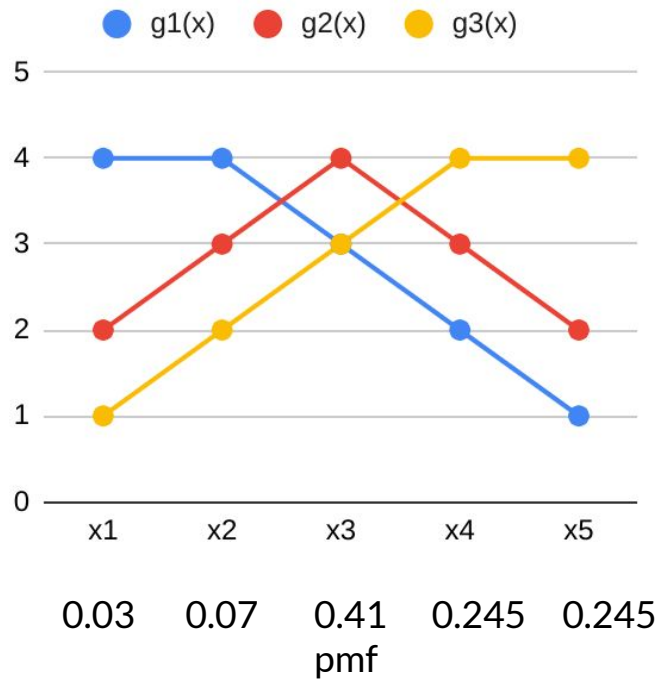


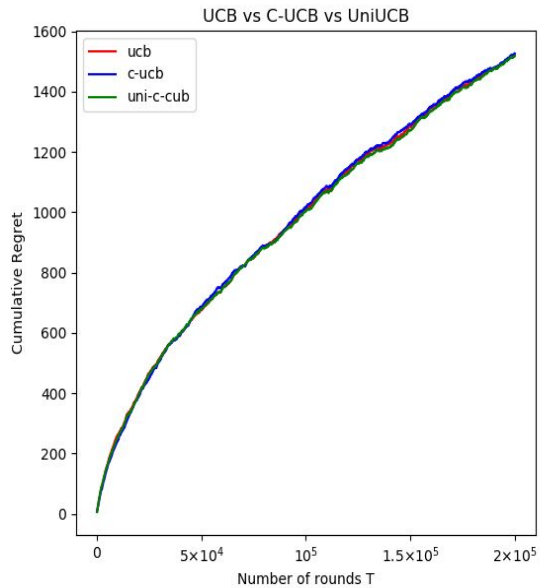
No Competitive arm



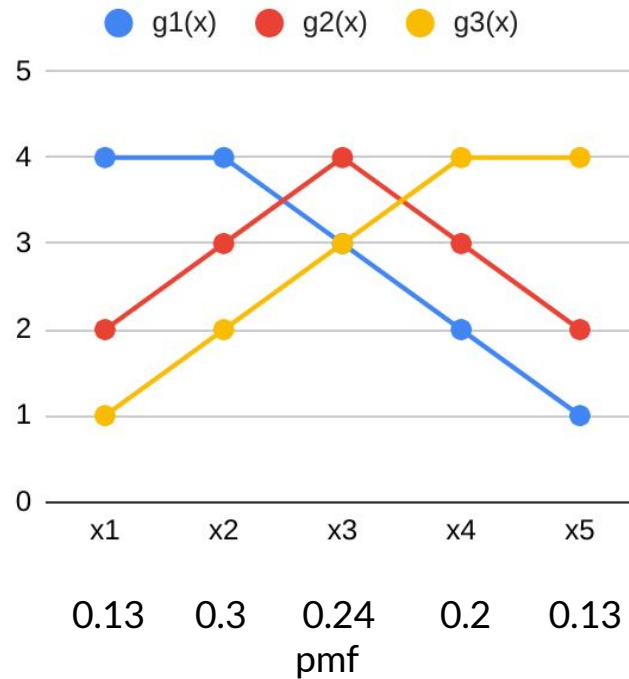


1 Competitive arm



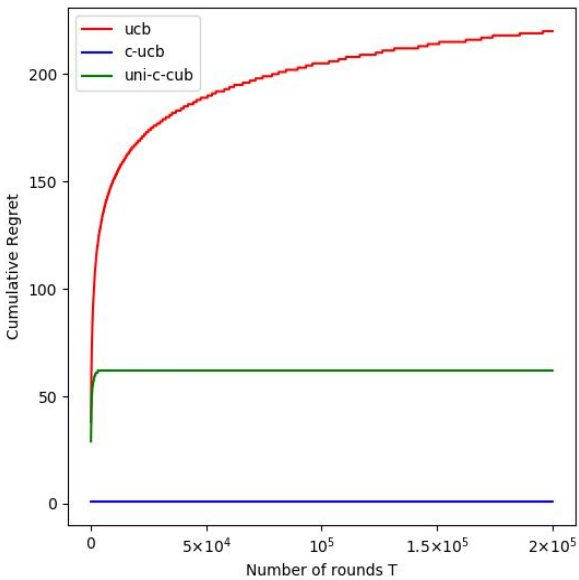


All arms are competitive



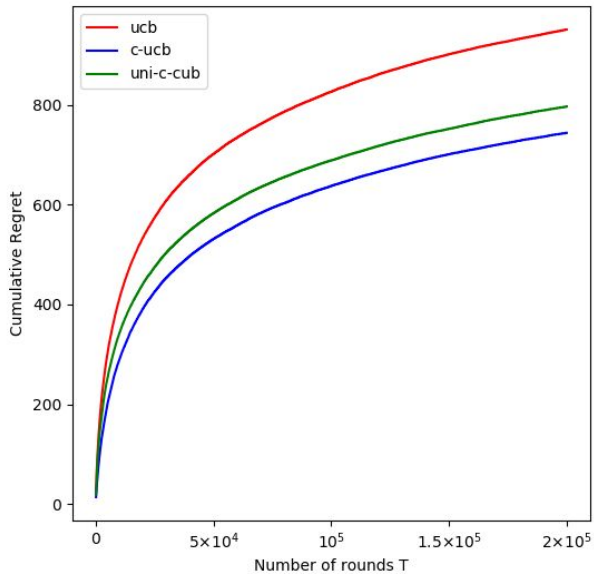


UCB vs CUCB vs UniUCB



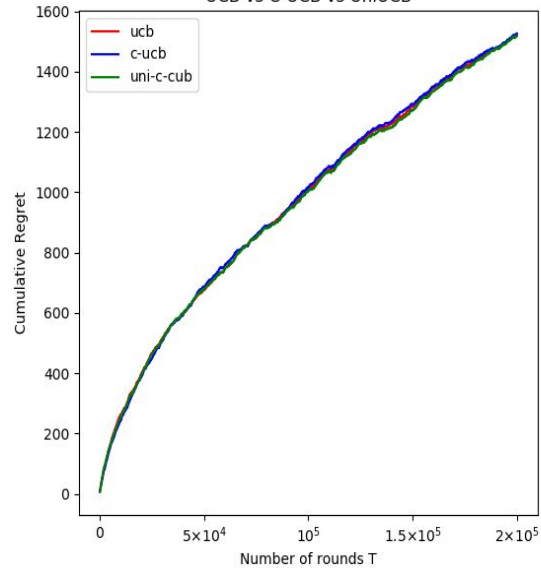
No Competitive arm

UCB vs CUCB vs UniUCB



1 Arm Competitive

UCB vs C-UCB vs UniUCB

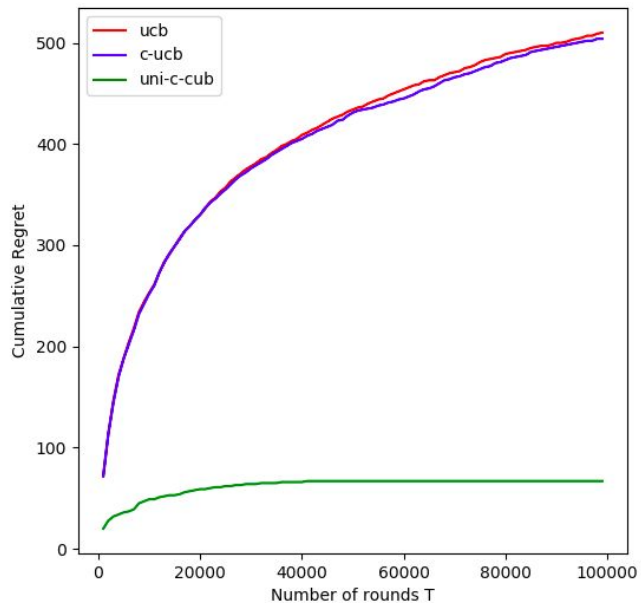


All Arms Competitive

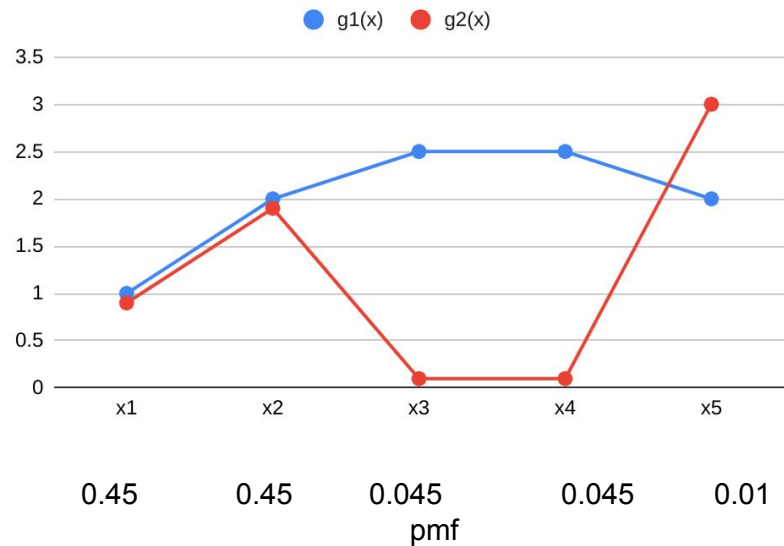


Can Uni-C-UCB outperform C-UCB ?

Yes, Indeed !



Uniform C-UCB works better than C-UCB





Concluding Remarks

- Two approaches presented
- Both perform better than UCB by exploiting correlation
