

1.2

$$V_{\text{opt}}^{(t)}(s) = \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [\text{Reward}(s, a, s') + \gamma V_{\text{opt}}^{(t-1)}(s')]$$

$$V_{\text{opt}}^{(0)}(-2) = 0; \quad V_{\text{opt}}^{(0)}(-1) = 0; \quad V_{\text{opt}}^{(0)}(0) = 0;$$

$$V_{\text{opt}}^{(0)}(1) = 0; \quad V_{\text{opt}}^{(0)}(2) = 0.$$

Iteration 1:

$$a = -1 \Rightarrow 80\% (20+0) + 20\% (-5+0) = 15$$

$$a = 1 \Rightarrow 70\% (20+0) + 30\% (-5+0) = 12.5$$

$$a = -1 \Rightarrow 80\% (-5+0) + 20\% (-5+0) = -5$$

$$a = 1 \Rightarrow 70\% (-5+0) + 30\% (-5+0) = -5$$

$$a = -1 \Rightarrow 80\% (-5+0) + 20\% (100+0) = 16$$

$$a = 1 \Rightarrow 70\% (-5+0) + 30\% (100+0) = 26.5$$

$$V_{\text{opt}}^{(1)}(-2) = 0; \quad V_{\text{opt}}^{(1)}(-1) = 15; \quad V_{\text{opt}}^{(1)}(0) = -5;$$

$$V_{\text{opt}}^{(1)}(1) = 26.5; \quad V_{\text{opt}}^{(1)}(2) = 0.$$

Iteration 2:

$$Q=1 \Rightarrow 70\% (20+0) + 30\% (-5+(-5)) = 11$$

$$Q=-1 \Rightarrow 80\% (20+0) + 20\% (-5+(-5)) = 14$$

$$Q=-1 \Rightarrow 80\% (-5+15) + 20\% (-5+26.5) = 12.3$$

$$Q=1 \Rightarrow 70\% (-5+15) + 30\% (-5+26.5) = 13.45$$

$$Q=-1 \Rightarrow 80\% (-5+(-5)) + 20\% (0+100) = 12$$

$$Q=1 \Rightarrow 70\% (-5+(-5)) + 30\% (0+100) = 23$$

$$V_{opt}^{(2)}(-2) = 0; V_{opt}^{(2)}(-1) = 14; V_{opt}^{(2)}(0) = 13.45;$$

$$V_{opt}^{(2)}(1) = 23; V_{opt}^{(2)}(2) = 0;$$

1. b.

The optimal policy $\pi_{opt} = \begin{cases} \text{po if isEnd(s)} \\ \max_{a \in \text{Actions}} Q_{opt}(s, a) \end{cases}$

Resulting policy for non terminal states:

$$\pi_{opt}(-1) = -1; \pi_{opt}(0) = 1; \pi_{opt}(1) = 1.$$

2. a

We have define a new MDP with states $S' = S \cup \{o\}$, where o is a new state.

Let's assume that the probability reaching that state from any other state is $1-\lambda$. The probability of not reaching this o state being in any other state is λ . So we have to adjust our transition probabilities by λ . It means that our transition probabilities $T'(s, a, s') = \lambda T(s, a, s')$

We also need to adjust rewards to rewards_λ(s, a, s'), if we want to have

λ

the same optimal value. Thus having an MDP solver that can only handle with discount of γ , while solving for another discount factor $\lambda \neq \gamma$ we need to adjust probabilities and rewards like:

$$T(s, a, s') = \begin{cases} p(1-\lambda) & \text{if } s' = 0 \\ \lambda T(s, a, s') & \text{otherwise} \end{cases}$$

$$L'(s, a, s') = \begin{cases} 0 & \text{if } s' = 0 \\ \frac{1}{\lambda} L(s, a, s') & \text{otherwise} \end{cases}$$

46.

On the small test case Q-learning shows great results. We've got only 12 states different out of 21 possible. We can assume that it can learn properly when to take, peek, quit.

On the other hand the large test wasn't so successful. We have got 878 different out of 2705 states.

It couldn't be able to learn properly when to quit, take or peek because it is quite easy to be over 40 with small multiplicity and big amount of cards.

4d.

FixedLLAlgorithm gets its policy from original MDP with average reward of 6.48. But when we using Q-learning, average reward is 9.43. Thus it has better results because it can easily adapt to the new problem, adjusting parameters to new circumstances. While FixedLLAlgorithm will follow that fixed policy and couldn't adjust to the new problem.

5.Q.

According to our MDP simulation we've got different results. For 40 year time horizon economically preferable would be holding off on investment and rearrange cities budget on other spheres, as ratio of investment to wait states is very low (0.0015). For 100 year period it would be opposite situation as our simulation shows almost 50% ratio on investment policy due to much higher probability of sea level rising.

5.6.

As to 50 years economic agenda our recommendations would be to invest in infrastructure. Taking into consideration not only isolated predictions on sea level rising, which is highly likely to occur in longer term period but also vulnerability of future generations. Their welfare depends on environmental decisions we are making today. Like with human health better to preserve than to cure.

5.G

Dots strategies recommend to invest in infrastructure. With almost 40% ratio based on discounted model and 100% based on non-discounted. Discounted model which take into account devastating weather events shows us that our efforts will be less paid off due to certain destruction of the infrastructure, but we will receive a reward after each time step. On the other hand non-discounted model shows that our investments will succeed.

5. d

Our original MDP predicted economic benefits despite negative flooding cost. As it still is a model it can't predict everything. But it can adjusted with discount factor which will allow our model to be more precise in its predictions in different scenarios.