

# HW2

Bucky Park

## 1 Introduction

I recently have been interested in finance and realized that the reinforcement learning can be applied to finance. While studying the textbook *Foundation of Reinforcement Learning with Applications in Finance* by Ashwin Rao and Tikhon Jelvis, I would like to summarize the methods for predicting and analyzing stock prices using RL from the first few chapters, mainly chapter 3 and 4. Codes are available at <https://github.com/TikhonJelvis/RL-book>. This textbook introduces some basic concepts in RL such as Markov processes, Markov reward processes and Markov decision processes with simple inventory examples. We first explore the setup of this example.

## 2 Setup

Let us focus on the inventory of a particular type of bicycle. Each day, there is random demand which is a non-negative integer. The probability of demand follows a Poisson distribution with parameter  $\lambda$ . In other words, for each  $i = 0, 1, 2, \dots$ , the probability of a demand  $i$  is

$$f(i) = \frac{e^{-\lambda} \lambda^i}{i!}.$$

Let  $C$  denote the maximum storage capacity. At the end of each day, you have the option to order a certain number of bicycles from your supplier. The ordered bicycles will arrive during the store's open hours in 2 days. Let the store's state at the end of each day be  $(\alpha, \beta)$ , where  $\alpha$  is the current inventory in the store and  $\beta$  is the inventory on a truck from the supplier, which was ordered the previous day and will arrive the next morning. Given the storage capacity constraint of at most  $C$ , your ordering policy for today is to order  $C - (\alpha + \beta)$  if  $\alpha + \beta < C$ ; and to not order otherwise.

## 3 Markov Processes with Inventory Example

In the setting above, we consider the set of states

$$S = \{(\alpha, \beta) \in \mathbb{Z}_{\geq 0}^2 \mid 0 \leq \alpha + \beta \leq C\}.$$

Let  $S_t = (\alpha, \beta)$  be the current state, then there are only  $\alpha + \beta + 1$  possible next states  $S_{t+1}$  for  $i = 0, 1, \dots, \alpha + \beta$  as follows:

$$(\alpha + \beta - i, C - (\alpha + \beta))$$

with transition probabilities by the Poisson distribution of demand as

$$\begin{cases} P((\alpha, \beta), (\alpha + \beta - i, C - (\alpha + \beta))) = f(i), & \text{for } 0 \leq i \leq \alpha + \beta - 1 \\ ((\alpha, \beta), (0, C - (\alpha + \beta))) = \sum_{j=\alpha+\beta}^{\infty} f(j), & \text{otherwise.} \end{cases}$$

