

Prof. Dr. Leif Döring
Sara Klein, Benedikt Wille

Reinforcement Learning

3. Exercise Sheet

1. Upper bound on $\hat{Q}_a(t)$ for many samples

Suppose ν is a bandit model with 1-sub-gaussian arms. Show that under the UCB Algorithm $\hat{Q}_a(t) < Q_a + \Delta_a$ with probability $1 - \delta$, given that $T_a(t) > \frac{2 \log(1/\delta)}{\Delta_a^2}$.

2. Best Baseline

The variance of a random vector X is defined by to be $\mathbb{V}[X] := \mathbb{E}[\|X\|_2^2] - \|E[X]\|_2^2$. Show by differentiation that

$$b_* = \frac{\mathbb{E}_{\pi_\theta}[X_A \|\nabla \log \pi_\theta(A)\|_2^2]}{\mathbb{E}_{\pi_\theta}[\|\nabla \log \pi_\theta(A)\|_2^2]}$$

is the baseline that minimises the variance of the unbiased estimators

$$(X_A - b) \nabla \log(\pi_\theta(A)), \quad A \sim \pi_\theta,$$

of $\nabla J(\theta)$.

3. Programming task*: Gradient Bandit Methods

Implement the Gradient Bandit algorithm from example 1.2.15 of the lecture. The probability that an arm is drawn is given by the soft-max distribution

$$\mathbb{P}_\pi(A_t = a) = \frac{\exp(\theta_t(a))}{\sum_{b=1}^k \exp(\theta_t(b))} =: \pi_t(a), a \in \mathcal{A}.$$

The weights are updated as follows

$$\theta_{t+1}(a) = \begin{cases} \theta_t(a) + \alpha(x_t - \bar{x}_t)(1 - \pi_t(a)) & , a = A_t \\ \theta_t(a) - \alpha(x_t - \bar{x}_t)\pi_t(a) & , \text{otherwise} \end{cases},$$

where $\alpha > 0$ is a step-size parameter and $\bar{x}_t := \frac{1}{t-1} \sum_{i=1}^{t-1} x_i$ is the mean reward until time $t - 1$.

Implement the Gradient Bandit algorithm with and without the baseline term $\bar{x}_t, t \in \{1, \dots, n\}$ and test both algorithm on a Gaussian Bandit. Play around with the mean parameters of the Gaussian Bandit. What do you notice?

4. Programming task*: Boltzmann Exploration

In this task we want to implement different variants of the Boltzmann Exploration Algorithm.

- (a) Use the implementation from Algorithm 1.
- (b) Use the Gumbel trick from Lemma 1.2.11 to implement the algorithm.
- (c) Use the implementation of Boltzmann exploration from algorithm 2. Source: paper section 4.
- (d) Test the different algorithms on a multi-armed bandit. Which algorithm with which parameter configuration leads to the best results (minimum reward, maximum best action probability)?

Input : Initialization $\hat{Q}_a(0), a \in \mathcal{A}$, number of total timesteps $n \in \mathbb{N}$, number of arms $k \in \mathbb{N}$ and parameter $\theta > 0$

Output: Trajectory of Rewards $(x_t)_{t \in \{1, \dots, n\}}$ and actions $(a_t)_{t \in \{1, \dots, n\}}$

begin

for $t \leftarrow 1$ **to** n **do**

Sample a_t from $\text{SM}(\theta, (\hat{Q}_a(t))_{a \in \mathcal{A}})$;

Obtain reward x_t by playing arm a_t ;

Update the estimated action value functions $(\hat{Q}_a(t), a \in \mathcal{A})$;

end

return $(x_t)_{t \in \{1, \dots, n\}}, (a_t)_{t \in \{1, \dots, n\}}$;

end

Algorithm 1: Boltzmann exploration algorithm

Input : Initialization $\hat{Q}_a(0), a \in \mathcal{A}$, number of total timesteps n , number of arms $k \in \mathbb{N}$ and parameter $C \in \mathbb{R}$

Output: Trajectory of Rewards $(X_t)_{t \in \{1, \dots, n\}}$ and actions $(A_t)_{t \in \{1, \dots, n\}}$

begin

for $t \leftarrow 1$ **to** n **do**

Simulate $z_a, a \in \mathcal{A}$ independently identically standard Gumbel;

Set $a_t = \arg \max_{a \in \mathcal{A}} \{ \hat{Q}_a(t) + \sqrt{\frac{C^2}{N_a}} z_a \}$;

Obtain reward x_t by playing arm a_t ;

Set $N_{a_t} = N_{a_t} + 1$;

Update the estimated action value functions $(\hat{Q}_a(t))$;

end

return $(X_t)_{t \in \{1, \dots, n\}}, (A_t)_{t \in \{1, \dots, n\}}$;

end

Algorithm 2: Boltzmann exploration algorithm modified