- Which approach is better ?
 - Depends on the exact nature of functions
- In the report, we would include our attempt on the regret analysis.



Thank You

Questions ?

Algorithm 1 C-UCB Correlated UCB Algorithm

- 1: **Input:** Reward Functions $\{g_1, g_2 \dots g_K\}$
- 2: **Initialize:** $n_k = 0, I_k = \infty$ for all $k \in \{1, 2, \dots, K\}$
- 3: **for** each round t **do**
- 4: Find $k^{\max} = \arg \max_k n_k(t-1)$, the arm that has been pulled most times until round $t-1$
- 5: Initialize the empirically competitive set $\mathcal{A} = \{1, 2, \dots, K\} \setminus \{k^{\max}\}$.
- 6: **for** $k \neq k^{\max}$ **do**
- 7: **if** $\hat{\mu}_{k^{\max}} > \hat{\phi}_{k, k^{\max}}$ **then**
- 8: Remove arm k from the empirically competitive set: $\mathcal{A} = \mathcal{A} \setminus \{k\}$
- 9: **end if**
- 10: **end for**
- 11: Apply UCB1 over arms in $\mathcal{A} \cup \{k^{\max}\}$ by pulling arm $k_t = \arg \max_{k \in \mathcal{A} \cup \{k^{\max}\}} I_k(t-1)$
- 12: Receive reward r_t , and update $n_{k_t} = n_{k_t} + 1$
- 13: Update Empirical reward: $\hat{\mu}_{k_t}(t) = \frac{\hat{\mu}_{k_t}(t-1)(n_{k_t}(t)-1) + r_t}{n_{k_t}(t)}$
- 14: Update the UCB Index: $I_{k_t}(t) = \hat{\mu}_{k_t} + B \sqrt{\frac{2 \log t}{n_{k_t}}}$
- 15: Compute pseudo-rewards for all arms $k \neq k_t$: $s_{k, k_t}(r_t) = \max_{x: g_{k_t}(x) = r_t} g_k(x)$.
- 16: Update empirical pseudo-rewards for all $k \neq k_t$: $\hat{\phi}_{k, k_t}(t) = \sum_{\tau: k_\tau = k_t} s_{k, k_\tau}(r_\tau) / n_{k_t}$
- 17: **end for**