We write code for this simple inventory example by defining transition map of the Markov processes above. Here, we cannot implement  $j$  going to infinity, in the second case where  $i > \alpha + \beta$ , we will store it as a single state.

```

1 from dataclasses import dataclass
2 from typing import Mapping, Dict
3 from rl.distribution import Categorical, FiniteDistribution
4 from rl.markov_process import FiniteMarkovProcess
5 from scipy.stats import poisson
6
7 @dataclass(frozen=True)
8 class InventoryState: # class for the inventory state
9     on_hand: int # represent alpha
10    on_order: int # represent beta
11
12    def inventory_position(self) -> int:
13        return self.on_hand + self.on_order # return alpha + beta
14
15 class SimpleInventoryMPFinite(FiniteMarkovProcess[InventoryState]):
16
17     def __init__(
18         self,
19         capacity: int,
20         poisson_lambda: float
21     ):
22         # initialization for the class from the given paramters
23         self.capacity: int = capacity
24         self.poisson_lambda: float = poisson_lambda
25
26         self.poisson_distr = poisson(poisson_lambda) # poisson distribution
27         super().__init__(self.get_transition_map()) # get_transition_map() function inherited from the superclass
28         # FinitMarkovProcess
29
30     def get_transition_map(self) -> \
31         Mapping[InventoryState, FiniteDistribution[InventoryState]]: # override get_transition_map method
32         d: Dict[InventoryState, Categorical[InventoryState]] = {} # dictionary representing the transition
33         # probabilities
34         for alpha in range(self.capacity + 1): # iterate all the states, i.e. all possible combinations of (alpha,
35             # beta)
36             for beta in range(self.capacity + 1 - alpha):
37                 state = InventoryState(alpha, beta) # extract the current state
38                 ip = state.inventory_position() # extract alpha + beta from the current state
39                 beta1 = self.capacity - ip # obtain the amount of order for today, will be beta for the next step
40                 state_probs_map: Mapping[InventoryState, float] = {
41                     InventoryState(ip - i, beta1): # for each i, create the next possible state
42                     (self.poisson_distr.pmf(i) if i < ip else # with probability f(i) if i < ip
43                      1 - self.poisson_distr.cdf(ip - 1)) # with 1 - cdf otherwise
44                     for i in range(ip + 1) # we do not iterate i to infinity as above.
45                     # Just iterate i until alpha + beta and gather all probability mass for i >
46                     # alpha + beta to i = alpha + beta
47                 } # calculate transition probabilities of the current state
48                 d[InventoryState(alpha, beta)] = Categorical(state_probs_map) # put transition probabilities of the
49                 # current state into the dictionary
50         return d # return obtained transition probabilities
51
52 if __name__ == '__main__':
53     user_capacity = 2 # set the initial number of bicycle
54     user_poisson_lambda = 1.0 # set the Poisson parameter
55
56     si_mp = SimpleInventoryMPFinite(
57         capacity=user_capacity,
58         poisson_lambda=user_poisson_lambda
59     ) # define Markov processes based on the initial paramters
60
61     print("Transition Map")
62     print("-----")
63     print(si_mp)

```

```

60 print("Stationary Distribution")
61 print("-----")
62 si_mp.display_stationary_distribution()

```

Code 1: rl/chapter2/simple\_inventory\_mp.py

## 4 Markov Reward Processes with Inventory Example

Let us assume that there are two types of costs:

- A holding cost of  $h$  for each bicycle that remains in your store overnight.
- A stockout cost of  $p$  for each unit of unmet demand.

After closing the store each day, the reward  $R_{t+1}$  is calculated as the negative sum of the overnight holding cost and the day's stockout cost. Since customer demand can take an infinite number of values, this results in an infinite set of possible next states and rewards from the current state. For the current state  $S_t = (\alpha, \beta)$ , and customer demand  $i$ ,  $S_{t+1}$  can be represented as

$$(\max(\alpha + \beta - i, 0), \max(C - (\alpha + \beta), 0))$$

and the reward  $R_{t+1}$  is

$$-h \cdot \alpha - p \cdot \max(i - (\alpha + \beta), 0)$$

Therefore, the transition probability for each  $i = 0, 1, 2, \dots$  can be represented as

$$P((\alpha, \beta), -h \cdot \alpha - p \cdot \max(i - (\alpha + \beta), 0), (\max(\alpha + \beta - i, 0), \max(C - (\alpha + \beta), 0))) = \frac{e^{-\lambda} \lambda^i}{i!}$$

We now work out the calculation of the reward function  $R_t$ . When  $\alpha$  of the next state  $S_{t+1}$  is positive, this correspond to no stockout cost and only an overnight holding cost  $h\alpha$ . Therefore, for all  $0 \leq \theta \leq C - (\alpha + \beta)$ ,

$$R_t((\alpha, \beta), \theta, (\alpha + \beta - i, \theta)) = -h\alpha$$

for  $0 \leq i \leq \alpha + \beta - 1$ . On the other hand, if  $\alpha$  of  $S_{t+1}$  is zero, the demand for the day was at least  $\alpha + \beta$  meaning there is some stockout cost and overnight holding cost. In this case, the exact stockout cost can be calculated by the expectation of the number of missed demand conditioned on the corresponding Poisson probabilities of demand exceeding  $\alpha + \beta$ . Therefore,

$$R_t((\alpha, \beta), \theta, (0, \theta)) = -h\alpha - p \frac{\sum_{j=\alpha+\beta+1}^{\infty} f(j) \cdot (j - (\alpha + \beta))}{\sum_{j=\alpha+\beta}^{\infty} f(j)} = -h\alpha - p \left( \lambda - (\alpha + \beta) \left( 1 - \frac{\alpha + \beta}{1 - F(\alpha + \beta - 1)} \right) \right).$$

