# Lecture 4: Model-Free Prediction

Tuesday, August 9, 2022    12:10 PM

Lecture 3: Planning w/ Dynamic Programming for a <u>KNOWN</u> MDP

This lecture: Model-Free prediction (w/ unknown MDP)

    — estimate value function of unknown MDP

Next lecture: Model-Free control

    — Optimise value function of unknown MDP

## Monte-Carlo Learning

  — Learn directly from episodes of experience

  — don't need knowledge of MDP transitions

  — only works for episodic MDPs (all episodes must terminate)

    Goal: learn $v_\pi$ from episodes of experience under policy $\pi$

$$v_\pi(s) = \mathbb{E}_\pi \left[ G_t \mid S_t = s \right]$$

    • Monte-Carlo evaluation uses empirical mean return instead of expected value

    — Version 1 of Monte Carlo learning

      First-Visit Monte-Carlo Policy Evaluation

      — For a state s

        on the first time-step t that s is visited

        — Increment counter $N(s) \leftarrow N(s) + 1$

        — Increment total return $S(s) \leftarrow S(s) + G_t$

$N(s) \rightarrow$ number of episodes in which

we have visited
s

Note that we are adding $G_t$ to $S(s)$, $G_t$ is TOTAL return from the time we saw the state

— Value is estimated by $V(s) = \dfrac{S(s)}{N(s)}$ mean return

— Law of Large Numbers

$$\Longrightarrow V(s) \to v_\pi(s) \text{ as}$$
$$N(s) \to \infty$$

## Version 2 of Monte Carlo Learning

Every - visit Monte-Carlo Evaluation

— Exact same algo as first-visit, except don't just update on first-visit, ==update every time you see a specific state==

## Incremental Mean Formula

$$\mu_k = \frac{1}{k} \sum_{j=1}^{k} x_j$$

$$= \frac{1}{k}\left( x_k + \sum_{j=1}^{k-1} x_j \right)$$

$$= \frac{1}{k}\left( x_k + (k-1)\mu_{k-1} \right) \quad \leftarrow \text{Previous mean}$$

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

$$\mu_k = \mu_{k-1} + \frac{1}{k}\left( \underbrace{x_k - \mu_{k-1}}_{\text{error term}} \right) \quad \left.\begin{array}{c}\\ \\ \end{array}\right\} \text{move mean in direction of error}$$

## Incremental Monte-Carlo Updates

- Update $V(s)$ incrementally after episode $S_1, A_1, R_2, \ldots S_T$
- For each state $S_t$ w/ return $G_t$

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

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$$

- For non-stationary problems, it can be helpful to track a running mean i.e; forget old episodes

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

## TD-Learning

Temporal-Difference: • Learn directly from episodes of experience
- model free: no knowledge of MDP transitions/
rewards
- learns from incomplete episodes, by bootstrapping
$\hookrightarrow$ updates guess before a guess

Goal: learn $V_\pi$ online from experience under policy $\pi$

Incremental, every-visit Monte-Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha (\underline{G_t} - V(S_t))$$

### Simplest T-D algorithm TD(0)

Update value $V(S_t)$ toward estimated return

$$R_{t+1} + \gamma V(S_{t+1})$$

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

$$R_{t+1} + \gamma V(s_{t+1}) \quad \text{is} \quad \text{TD-target}$$

$$\delta_t = \boxed{R_{t+1} + \gamma V(s_{t+1})} - V(s_t) \quad \text{is} \quad \text{TD-error}$$

- Estimates are grounded by end-of-episode

## Advantages and Disadvantages

- TD can learn before knowing final outcome
  MC must wait till end of episode

- TD can learn w/o final outcome

## Bias/Variance Trade-Off

- Return $G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T$ is unbiased
  estimate of $v_\pi(s_t)$

- True TD target $R_{t+1} + \gamma v_\pi(s_{t+1})$ is unbiased estimate
  of $v_\pi(s_t)$
  of
  oracle value (true value function)

