

Abstract

Introduction

Explore/Exploit

Many decisions that people are faced with require finding a balance between exploiting known information for short-term gain and exploring new sources of information that may lead to future reward. For example, a person might choose to go to a restaurant that they know well and have had multiple satisfying meals at, or choose to try a new restaurant that could result in a better or worse experience. This is known as the explore-exploit trade-off.

MAB

Many tasks requiring this kind of trade-off can be characterized as multi-arm bandit (MAB) problems. A MAB problem consists of A possible actions, where the a -th action yields the random reward r_t drawn from an unknown distribution. Over T trials, the agent must choose $a_{1:T}$ such that $\sum r_{1:T}$ is maximized. The policy $\tilde{\pi}$ maps the observed history of actions and rewards $h_{1:t-1} = [(a_1, r_1), \dots, (a_{t-1}, r_{t-1})]$ to the next action a_t for all possible histories $h_{1:t-1} \in \mathcal{H}$. The goal in a MAB problem is to find the optimal policy such that the expected reward over T trials, $\mathbf{E}[r_{1:T} | \tilde{\pi}(h_{1:t-1})]$, is maximized. This is the solution to the Bellman equation

$$\mathbf{E}[r_{1:T} | \tilde{\pi}(h_{1:t-1})] = \mathbf{E}[r_t | \tilde{\pi}(h_{1:t-1})] + \mathbf{E}[r_{t+1:T} | \tilde{\pi}(h_{1:t})]$$

where the first term yields the expected reward on trial t and the second term recursively defines the expected reward on all subsequent trials up to T . This second term can be evaluated for trial $t+1$ by weighing the expected reward given all possible histories $\mathcal{H}_{1:t} = [(a_1, r_1), \dots, (a_{t-1}, r_{t-1}), (a_t, r_t)]$ for $r_t \in \mathcal{R}$ by the probability of that history occurring:

$$\mathbf{E}[r_{t+1:T} | \tilde{\pi}(h_{1:t})] = \sum_{r_t \in \mathcal{R}} p(r_t | a_t, h_{1:t-1}) \mathbf{E}[r_{1:T} | \tilde{\pi}(h_{1:t})]$$

If the reward r observed after choosing action a is drawn from the distribution $H(\theta_a)$, an agent must employ a strategy that balances choosing the action that maximizes the reward on trial t , $a_t^* = \operatorname{argmax}_{a_t \in A} \mathbf{E}[H(\theta_{a_t})]$, and learning more about θ_A by choosing novel actions to maximize reward on future trials. As the problem of evaluating all possible histories is intractable, one of two general classes of approximations can be used. The first, sparse sampling, involves using a subset of possible histories to evaluate the long-term reward of an action, and includes Monte Carlos tree search (?, ?) and its variants (e.g., 5035667, Gelly:2007:COO:1273496.1273531). Rather than relying on simulations, myopic strategies define a value function for approximate long-term reward. Common myopic strategies include upper confidence bound (UCB) (?, ?) and epsilon-greedy (?, ?) algorithms.

CMAB/Function Learning

Contextual multi-armed bandit (CMAB) problems introduce additional information into the standard MAB problem in the form of a set of features associated with each option. By learning a map between these features and reward, an agent can gain information about the expected reward of actions besides the one that it has chosen. (Function Learning)

Basic CMAB Strategies

(eps greedy, mean/var greedy, ucb/pmi/mi, entropy search)

Once an agent has a mapping between features and reward, the policy $\tilde{\pi}$ must be selected that chooses the action a_t given the history of observed actions, features, and rewards $h_{1:t-1}$. Given that this policy does not take into account all possible trajectories, it must include a value function that approximates the long-term reward of any given action.

Alternative Goals

In addition to CMABs, the explore/exploit tradeoff can also be observed in other common context-dependent problems. In active learning (AL) (?, ?), participants assume some degree of control over the contexts that they observe, rather than passively observing a predetermined set of examples. Optimization problems (?, ?) require finding the set of features x that maximizes some reward $f(x)$; for example, finding the number of hours to dedicate to work and to leisure respectively that maximizes satisfaction.

While these tasks might all share the same mapping between action and reward, their goal-specific reward, and thus the function used to approximate long-term reward, are distinct. For any particular action reward pair (a_t, r_t) , there exists the goal-specific action reward pair (a_t, r'_t) . Since the goal in the CMAB task is to maximize cumulative reward over time, the goal-specific reward is the same as the reward, that is:

$$\mathbf{E}_{cmab}[r'_t | a_t, h_{1:t-1}] = \mathbf{E}[r_t | a_t, h_{1:t-1}]$$

In contrast, the goal of the optimization problem is simply to find the maximum possible reward withing T steps. As such, the goal specific reward is defined as:

$$\mathbf{E}_{opt}[r'_t | a_t, h_{1:t-1}] = \max(\mathbf{E}[r_t | a_t, h_{1:t-1}] - \sum_{i=1}^{t-1} r'_i, 0)$$

That is, on trial t if trial t yields an increase in reward over the previous maximum reward, r'_t is the difference. If not, r'_t is 0. Since active learning is concerned with learning the reward function rather than the magnitude of the rewards themselves, its goal-specific reward can be described as the sum of the decrease in variance across all possible actions:

$$\mathbf{E}_{al}[r'_t | a_t, h_{1:t-1}] = \sum_{a \in A} \mathbf{Var}[r | a, h_{1:t-1}] - \mathbf{Var}[r | a, h_{1:t}]$$