

# Reinforcement Learning and Blackjack

Benjamin Allévius<sup>1</sup>, Sebastian Rosengren<sup>1</sup>, and Erik Thorsén<sup>1</sup>

<sup>1</sup>Department of Mathematics, Stockholm University, Sweden

**Abstract**

# 1 Introduction

The purpose of this project is to use reinforcement learning to train an *agent* to play blackjack and investigate the learning capabilities of the framework. We implement two different representations of the state space, one of which, albeit intuitive, does not satisfy the model's stationarity assumptions. In comparing the two different representations we investigate the importance of the stationarity assumption. We apply these state representations and reinforcement learning to different versions of blackjack—namely difference in how many decks of cards that are used, from a finite number to infinite.

## 2 Reinforcement Learning

Reinforcement learning is, as many things are in machine learning, both the problems and solutions pertaining to a specific domain. The problem is how an *agent* ought to act in an *environment* with imperfect information and randomness, but with feedback. The solutions are many, but most common is the Markov decision process approach.

### 2.1 Markov Decision Process

Markov decision processes (MDP) is the formalism which allows us to reason about reinforcement learning in a mathematical way, and all standard reinforcement learning problems are formulated in a Markov decision process framework. We start by giving the definition to the simplest case of a *stationary* Markov decision process.

The building blocks of a Markov decision process consist of

- (i)  $S$  a finite space of *states*
- (ii)  $A$  a finite space of *actions*
- (iii)  $R$  a finite space of *rewards*
- (iv)  $P_a(s, s')$  a transition probability function defined for all  $(s, a, s') \in S \times A \times S$
- (v)  $r(s, a)$  the immediate or expected immediate reward of taking action  $a$  in state  $s$ .
- (vi)  $\gamma \in [0, 1]$ , the discount factor (here always equal to 1).

Let  $\{(S_t, A_t, R_{t+1}), t \geq 0\}$  be a stochastic process, where  $S_t \in S$  is the state of the system at time  $t$ ;  $A_t \in A$  is the (potentially random) action taken at time  $t$ , and  $R_t \in R$  is the immediate reward at time  $t$ . The *history*  $H_t$  of the system up to time  $t$  is the random vector given by  $(S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}, S_t)$ . An important concept going further is that of a policy: a policy  $\pi = (\pi_0, \pi_1, \dots)$  is a sequence of decision rules, where  $\pi_n$  is a function of the history of the process up to time  $n$  (i.e.  $H_n$ ) mapping to a probability measure over  $A$ .

Now, we call the process  $\{(S_t, A_t, R_{t+1}), t \geq 0\}$  a Markov decision process if

$$\begin{aligned} \mathbb{P}(S_t = s | S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}, \dots, S_0 = s_0, A_0 = a_0) \\ = \mathbb{P}(S_t = s | S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}) = P_a(s_{t-1}, s) \end{aligned} \quad (1)$$

and  $R_{t+1}$  given  $s_t, a_t$  is independent of  $H_t$ , with  $\mathbb{E}(R_{t+1} | H_t) = r(S_t, A_t)$ . The process is called stationary since the transition probabilities nor the reward kernel depends on  $t$ .

At first glance expression (1) makes sense, but there are subtle details which need to be addressed. For instance, we have said nothing about the process  $\{A_t\}$  and therefore we should be careful defining conditional probabilities containing it. However, if  $\{A_t\}$  is governed by a policy  $\pi = (\pi_0, \pi_1, \dots)$  then we can define the above probability. In order for  $\{(S_t, A_t, R_{t+1}), t \geq 0\}$  to be a Markov decision process we formally require that (1) holds for all policies. For a thorough treatment of Markov decision processes see Puterman (1994).

The goal of reinforcement learning is to take a MDP and choose a policy which maximizes expected (discounted) rewards. This can seem hopeless since policy decisions in general depends on the whole history of the chain. However, it can be shown that for finite stationary MDPs this is equivalent to maximizing rewards over policies which take into account only the current state of the process and which maps this to a single action, i.e.  $\pi_n(H_n) = \pi(S_t) \in A$  — this is known as a *stationary Markov policy*. This is something that make the optimizing problem considerably easier, and allows e.g. the *Q-learning* algorithm to give estimates that converges to the optimal solution. Hence, from here on out we consider only stationary Markov policies.

