

Reinforcement Learning and Dynamic Optimization (ΠΛΗ423/ΠΛΗ723)

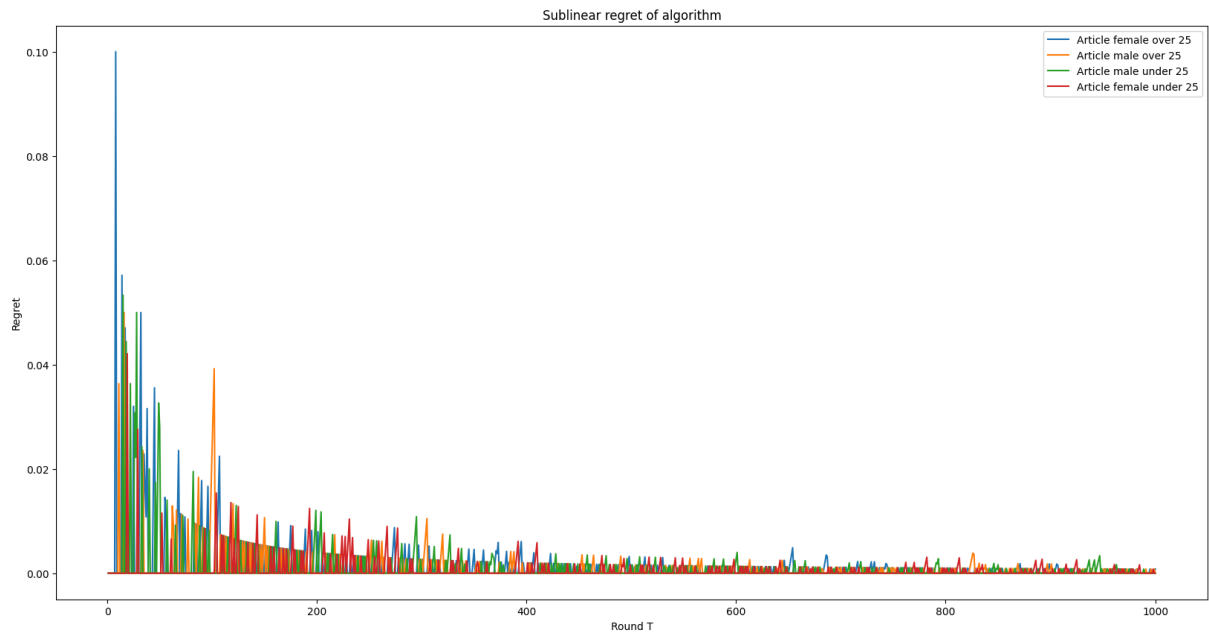
Assignment 1

Recommending News Articles to Unknown Users



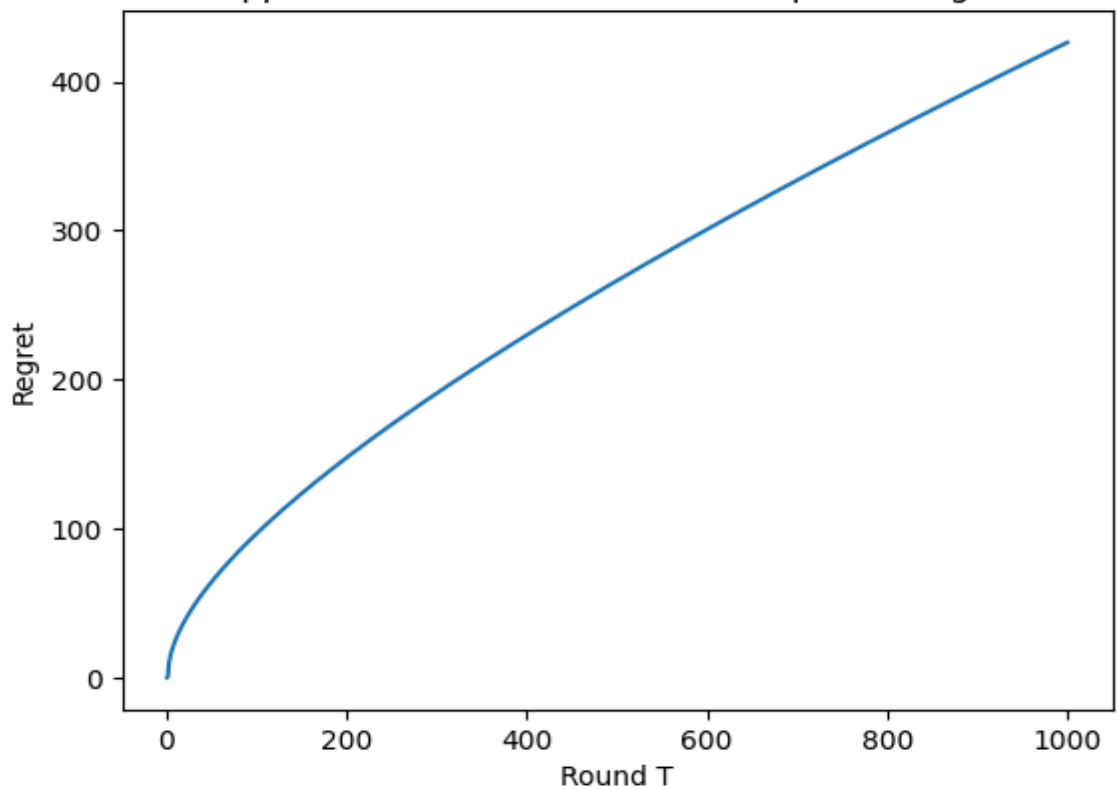
