# Chapter 6: Temporal Difference Learning

Objectives of this chapter:

- Introduce Temporal Difference (TD) learning
- Focus first on policy evaluation, or prediction, methods
- Then extend to control methods
  - i.e. policy improvement.

# TD Prediction

**Policy Evaluation (the prediction problem)**:
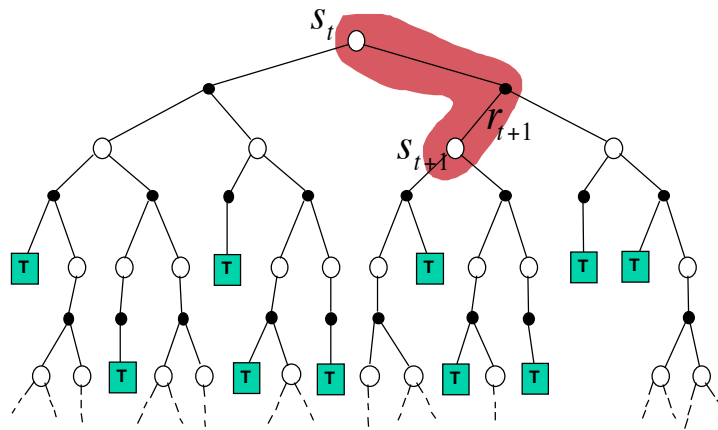for a given policy $\pi$, compute the state-value function $V^{\pi}$

The simplest TD method, TD(0):

$$V(s_t) \leftarrow V(s_t) + \alpha\left[r_{t+1} + \gamma\, V(s_{t+1}) - V(s_t)\right]$$

**target**: an estimate of the return
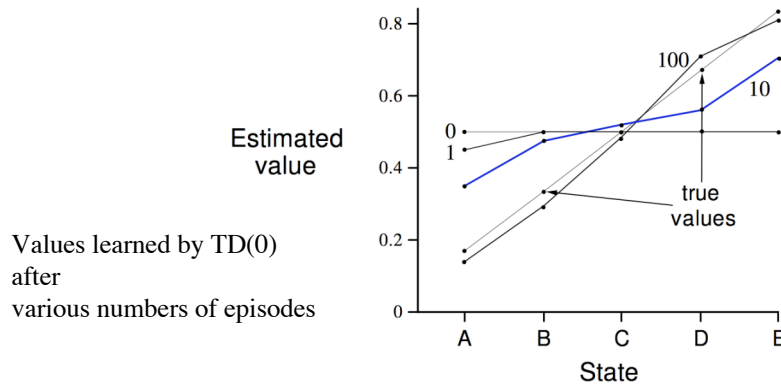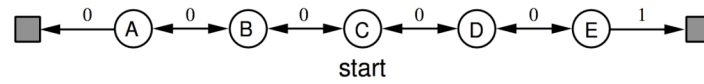
# Simplest TD Method

$$V(s_t) \leftarrow V(s_t) + \alpha\left[r_{t+1} + \gamma\, V(s_{t+1}) - V(s_t)\right]$$



# Advantages of TD Learning

- TD methods do not require a model of the environment, only experience
- TD methods can be fully incremental
  - You can learn before knowing the final outcome
    - Less memory
    - Less peak computation
  - You can learn without the final outcome
    - From incomplete sequences
- TD converges to an optimal policy (under certain assumptions to be detailed later)

# Random Walk Example



Values learned by TD(0)
after
various numbers of episodes

---

# Optimality of TD(0)

Batch Updating: train completely on a finite amount of data,
e.g., train repeatedly on 10 episodes until convergence.

Compute updates according to TD(0), but only update
estimates after each complete pass through the data.

For any finite Markov prediction task, under batch updating,
TD(0) converges for sufficiently small $\alpha$.

# You are the Predictor

Suppose you observe the following 8 episodes:
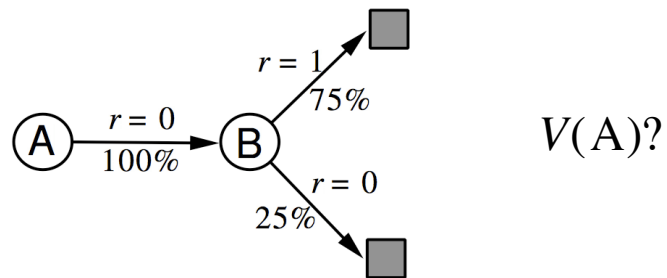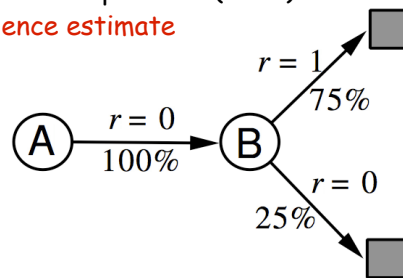
A, 0, B, 0
B, 1
B, 1
B, 1
B, 1
B, 1
B, 1
B, 0
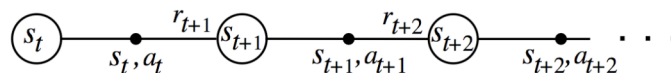
$V(A)?$

