

**NANYANG  
TECHNOLOGICAL  
UNIVERSITY**  
SINGAPORE

# CZ3005 Tutorial 3

*Yu Han*

*Nanyang Assistant Professor*

*[han.yu@ntu.edu.sg](mailto:han.yu@ntu.edu.sg)*

*N4-02c-109*



# Question 3.1

Assume that you are a manager of a warehouse (with a maximum capacity of  $W$  items). Each month  $t$ , you know the current inventory  $i$  (how many items left) in your warehouse. You might have a guess of the external demand in the next month  $(t + 1)$  with a distribution  $p$  (the probability that the external demand are  $d$  items is  $p(D_t = d)$ ,  $d = 0, 1, 2, \dots$ ). Based on this information, you decide to order additional items from a supplier. The cost might come from the storing cost of items in warehouse. Your objective is to maximize the profit. Use your own parameters for **fixed costs to buy** and store for each item and a **fixed selling price**.

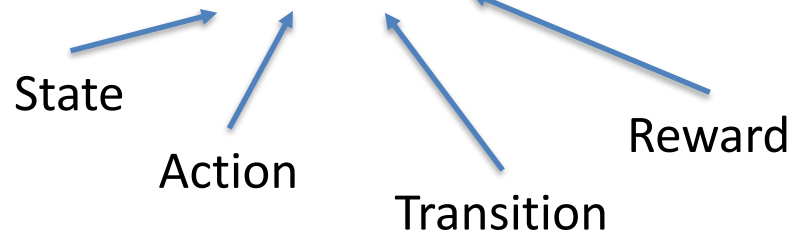
Please **write an MDP formulation** for the above problem.

**Hint:** Decision epochs are made at the beginning of each month, hence all events (more items arrive, fill external orders) would make states change. Actions are the amount of an order.



# Question 3.1

One form of MDP formulation  $\{S, A, T, R\}$



State space is  $S = \{0, 1, 2, \dots, W\}$   
Action space is  $A = \{0, 1, 2, \dots, W\}$

The reward term  $R(s_t, a_t)$  consists of three components:

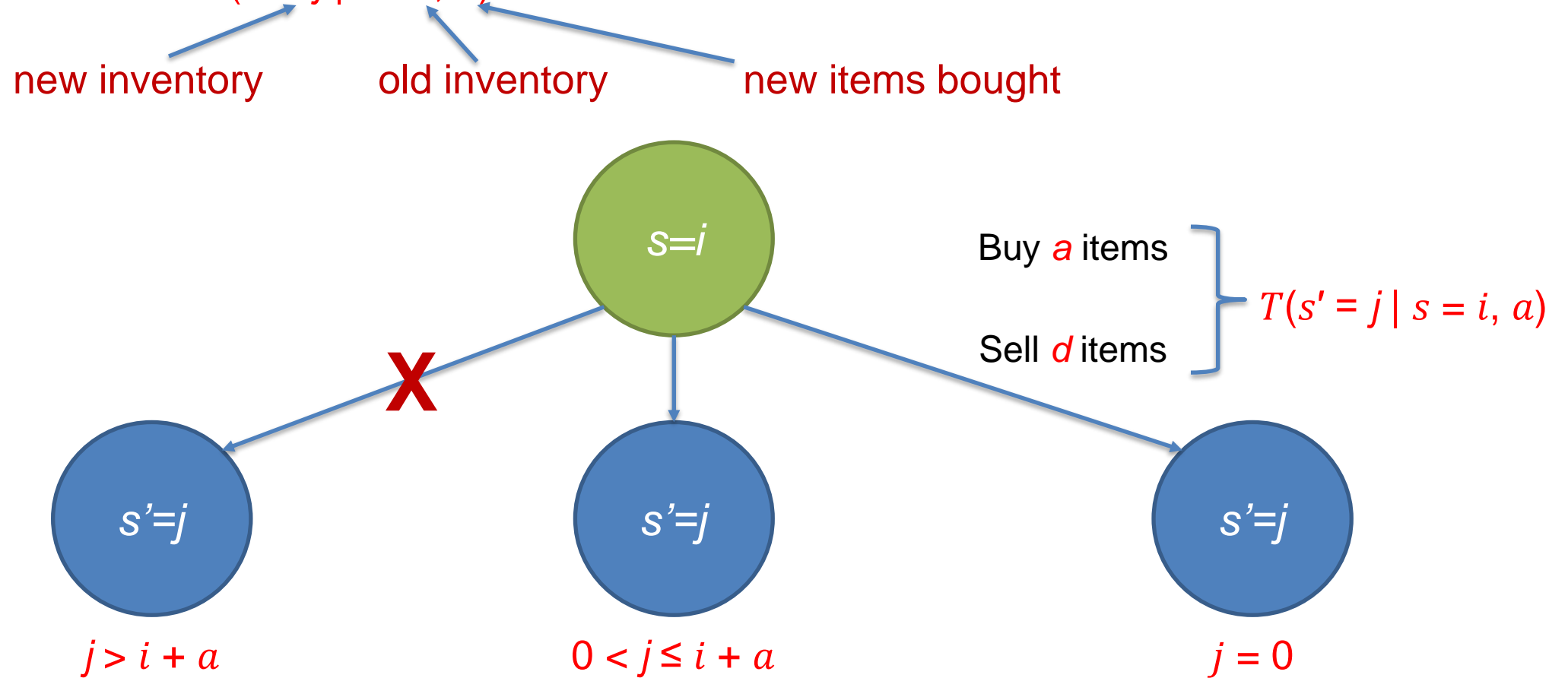
- Cost of **buying**  $a_t$  items are  $buy(a_t)$
- Cost of **storing**  $(s_t + a_t)$ . This cost is fixed and denoted as:  $Store(s_t + a_t)$ .
- Assume the **selling** price of  $D_t$  items is  $f(D_t)$ . The total sale price is:

$$Sell(s_t + a_t) = \sum_{d=0}^{s_t + a_t} p(D_t = d)f(d)$$

In summary, the reward function is:  $R(s_t, a_t) = Sell(s_t + a_t) - buy(a_t) - Store(s_t + a_t)$

# Question 3.1

The transition function  $T(s' = j \mid s = i, a)$  has three cases:





# Question 3.1

The transition function  $T(s' = j \mid s = i, a)$  has three cases:

new inventory

old inventory

new items bought

- If  $j > i + a$ , that means even after we sold  $d$  items, the remaining inventory is still larger than what we had before,  $i$ , plus what we bought,  $a$ , i.e.  $(i+a)$ . This is physically impossible. Thus, we will have 0 probability of transiting to this state:

$$T(j \mid i, a) = 0$$

- If  $0 < j \leq i + a$ , that means after we made the sales, the remaining inventory is more than 0, but fewer than  $(i+a)$ . This means the external demand  $d$  in this time step is  $(i + a - j)$ :

$$T(j \mid i, a) = p(D_t = i + a - j)$$

- If  $j = 0$ , that means we sold out entire inventory. Hence, the external demand  $d$  is equal to or more than  $(i+a)$ :

$$T(j \mid i, a) = p(D_t \geq i + a) = \sum_{d=i+a}^{\infty} p(D_t = d)$$



## Question 3.2

This Gridworld MDP operates like to the one we saw in class. The states are grid squares, identified by their row and column number (row first). The agent always starts in state  $(1,1)$ , marked with the letter **S**. There are **two terminal goal states**,  $(2,3)$  with reward **+5** and  $(1,3)$  with reward **-5**. Rewards are **0** in non-terminal states (The reward for a state is received as the agent moves into the state). The transition function is such that the **intended** agent movement (North, South, West, or East) happens with probability **0.8**. With probability **0.1** each, the agent ends up in one of the states perpendicular to the intended direction. If a **collision** with a wall happens, the agent stays in the same state.