For a given policy  $\pi$  define the *value function*

$$V^\pi(s) = \mathbb{E}_\pi\left(\sum_{n=0}^{\infty} R_{t+1} | S_0 = s\right) \quad (2)$$

i.e. the expected reward obtained when following policy  $\pi$ . If  $s_t$  is an *terminal*

state — a state from which the process never leaves — we take  $R_{t+1} = 0$ , which implies that  $V^\pi(s) = 0$  if  $s$  is terminal.

We say that a policy  $\pi^*$  is optimal if

$$V^{\pi^*}(s) \geq V^\pi(s) \quad \forall s \text{ and policies } \pi.$$

It can be shown that for finite Markov decision processes such a policy exists uniquely. Under the optimal policy  $\pi^*$  let  $V^*(s) = V^{\pi^*}(s)$ . Say we know  $V^*(s)$  for all  $s$ , can this be used to calculate an optimal policy? The answer is yes. One can prove that

$$V^*(s) = r(s, \pi^*(s)) + \sum_{s'} P_{\pi^*(s)}(s, s') V^*(s') \quad (3)$$

$$= \max_{a \in A} (r(s, a) + \sum_{s'} P_a(s, s') V^*(s')). \quad (4)$$

The equations are known as the *Bellman equations* and they can be used to show that  $\pi^*(s)$  is given by the action  $a$  which achieves the above maximum. So if we can find a way to calculate the optimal value function we get the optimal strategy. In small state spaces the value function can actually be calculated but in large spaces approximation is the best one can hope for.

## 2.2 Q-Learning

Q-learning is an algorithm which can be used to approximate the optimal value function and corresponding strategy. Here we will assume that the underlying MDP is terminal.

In order to specify the algorithm we need some definitions. Define the *action-value function* as

$$Q^\pi(s, a) = r(s, a) + \sum_{s'} P_a(s, s') V^\pi(s'). \quad (5)$$

That is the expected reward of starting in state  $s$  taking action  $a$  and then following the policy  $\pi$ . For terminal states  $s$  we set  $Q(s, a) = 0$  for all  $a \in A$ . Furthermore, define the optimal action-value function as

$$Q^*(s, a) = (r(s, a) + \sum_{s'} P_a(s, s') V^*(s')) \quad (6)$$

i.e. the expected reward of being in state  $s$  taking action  $a$  and then following the optimal policy  $\pi^*$ . We see that  $\max_{a \in A} Q^*(s, a) = V^*(s)$ , and with the same argument as above  $\pi^*(s) = \arg \max_{a \in A} Q^*(s, a)$ .

We can now specify the Q-learning algorithm:

- (i) Initialize  $Q(s, a)$  for all  $s, a$ , learning rate  $\alpha_t(s, a) \in (0, 1]$ , and decide on some policy e.g.  $\epsilon$ -greedy.
- (ii) Repeat for each  $N$  episodes (start  $\rightarrow$  final state):  
While  $s_t$  not terminal:

$$Q(s_t, a_t) \leftarrow (1 - \alpha_t(s_t, a_t))Q(s_t, a_t) + \alpha_t(s_t, a_t)(r_{t+1} + \max_{a \in A} Q(s_{t+1}, a)) \quad (7)$$

It is known that for stationary finite Markov decision process (and some additional assumptions) the Q-learning produces estimates converging to the optimal action-value function.

### 3 Blackjack as a Markov Decision Process

In order to test this reinforcement learning framework we have chosen to work with the simplest form of Black Jack. We assume the following setup:

- One player against the dealer, with player staking one unit on each hand;
- two actions possible: ask for another card, or stay;
- cards 2–10 counts as their numerical value, suites counts as 10, and ace counts as either 1 or 11 depending on whichever is best. If an ace can be counted as 11 without player going bust it is known as a *usable* ace, the same goes for the dealer.

The goal of the player is to beat the dealer in one of the following ways

- Get 21 points on the first two cards, knows as a blackjack, without a dealer blackjack. **Net profit:** 1.5 times stake;
- Reach a final score higher than the dealer without exceeding 21. **Net profit:** stake;
- Dealer gets points exceeding 21 and player does not. **Net profit:** stake.

