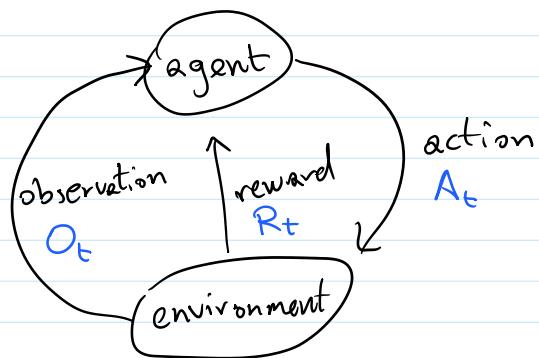


► Intro



- Markov property \rightarrow the next state is only dependent on the current state, not past states

$$P(S_{t+1} | S_t) = P(S_{t+1} | S_1, \dots, S_t)$$

- fully observable : $O_t = S_t$

- POMDPs \rightarrow partially observable states

$$S_t = \sigma(S_{t-1}W_s + O_tW_s)$$

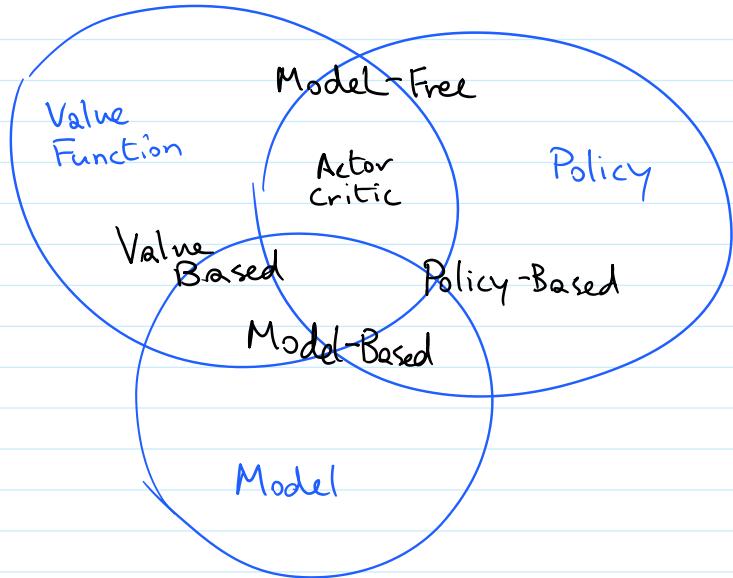
↳ use RNN

- Dynamics Model

$$M(s, a, s') = P(S_{t+1} = s' | S_t = s, A_t = a)$$

$$R(s, a) = E(R_{t+1} | S_t = s, A_t = a)$$

- RL Agent Taxonomy



- Fundamental Problems
 - Learning → learn unknown dynamics model
 - improve policy based on experience
 - Planning → given known dynamics, find optimal policy
- Exploration vs Exploitation → exploration may seem suboptimal but it might lead to finding even better states
 - Exploration → gain new info about environment
 - Exploitation → maximize reward w/ known info
- Prediction vs Control
 - Prediction → evaluate given policy
 - find true value of states
 - Control → find the best policy

► MDPs

- gain → total reward w/ discount

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$$

- Markov Reward Process (MRP) → no actions
 - value function

$$v(s) = E(A_t \mid S_t = s)$$

$$= E(R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s)$$

- Bellman Equation

$$v(s) = R_s + \gamma \sum_{s' \in S} M(s, s') v(s')$$

prob going from s to s'

↳ matrix form

$$\bar{v} = \bar{R} + \gamma \hat{M} \bar{v}$$

$$\bar{v} = (I - \gamma \hat{M})^{-1} \bar{R} \rightarrow \text{solution}$$

- Markov Decision Process (MDP)

