# Markov Decision Problems (MDP) [17.1]

☐ <u>State</u>
- → general representation
- → $S_0$ Initial State
- → IS-TERMINAL(S) - Terminal state test

☐ Actions(s) · state ⟶ {actions}
- → get legal actions from S

☐ <u>Transition Model</u>

· Deterministic: Result(s,a) → s'
  STATE × ACTION ⟶ STATE

· Stochastic: $P(s'|s,a)$ PROBABILITY DISTRIBUTION
  STATE × ACTION × STATE ⟶ p    $0 \leq p \leq 1$

<u>Markov Assumption</u>: Transitions only depend on the
current state & no other history.

☐ <u>Reward Function</u>
- → per-action reward

· $R(s,a,s')$ · STATE × ACTION × STATE ⟶ r    $-R_{MAX} \leq r \leq R_{MAX}$
Reward for taking action a in state s to reach state s'
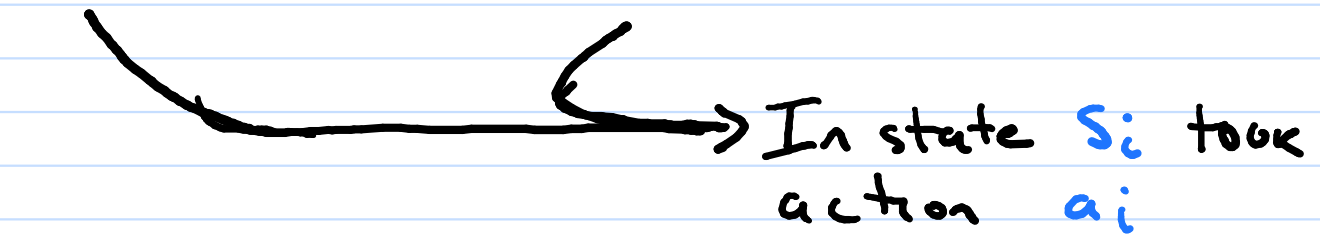
---

<u>Sequential Decision Problem</u>: Utility is determined
by a sequence of decisions.

# From Rewards To Utility [17.1.1]

**Goal:** $U(s)$    Utility value for state $s$.

Sequential Decision Utility: Sequence of state, action $\rightarrow$ u

$$U_h([\underbrace{s_0, a_0}_{\substack{\text{Next} \\ \text{Decision}}}, s_1, a_1, \dots, \underbrace{s_n, a_n}_{\substack{\text{Future} \\ \text{Decision}}}])$$

Utility of decision history.

$\longrightarrow$ In state $s_i$ took action $a_i$

## Additive Discount Rewards

Utility is the <u>sum of rewards</u>, <u>weighed for their age</u>, with more recent rewards carrying more weight then potential future rewards

□ Discount Factor $\gamma$   with   $0 \le \gamma \le 1$

$$U_h([\dots s_i, a_i \dots]) = \sum_{i=0}^{n} \underbrace{\gamma^i}_{\text{Discount}} \times \underbrace{R(s_i, a_i, s_{i+1})}_{\substack{\text{Reward for} \\ i^{th} \text{ Decision}}}$$

$\gamma = 1 \longrightarrow$ <u>No discount. Past $\equiv$ Future</u>

$\qquad 1^n = 1$ for all $n \in \mathbb{N} \longrightarrow \{0,1,2,3,\dots\}$

$\gamma = 0 \longrightarrow$ Ignore future. Only $s_0, a_0$ count.

$\qquad 0^0 = 1$ , $0^n = 0$ for $n > 0$.

# From $U_h$ to Utility of States $[17.1.1 - 17.1.2]$

Do the "first & rest" trick to isolate one decision.

$$U_h([\cdots s_i, a_i \cdots]) = \sum_{i=0}^{\infty} \gamma^i R(s_i, a_i, s_{i+1})$$

$$= \underbrace{R(s_0, a_0, s_1)}_{\substack{\text{current} \\ \text{reward}}} + \underbrace{\gamma}_{\substack{\uparrow \\ \text{Discount}}} \times \underbrace{U_h([s_1, a_1, \cdots s_n, a_n])}_{\text{reward-to-go}}$$

But how should we isolate the state from its associated action?

↳ State utility $\longrightarrow$ Utility of best action in state!

## Bellman Equation (With Deterministic Transitions)

$$U(s) = \underset{a \in \text{Action}(s)}{\text{MAX}} \left\{ \underbrace{R(s, a, \text{Result}(s,a))}_{\text{Reward for } s} + \gamma \underbrace{U(\text{Result}(s,a))}_{\text{reward-to-go}} \right\}$$

⚹ This is really a SYSTEM OF EQUATIONS:
- One equation per state
- Per-State equations refering to other states.

⚹ Methods exist that let us solve the system (see 17.2).
They generally work through iterative refinement.

We're intrested in LEARNING from experience.

# From Utility to Policy [17.1.2]

A policy $\pi(s)$ determines what action to take in state $s$.

Given $U(s)$,

$$\pi(s) = \underset{a \in \text{Action}(s)}{\text{argmax}} \left\{ U(\text{RESULT}(s,a)) \right\}$$

$\Rightarrow$ action that results in the state with the highest utility.