We write code for this scenario by defining transition reward map of the Markov processes above. As in the case of Markov Processes, since we cannot implement  $i$  going to infinity, we will handle it as a single state.

```

1 from dataclasses import dataclass
2 from typing import Tuple, Dict, Mapping
3 from rl.markov_process import MarkovRewardProcess
4 from rl.markov_process import FiniteMarkovRewardProcess
5 from rl.markov_process import State, NonTerminal
6 from scipy.stats import poisson
7 from rl.distribution import SampledDistribution, Categorical, \
8     FiniteDistribution
9 import numpy as np
10
11
12 @dataclass(frozen=True)
13 class InventoryState: # same as the above
14     on_hand: int

```

```

15     on_order: int
16
17     def inventory_position(self) -> int:
18         return self.on_hand + self.on_order
19
20
21 class SimpleInventoryMRP(MarkovRewardProcess[InventoryState]):
22
23     def __init__(
24         self,
25         capacity: int,
26         poisson_lambda: float,
27         holding_cost: float,
28         stockout_cost: float
29     ):
30         # initialize the class from the given parameters
31         self.capacity = capacity
32         self.poisson_lambda: float = poisson_lambda
33         self.holding_cost: float = holding_cost
34         self.stockout_cost: float = stockout_cost
35
36     def transition_reward(
37         self,
38         state: NonTerminal[InventoryState]
39     ) -> SampledDistribution[Tuple[State[InventoryState], float]]:
40         # implement the abstract method transition_reward
41
42         def sample_next_state_reward(state=state) -> \
43             Tuple[State[InventoryState], float]:
44             # sample pair of next state and reward, return an instance of SampledDistribution
45
46             demand_sample: int = np.random.poisson(self.poisson_lambda) # sample the demand from Poisson distribution
47             ip: int = state.state.inventory_position() # get alpha + beta
48             next_state: InventoryState = InventoryState(
49                 max(ip - demand_sample, 0),
50                 max(self.capacity - ip, 0)
51             ) # next state based on (alpha + beta) and the demand
52             reward: float = - self.holding_cost * state.state.on_hand \
53                 - self.stockout_cost * max(demand_sample - ip, 0) # calculate the reward
54             return NonTerminal(next_state), reward # return next state and reward
55
56         return SampledDistribution(sample_next_state_reward) # return an instance of SampledDistribution with
57             sample_next_state_reward
58
59 class SimpleInventoryMRPFinite(FiniteMarkovRewardProcess[InventoryState]):
60
61     def __init__(
62         self,
63         capacity: int,
64         poisson_lambda: float,
65         holding_cost: float,
66         stockout_cost: float
67     ):
68         # initialize the class from the given parameters
69         self.capacity: int = capacity
70         self.poisson_lambda: float = poisson_lambda
71         self.holding_cost: float = holding_cost
72         self.stockout_cost: float = stockout_cost
73
74         self.poisson_distr = poisson(poisson_lambda)
75         super().__init__(self.get_transition_reward_map())
76         # get_transition_reward_map function inherited from the superclass FiniteMarkovRewardProcess
77
78     def get_transition_reward_map(self) -> \