- policy

$$\pi(a \mid s) = P(A_t = a \mid S_t = s)$$

↳ MRP given a policy

$$M^\pi(s, s') = \sum_{a \in A} \pi(a \mid s) M(s, a, s')$$

$$R^\pi(s) = \sum_{a \in A} \pi(a \mid s) R(s, a)$$

- Q-function → action-value

$$Q_\pi(s, a) = E_\pi(A_t \mid S_t = s, A_t = a)$$

- Bellman Expectation Equation

$$v_\pi(s) = \sum_{a \in A} \pi(a \mid s) Q_\pi(s, a)$$

$$Q_\pi(s, a) = R(s, a) + \gamma \sum_{s' \in S} M(s, a, s') v_\pi(s')$$

$s \in S$

- Bellman Optimality Equation

$$v^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \sum_{s' \in S} M(s, a, s') v^*(s')$$

- Additional Topics \rightarrow
 - \rightarrow MDPs
 - \rightarrow POMDPs
 - \rightarrow Ergodic MDPs

► Dynamic Programming

- DP \rightarrow the problem must have a sequential component
 \rightarrow breakdown problem into subproblems & combine solutions
- this is all assuming agent knows $M(s, a, s')$ & $R(s, a)$
 \hookrightarrow how to solve for $v^*(s)$ & $Q^*(s, a)$ when S is large
- Policy Evaluation

$$v_{k+1}(s) = \sum_{a \in A} \pi(a|s) \left(R(s, a) + \gamma \sum_{s' \in S} M(s, a, s') v_k(s') \right)$$

$$\bar{v}^{k+1} = \bar{R}^\pi + \gamma \hat{M}^\pi \bar{v}^k$$

- Improving Policy \rightarrow greedily

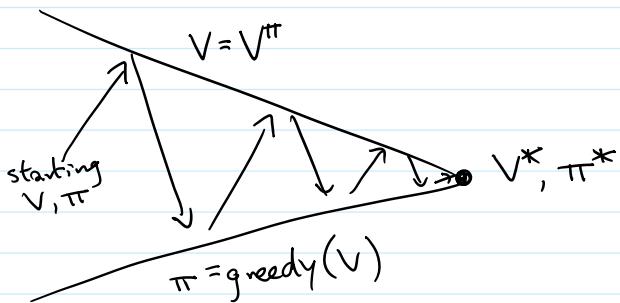
$$\pi'(s|a) = \operatorname{argmax}_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} M(s, a, s') v_\pi(s') \right]$$

- Policy Iteration

- algorithm :

- initialize π (arbitrarily)
- repeat until convergence ($\pi = \pi'$)
- estimate v^π by evaluating policy

improve policy $\pi' = \text{greedy}(v^\pi)$
 $\pi' \leftarrow \pi$



- Value Iteration \rightarrow skip policy eval

$$v_{k+1}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} M(s, a, s') v_k(s') \right]$$

- DP Algorithms

<u>Problem</u>	<u>Bellman Eqn.</u>	<u>Algorithm</u>
Prediction	Expectation	Policy Eval
Control	Expectation + Greedy Improvement	Policy Iteration
Control	Optimality	Value Iteration

- Async DP

- in-place DP \rightarrow analogue to SGD
- prioritized sweeping \rightarrow update based on Bellman error
- real-time DP \rightarrow update only those states agent visits
- Problem Size \rightarrow only works up to $|S| \sim 10^6$
- For larger problems use sample backups instead full-backups

► Prediction

- Monte Carlo \rightarrow update $V(s)$ after episode

$$V(S_t) \leftarrow V(S_t) + \alpha (C_t - V(S_t))$$

$$C_t = \sum_{\text{episode}} R_t$$

- Temporal Difference \rightarrow adjust $V(s)$ after each action \rightarrow bootstrapping

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

- MC vs TD

TD

bootstraps

online or episodic

learn from incomplete sequences

works in non-terminating systems

low variance

some bias

converges to Markov model solution

uses Markov property

MC

learns at the end of an episode

episodic only

high variance

no bias

Converges to min mean squared error of

$$\sum \sum (C_t^k - V(S_t^k))^2$$

change depends on many random actions

starting state doesn't matter

- $TD(\lambda) \rightarrow$ generalizing MC & TD

- use eligibility function to rank how much you can learn about other states based on your current experience

