

UCB bandits w/independent arms

a^* = optimal action

a_t = action selected at time step t

$f^*(a)$ = mean reward of a under true reward

$$\text{regret}(\text{UCB}, T) = \sum_{t=1}^T f^*(a^*) - f^*(a_t)$$

time steps

difference in
expected best reward
& expected reward of action selected

$$= \sum_{t=1}^T [U_t(a_t) - f^*(a_t)] + [f^*(a^*) - U_t(a_t)]$$

where $U_t(a_t)$ is an upper confidence bound on
the expected reward of action a_t at time t

can use Hoeffding or Azuma inequalities to compute an
upper bound on expected reward of each action at
each time step

e.g. $U_t(a_t) = \underbrace{\hat{f}_t(a_t)}_{\text{empirical}} + \sqrt{\frac{4 \log(1/\delta)}{n_t(a)}} \quad \text{holds w/prob } \geq 1 - \delta t^{-2}$

$n_t(a) \leftarrow \# \text{times tried action } a_t$

what is the prob that on some step (in T steps) one of the
upper confidence bounds will fail to hold?

$$\begin{aligned} \sum_{t=1}^T P(f^*(a^*) > U_t(a_t)) &\leq \sum_{t=1}^T \sum_{i=1}^m P(|f^*(a_i) - \hat{f}(a_i)| > \sqrt{\frac{4 \log(1/\delta)}{n_t(a)}}) \\ &\leq \sum_{t=1}^T \sum_{i=1}^m \delta t^{-2} \quad \text{aside: } \sum_{i=1}^{\infty} i^{-2} = \frac{\pi^2}{6} < 2 \\ &\leq 2m\delta \end{aligned}$$

\Rightarrow w/prob at least $1 - 2m\delta$ UCB will overestimate the
true rewards on all T timesteps

$$\Rightarrow f^*(a^*) - U_t(a_t) \leq 0$$

$$\Rightarrow \text{regret}(\text{UCB}, T) \leq \sum_{t=1}^T U_t(a_t) - f^*(a_t) \quad (\text{dropped 2nd term})$$

* This is important because now can compare our bounds
 U_t (which are calculated) to observed reward $f^*(a_t)$
instead of unobserved $f^*(a^*)$

therefore with prob $\geq 1 - 2m\delta$

$$\text{regret}(\text{UCB}, T) \leq \sum_{t=1}^T u_t(a_t) - f^*(a_t)$$

$$\leq \sum_{t=1}^T \sqrt{\frac{4 \log(T\delta)}{n_t(a_t)}}$$

$$\leq \sqrt{4 \log(T\delta)} \sum_{t=1}^T \sqrt{\frac{1}{n_t(a_t)}}$$

$$(*) = \sqrt{4 \log(T\delta)} \sum_{i=1}^m \sum_{n=1}^{n(i)} \sqrt{\frac{1}{n}} \quad \begin{matrix} \text{split by each} \\ \text{of } m \text{ actions} \end{matrix}$$

note $\sum_{n=1}^m \sqrt{\frac{1}{n}} \leq 2\sqrt{T}$ by an integral comparison, then

$\sum_{i=1}^m \sum_{n=1}^{n(i)} \frac{1}{n}$ is maximized if all arms pulled equally

$$\text{so } \sum_{i=1}^m \sum_{n=1}^{n(i)} \frac{1}{n} \leq \sum_{i=1}^m \sum_{n=1}^{T/m} \frac{1}{n} \\ \leq \sum_{i=1}^m 2\sqrt{T/m} \\ = 2m\sqrt{T/m} \\ = 2\sqrt{Tm}$$

substitute this back into (*)

$$\text{regret}(\text{UCB}, T) \leq 4\sqrt{mT \log(T\delta)} \quad \text{with prob at least } 1 - 2m\delta$$

can be made tighter in terms of log & constants

- Regret vs Bayesian regret

$$\text{regret}(\pi, T, \theta) = \sum_{t=1}^T \mathbb{E} [f^*(a^*) - f^*(a_t) | \theta] \quad \leftarrow \text{parameters}$$

$$\text{Bayes regret}(\pi, T) = \sum_{t=1}^T \mathbb{E} [f^*(a^*) - f^*(a_t)] \quad \text{with respect to prior } \mu \text{ over } \theta$$

can be related

$$\text{let Bayes regret}(\pi, T) = O(g(T)) \text{ for some non-neg func } g$$

then $\forall \epsilon > 0$

$$P(\text{Regret}(\pi, T, \theta) \geq g(T)/\epsilon) \leq \epsilon \cdot HT \quad \text{by Markov's inequality}$$

can also use to help bound impact of using wrong prior

- Lemma 2 (Posterior sampling). Let ϕ be the distib of f^* & consider a function g (must satisfy a few properties - see references)

$$E[g(f^*) | H_t] = E[g(f_t) | H_t] \quad \text{where } H_t \text{ is the history of observed actions & rewards}$$

$$\text{regret}(\text{Posterior Sampling}, T) = \sum_{t=1}^T f^*(a^*) - f^*(a_t) \quad f_t(a_t) \text{ is our imagined reward}$$

$$= \sum_{t=1}^T |f_t(a_t) - f^*(a_t)| + |f^*(a^*) - f_t(a_t)|$$

expct. wrt prior

$= 0$ in expct due to lemma

Posterior sampling bound cont.

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T |f_t(x_t) - f^*(x_t)| \right] &= \mathbb{E} \left[\sum_{t=1}^T [f_t(x_t) - U_t(x_t)] + \sum_{t=1}^T [U_t(x_t) - f^*(x_t)] \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T |f_t(x_t) - U_t(x_t)| + \sum_{t=1}^T |U_t(x_t) - f^*(x_t)| \right] \\ &\leq 2 \mathbb{E} \left[\underbrace{\sum_{t=1}^T |U_t(x_t) - L_t(x_t)|}_{\text{so now related to upper & lower confidence bounds}} \right] \end{aligned}$$

$$\begin{aligned} \text{regret (posterior sampling, T)} &\leq 2 \mathbb{E} \left[\sum_{t=1}^T |U_t(x_t) - L_t(x_t)| \right] \\ &\leq 2 \mathbb{E} \left[\sum_{t=1}^T 2 \sqrt{\frac{4 \log(\delta)}{n_t(x_t)}} \right] \leftarrow \begin{array}{l} \text{using same bounds} \\ \text{as UCB} \end{array} \\ &\leq 16 \sqrt{m T \log(\delta)} \quad w/\text{prob} \geq 1 - 2m\delta \end{aligned}$$

so posterior sampling has same (ignoring constants, logs)
regret bounds!