- TD target $R_{t+1} + \gamma V(s_{t+1})$ is biased estimate of $v_\pi(s_t)$
  — has much less variance than the return
  — Return depends on many random actions, transitions,
                                              rewards
  → TD target depends on one random action, transition,
    reward

## MC

- high variance, zero bias
- good convergence
- not sensitive to initial value
  very simple

## TD

- Usually more efficient than MC
- TD(0) converges to $v_\pi(s)$
  - (but not always w/ function
    approximation)
- more sensitive to initial value

- MC and TD converge: $V(s) \rightarrow V_\pi(s)$ as experience $\rightarrow \infty$
- But what about batch solutions

  - $s_1^1, a_1^1, r_2^1, \ldots s_{T_1}^1$   $\leftarrow$ episode 1
    $\vdots$
  - $s_1^k, a_1^k, r_2^k \ldots s_{T_k}^k$   $\leftarrow$ episode $k$

- MC converges to solution w/ minimum mean-squared error
  - Best fit to observed returns
    $$\sum_{k=1}^{K} \sum_{t=1}^{T_k} \left( g_t^k - V(s_t^k) \right)^2$$

- TD(0) converges to solution of max-likelihood Markov
  model

- TD exploits Markov property
  - Usually more efficient in Markov environments
- MC does not exploit Markov property
  - more efficient in non-markov environments

## Recap

$$V(S_t) \longleftarrow V(S_t) + \alpha\left(G_t - V(S_t)\right) \qquad \text{Monte-Carlo}$$

$$V(S_t) \longleftarrow V(S_t) + \alpha\left(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)\right) \qquad \text{TD-Learning}$$

$$V(S_t) \longleftarrow E_\pi\left[R_{t+1} + \gamma V(S_{t+1})\right] \qquad \begin{array}{l}\text{Dynamic} \\ \text{Programming}\end{array} \quad \left(\begin{array}{l}\text{Need to know} \\ \text{environment} \\ \text{dynamics}\end{array}\right)$$
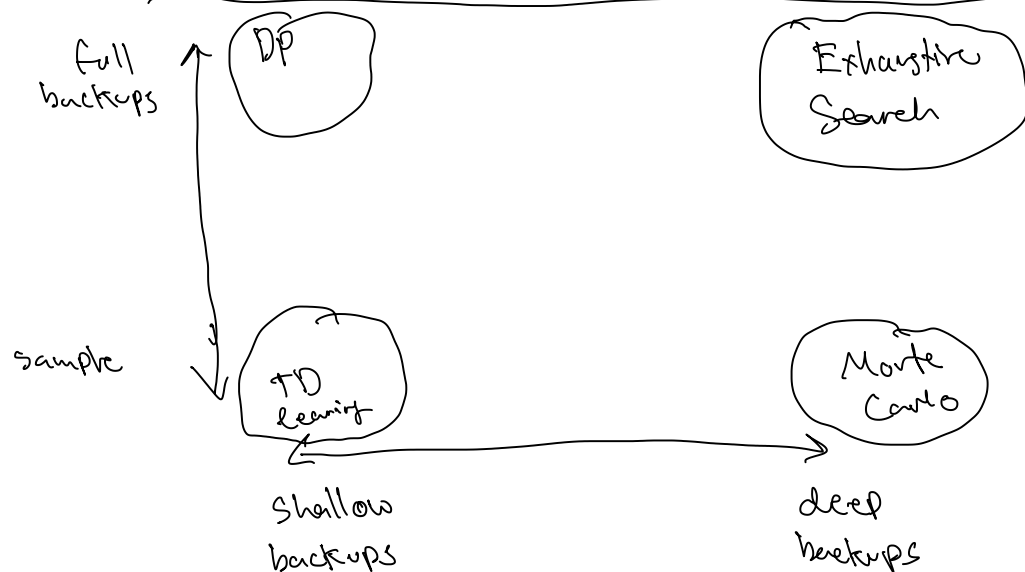
Categorization