**ΠΟΛΥΤΕΧΝΕΙΟ
ΚΡΗΤΗΣ**

Νικόλαος Αγγελίδης 2019030190

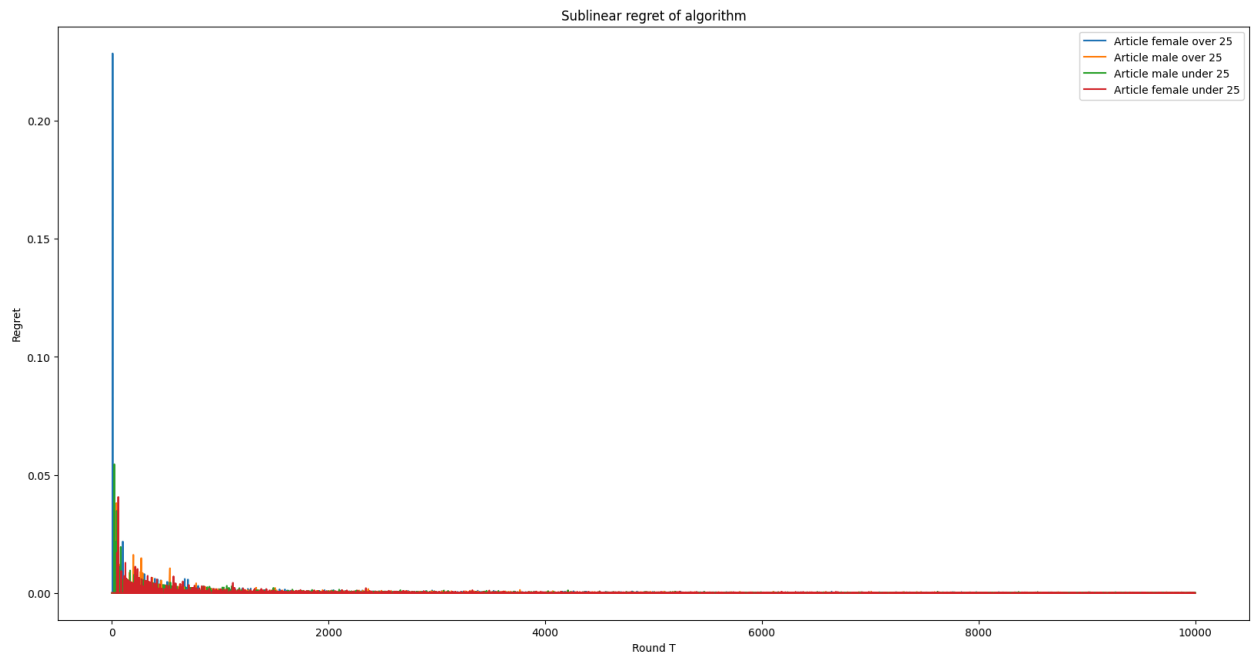


As shown in the measurement plot the regret goes to zero

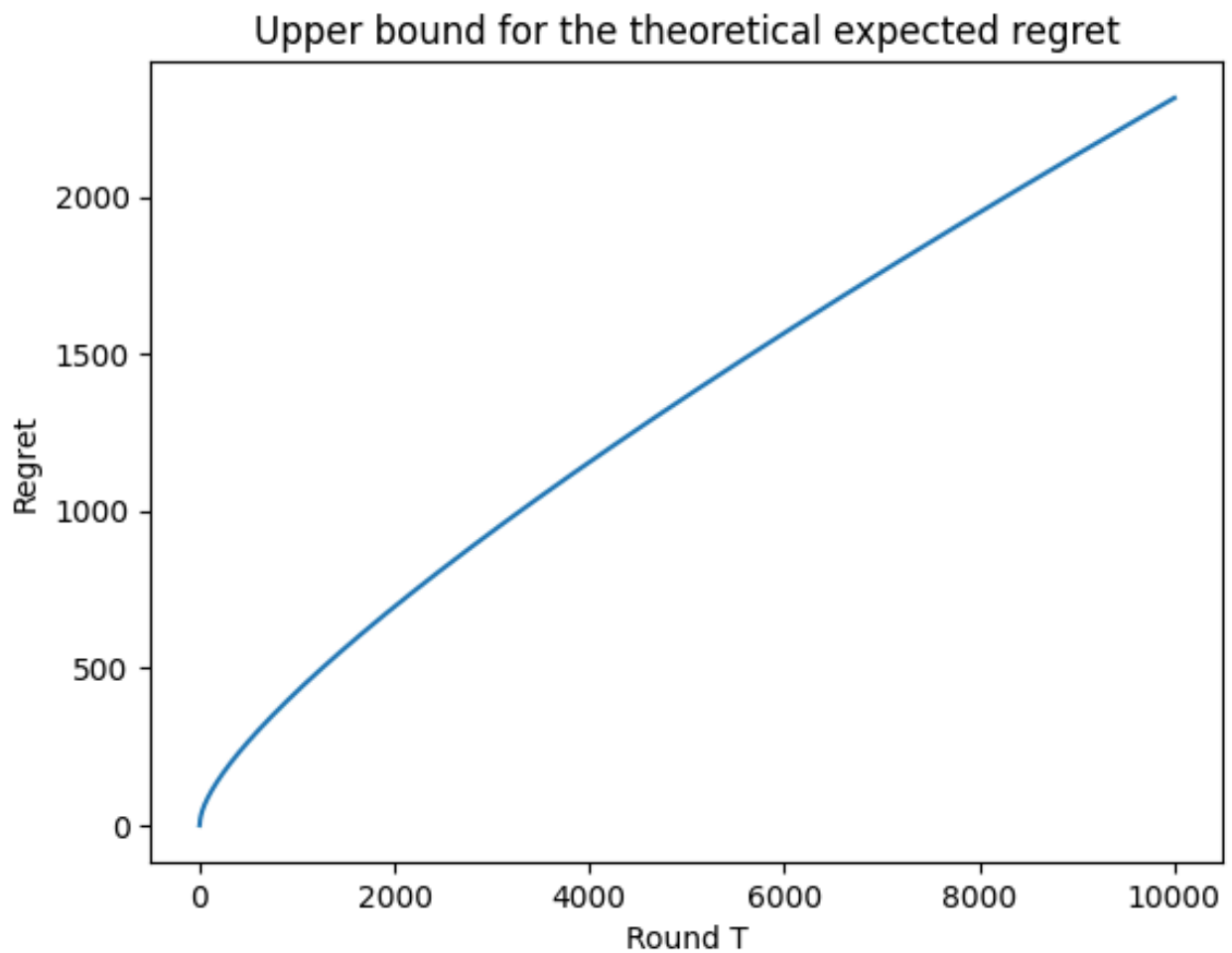
Upper bound for the theoretical expected regret



Upper bound for the theoretical expected regret of the algorithm ($T = 1000$)



Same as before the regret inclines to zero, much faster than before.