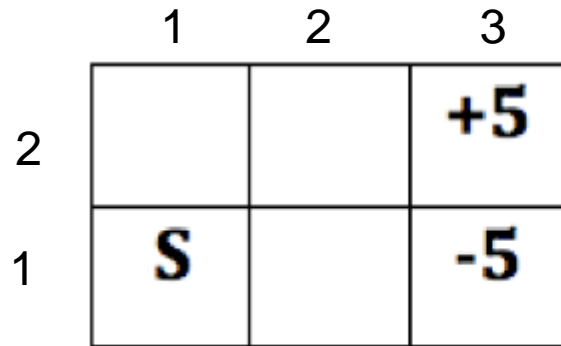


Figure (a) Gridworld MDP.

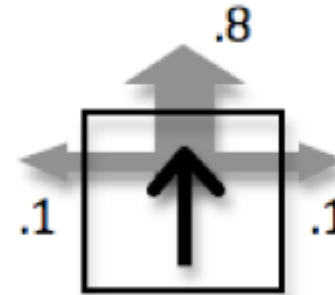


Figure (b) Transition function.

## Question 3.2(a)

Value Grid:

|   | 1 | 2 | 3 |
|---|---|---|---|
| 2 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 |

Suppose the agent knows the transition probabilities. Give the first two rounds of value iteration updates for each state, with a discount of **0.9**. (Assume  $V_0$  is 0 everywhere and compute  $V_i$  for times  $i = 1, 2$ ). (Assume values of termination states  $((1, 3)$  and  $(2, 3))$  are always 0)



# Question 3.2(a)

Value Grid:

|   |   |   |   |
|---|---|---|---|
|   | 1 | 2 | 3 |
| 2 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 |

Apply the Bellman backups:  $V_{i+1}(s) = \max_a \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V_i(s')]$

Reward Grid:

|   |          |            |                  |
|---|----------|------------|------------------|
|   | 1        | 2          | 3                |
| 2 |          | 0.1<br>0.1 | 0.8<br><b>+5</b> |
| 1 | <b>S</b> |            | <b>-5</b>        |

$$\begin{aligned}
 V_1(2,2) &= 0.8 \times (5 + 0.9 \times 0) && \text{[moved to (2,3)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[moved to (1,2)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[hit the wall and remained in (2,2)]} \\
 &= 4 + 0 + 0 = \mathbf{4} && \text{max} \\
 V_1(2,2) &= 0.8 \times (0 + 0.9 \times 0) && \text{[moved to (1,2)]} \\
 &+ 0.1 \times (5 + 0.9 \times 0) && \text{[moved to (2,3)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[moved to (2,1)]} \\
 &= 0 + 0.5 + 0 = \mathbf{0.5} \\
 V_1(2,2) &= 0.8 \times (0 + 0.9 \times 0) && \text{[moved to (2,1)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[moved to (1,2)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[hit the wall and remained in (2,2)]} \\
 &= 0 + 0 + 0 = \mathbf{0}
 \end{aligned}$$



# Question 3.2(a)

Value Grid:

|   |   |   |   |
|---|---|---|---|
|   | 1 | 2 | 3 |
| 2 | 0 | 4 | 0 |
| 1 | 0 | 0 | 0 |

Reward Grid:

|   |          |            |                  |
|---|----------|------------|------------------|
|   | 1        | 2          | 3                |
| 2 |          |            | <b>+5</b>        |
| 1 | <b>S</b> | 0.8<br>0.1 | <b>-5</b><br>0.1 |

$$\begin{aligned} V_2(1,2) &= 0.8 \times (0 + 0.9 \times 4) && [\text{moved to } (2,3)] \\ &+ 0.1 \times (-5 + 0.9 \times 0) && [\text{moved to } (1,3)] \\ &+ 0.1 \times (0 + 0.9 \times 0) && [\text{moved to } (1,1)] \\ &= 2.88 - 0.5 + 0 \\ &= \mathbf{2.38} \end{aligned}$$

# Question 3.2(a)

Value Grid:

|   | 1 | 2    | 3 |
|---|---|------|---|
| 2 | 0 | 4    | 0 |
| 1 | 0 | 2.38 | 0 |