Bootstrap : update involves an estimate

- MC does not bootstrap
- DP bootstraps
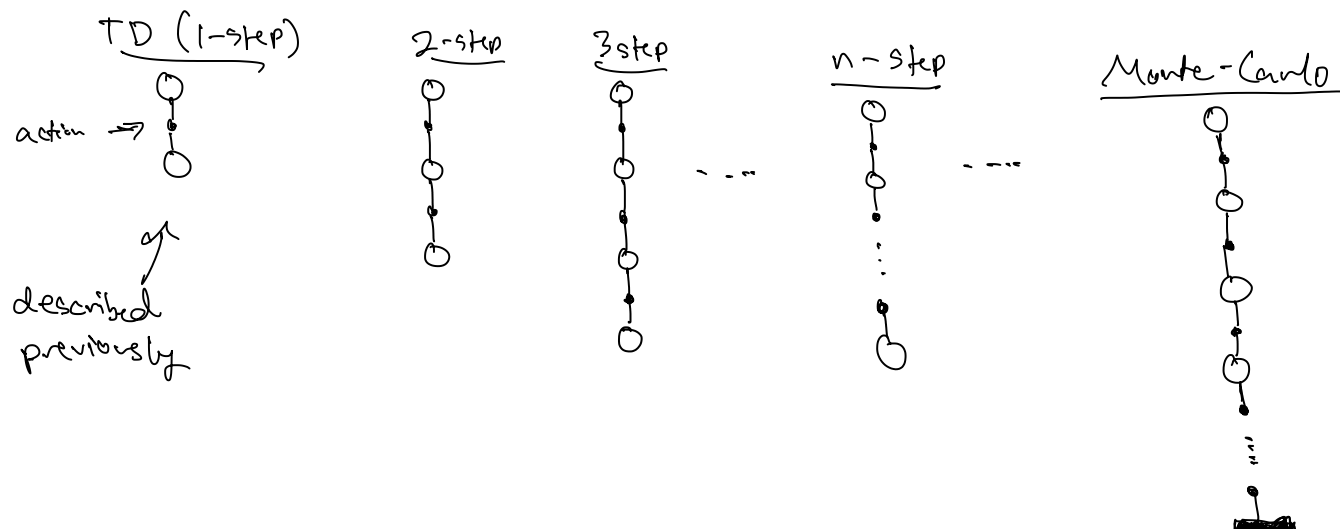- TD bootstraps

Sampling : update samples an expectation

- MC samples
- DP does not sample
- TD samples

# Unified View of Policy Evaluation

full
backups

sample

DP

Exhaustive
Search

TD
learning

Monte
Carlo

Shallow
backups

deep
backups

$TD(\lambda)$ — mix of TD-learning and Monte-Carlo
updates

• TD target looks n steps into the future

TD (1-step)    2-step    3step    n-step    Monte-Carlo

action →

described
previously

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

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

$$\vdots$$

$$n=\infty \qquad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T$$

$$\text{n-step return} \qquad \overset{\text{real}}{\overbrace{\phantom{xxxxxxxxxxxxxxxxxxxx}}} \qquad \overset{\text{estimate}}{\overbrace{\phantom{xxxxx}}}$$

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

**n-step temporal-difference learning**

$$V(S_t) \leftarrow V(S_t) + \alpha\left(G_t^{(n)} - V(S_t)\right)$$
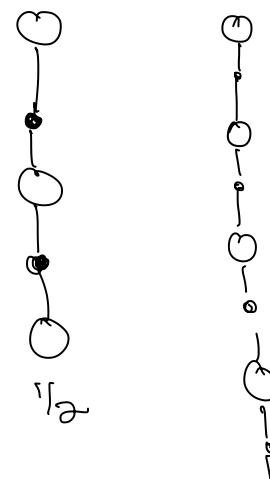
## Averaging n-step returns