```

```

79         Mapping[
80             InventoryState,
81             FiniteDistribution[Tuple[InventoryState, float]]
82         ]: # override get_transition_reward_map method
83         d: Dict[InventoryState, Categorical[Tuple[InventoryState, float]]] = {} # dictionary representing the
            transition reward map
84         for alpha in range(self.capacity + 1):
85             for beta in range(self.capacity + 1 - alpha): # iterate all the possible states
86                 state = InventoryState(alpha, beta) # extract the current state
87                 ip = state.inventory_position() # extract alpha + beta from the current state
88                 beta1 = self.capacity - ip # obtain the amount of order for today, will be candidate of beta for the
                    next step
89                 base_reward = - self.holding_cost * state.on_hand # precalculate h * alpha, will store the reward for
                    the given i
90
91                 sr_probs_map: Dict[Tuple[InventoryState, float], float] = \
92                     {(InventoryState(ip - i, beta1), base_reward):
93                      self.poisson_distr.pmf(i) for i in range(ip)}
94                 # construct transition reward map when (alpha + beta) > i, in this case the reward = - h * alpha
95
96                 # construct transition reward map when (alpha + beta) < i, but i can be infinite so gather all
                    probability mass to one
97                 probability = 1 - self.poisson_distr.cdf(ip - 1)
98                 # in this case, calculate the expected reward assuming i >= alpha + beta
99                 reward = base_reward - self.stockout_cost * \
100                     (self.poisson_lambda - ip *
101                      (1 - self.poisson_distr.pmf(ip) / probability))
102
103                 # put this (state, reward) and probability into the transition map
104                 sr_probs_map[(InventoryState(0, beta1), reward)] = probability
105                 d[state] = Categorical(sr_probs_map) # put transition map of the current state into the dictionary
106         return d
107
108
109 if __name__ == '__main__':
110     # set the initial parameters
111     user_capacity = 2 # represents C
112     user_poisson_lambda = 1.0 # represents lambda
113     user_holding_cost = 1.0 # represents h
114     user_stockout_cost = 10.0 # represents p
115
116     user_gamma = 0.9 # represents discount factor
117
118     si_mrp = SimpleInventoryMRPFinite(
119         capacity=user_capacity,
120         poisson_lambda=user_poisson_lambda,
121         holding_cost=user_holding_cost,
122         stockout_cost=user_stockout_cost
123     ) # define Makov Reward Processes based on the parameters
124
125     from rl.markov_process import FiniteMarkovProcess
126     print("Transition Map")
127     print("-----")
128     print(FiniteMarkovProcess(
129         {s.state: Categorical({s1.state: p for s1, p in v.table().items()})
130          for s, v in si_mrp.transition_map.items()})
131     ))
132
133     print("Transition Reward Map")
134     print("-----")
135     print(si_mrp)
136
137     print("Stationary Distribution")
138     print("-----")
139     si_mrp.display_stationary_distribution()

```

```

140     print()
141
142     print("Reward Function")
143     print("-----")
144     si_mrp.display_reward_function()
145     print()
146
147     print("Value Function")
148     print("-----")
149     si_mrp.display_value_function(gamma=user_gamma)
150     print()

```

Code 2: rl/chapter2/simple\_inventory\_mrp.py

## 5 Markov Decision Processes with Inventory Example

We have assumed that the policy is fixed to order  $C - (\alpha + \beta)$  if  $\alpha + \beta < C$ , and not to order otherwise. Now, we find the optimal value function and the optimal policy using dynamic programming. With the Markov reward process we have constructed above. First, we create a class inherited from `FiniteMarkovDecisionProcess` and use some classes methods to RL in where `rl.policy` and `rl.dynamic_programming`.

```

1  from dataclasses import dataclass
2  from typing import Tuple, Dict, Mapping
3  from rl.markov_decision_process import FiniteMarkovDecisionProcess
4  from rl.policy import FiniteDeterministicPolicy
5  from rl.markov_process import FiniteMarkovProcess, FiniteMarkovRewardProcess
6  from rl.distribution import Categorical
7  from scipy.stats import poisson
8
9
10 @dataclass(frozen=True)
11 class InventoryState: # same as the above
12     on_hand: int
13     on_order: int
14
15     def inventory_position(self) -> int:
16         return self.on_hand + self.on_order
17
18
19 InvOrderMapping = Mapping[
20     InventoryState,
21     Mapping[int, Categorical[Tuple[InventoryState, float]]]
22 ]
23
24
25 class SimpleInventoryMDPCap(FiniteMarkovDecisionProcess[InventoryState, int]):
26
27     def __init__(
28         self,
29         capacity: int,
30         poisson_lambda: float,
31         holding_cost: float,
32         stockout_cost: float
33     ):
34         self.capacity: int = capacity
35         self.poisson_lambda: float = poisson_lambda
36         self.holding_cost: float = holding_cost
37         self.stockout_cost: float = stockout_cost
38
39         self.poisson_distr = poisson(poisson_lambda)
40         super().__init__(self.get_action_transition_reward_map())
41