The house always plays according to the same strategy: draw cards until it has a card sum greater than or equal to 17. The game starts with the dealer giving the player two cards (visible) and himself two cards (one visible, one hidden). The player is then allowed to take actions until done, after which the dealer follows his strategy until done. If player's points exceeds 21 his stake is

lost regardless of dealers outcome (this is where the house edge comes from); if player's points equals dealer's points this is a push and stake is returned to player; otherwise payout is made according to above rules. Standard practice is that the cards are drawn from 6–10 decks, and we implement these situations as well as when the cards are drawn from an infinite deck so that every card has same probability of being drawn at all times.

### 3.1 As a Markov Decision Process

Next we fit the above situation in to the Markov decision process framework. We do it in two ways: one with a large state space such that all MDP model assumptions are satisfied; one with a smaller state space, where the resulting process is no longer stationary. In both cases actions and rewards are as specified above, and we need only specify the state space in order to have the model.

### 3.2 Stationary Markov Decision Process

We begin by noting that the color and suite of the cards does not matter, only their numerical values. We use this as the cards identifier, with aces equal to 1. We use the following representation of the state space

$$S = \{(s_{p1}, \dots, s_{p10}, \Sigma_d) \mid i(s_{pi} - 1) \leq 21 \text{ for } i = 1, \dots, 10, \text{ and } \Sigma_d \leq 26\}.$$

Hence, an element of the state space is a vector of length 11, where the 10 first elements indicates the number of each card the player is holding, and the last element indicates the dealer's *visible* card sum.

Note that for a card of value  $i$ , the maximum number  $s$  of cards with this value the player can have is such that  $n - 1$  keeps the player's card sum less than or equal to 21, but drawing one more card will lead to a bust (assuming the player has no cards of other values). At first glance this state space seems enormous, and while that is true the *effective* state space is much smaller, i.e. the states which the system realistically will visit is relatively small. Also, with this state representation it is always possible to determine if a state is terminal—we need only check if the player or dealer is bust, or if dealer's card sum exceeds 17. The rewards  $R_1, R_2, \dots$  follows the above payout specification, with  $R_t = 0$  if  $s_t$  is not terminal. It is also possible, albeit tedious, to calculate the transition probabilities and to show that the resulting process is a Markov decision process.

### 3.3 Non-stationary Markov Decision Process

Most strategies in blackjack involves keeping track only of your own card sum (with some possible extensions) and not every card you are holding. Hence, it seems natural to base the state space on this observation. It is also practical to know if the player has an usable ace, so we include this in the state representation. Let

$$S = \{(\Sigma_p, a_p, \Sigma_d)\}$$

where  $\Sigma_p, \Sigma_d$  is the card sum of the player and the dealer, and  $a_p$  indicates if the player is holding a usable ace or not. This seem to be a popular representation of the state space for blackjack across the RL-community. Although natural and popular, this representation does not yield a stationary Markov process. Consider,

$$s_0 = (21, 1, x), \quad s_5 = (21, 1, x).$$

The first situation is a blackjack and the second is not. Hence, the expected rewards of the action “staying” will differ for the two states

$$r(s_0, \text{stay}) = 1.5 \neq r(s_5, \text{stay}).$$

For finite decks this state representation also yields non-stationary transition probabilities, however they are stationary in the infinite deck case.

## 4 Implementation

In implementing a Q-learning algorithm for Blackjack, we face a number of choices. First, as discussed in Sections 3.2 and 3.3, is the choice of state space. Second, the parameters  $\alpha$  and  $\epsilon$  of the algorithm, which may be functions of the current state  $s$  and the action  $a$ . Third, the implementation itself—how do we implement the algorithm in code? The first and third points are interconnected here, so we deal with the second point first, then the other two.