↳ using both a frequency & recency heuristic

$$E_o(s) \leftarrow 0$$

$$E_{t+1}(s) \leftarrow \gamma \lambda E_t(s) + 1(S_t = s)$$

- now adjust all states each step according to eligibility

$$\delta_t = R_t + \gamma V(S_{t+1}) - V(S_t)$$

TD error

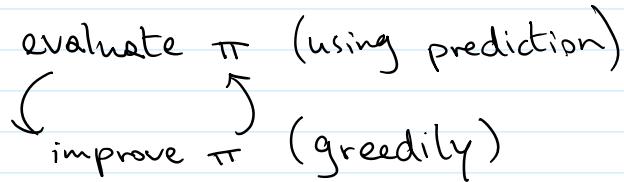
$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

- MC is equivalent to TD(1)
- TD is equivalent to TD(0)

Control

- now learn best policy

↳ naive approach



- greedy policy improvement \rightarrow given dynamics & value function

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} \quad R(s, a) + M(s, a, s') V(s')$$

↳ model-free

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} \quad Q(s, a)$$

- Exploration vs Exploitation $\rightarrow \varepsilon$ -greedy action selection

- w/ prob. $1-\varepsilon$ choose greedy action
- w/ prob. ε choose random action

- Monte-Carlo Control \rightarrow policy iteration

- policy eval \rightarrow MC $Q \leftarrow Q^\pi$

- policy improvement $\rightarrow \varepsilon$ -greedy

- gradually $\downarrow \varepsilon$ in ∞ limit \rightarrow GLIE

↳ in k^{th} episode w/ π

↳ for each S_t & A_t in episode

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(R_t - Q(S_t, A_t))$$

$$\varepsilon \leftarrow \frac{1}{k}$$

$$\pi(s) \leftarrow \underset{a \in A}{\operatorname{argmax}} \quad Q(s, a)$$

- SARSA \rightarrow TD(0) control
 \rightarrow on policy \rightarrow choose future actions using π

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ ;
    until  $S$  is terminal
  
```

- SARSA(λ) \rightarrow analogue of TD(λ)
 \rightarrow eligibility list to learn in advance

```

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
Repeat (for each episode):
     $E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
    Initialize  $S, A$ 
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
         $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$ 
         $E(S, A) \leftarrow E(S, A) + \delta$ 
        For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :
             $Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$ 
             $E(s, a) \leftarrow \gamma \lambda E(s, a)$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
  
```

- Q-Learning \rightarrow on or off policy

$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \max_{a'} \gamma Q(S', a') - Q(S, A))$$

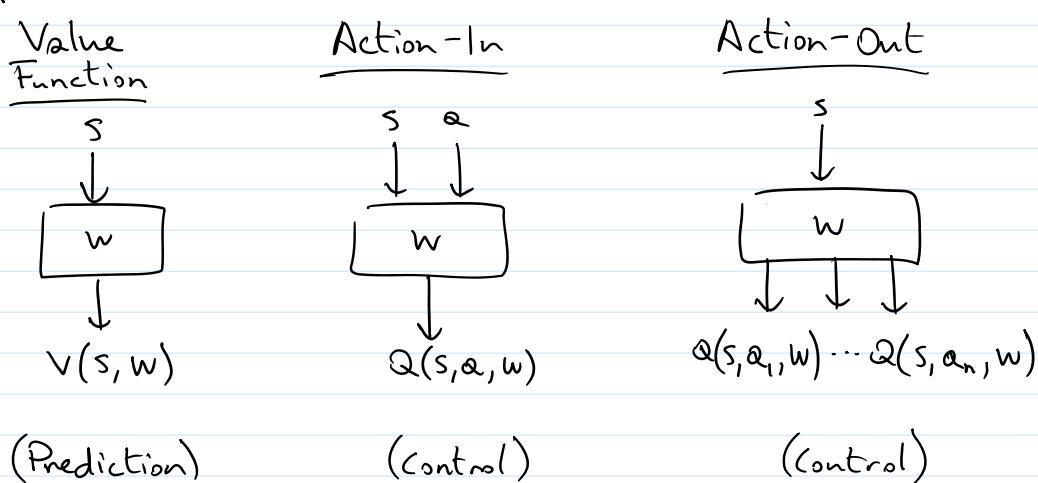
```

