

Cheat Sheet

Step breakdown – Q learning

$$Q(s, a) = w_1 f_1 + w_2 f_2$$

find the max of multiple Q values

$$Q(s, \text{west}) = 1(\text{weight}) \mathbf{(1/1)} + \mathbf{10(3)}$$

$$Q(s, \text{south}) = 1*(1/1) + 10*(1)$$

$$Q(s, \text{west})(\text{current}) = 31$$

$$Q(s, \text{south})(\text{current}) = 11$$

Q – learning still

$$Q(s, \text{west}) = 11$$

$$Q(s, \text{south}) = 11$$

$$\text{sample} = r + \Gamma * \max Q(s', a')$$

$$\text{sample} = 9 + 1 * 11 = 20$$

$$\text{difference} = \text{sample} - \text{current}$$

$$\text{difference} = 20 - 31 = -11$$

update weights

$$Q(s, a) = w_1 f_1 + w_2 f_2$$

$$w_1 = w_1 + \alpha * \text{difference} * f_1$$

$$w_2 = w_2 + \alpha * \text{difference} * f_2$$

$$w_1 = 1 + 0.5 * -11 * 1$$

$$w_2 = 10 + 0.5 * -11 * 3$$

----- value iteration

$$(\text{Max}) \sum T(s, a, s') * ((R(s, a, s') + \Gamma * V(s'))$$

$$0.8(0 + 1 * (-2)) + 0.2(0 + 1 * 0) = -1.6 // \text{Going Right}$$

$$0.8(0 + 1 * (0)) + 0.2(0 + 1 * (-2)) = -0.4 // \text{Going Down}$$

Direct Evaluation:

$$\text{Episode 1 : } (0 + 0 + 10) = 10$$

$$\text{Episode 2: } (0 - 10) = -10$$

$$\text{Episode 3: } (0+0-10) = -10$$

$$\text{Episode 4: } (0-10) = -10$$

$$(10-10-10-10) / 4 \text{ (episodes given)}$$

$$= -5$$

$$(1,2)$$

only episodes 2, 4

$$\text{Episode 2: } (0-10) = -10$$

$$\text{Episode 4: } (0-10) = -10$$

$$(-10-10)/2$$

$$= -10$$

$$\text{prob } ((1,2) | (1,1), \text{down})$$

$$-5/-10 = 1/2 = 0.5$$

temporal difference learning:

$$1 - \alpha * V(s) + \alpha * \text{reward} + \Gamma * V(s')$$

$$V(s) = \text{direct evaluation}$$

$$\text{gamma} + \alpha = \text{given}$$

Probability:

$$P(A) = \text{summation of all values that equals A}$$

$$P(B | A) = P(A) / P(A, B)$$

$$P(A, B | C) = P(C) / P(A, B, C)$$

$$P(C | A, B) = P(A, B, C) / P(A, B)$$