$V(B)?$

---

# You are the Predictor



$V(A)?$

# You are the Predictor

- The prediction that best matches the training data is V(A)=0
    - This minimizes the mean-square-error on the training set
- If we consider the sequentiality of the problem, then we would set V(A)=.75
    - This is correct for the maximum likelihood estimate of a Markov model generating the data
    - i.e, if we do a best fit Markov model, and assume it is exactly correct, and then compute what it predicts (how?)
    - This is called the certainty-equivalence estimate
    - This is what TD(0) gets

- Thought from Dan: If P(start at A) is so low, apparently, who cares?

$r = 0$ 100% (A) → (B)

$r = 1$ 75%

$r = 0$ 25%

---

# Learning An Action-Value Function

Estimate $Q^{\pi}$ for the current behavior policy $\pi$.

$(s_t)$ —•— $r_{t+1}$ —$(s_{t+1})$ —•— $r_{t+2}$ —$(s_{t+2})$ —•— $\cdot$ $\cdot$ $\cdot$
$s_t, a_t$ ⠀⠀⠀ $s_{t+1}, a_{t+1}$ ⠀⠀⠀ $s_{t+2}, a_{t+2}$

After every transition from a nonterminal state $s_t$, do this:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \, Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right]$$

If $s_{t+1}$ is terminal, then $Q(s_{t+1}, a_{t+1}) = 0$.

# Sarsa: On-Policy TD Control

Turn this into a control method by using the current greedy policy:

One - step Sarsa :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right]$$

Initialize $Q(s, a)$ arbitrarily
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)
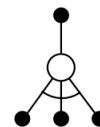      $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$
      $s \leftarrow s'; a \leftarrow a';$
   until $s$ is terminal

# Q-Learning: Off-Policy TD Control

One - step Q - learning :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

Initialize $Q(s, a)$ arbitrarily
Repeat (for each episode):
   Initialize $s$
   Repeat (for each step of episode):
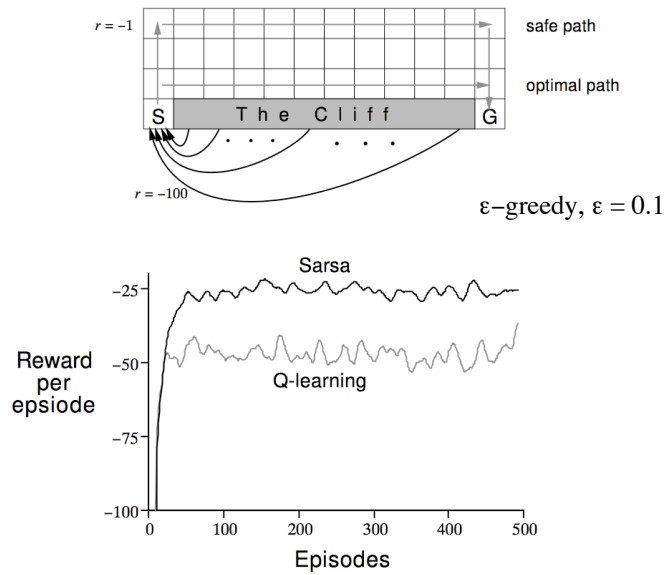      Choose $a$ from $s$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
      Take action $a$, observe $r$, $s'$
      $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$
      $s \leftarrow s';$
   until $s$ is terminal

# Cliffwalking

r = -1     safe path

optimal path

S    T h e   C l i f f    G

r = -100

$\varepsilon$−greedy, $\varepsilon = 0.1$

Sarsa

−25

Reward per epsiode

−50

Q-learning

−75

−100

0   100   200   300   400   500

Episodes

---

# Summary

- TD prediction
- Introduced *one-step tabular model-free TD methods*
- Extend prediction to control by employing some form of GPI
    - On-policy control: Sarsa
    - Off-policy control: Q-learning