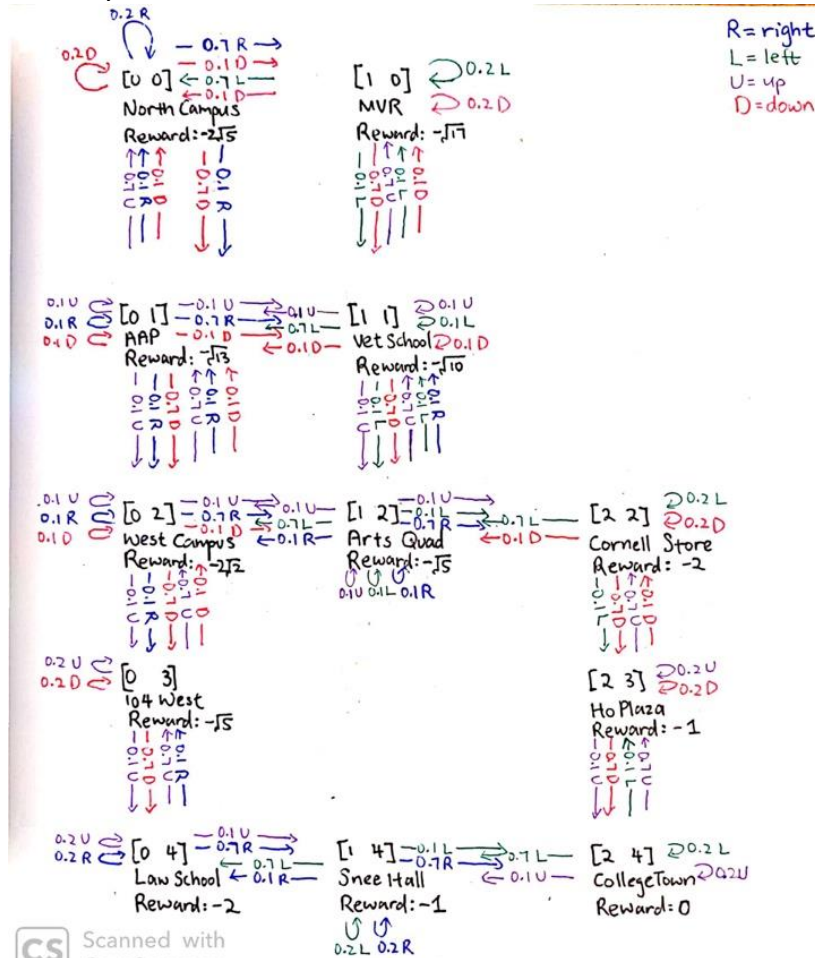


1. MPD and Utility Function



- a. Scanned with CamScanner
- b. Let $N = U^\pi(\text{North Campus})$, $V = U^\pi(\text{Vet School})$, $H = U^\pi(\text{Ho Plaza})$,
 $W = U^\pi(\text{104 West})$

$$\begin{aligned}
 N &= 0.7 \left(-\sqrt{13} + 0.5U^\pi([0\ 1]) \right) + 0.1 \left(-\sqrt{17} + 0.5U^\pi([1\ 0]) \right) \\
 &\quad + 0.2 \left(-2\sqrt{5} + 0.5N \right) \\
 V &= 0.7 \left(-\sqrt{5} + 0.5U^\pi([1\ 2]) \right) + 0.1 \left(-\sqrt{17} + 0.5U^\pi([1\ 0]) \right) + 0.1 \left(-\sqrt{13} \right. \\
 &\quad \left. + 0.5U^\pi([0\ 1]) \right) + 0.1 \left(-\sqrt{10} + 0.5U^\pi(V) \right) \\
 H &= 0.7(0 + 0.5U^\pi([2\ 4])) + 0.1(-2 + 0.5U^\pi([2\ 2])) + 0.2(-1 + 0.5H) \\
 W &= 0.7(-2 + 0.5U^\pi([0\ 4])) + 0.1(-2\sqrt{2} + 0.5U^\pi([0\ 2])) + 0.2(-\sqrt{5} \\
 &\quad + 0.5W)
 \end{aligned}$$

2. Discount Rewards

- a. 1) If the environment does not contain a terminal state or if the agent never reaches one, then all environment histories will be infinitely long, and utilities will be infinite, so you can't compare state sequences then.
- 2) Discounting favors near-term rewards, which is what we want to do, but you can't do that if you're just adding up all future rewards as is.

