

- Model-Free Prediction - See slide 4F

- Estimate Value Function using:

- ① Monte Carlo
- ② Temporal Difference
- ③ TD( $\lambda$ )

- This lecture just looks at predicting reward for a policy w/o models
- Next lecture looks at control
- Two Dichotomies:

- Planning vs. Model Free / Prediction vs. Control

- Monte Carlo Reinforcement Learning:

- Monte Carlo methods learn directly from episodes of experiences

→ Requires that experiences terminate

- Goal: Learn  $V_\pi$  from episodes of experiences under policy  $\pi$

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$$

→ In Monte Carlo, we use empirical mean return

- ① First-visit MC Evaluation: In each episode, obtain the value function from each state using trajectories obtained after that state is reached for the first time. Average value function from each state using many episodes.

$$N(s) = \# \text{ of first visits to } s \text{ in the episode}$$

$$\rightarrow N(s) \leftarrow N(s) + 1 \text{ upon first visit}$$

$$S(i) = \text{Total return from visits to } s$$

$$\rightarrow S(i) \leftarrow S(i) + G_i \text{ from first visit}$$

$$\therefore V(i) = S(i) / N(s)$$

- By Law of large numbers,  $V(s) \rightarrow v\pi(s) \Rightarrow N(s) \rightarrow \infty$

② Every Visit MC Evaluation: Just like FVMC, but increment counter  $N(s)$  and store  $S(s)$  for Every visit to each state.

- Blackjack as MDP

↳ State space has 3 different variables:

① Current sum (only consider sums with decisions  $\Rightarrow$  sum between 12 - 21).

② Dealer's showing card

③ Do I have "hole" ace (yes, no)  $\leftarrow$  can take values 1 or 11

↳ Actions: Hit, stand

↳ Reward: +1 for win; 0 for tie; -1 for miss

↳ Policy evaluated using every visit MC by simulating some # of episodes

↳ Noise is observed around rare states, since you don't encounter rare states in MC as frequently

↳ Sample returns and model VF without a model.

- Incremental Means: Mean can be computed incrementally

$$\mu_k = \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

↳ New mean = old mean + step = error

↳ Current prediction in direction of error

- Incremental MC Updates:

For each state  $s_t$  with return  $G_t$ :

$$N(s_t) \leftarrow N(s_t) + 1$$

$$v(s_t) \leftarrow v(s_t) + \frac{1}{N(s_t)} (G_t - v(s_t))$$

Always forward-looking, using reward from state obtained at the end of the episode.

For non-stationary problems, track a running mean and forget old episodes

$$V(s_t) \leftarrow v(s_t) + \alpha (G_t - V(s_t))$$

Constant Learning Rate

↪ Don't store  $N$

↪ Memoryless algorithms

### • Temporal Difference Learning (TD)

↪ Learn from actual atomic experiences (even incomplete episodes)

↪ Model-Free Approach

↪ Bootstrap = Update guess of value function as you get more data.

Goal: Learn  $V\pi$  online from experiences under  $\pi$

↪ TD(0) = simplest TD

$$V(s_t) \leftarrow V(s_t) + \alpha \underbrace{(R_{t+1} + \gamma V(s_{t+1}) - V(s_t))}_{\text{TD Error}}$$

Estimated return from one step.

"TD Target" = Random target depending on next time step

↪ Estimated, biased algorithm

↪ TD learning reflects how we learn in real life

- TD can learn before you see the final outcome, while MC requires complete episodes
- TD can learn for non-terminal processes
  - ↳ TD will eventually find the correct value function
- Bias/Variance Tradeoff: Return  $G_t$  is an unbiased estimate of  $V_\pi(s_t)$ 
  - ↳ TD Target contains  $V(s_{t+1})$  which is not necessarily the true value function so introduces bias
  - ↳  $G_t$  incurs noise from many steps  $\Rightarrow$  High variance
  - ↳ TD looks at noise from a single step  $\Rightarrow$  Low variance
  - ↳ MC has very good convergence even w/ function approx.
  - ↳ MC is not very sensitive to initialization
  - ↳ MC is simple to understand and use
  - ↳ TD has low variance, some bias
  - ↳ Usually more efficient than MC
  - ↳  $TD(\alpha)$  converges to  $V_\pi(s)$ , but not always w/ function approx.
  - ↳ More sensitive to initial values
- Random Walk TD Example:

