

Regret Bounds for (α, γ) -UCB

Finnian Lattimore

April 1, 2015

The basic setting is the stochastic bandit problem:

- We have a set of arms $i \in \{1 \dots K\}$
- For each arm i , there is an unknown distribution of rewards $P_i(X)$
- Each time we select a given arm, i , the reward is sampled i.i.d from $P_i(X)$. This is a big assumption - it states that the reward at a given timestep depends only on the action selected at that timestep, not on the sequence of previous actions.
- Can the reward distributions be (fixed) functions of the timestep? (No I think - otherwise the idea of an optimal arm independent of t doesn't make sense)

Some notation:

- K the number of arms
- i identifies an arm, $i \in \{1 \dots K\}$
- i^* the arm with the highest true expected reward
- I_t the arm selected by the algorithm at timestep t
- $\hat{\mu}_i$ an estimator for the expected reward of arm i based on the sample mean
- μ_i the true expected reward for an arm i
- μ^* the true expected reward of i^* (the best arm)
- $\Delta_i = \mu^* - \mu_i$ how much worse arm i is than the best arm
- $T_i(s) = \sum_{t=1}^s \mathbb{1}\{I_t = i\}$ the number of times arm i was selected upto timestep s
- $\hat{\epsilon}_{it}$ an estimate of the uncertainty in the empirical estimator for the expected reward

The goal is to get a bound on the pseudo-regret, defined as:

$$R_n = n\mu^* - E \left[\sum_{t=1}^n \mu_{I_t} \right] \quad (1)$$

$$= \sum_{i=1}^K \Delta_i E[T_i(n)] \quad (2)$$

UCB-1

For simplicity just consider Bernoulli bandits: $P(X_i) \sim \text{Bernoulli}(p_i)$. Hoeffding's inequality gives us a high probability bound on the how much our sample based estimate of the expected reward can be below the true expected reward. Let $\hat{\mu}_{is} = \frac{1}{s} \sum_{t=1}^s X_t$ be the sample average and $\mu_i = E[P_i(X)]$ if we select arm i a fixed number of times s :

$$P(\mu_i - \hat{\mu}_{is} > \epsilon) \leq e^{-2s\epsilon^2} \quad (3)$$

$$\implies P\left(\mu_i - \hat{\mu}_{is} > \sqrt{\frac{\log(1/\delta)}{2s}}\right) \leq \delta \quad (4)$$

The UCB Algorithm

Define the upper confidence bound for each arm i at timestep t as:

$$ucb_{it} = \hat{\mu}_{it} + \sqrt{\frac{2\log(t)}{T_i(t)}} \quad (5)$$

$$= \hat{\mu}_{it} + \hat{\epsilon}_{it} \quad (6)$$

At time t select arm I_t with the highest upper confidence bound:

$$I_t = \operatorname{argmax}_{i=1\dots K} (ucb_{it}) \quad (7)$$

$$R_n \leq \sum_{i:\Delta_i > 0} \left(\frac{8\log(n)}{\Delta_i} + (1 + \frac{\pi}{3})\Delta_i \right) \quad (8)$$

$$R_n < O(\sqrt{Kn\log(n)}), \text{ provided } n \gg K \quad (9)$$

The confidence bound we use in the UCB algorithm is clearly related but not identical to the bound in equation 4. The difference is due to the fact that when we use the UCB algorithm, the number of times each arm is selected is not fixed in advance but depends on the results of previous actions.

We now want to prove the bound in equation 8 holds.

Theorem 1. *If $I_t = i \neq i^*$ at least one of the following statements is true:*

1. *The estimated UB on the best arm, i^* , is less than or equal to the actual reward for that arm:*

$$\hat{\mu}_{i^*t} + \hat{\epsilon}_{i^*t} \leq \mu_{i^*} + \Delta_{i^*}$$

2. *The UB on the non-optimal arm is too high:*

$$\hat{\mu}_{it} + \hat{\epsilon}_{it} \geq \mu_i + 2\hat{\epsilon}_{it} \quad (10)$$

3. *The uncertainty bound is too wide compared the difference between the payoff of this arm and the optimal one. Since $\hat{\epsilon}_{it}$ is a function of the number of times we have selected arm i , this can also be thought of as 'we have not selected arm i enough yet'.*

$$\Delta_i < 2\hat{\epsilon}_{it} \quad (11)$$

Proof. Assume statements 1-3 are all false.

$$\begin{aligned}
UB(i^*) &\equiv \hat{\mu}_{i^*t} + \hat{\epsilon}_{i^*t} > \mu_i + \Delta_i > \mu_i + 2\hat{\epsilon}_{it} \\
UB(i) &\equiv \hat{\mu}_{it} + \hat{\epsilon}_{it} < \mu_i + 2\hat{\epsilon}_{it} \\
\implies UB(i) &< UB(i^*) \text{ which contradicts the statement that we will play arm } i, I_t = i
\end{aligned}$$

□

Suppose we had selected a non-optimal arm i in all timesteps until γ , where we choose a value γ that ensures statement 3 is false. In the remaining time-steps, we can then only select arm i if either statement 1 or statement 2 is true.

$$\begin{aligned}
E[T_i(n)] &\leq \gamma + E\left[\sum_{t=\gamma+1}^n \mathbb{1}\{(1) \text{ or } (2) \text{ is true}\}\right] \leftarrow \text{since if (3) is false, (1) or (2) must be true} \\
&\leq \gamma + \sum_{t=\gamma+1}^n [\mathbb{P}((1) \text{ is true}) + \mathbb{P}((2) \text{ is true})]
\end{aligned}$$

$$\begin{aligned}
P((1) \text{ is true}) &= P(\hat{\mu}_{i^*t} + \sqrt{\frac{2\log(t)}{T_i(t)}} \leq \mu^*) \\
&\leq P(\exists s \in \{1 \dots t\} : \hat{\mu}_{i^*s} + \sqrt{\frac{2\log(t)}{s}} \leq \mu^*) \leftarrow \text{to get around the problem that } T_i(t) \text{ is random} \\
&\leq \sum_{s=1}^t P\left(\hat{\mu}_{i^*s} + \sqrt{\frac{2\log(t)}{s}} \leq \mu^*\right) \leftarrow \text{union bound}
\end{aligned}$$

From equation (4) we have:

$$\begin{aligned}
&P\left(\mu_i - \hat{\mu}_{is} > \sqrt{\frac{\log(1/\delta)}{2s}}\right) \leq \delta, \text{ letting } \delta = t^{-4}, \\
\implies P\left(\hat{\mu}_{i^*s} + \sqrt{\frac{2\log(t)}{s}} \leq \mu^*\right) &< t^{-4} \\
\implies P((1) \text{ is true}) &\leq \sum_{s=1}^t t^{-4} = t * t^{-4} = t^{-3}
\end{aligned}$$

Similarly, $P((2) \text{ is true}) \leq t^{-3}$ so we have proved:

$$E[T_i(n)] \leq \gamma + \sum_{t=\gamma+1}^n 2t^{-3} \quad (12)$$

Plugging this into the definition of pseudo-regret in equation 2 gives:

$$R_n \leq \sum_{i=1}^K \Delta_i \left(\gamma + \sum_{t=\gamma+1}^n 2t^{-3} \right) \quad (13)$$