Reward Grid:

|   | 1                 | 2 | 3         |
|---|-------------------|---|-----------|
| 2 | 0.1<br>0.8<br>0.1 |   | <b>+5</b> |
| 1 | <b>S</b>          |   | <b>-5</b> |

$$\begin{aligned}
 V_2(2,1) &= 0.8 \times (0 + 0.9 \times 4) && \text{[moved to (2,2)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[moved to (1,1)]} \\
 &+ 0.1 \times (0 + 0.9 \times 0) && \text{[hit the wall and remained in (2,1)]} \\
 &= 2.88 + 0 + 0 \\
 &= \mathbf{2.88}
 \end{aligned}$$

# Question 3.2(a)

Value Grid:

|   |      |      |   |   |
|---|------|------|---|---|
|   |      | 1    | 2 | 3 |
| 2 | 2.88 | 4    | 0 |   |
| 1 | 0    | 2.38 | 0 |   |

Optimal Route:

|   |          |   |           |
|---|----------|---|-----------|
|   | 1        | 2 | 3         |
| 2 |          |   | <b>+5</b> |
| 1 | <b>S</b> |   | <b>-5</b> |

|            |       |       |       |       |       |       |
|------------|-------|-------|-------|-------|-------|-------|
| S =        | (1,1) | (1,2) | (1,3) | (2,1) | (2,2) | (2,3) |
| $V_0(S) =$ | 0     | 0     | 0     | 0     | 0     | 0     |
| $V_1(S) =$ | 0     | 0     | 0     | 0     | 4     | 0     |
| $V_2(S) =$ | 0     | 2.38  | 0     | 2.88  | 4     | 0     |

## Question 3.2(b)

---

Suppose the agent does not know the transition probabilities. What does it need to be able to do (or have available) in order to learn the optimal policy?

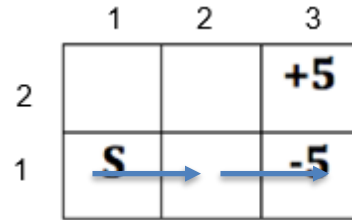
**Ans:** The agent must be able to explore the world by taking actions and observing the effects.



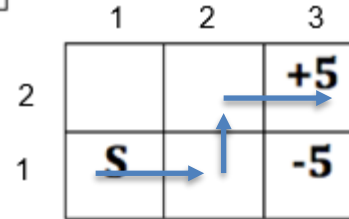
## Question 3.2(c)

The agent starts with the policy that always chooses to **go right**, and executes the following three trials:

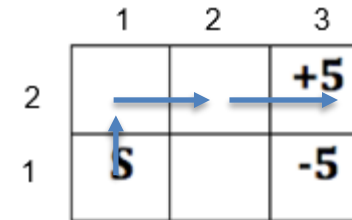
1)  $(1,1) - (1,2) - (1,3)$ ,



2)  $(1,1) - (1,2) - (2,2) - (2,3)$ , and



3)  $(1,1) - (2,1) - (2,2) - (2,3)$ .



What are the Monte Carlo estimates for states  $(1,1)$  and  $(2,2)$ , given these traces (assuming that the discount factor is 1)?



## Question 3.2(c)

**Ans:** To compute the estimates, average the rewards received in the trajectories that went through the indicates states:

$$V(1,1)=(-5+5+5)/3=5/3= 1.666,$$

$$V(1,2)=(-5+5)/2=0/2= 0,$$

$$V(2,1)=5/1= 5$$

$$V(2,2)=(5+5)/2= 5$$

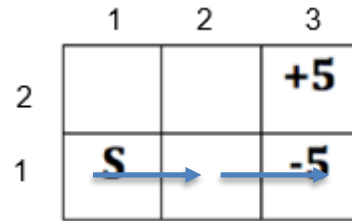
|   | 1                 | 2 | 3         |
|---|-------------------|---|-----------|
| 2 | 5                 | 5 | <b>+5</b> |
| 1 | <b>S</b><br>1.666 | 0 | <b>-5</b> |

# Question 3.2(d)

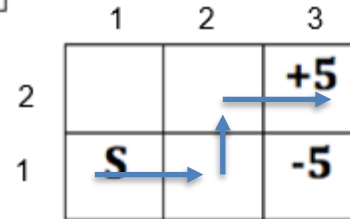
Using a learning rate of **0.1** and assuming initial values of **0**, what updates does the Q-learning agent make after trials 1 and 2, above?