## $Q(s,a)$ : The Action-Utility Function.

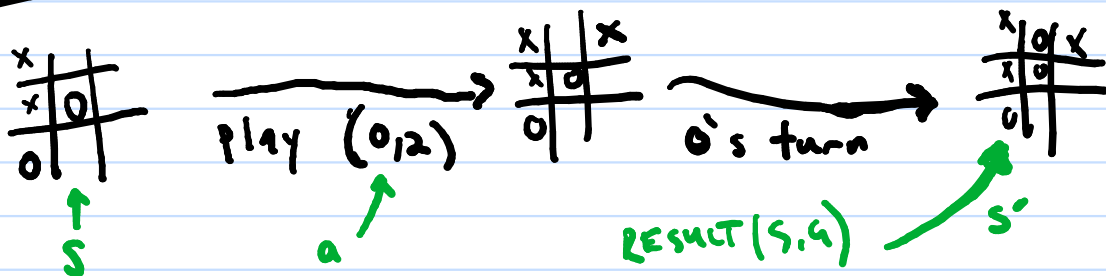Utility of ACTION $a$ in state $s$. More fine-grained than $U(s)$.

$$U(s) = \underset{a \in \text{Action}(s)}{\text{MAX}} \; Q(s,a) \qquad \begin{array}{c} \text{State Utility} \\ \Updownarrow \\ \text{BEST ACTION UTILITY} \end{array}$$

$$U(s) = \underset{a \in \text{Action}(s)}{\text{MAX}} \left\{ R(s,a,s') + \gamma \underset{a'}{\text{Max}} \, Q(s',a') \right\}$$

Can be observed. No RESULT NEEDED!

## Why?



$\Rightarrow$ No need to predict/model opponent. Just observe

# Temporal Difference Learning [22.3.3]

Say we have an estimate for the Action-Utility function $Q$. We then make action $a$ in state $s$ to result in state $s'$ and receive reward $R(s,a,s')$. We can calculate $Q(s,a)$ as.

This should be equal to our estimate for $Q(s,a)$. If not we should adjust our estimate based on the error.

$$Q(s,a) \longleftarrow Q(s,a) + \alpha \left[ \underbrace{R(s,a,s') + \gamma \max_{a'} Q(s',a')}_{\substack{\text{Calculated} \\ Q(s,a)}} - \underbrace{Q(s,a)}_{\substack{\text{expected} \\ Q(s,a)}} \right]$$
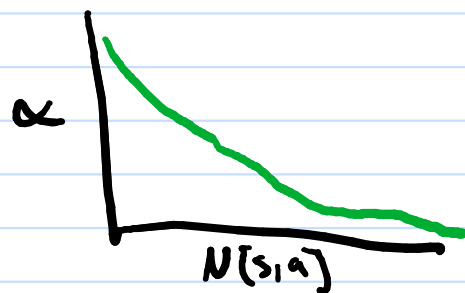
Learning rate parameter

Error of expected & observed values

Do this enough and our $Q$ values will converge to the actual values. The learning rate parameter $\alpha$ determines the size of the adjustment. Ideally, it's a function of time and decreases the more you update the value for $(s,a)$

$$\alpha \left( \underbrace{N(s,a)}_{\substack{\text{number of updates} \\ \text{to } Q(s,a)}} \right)$$

# ACTIVE LEARNING [17.3, 22.3-22.3.1]

But ... actions taken are based on the policy $\pi$, which is ideally based on the utility function $Q$.

$\Rightarrow$ While learning, we want to balance EXPLOITING moves we think are good (high utility) while EXPLORING moves of unknown quality. After sufficient exploration, we can stick to "good" moves.

The EXPLORATION FUNCTION $f$ determines our preference for Exploration vs exploitation.

$$f : \underbrace{\text{STATE} \times \text{ACTION}}_{\substack{\text{given state} \\ \text{action}}} \longrightarrow \underbrace{V}_{\substack{\text{optimistic} \\ \text{utility} \quad \text{or "RANK"}}}$$

○ TRY AT LEAST
$N_e$

$$f(s,a) = \begin{cases} R_{max} & N[s,a] < N_e \\ Q(s,a) & \text{otherwise} \end{cases}$$

□ UCB1 [5.4]

$$f(s,a) = Q(s,a) + C \times \sqrt{\frac{\log\left(\sum_{a'} N[s,a']\right)}{N[s,a]}} \longleftarrow \log(\text{Times in } s)$$

# TEMPORAL DIFFERENCE Q-LEARNING

□ Specify Game as MDP. [Transition Model Optional!]

□ Specify Rewards Function $R$

□ Specify learning rate function $\alpha$ $( \mathbb{N} \to r, \; 0 \leq r \leq 1 )$

□ Specify discount rate $\gamma$

□ Specify Exploration function $f$

## While Learning
### Maintain

$Q[s,a] \longrightarrow$ Action-Utility Estimates. Initially $\varnothing$.

$N[s,a] \longrightarrow$ Frequency Counts for state $s$, action $a$. Init. $\varnothing$.

→ **Update $Q$ values**

Given states $s$ and $s'$ and action $a$ such that taking action $a$ in $s$ resulted in $s'$.

$$Q[s,a] \mathrel{+}= \alpha(N[s,a])\left[ R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

→ **Choose Action**

In state $s$,

$a \longleftarrow \underset{a' \in Actions(s)}{argmax} \; f(s,a')$