Upper bound for the theoretical expected regret of the algorithm ($T = 10000$)

$$\begin{aligned}
E[R(T)] &\triangleq E[\text{Regret}(T)] \\
E[R(T)] &= E[R(T)|\text{explore}] + E[R(T)|\text{exploit}] \\
\cdot E[R(T)|\text{explore}] &= \sum_{t=1}^{NKU} \mu^* - \mu_{i_t} = \sum_{i=1}^K N \cdot U \cdot \Delta_i \leq N \cdot U \cdot K (\Delta_i \leq 1) \\
\cdot E[R(T)|\text{exploit}] &= (T - NKU) \cdot E[\Delta_{i_{\text{ALG}}}] \\
E[R(T)|\text{exploit}] &\leq T \cdot E[\Delta_{i_{\text{ALG}}}] \text{ (since } T - NKU \leq T \text{)} \\
&\leq T \cdot (1 - P(\text{Bad})) \Delta_{i_{\text{ALG}}|\text{Good}} + T \cdot P(\text{Bad})
\end{aligned}$$

$$\begin{aligned}
E[R(T)] &\leq N \cdot K \cdot U + T \cdot E[\Delta_{i_{\text{ALG}}}] \\
* \epsilon &= \sqrt{2 \log(T) / NU} \\
* P(\text{Bad}) &= P(\exists i : |\hat{\mu}_i - \mu_i| > \sqrt{2 \log(T) / NU}) \leq KT^{-4} \\
* P(\text{Good}) &= P(\forall i : |\hat{\mu}_i - \mu_i| \leq \sqrt{2 \log(T) / NU}) \geq 1 - KT^{-4} \\
* \Delta_{i_{\text{ALG}}} &= \mu^* - \mu_i \leq \mu^* - \hat{\mu}_i + \epsilon \text{ (since } \mu_i \geq \hat{\mu}_i - \epsilon \text{)} \\
&\leq \hat{\mu}^* + \epsilon - \hat{\mu}_i + \epsilon \text{ (since } \mu^* \leq \hat{\mu}_i + \epsilon \text{)} \Rightarrow \Delta_{i_{\text{ALG}}|\text{Good}} \leq 2\epsilon
\end{aligned}$$

$$\begin{aligned}
E[R(T)] &\leq N \cdot K \cdot U + T \cdot (1 - P(\text{Bad})) \cdot \Delta_{i_{\text{ALG}}|\text{Good}} + T \cdot P(\text{Bad}) \\
E[R(T)] &\leq N \cdot K \cdot U + T \cdot (1 - KT^{-4}) \cdot \Delta_{i_{\text{ALG}}|\text{Good}} + T \cdot KT^{-4} \\
E[R(T)] &\leq N \cdot K \cdot U + T \cdot \Delta_{i_{\text{ALG}}|\text{Good}} \text{ (since } KT^{-3} \rightarrow 0, \text{ as } T \text{ grows)} \\
E[R(T)] &\leq N \cdot K \cdot U + T \cdot 2 \sqrt{2 \log(T) / NU}
\end{aligned}$$

$$\text{We want : } N \cdot K \cdot U \approx T \cdot 2 \sqrt{2 \log(T) / NU}$$

$$N = O\left(\left(\frac{TK^2 \log(T)}{U^2}\right)^{1/3}\right) \quad (1)$$

Insert (1) into the regret formula :

$$E[R(T)] \leq O\left(\left(\frac{TK^2 \log(T)}{U^2}\right)^{1/3}\right) \cdot K \cdot U + T \cdot 2 \sqrt{\frac{2 \log(T)}{O\left(\left(\frac{TK^2 \log(T)}{U^2}\right)^{1/3}\right) U}}$$

* First term :

$$\begin{aligned}
O\left(\left(\frac{TK^2 \log(T)}{U^2}\right)^{1/3}\right) \cdot K \cdot U &= O((T^{1/3} K^{2/3} \log(T)^{1/3}) \cdot K / U^{1/3}) = \\
O\left(\frac{T^{1/3} K^{5/3} \log(T)^{1/3}}{U^{1/3}}\right)
\end{aligned}$$

* Second term :

$$T \cdot 2 \sqrt{\frac{2 \log(T)}{O\left(\left(\frac{TK^2 \log(T)}{U^2}\right)^{1/3}\right) U}} = O\left(T \cdot \sqrt{\frac{\log(T)}{T^{1/3} K^{2/3} \log(T)^{1/3} / U^{2/3}}}\right) = O\left(T \cdot \sqrt{\frac{U^{2/3}}{T^{1/3} K^{2/3}}}\right)$$

So the theoretical expected regret of the algorithm :

$$E[R(T)] \leq O\left(\frac{T^{1/3} K^{5/3} \log(T)^{1/3}}{U^{1/3}}\right) + O\left(T \cdot \sqrt{\frac{U^{2/3}}{T^{1/3} K^{2/3}}}\right)$$