# 11/9/2023: Reinforcement Learning

| Supervised Learning | Unsupervised Learning | Reinforcement Learning |
|---|---|---|

Dataset =

$\{ (x^{(1)}, y^{(1)}), \dots \}$

↑ input   ↑ desired output

Dataset =

$\{ x^{(1)}, x^{(2)}, \dots \}$

↑ No Supervision at all

① Algorithm creates dataset as it runs

② Learning signal = Rewards for taking actions (ie. did a good thing or bad thing happen)?

↙ more than unsupervised learning   ↘ weaker than supervised learning

Someone hands us a dataset.
Algorithm takes dataset as input, doesn't influence what's in the dataset
   "passive"

Example Students selecting classes
- Action: Take some classes, not others
- State: What you know / what prereqs you satisfy
- Reward: Enjoyment, job

Other examples: Robotics, video games

→ ① Defining the world   (No Learning)
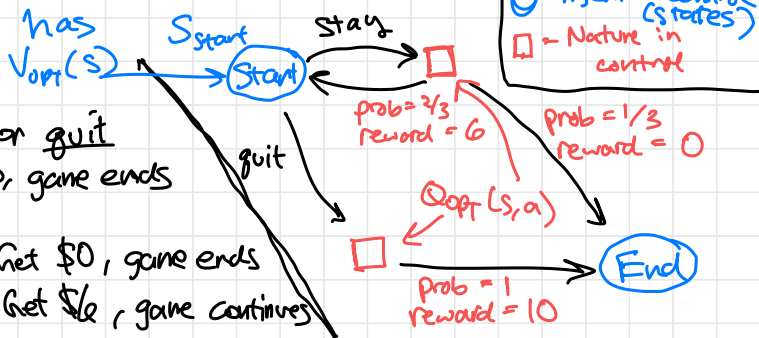→ ② Learning b/c the agent doesn't know how the world works

① Markov Decision Process (MDP)
   Formal description of a world with states, actions, rewards, etc...

Example MDP:
At each timestep:
- Agent can **stay** or **quit**
- If quit: receive $10, game ends
- If stay:
  - Probability $1/3$: Get $0, game ends
  - Probability $2/3$: Get $4, game continues

has $V_{opt}(s)$

$S_{start}$

stay

○ = Agent is control (states)
□ = Nature in control

Start

prob = $2/3$ reward = 4

prob = $1/3$ reward = 0

quit

$Q_{opt}(s,a)$

prob = 1 reward = 10

End

Formal ingredients of MDP:
- Set of states $S$ (e.g. possible configurations/ locations of robot)
- Starting state $S_{start}$ (OR distribution over states)
- Actions ($s$): Possible actions at state $s$
- $T(s, a, s')$: Probability of transitioning to state $s'$
  starting at state $s$ & taking action $a$
- Reward $(s, a, s')$: Reward received when transitioning from
  $s$ to $s'$ after taking action $a$

unknown In the example MDP:
in          $T(Start, stay, End) = 1/3$
RL          Reward$(Start, stay, End) = 0$
- Is End ($s$): Is $s$ an end state?
  Game ends when reaching an end state

_____

Given a __Known__ MDP, what should the agent do?

| Policy |  Strategy used by an agent
     Formally: mapping $\pi(s) \to a \in Actions(s)$
                            ↑                  ↖
                      Current state        chosen action

Why? Visiting $s$ multiple times has exact same
      transitions & rewards ⇒ best action is same

Value function: The value $V_\pi(s)$ for policy $\pi$
         and state $s$ is:  ___(discounted)
         the expected sum of rewards when
         Starting at $s$, use policy $\pi$

Discounting: Future rewards are less valuable
         - At each timestep, probability of survival $< 1$
      we introduce a  discount factor $\gamma \in [0, 1]$
              = probability of survival at each timestep  e.g. $\gamma = .99$

we care about discounted sum of rewards
For sequence $r_1, r_2, r_3$, discounted sum $= r_1 + \gamma r_2 + \gamma^2 r_3 + \cdots$

## Optimal Policies

$$V_{OPT}(s) = \text{maximum possible expected discounted sum of rewards starting at } s \text{ for any policy}$$

"optimal value"

$$Q_{opt}(s,a) = \text{maximum possible expected discounted sum of rewards starting at } s \text{ and forced to take action } a$$

"Q-value"

$$V_{OPT}(s) = \begin{cases} 0 & \text{if } \text{IsEnd}(s) \\ \max\limits_{a \in Actions(s)} Q_{OPT}(s,a) & \text{else} \end{cases}$$

$$Q_{OPT}(s,a) = \sum_{s' \in S} \underbrace{T(s,a,s')}_{\substack{\text{Prob of} \\ \text{going to } s'}} \left[ \underbrace{Reward(s,a,s')}_{\substack{\text{Reward now} \\ \text{(no discount)}}} + \underbrace{\gamma}_{\substack{\text{discount}}} \underbrace{V_{OPT}(s')}_{\substack{\gamma \\ \text{future} \\ \text{reward}}} \right]$$

Optimal policy: $\pi^{*}(s) = \underset{a \in Actions(s)}{argmax} \ Q_{OPT}(s,a)$

Lesson: If we can estimate $Q_{OPT}(s,a)$ well, ⭐ we can deduce the optimal policy

### Reinforcement Learning time
- Believe world is an MDP
- But $T(s,a,s')$, $Rewards(s,a,s')$ unknown to us

RL training pseudocode:
For episode = 1, 2, 3, ... :
  $s_1 \leftarrow s_{start}$   // or sample distribution over start state
  for $t = 1, ...$
    • agent chooses action $a_t$ · $\pi_{act}(s_t)$
      policy used to act by agent

    • Agent receives:

Reward $r_t$
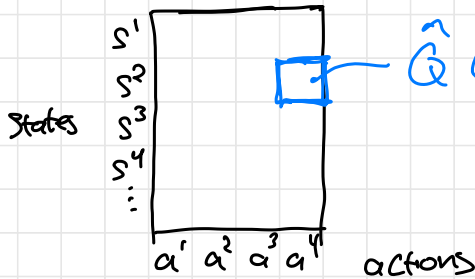
New state $s_{t+1}$ <span style="color:red">New data about the world</span>

- Update agent's parameters
  AKA Learning

- If IsEnd($s_{t+1}$): break

## Learning algorithm: Q-Learning

Goal: Learn $Q_{opt}(s,a)$ for every $(s,a)$

We will maintain best guess $\hat{Q}(s,a)$ for each $(s,a)$



States $s^1$ $s^2$ $s^3$ $s^4$ ...

$a^1$ $a^2$ $a^3$ $a^4$   actions

$\hat{Q}(s^2, a^4) = $ our estimate of
$Q_{opt}(s^2, a^4)$

each $\hat{Q}(s,a)$ is one parameter of model