```

```

42 # implement get_action_transition_reward_map method for MDP
43 def get_action_transition_reward_map(self) -> InvOrderMapping:
44     d: Dict[InventoryState, Dict[int, Categorical[Tuple[InventoryState,
45                                                     float]]]] = {}
46
47     for alpha in range(self.capacity + 1):
48         for beta in range(self.capacity + 1 - alpha): # iterate all the possible states
49             state: InventoryState = InventoryState(alpha, beta) # extract the current state
50             ip: int = state.inventory_position() # extract alpha + beta from the current state
51             base_reward: float = - self.holding_cost * alpha # precalculate h * halpha, will store the reward for
                    the given i
52             d1: Dict[int, Categorical[Tuple[InventoryState, float]]] = {} # dictionary for transition maps for
                    different order policies
53
54             for order in range(self.capacity - ip + 1): # iterate all the possible order from 0 to C - alpha + beta
55
56                 # this part is the same with above code
57                 sr_probs_dict: Dict[Tuple[InventoryState, float], float] = \
58                     {(InventoryState(ip - i, order), base_reward):
59                      self.poisson_distr.pmf(i) for i in range(ip)}
60
61                 probability: float = 1 - self.poisson_distr.cdf(ip - 1)
62                 reward: float = base_reward - self.stockout_cost * \
63                     (self.poisson_lambda - ip *
64                      (1 - self.poisson_distr.pmf(ip) / probability))
65                 sr_probs_dict[(InventoryState(0, order), reward)] = \
66                     probability
67                 d1[order] = Categorical(sr_probs_dict)
68
69             # store the obtained dictionary of transition reward maps into the dictionaries
70             d[state] = d1
71     return d
72
73
74 if __name__ == '__main__':
75     from pprint import pprint
76     #set the initial parameter
77     user_capacity = 2
78     user_poisson_lambda = 1.0
79     user_holding_cost = 1.0
80     user_stockout_cost = 10.0
81
82     user_gamma = 0.9
83
84     si_mdp: FiniteMarkovDecisionProcess[InventoryState, int] = \
85         SimpleInventoryMDPCap(
86             capacity=user_capacity,
87             poisson_lambda=user_poisson_lambda,
88             holding_cost=user_holding_cost,
89             stockout_cost=user_stockout_cost
90         ) # define MDP based on the paramters
91
92     print("MDP Transition Map")
93     print("-----")
94     print(si_mdp)
95
96     fdp: FiniteDeterministicPolicy[InventoryState, int] = \
97         FiniteDeterministicPolicy(
98             {InventoryState(alpha, beta): user_capacity - (alpha + beta)
99              for alpha in range(user_capacity + 1)
100              for beta in range(user_capacity + 1 - alpha)}
101         ) # Define policies for all the possible (alpha, beta)
102
103     print("Deterministic Policy Map")
104     print("-----")

