

Carissa Ying Geok Teng (A0205190R/E0425113)

1a)

State

On(Start), Path(s_0, s_1) ...

Action

Travel(from, to)

Precond: On(from) \wedge \sim Visited(to) \wedge Path(from, to)

Effect: \sim On(from) \wedge On(to) \wedge Visited(to)

Transition

$P(s_1 | s_0, \text{Travel}(s_0, s_1))$...

Reward

$R(s) = -0.4$

$R(\text{Visited}(s_0) \wedge \text{Visited}(s_1) \wedge \dots) = 1$

1b)

State

ItemsInStorage(n), NumOfOrders(o), \sim Ordered, Backordered(u)

Action

Order(n, m): Order 1 time such that inventory = M

Precond: ItemsInStorage(n) \wedge (n + m = M) \wedge \sim Ordered

Effect: \sim ItemsInStorage(n) \wedge ItemsInStorage(n + m) \wedge Ordered

Backorder(n, u, o): Backorder up to B units

Precond: ItemsInStorage(n) \wedge Backordered(u) \wedge (n + o < N) \wedge (u + o < B)

Effect: \sim ItemsInStorage(n) \wedge \sim Backordered(u) \wedge ItemsInStorage(n + o) \wedge Backordered(u + o)

FulfillOrder(n, o, f)

Precond: ItemsInStorage(n) \wedge NumOfOrders(o) \wedge (n - f \geq 0) \wedge (o - f \geq 0)

Effect: \sim ItemsInStorage(n) \wedge \sim NumOfOrders(o) \wedge ItemsInStorage(n - f) \wedge NumOfOrders(o - f)

NextDay(n, o, u):

Precond: ItemsInStorage(n) \wedge NumOfOrders(o) \wedge Backordered(u) \wedge (o == 0)

Effect: \sim Ordered \wedge \sim Backordered(u)

Transition

P(New inventory and orders | Old inventory and orders, Action to order more or fulfil orders)

Reward

$$R(s, \text{Order}(n, m)) = -c$$

$$R(s, \text{Backorder}(n, u, o)) = -b * o$$

$$R(s, \text{FulfillOrder}(n, o, f)) = f$$

$$R(s, \text{NextDay}(n, o, u)) = -n$$

1c)

State

Screen display where each screen pixel with one of its values from 0-127

Action

One of the 18 actions

Transition

P(display after an action is taken | current screen display, one of the 18 actions)

Reward

$$\Delta \text{Score} - \Delta \text{Time}$$

2a)

Policy

$$\pi^*(s_1) = a_2$$

$$\pi^*(s_2) = a_1$$

Value Function

$$U^*(s_1) = P(s_1 | s_1, a_2)R(s_1 | a_2) + P(s_2 | s_1, a_2)R(s_2 | a_2) = 0.1(0) + 0.9(3) = 2.7$$

$$U^*(s_2) = P(s_1 | s_2, a_1)R(s_1 | a_1) + P(s_2 | s_2, a_1)R(s_2 | a_1) = 0(1) + 1(3) = 3$$

2b)

Policy

$$\pi^*(s_1, t_1) = a_2$$

$$\pi^*(s_2, t_1) = a_1$$

$$\pi^*(s_1, t_2) = a_2$$

$$\pi^*(s_2, t_2) = a_1$$

Value Function

$$\begin{aligned}U^*(s_1) &= P(s_1 | s_1, a_2)R(s_1 | a_2) + P(s_2 | s_1, a_2)R(s_2 | a_2) \\&\quad + P(s_1 | s_1, a_2) P(s_1 | s_1, a_2)R(s_1 | a_2) + P(s_2 | s_1, a_2) P(s_1 | s_1, a_2)R(s_2 | a_2) \text{ (s}_1 \text{ on first move)} \\&\quad + P(s_1 | s_2, a_1)P(s_2 | s_1, a_2)R(s_1 | a_2) + P(s_2 | s_2, a_1)P(s_2 | s_1, a_2)R(s_2 | a_2) \text{ (s}_2 \text{ on first move)} \\&= 0.1(0) + 0.9(3) + 0.1(0.1)(0) + 0.9(0.1)(3) + 0(0.9)(0) + 1(0.9)(3) = 5.67 \\U^*(s_2) &= 3 + 3 = 6\end{aligned}$$

2c)

Policy

$$\pi^*(s_1) = a_2$$

$$\pi^*(s_2) = a_1$$

Value Function

$$U^*(s_2) = \frac{3}{1-\gamma} = 3 / 0.1 = 30$$

$$U^*(s_1) = \gamma P(s_1 | s_1, a_2)U^*(s_1) + \gamma P(s_2 | s_2, a_2)U^*(s_2)$$

$$= 0.9(0.1)U^*(s_1) + 0.9(0.9)(30)$$

$$0.91U^*(s_1) = 24.3$$

$$U^*(s_1) = 26.703$$