- ↳ 1D Random Walk with uniform left/right probability; reward 1 at RHS and 0 elsewhere
- ↳ TD learning converges faster than MC.
- ↳ By backtracking, you learn more efficiently than when incurring full variance of episode using MC

- Batch MC and TD:

↳ VF is average of  $G_t$  observed from each state in the batch  $\leftarrow$  MC estimate

↳ Alternatively, looks at empirical transition probabilities and returns associated with each state in the transitions  $\leftarrow$  TD Estimate

- MC converges to the solution with the minimum Mean Squared Error: 
$$\sum_{k=1}^K \sum_{t=1}^T (g_t^k - v(s_t^k))^2$$
  
Episodes

- $TD(0)$  converges to solution of maximum likelihood Markov Model

↳ First get empirical transition probabilities

↳ Calculate mean reward from each state

- Markov Decision Process has Markov Property, which is exploited by TD, which makes use of this property to be more efficient

- MC ignores Markov Property, which can be inefficient if process exhibits non-Markov behavior

- Monte Carlo Samples trajectories from states and uses those samples to estimate VF

↳ TD learning looks just 1 step ahead

↳ DP looks 1 step look-ahead, but did not sample, instead did full look-ahead

- Bootstrapping vs. Sampling:

↳ Bootstrapping uses estimates instead of real returns for updates

↳ MC = No bootstrapping

↳ TD = Bootstrapping

↳ DP = Bootstrapping

↳ Sampling:

↳ MC, TD sample

↳ DP does not sample, instead using full-width lookahead

\* TD( $\lambda$ ) lets you specify the depth of your backups

↳ As you take steps, increase depth, you ground yourself in real dynamics and rewards.  
Real dynamics always brings you closer to ground truth.

↳ What if you use multiple steps in TD learning.

↳ TD can be generalized to n-step experiences

$$G_t^{(1)} \Rightarrow n=1, 1\text{-step return} \therefore G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$$

$$G_t^{(2)} \Rightarrow n=2, 2\text{-step return} \therefore G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$$

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$G_t^{(\infty)} = \text{MC}$$

↳ n-step TD Update:

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^{(n)} - V(s_t))$$

↳ As  $n \rightarrow \infty$ , you get very high errors for short training times, like MC

↳ There is a sweet spot around  $n=3$  or  $5$  which minimizes error.

\* We want an algorithm that efficiently considers all  $n$  at once

↳ We can average over n-step returns and call that the target

↳ How can we efficiently combine all  $n$ ?

↳  $\lambda$  return!

↳  $G_t^\lambda$  combines n-step returns  $G_t^{(n)}$  using weights  $(1-\lambda) \lambda^{n-1}$

$$\therefore G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \quad \leftarrow \text{Geometric Weighting}$$

↳ Forward-viewing  $\lambda$ :

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^\lambda - V(s_t))$$

\* Weighting decays with the number of steps you look into the future

↳ With geometric weightings, you can run TD( $\lambda$ ) with same time as TD(0).

↳ Forward-view TD( $\lambda$ ) suffers the same consequences as MC

↳  $\lambda = 1 \Rightarrow$  MC

$\lambda = 0 \Rightarrow$  TD(0)

↳ Backward view TD( $\lambda$ ):

↳ Eligibility Trace: Assigning credit for recency and frequency

↳ Combines frequency and recency heuristics

↳ Backward view TD( $\lambda$ ):

① Keep an eligibility trace for each state  $s$

② Update value  $V(s)$  for each state  $s$  in proportion to TD Error  $\delta_t$  on eligibility trace  $E_t(s)$

$$\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

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

↳ VF updated on previous states

↳ When  $\lambda = 0$ , only current state is updated

↳ When  $\lambda = 1$ , there is no recency decay  $\Leftrightarrow$  MC

\* The sum of offline updates is the same for forward-view and backward-view  
 $TD(\alpha)$

↑  
Theorem.

$TD(\alpha)$  gives you spectrum between  $TD(0)$  and  $MC$