- b. Lower bound is when $R(s_t, a_t, s_{t+1}) = R_{\min}$.

$$\sum_{t=0}^{\infty} \gamma^t R_{\min.} = \frac{R_{\min.}}{1-\gamma}$$

Upper bound is when $R(s_t, a_t, s_{t+1}) = R_{\max}$.

$$\sum_{t=0}^{\infty} \gamma^t R_{\max.} = \frac{R_{\max}}{1-\gamma}$$

Anything between the two bounds is also finite.

- c. $\gamma^5(1) = 2\gamma^{10}$

$$\frac{1}{2} = \gamma^5$$

$$\gamma = 0.871$$

- d. You should **decrease** γ because now you have $0.871 > \gamma$

3. Policy Iteration

- a. There are 2 states and 3 actions in this environment. $2 \times 3 = 6$

- b. $s \times a$

- c. :

$$\hat{\pi}_0(low) = searching$$

$$\hat{\pi}_0(high) = searching$$

$$\hat{U}_0(low) = 0$$

$$\hat{U}_0(high) = 0$$

$$i = 1$$

$$\hat{U}_1(low) = 0.4(-10 + 0.5\hat{U}_0(high)) + 0.6(4 + 0.5\hat{U}_0(low)) = -1.6$$

$$\hat{U}_1(high) = 0.7(4 + 0.5\hat{U}_0(high)) + 0.3(4 + 0.5\hat{U}_0(low)) = 4$$

$$\begin{aligned} \hat{\pi}_1(low) &= \operatorname{argmax}_{a \in A}(\hat{U}_1(search), \hat{U}_1(wait), \hat{U}_1(recharge)) \\ &= \operatorname{argmax}_{a \in A}(-1.28, 0.2, 2) \end{aligned}$$

$$\hat{\pi}_1(high) = \operatorname{argmax}_{a \in A}(\hat{U}_1(search), \hat{U}_1(wait), \hat{U}_1(recharge)) = \operatorname{argmax}_{a \in A}(4, 3, 2)$$

$$i = 2$$

$$\hat{U}_2(low) = 1(0 + 0.5\hat{U}_1(high)) = 2$$

$$\hat{U}_2(high) = 0.7(4 + 0.5\hat{U}_1(high)) + 0.3(4 + 0.5\hat{U}_1(low)) = 5.16$$

$$\begin{aligned} \hat{\pi}_2(low) &= \operatorname{argmax}_{a \in A}(\hat{U}_2(search), \hat{U}_2(wait), \hat{U}_2(recharge)) \\ &= \operatorname{argmax}_{a \in A}(-0.6, 2, 2.58) \end{aligned}$$

$$\begin{aligned} \hat{\pi}_2(high) &= \operatorname{argmax}_{a \in A}(\hat{U}_2(search), \hat{U}_2(wait), \hat{U}_2(recharge)) \\ &= \operatorname{argmax}_{a \in A}(5.16, 3.58, 2.58) \end{aligned}$$

Converge!

4. Q-learning

- a. The Q-value is the expected reward from taking an action and the expected utility of the state we end up in from taking that action.

- b. If $\alpha = 1$, then $Q(s, a) = R(s, a, s') + \gamma \max_{a'} Q(s', a')$, i.e. the Q-value we learned for the state-action pair will be gone; the Q-function will also be deterministic.
 If $\alpha = 0$, then $Q(s, a)$ is not changing.
 If $\alpha > 1$, then we are overestimating the reward.
 If $\alpha < 0$, then reward will be negative
- c. It ensures we always take the best action, thus learning the best experience
- d. $Q([0\ 0], right) = 0.7(-\sqrt{17} + 0.5U^*([1\ 0])) + 0.1(-\sqrt{13} + 0.5U^*([0\ 1])) + 0.2(-2\sqrt{5} + 0.5U^*([0\ 0]))$

$$Q([0\ 0], down) = 0.7(-\sqrt{13} + 0.5U^*([0\ 1])) + 0.1(-\sqrt{17} + 0.5U^*([1\ 0])) + 0.2(-2\sqrt{5} + 0.5U^*([0\ 0]))$$