Last Time:

☐ decision-making

☐ MDPs

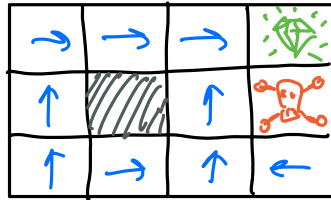Today:

☐ values & Q-values

☐ hands-on exercise

- Today: MDPs
- Thurs: Probability

- Next Tues: Paper Reading (Fundamentals)
- Next Thurs: —"—  (Math. Human Models)

<u>RECAP</u> : We left off having defined a Markov Decision Process:

- set of states $S$

- set of actions $A$

- transition function $P(s'|s,a)$ ← {or $T(s', a, s)$}

- reward function $r(s,a)$   // sometimes $r(s)$, $r(s,a,s')$, $r(s,a)$

- discount factor $\gamma \in [0,1]$

MDP quantities we have seen so far:

- policy $[\pi]$ : choice of action for each state



- cumulative reward : sum of (discounted) rewards

Our goal is to find a policy $\pi$ that maximizes the discounted sum of rewards:

$$\gamma^0 r(s^0, a^0) + \gamma^1 r(s', a') + \gamma^2 r(s^2, a^2) + \ldots$$

this trades off short & long-horizon reward

## VALUE FUNCTION

lets quantify the best cumulative reward we could ever get starting from a state $s$. This will precisely be our value function: $V : S \rightarrow \mathbb{R}$

Simpler Case: <u>Deterministic</u> Transition.

$$V(s) = \max_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^t, a^t) \right]$$

② try every possible policy and get the maximum possible value

① discounted sum of rewards

It turns out that we can rewrite this equation <u>recursively</u> because of Bellman's "principle of optimality." It says that you can decompose a complex problem into smaller sub-problems. Mathematically, it lets us redefine our value:
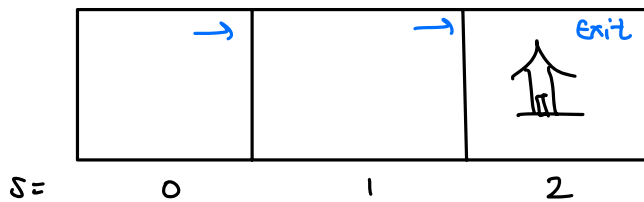
$$V(s) = \max_{a \in A} \left[ r(s,a) + \gamma \cdot V(s') \right]$$

reward right now!

s' is the next state you get to after taking a from state s

maximum possible discounted value of next state <u>plus</u> current reward.

<u>ex.</u>



$A = \{ \rightarrow, \text{Exit} \}$
you can take Exit @ home
(s=2)

deterministic $\rightarrow P(s'|s,a) = \begin{cases} 1 & \text{if } s'=s+a \\ 0 & \text{else} \end{cases}$

$r(s,a) = 0$ UNLESS you $r(s=2, a=\text{Exit}) = +1$

① $V(s=2) = r(s=2, a=\text{Exit}) = \boxed{+1}$

② $V(s=1) = r(s=1, \underset{\rightarrow = 0}{a = \rightarrow}) + \gamma \cdot V(s=2) = \boxed{\gamma}$

③ $V(s=0) = r(s=0, \underset{\rightarrow = 0}{a = \rightarrow}) + \gamma \cdot V(s=1) = \boxed{\gamma^2}$

## Stochastic Transitions

$$V(s) = \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V(s') \right]$$

<span style="color:green">when transitions are stochastic, we get <u>expected value</u> across next states we may visit.</span>

## Q-function

One more helpful quantity is the Q-function (Q value). You'll see this alot in papers ∵. Expected value of taking <span style="color:blue">some action $a$ from state $s$,</span> and then acting optimally from there on.

$\boxed{f(s,a) = s'}$

Determistic: $\quad Q(s,a) = r(s,a) + \gamma \cdot V(s')$ <span style="color:blue">← taking $a$</span>

Stochastic: $\quad Q(s,a) = r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V(s')$

Intuition: How good is action $a$ in the short & long-term?