Initialize  $Q(s, a), \forall s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
    Initialize  $S$ 
    Repeat (for each step of episode):
        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
        Take action  $A$ , observe  $R, S'$ 
         $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
         $S \leftarrow S'$ ;
    until  $S$  is terminal

```

Value Function Approximation

- for problems w/ too many states for table look-up \rightarrow eg. continuous
 - Types



- Types of Approximators

- linear combinations of features
 - neural network
 - decision trees
 - nearest neighbor
 - fourier/wavelet bases

} differentiable, non-stationary, non-iid

- Linear Function Approximation

$$v(s, \bar{w}) = \bar{x}(s) \cdot \bar{w}$$

- loss

$$L(\bar{w}) = E_{\pi} \left[\left(v_{\pi}(s) - \bar{x}(s) \cdot \bar{w} \right)^2 \right]$$

- update rule

$$\nabla_w v(s, w) = \bar{x}(s)$$

$$\Delta \bar{w} = \alpha (v_{\pi}(s) - v(s, \bar{w})) \bar{x}(s)$$

- incremental method \rightarrow approximate v_{π} by experience

$$\Delta \bar{w} = \alpha (R_{t+1} + \gamma v(S_{t+1}, \bar{w}) - v(S_t, \bar{w})) \bar{\nabla}_w v(s, \bar{w})$$

- same idea w/ Q function

- Convergence Issues

- prediction

gradient follows TD
 of true Bellman eqn.
 ↓

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	TD	✓	✓	✗
Off-Policy	Gradient TD	✓	✓	✓
	MC	✓	✓	✓
	TD	✓	✗	✗
	Gradient TD	✓	✓	✓

- control

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-learning	✓	✗	✗
Gradient Q-learning	✓	✓	✗

(✓) = chatters around near-optimal value function

- Experience Replay DQN \rightarrow batch method \rightarrow reuse training data

- take action $\rightarrow \epsilon$ -greedy

- store transition in replay memory $(s_t, a_t, r_{t+1}, s_{t+1})$

- sample random mini batches from memory

- use 2 networks \rightarrow fix one to \uparrow stability
 \rightarrow train the other w/ Q values from fixed network

- optimize MSL loss w/ SGD

— Fixed

- optimize MSL loss w/ SGD

$$L_i(\bar{w}_i) = E_{(s, a, r, s') \sim D_i} \left[(r + \gamma \max_{a'} Q(s', a', \bar{w}_i) - Q(s; a', w_i))^2 \right]$$

fixed

- Linear Analytical Solution

$$\bar{w} = (\bar{x}(s_i) \bar{x}(s_j))^{-1} \bar{x}(s_i) v_i^\pi$$

$$\begin{aligned} \hookrightarrow \text{don't know } v_t^\pi &\rightarrow \text{MC} \rightarrow v_t^\pi \approx c_i \\ &\rightarrow \text{TD}(0) \rightarrow v_t^\pi \approx R_{t+1} + \gamma v(S_{t+1}, \bar{w}) \\ &\rightarrow \text{TD}(\lambda) \rightarrow c_t^\lambda \end{aligned}$$

- convergence

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✓	✗
	LSTD	✓	✓	-
Off-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✗	✗
	LSTD	✓	✓	-

- LSTDQ \rightarrow least squares q-learning

LSTDQ algorithm: solve for total update = zero

$$0 = \sum_{t=1}^T \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, \pi(S_{t+1}), \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \mathbf{x}(S_t, A_t)$$

$$\mathbf{w} = \left(\sum_{t=1}^T \mathbf{x}(S_t, A_t) (\mathbf{x}(S_t, A_t) - \gamma \mathbf{x}(S_{t+1}, \pi(S_{t+1})))^\top \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t, A_t) R_{t+1}$$

```
function LSPI-TD( $\mathcal{D}, \pi_0$ )
     $\pi' \leftarrow \pi_0$ 
    repeat
         $\pi \leftarrow \pi'$ 
         $Q \leftarrow \text{LSTDQ}(\pi, \mathcal{D})$ 
        for all  $s \in \mathcal{S}$  do
             $\pi'(s) \leftarrow \underset{a \in \mathcal{A}}{\text{argmax}} Q(s, a)$ 
        end for
    until ( $\pi \approx \pi'$ )
    return  $\pi$ 
end function
```