

Chapter-4

① $q_{\pi}(11, \text{down}) = -1 + V_{\pi}(T)$
 $= -1$

$q_{\pi}(7, \text{down}) = -1 + V_{\pi}(11)$
 $= -1 - 14 = -15$

② Let the states look like this for consistency in code.

For 2nd case:

0	1 $\sqrt{-14}$	2 $\sqrt{-20}$	3 $\sqrt{-22}$
4 $\sqrt{-14}$	5 $\sqrt{-18}$	6 $\sqrt{-20}$	7 $\sqrt{-20}$
8 $\sqrt{-20}$	9 $\sqrt{-20}$	10 $\sqrt{-18}$	11 $\sqrt{-14}$
12 $\sqrt{-22}$	13 $\sqrt{-20}$	14 $\sqrt{-14}$	15
	16 $\sqrt{-20}$		

$V_{\pi}(15) = -20$

For 1st case:

$V_{\pi}(15) = -20$

(Please refer to code)

Ex 4.3

$$q_{\pi}(s) = E_{\pi}[R_{t+1}]$$

$$q_{\pi}(s, a) = \sum_{s'} p(s'|s, a) \{ E[r|s', s, a] + \gamma V_{\pi}(s') \}$$

$$= \sum_{s'} p(s'|s, a) \{ E[r|s', s, a] + \gamma \sum_{a'} \pi(a'|s') q(s', a') \}$$

$$q_{k+1}(s, a) = \sum_{s'} p(s'|s, a) \{ E[r|s', s, a] + \gamma \sum_{a'} \pi(a'|s') q_k(s', a') \}$$

Ex 4.8

Because betting \$1 is better when you have \$51. As if you loose then you could bet \$50. rather than betting all \$51 at once.
(Note: Agent gets reward when it has total of \$100).

Ex 4.9

please refer to the code

Ex 4.10

$$q_{k+1} = \sum p(s'|s, a) [r + \gamma \max_{a'} q_k(s', a')]$$

$$q_{k+1} = \sum p(s'|s, a) [r + \gamma \max_{a'} q_k(s', a')]$$