### 4.1 Parameter choices for Blackjack Q-learning

The (possibly) state- and action-dependent learning rate  $\alpha_t(s, a) \in (0, 1]$  determines to what extent we update the  $Q$  function each time we visit state  $s$  and take action  $a$ . A small  $\alpha_t$  means that  $Q$  changes slowly from its initial value (which is an additional parameter to choose), and a larger  $\alpha_t$  means

that we put more emphasis on the immediate (well, at  $t + 1$ ) reward and currently optimal action of the state  $s'$  we transition to. Intuitively, our estimate of  $Q(s, a)$  should improve with the number of times we have visited  $s$  and taken action  $a$ , meaning that it becomes less important to update the function in this state-action pair. Vice versa, in state-action pairs we have not encountered many times before, our estimate of the function value is uncertain, and we should put more emphasis on the reward we received for taking action  $a$ , as well as the currently known best action out of the state  $s'$  we transition to. For our Blackjack setup, we may note again that the immediate rewards are 0 in all states except for a terminal state, in which case the function  $Q$  is defined to be zero. Thus, the update rule 7 becomes

$$Q(s_t, a_t) \leftarrow \begin{cases} (1 - \alpha_t(s_t, a_t))Q(s_t, a_t) + \alpha_t(s_t, a_t)r_{t+1}, & s_{t+1} \text{ terminal} \\ (1 - \alpha_t(s_t, a_t))Q(s_t, a_t) + \alpha_t(s_t, a_t) \max_{a \in A} Q(s_{t+1}, a), & s_{t+1} \text{ non-terminal.} \end{cases} \quad (8)$$

A choice of  $\alpha_t(s_t, a_t)$  that has the intuitive properties mentioned above is  $\alpha_t(s_t, a_t) = \#[(s = s_t, a = a_t)]^{-1}$ , i.e. the reciprocal of the number of times action  $a = a_t$  has been taken in state  $s = s_t$ . This choice of  $\alpha_t$  satisfies the convergence criteria  $\sum_t \alpha_t(s_t, a_t) = \infty$  and  $\sum_t \alpha_t^2(s_t, a_t) < \infty$  outlined in e.g. Watkins and Dayan (1992). Further investigation shows that one can do better than this when the function  $Q$  is updated asynchronously, i.e. one  $(s, a)$ -pair at a time. Even-Dar and Mansour (2003) show that in this case, the optimal choice is  $\alpha_t(s_t, a_t) = \#[(s = s_t, a = a_t)]^{-\omega}$ , where  $\omega \approx 0.77$ . This is the choice of  $\alpha_t$  we choose, although it should be noted that the scenario considered by Even-Dar and Mansour (2003) included a discount factor  $\gamma < 1$ .

Another important choice in implementing Q-learning is the choice of policy—what action should we take in a given state? Under certain conditions, our learned  $Q$  will eventually guide us to the optimal policy, but this is not the case initially. If we used the greedy policy of always choosing the action  $a$  with the highest value of  $Q$  when in state  $s$ , an unlucky start may lead us to never explore some states and find (approximately) the true value  $Q^*(s, a)$  for all actions  $a$  available when in state  $s$ —some perhaps more valuable than the greedy action chosen. Indeed, the second condition for convergence to  $Q^*$  is that all pairs  $(s, a)$  are visited infinitely often, asymptotically (Buşoniu et al., 2010). Commonly, this is done by choosing a random action with some small probability  $\epsilon$ , rather than choosing the greedy action. Again, it is reasonable that if we have visited a given state  $s$  many times, and also taken all actions out of this state many times, our estimate  $Q(s, a)$  for all actions  $a$  available in  $s$  should be fairly certain. Thus, there is less



need to choose a random action; we may go with the greedy choice with a high probability. There are several ways of implementing such  $\epsilon$ -decay. We choose the rather simple option of choosing a random action when in state  $s_t = s$  with probability  $\epsilon_t(s_t) = \frac{c}{\#[(s=s_t)]}$  for some positive constant  $c \leq 1$ , though it should be noted that more advanced options such as Boltzmann exploration exist (Buşoniu et al., 2010).

## 4.2 Software implementation

# 5 Results

## References

- Buşoniu, L., Babuška, R., De Schutter, B., , and Ernst, D. (2010). *Reinforcement learning and dynamic programming using function approximators*. CRC Press.
- Even-Dar, E. and Mansour, Y. (2003). Learning rates for q-learning. *Journal of Machine Learning Research*, 5:1–25.
- Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition.
- Watkins, C. J. and Dayan, P. (1992). Technical note: Q-learning. *Machine Learning*, 8:279–292.