```

```

105 print(fdp)
106
107 implied_mrp: FiniteMarkovRewardProcess[InventoryState] =\
108     si_mdp.apply_finite_policy(fdp) # apply policies to mdp, it will return MRP based on the given policies
109 print("Implied MP Transition Map")
110 print("-----")
111 print(FiniteMarkovProcess(
112     {s.state: Categorical({s1.state: p for s1, p in v.table().items()})
113     for s, v in implied_mrp.transition_map.items()})
114 )) # print transition map from MRP
115
116 print("Implied MRP Transition Reward Map")
117 print("-----")
118 print(implied_mrp) # print transition reward map from MRP
119
120 print("Implied MP Stationary Distribution")
121 print("-----")
122 implied_mrp.display_stationary_distribution() # print stationary distribution from MRP
123 print()
124
125 print("Implied MRP Reward Function")
126 print("-----")
127 implied_mrp.display_reward_function() # print reward function from MRP
128 print()
129
130 print("Implied MRP Value Function")
131 print("-----")
132 implied_mrp.display_value_function(gamma=user_gamma) # print value function from MRP
133 print()
134
135 from rl.dynamic_programming import evaluate_mrp_result
136 from rl.dynamic_programming import policy_iteration_result
137 from rl.dynamic_programming import value_iteration_result
138
139 print("Implied MRP Policy Evaluation Value Function")
140 print("-----")
141 pprint(evaluate_mrp_result(implied_mrp, gamma=user_gamma)) # print policy evaluation value function
142 print()
143
144 print("MDP Policy Iteration Optimal Value Function and Optimal Policy")
145 print("-----")
146 opt_vf_pi, opt_policy_pi = policy_iteration_result(
147     si_mdp,
148     gamma=user_gamma
149 )
150 pprint(opt_vf_pi)
151 print(opt_policy_pi) # print optimal value function and optimal policy from policy iteration
152 print()
153
154 print("MDP Value Iteration Optimal Value Function and Optimal Policy")
155 print("-----")
156 opt_vf_vi, opt_policy_vi = value_iteration_result(si_mdp, gamma=user_gamma)
157 pprint(opt_vf_vi)
158 print(opt_policy_vi) # print optimal value function and optimal policy from value iteration
159 print()

```

Code 3: rl/chapter3/simple\_inventory\_mdp\_cap.py

Figure 1 shows the value functions and policies from policy iteration and value iteration, respectively.



#### MDP Policy Iteration Optimal Value Function and Optimal Policy

-----

```
{NonTerminal(state=InventoryState(on_hand=0, on_order=0)): -43.59563313047815,  
  NonTerminal(state=InventoryState(on_hand=0, on_order=1)): -37.97111179441265,  
  NonTerminal(state=InventoryState(on_hand=0, on_order=2)): -37.3284904356655,  
  NonTerminal(state=InventoryState(on_hand=1, on_order=0)): -38.97111179441265,  
  NonTerminal(state=InventoryState(on_hand=1, on_order=1)): -38.3284904356655,  
  NonTerminal(state=InventoryState(on_hand=2, on_order=0)): -39.3284904356655}
```

For State InventoryState(on\_hand=0, on\_order=0): Do Action 2

For State InventoryState(on\_hand=0, on\_order=1): Do Action 1

For State InventoryState(on\_hand=0, on\_order=2): Do Action 0

For State InventoryState(on\_hand=1, on\_order=0): Do Action 1

For State InventoryState(on\_hand=1, on\_order=1): Do Action 0

For State InventoryState(on\_hand=2, on\_order=0): Do Action 0

#### MDP Value Iteration Optimal Value Function and Optimal Policy

-----

```
{NonTerminal(state=InventoryState(on_hand=0, on_order=1)): -37.97111179441265,  
  NonTerminal(state=InventoryState(on_hand=0, on_order=2)): -37.3284904356655,  
  NonTerminal(state=InventoryState(on_hand=1, on_order=0)): -38.97111179441265,  
  NonTerminal(state=InventoryState(on_hand=1, on_order=1)): -38.3284904356655,  
  NonTerminal(state=InventoryState(on_hand=2, on_order=0)): -39.3284904356655,  
  NonTerminal(state=InventoryState(on_hand=0, on_order=0)): -43.59563313047815}
```

For State InventoryState(on\_hand=0, on\_order=0): Do Action 2

For State InventoryState(on\_hand=0, on\_order=1): Do Action 1

For State InventoryState(on\_hand=0, on\_order=2): Do Action 0

For State InventoryState(on\_hand=1, on\_order=0): Do Action 1

For State InventoryState(on\_hand=1, on\_order=1): Do Action 0

For State InventoryState(on\_hand=2, on\_order=0): Do Action 0

Figure 1: Optimal value function and policy from policy iteration and value iteration