Trials:

1) (1,1) – (1,2) – (1,3),



2) (1,1) – (1,2) – (2,2) – (2,3)



$$Q_{new}(s, a) = Q_{old}(s, a) + \alpha [r + \gamma \max_{a'} Q_{old}(s, a') - Q_{old}(s, a)]$$

Learning Rate
Discount Factor

## Question 3.2(d)

How fast to forget

Remembering Past Experience at  $s$

How fast to learn

Learning from New Experience  $a'$  at the same  $s$

Discount Factor

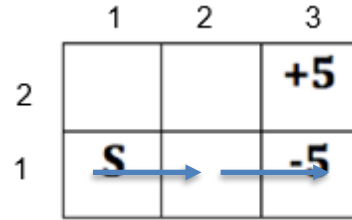
$$Q_{new}(s, a) = (1 - \alpha) Q_{old}(s, a) + \alpha [r + \gamma \max_{a'} Q_{old}(s, a')]$$

↓

$$Q_{new}(s, a) = Q_{old}(s, a) + \alpha [r + \gamma \max_{a'} Q_{old}(s, a') - Q_{old}(s, a)]$$

## Question 3.2(d)

1)  $(1,1) - (1,2) - (1,3)$



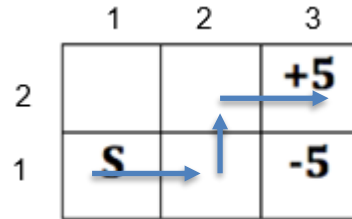
After trial 1, all of the updates will be zero, expect for:

$$Q((1,2), right) = 0 + .1 (-5 + 0.9 \times 0 - 0) = -0.5$$

|       | (1,1) | (1,2) | (1,3) | (2,1) | (2,2) | (2,3) |
|-------|-------|-------|-------|-------|-------|-------|
| Left  | 0     | 0     | 0     | 0     | 0     | 0     |
| Right | 0     | -0.5  | 0     | 0     | 0     | 0     |
| Up    | 0     | 0     | 0     | 0     | 0     | 0     |
| Down  | 0     | 0     | 0     | 0     | 0     | 0     |

## Question 3.2(d)

2) (1,1) – (1,2) – (2,2) – (2,3)



After trial 2, all of the updates will be zero, expect for:

Remember the probability of 0.1 for moving sideways?

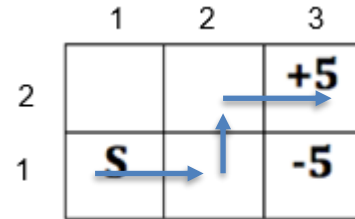
$$Q((1,2), right) = -0.5 + .1(0 + 0.9 \times 0 - (-0.5)) = -0.45$$

|       | (1,1) | (1,2) | (1,3) | (2,1) | (2,2) | (2,3) |
|-------|-------|-------|-------|-------|-------|-------|
| Left  | 0     | 0     | 0     | 0     | 0     | 0     |
| Right | 0     | -0.45 | 0     | 0     | 0     | 0     |
| Up    | 0     | 0     | 0     | 0     | 0     | 0     |
| Down  | 0     | 0     | 0     | 0     | 0     | 0     |



## Question 3.2(d)

2)  $(1,1) - (1,2) - (2,2) - (2,3)$



After trial 2, all of the updates will be zero, expect for:

$$Q((2,2), right) = 0 + .1(5 + 0.9 \times 0 - 0) = 0.5$$

|       | (1,1) | (1,2) | (1,3) | (2,1) | (2,2) | (2,3) |
|-------|-------|-------|-------|-------|-------|-------|
| Left  | 0     | 0     | 0     | 0     | 0     | 0     |
| Right | 0     | -0.45 | 0     | 0     | 0.5   | 0     |
| Up    | 0     | 0     | 0     | 0     | 0     | 0     |
| Down  | 0     | 0     | 0     | 0     | 0     | 0     |

# Thank you!

---