Optimal policy: $\quad \pi(s) = \arg\max_{a \in A} Q(s,a)$

But, we were staring at our value function, we need to deal w/ the recursion!

$$V(s) = \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V(s') \right]$$

$$V(s') = \max_{a \in A} \left[ r(s',a) + \gamma \sum_{s' \in S} P(s''|s',a) V(s') \right]$$

recursive ("dynamic programming")

# VALUE ITERATION

$V[s] \leftarrow 0$  // current estimate of value (zero $\forall s \in S!$)

$V'[s] \leftarrow 0$  // next / updated estimate of value

$\delta \leftarrow 0$  // maximum change btwn. $V$ and $V'$

while True:

    $V \leftarrow V'$ and $\delta \leftarrow 0$  // we will update $V'$ so remember prior values

    for each state $s \in S$:

        $V'[s] \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V[s'] \right)$

        if $|V'[s] - V[s]| > \delta$

            $\delta \leftarrow |V'[s] - V[s]|$

  if $\delta \leq \varepsilon$
      return $V'[s]$

// stopping condition: $\varepsilon$ = max error
// if we have less error than $\varepsilon$, return!

# Exercise [from CS188 Berkeley]

Pacman is using MDPs to maximize his expected utility. In each environment:

- **deterministic transition**  ↑ ← ↓ →
- Pacman has the standard actions {North, East, South, West} unless blocked by an outer wall
- There is a reward of 1 point when eating the dot (for example, in the grid below, $R(C, South, F) = 1$)
  ↳ i.e. transition to state F
- The game ends when the dot is eaten
  ↳ i.e. zero reward at F state

**(a)** Consider a the following grid where there is a single food pellet in the bottom right corner $(F)$. The **discount** factor is 0.5. There is <u>no living reward</u>. The states are simply the grid locations.

i.e. you only get ↙ reward at the dot

| A | B | C |
|---|---|---|
| D | E | F ○ |

**(i)** What is the optimal policy for each state?

**(ii)** What is the optimal value for the state of being in the upper left corner $(A)$? Reminder: the discount factor is 0.5.

Run 4 iterations of value iteration by filling out the table:

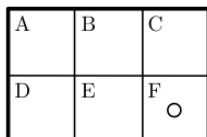| iter | V(A) | V(B) | r(C) | V(D) | V(E) | V(F) |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 |  |  |  |  |  |  |
| 2 |  |  |  |  |  |  |
| 3 |  |  |  |  |  |  |
| 4 |  |  |  |  |  |  |

# Solution:

Pacman is using MDPs to maximize his expected utility. In each environment:

- Pacman has the standard actions {North, East, South, West} unless blocked by an outer wall
- There is a reward of 1 point when eating the dot (for example, in the grid below, $R(C, South, F) = 1$)
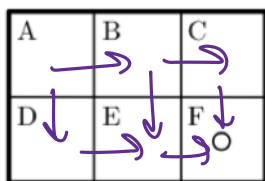- The game ends when the dot is eaten ⟶ i.e. transition to state F

(a) Consider a the following grid where there is a single food pellet in the bottom right corner ($F$). The **discount** factor is 0.5. There is no living reward. The states are simply the grid locations.

i.e. you only get ↙ reward at the dot

| A | B | C |
|---|---|---|
| D | E | F ○ |

(i) What is the optimal policy for each state?

| A | B | C |
|---|---|---|
| D | E | F ○ |

(ii) What is the optimal value for the state of being in the upper left corner ($A$)? Reminder: the discount factor is 0.5.

$V^*(A) = 0.25$

| k | V(A) | V(B) | V(C) | V(D) | V(E) | V(F) |
|---|------|------|------|------|------|------|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 | 0 | 1 | 0 |
| 2 | 0 | 0.5 | 1 | 0.5 | 1 | 0 |
| 3 | 0.25 | 0.5 | 1 | 0.5 | 1 | 0 |
| 4 | 0.25 | 0.5 | 1 | 0.5 | 1 | 0 |

$k=1:$ $V(E) = max\left[ r(E, →) + 0.5 V(F),\right.$
$r(E, ↑) + 0.5 V(B),$
$\left. r(E, ←) + 0.5 V(D)\right] = max[1, 0, 0] = \boxed{1}$

$k=2:$ $V(B) = max\left[ r(B, →) + 0.5 V(C),\right.$
$r(B, ↓) + 0.5 V(E),$
$\left. r(B, ←) + 0.5 V(A)\right] = max[0.5, 0.5, 0] = \boxed{0.5}$

$k=3:$ $V(A) = max\left[ r(A, →) + 0.5 V(B),\right.$
$\left. r(A, ↓) + 0.5 V(D)\right] = max[0.25, 0.25] = \boxed{0.25}$

If I give you $V^*(s)$ optimal value, then I can extract the optimal <u>policy</u> via a <u>ONE-STEP DECISION-MAKING PROBLEM!</u>

$$\pi^*(s) = \arg \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s'} P(s'|s,a) V^*(s') \right]$$