- average n-step returns over different n
- e.g. average 2-step and 4-step return $\longrightarrow$

### TD($\lambda$), $\lambda$-Return

- $\lambda$-Return $G_t^{\lambda}$ combines all
  n-step        $\quad$ $G_t^{(n)}$

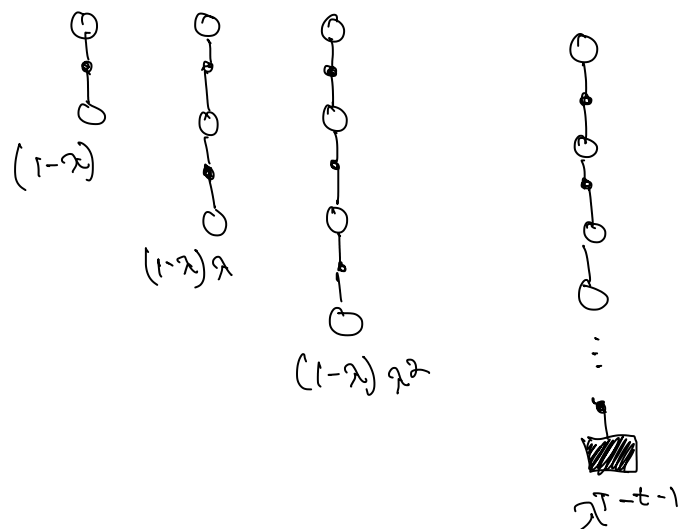## One backup



$^{1}/_{2}$

... of n returns $G_t^{(n)}$

- Usig weight $(1-\lambda)\lambda^{n-1}$

$$G_T^\lambda = (1-\lambda) \sum_{n=1}^\infty \lambda^{n-1} G_t^{(n)}$$

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

⋆ initial weight of $(1-\lambda)$ is used to ensure weights sum to one

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^\infty \lambda^{n-1} G_t^{(n)}$$

$1/2$

$(1-\lambda)$

$(1-\lambda)\lambda$

$(1-\lambda)\lambda^2$

$\lambda^{T-t-1}$

- Forward-view looks into future to compute $G_t^\lambda$
- Like MC, can only be computed from complete episodes

## BUT !!!

- Backward View TD($\lambda$)
  - Forward view provides theory
  - Backward provides mechanism
  - Update online, every step, from incomplete sequences

# Credit assignment Problem

Frequency Heuristic: assign credit to most frequent states

Recency heuristic: assign credit to most recent states

## Eligibility Trace

$$E_0(s) = 0$$

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

$\underset{\text{decay factor}}{\uparrow}$



accumulating eligibility trace

times of visit to state

- Keep eligibility trace for every state $s$
- Update value $V(s)$ for every state $s$ in proportion to TD-Error and eligibility trace

TD-Error $\delta_t = \left(R_{t+1} + \gamma V(S_{t+1})\right) - V(S_t)$

$$V(s) \leftarrow V(s) + \alpha \, \delta_t \, E_t(s) \qquad \alpha = \text{learning\_rate}$$

when $\lambda = 0$, only current state is updated

$$E_t(s) = \mathbf{1}(S_t = s)$$

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

$$V(s_t) \xleftarrow{=} V(s_t) + \alpha \delta_t$$

when $\lambda = 1$, credit is deferred until end of episode
- same updates at Monte Carlo

Theorem: Sum of offline updates is identical for forward and backward view $TD(\lambda)$

$$\underbrace{\sum_{t=1}^{T} \alpha \delta_t E_t(s)}_{\text{forward view}} = \underbrace{\sum_{t=1}^{T} \alpha \left( G_t^{\lambda} - V(s_t) \right) \mathbf{1}(s_t = s)}_{\text{backward view}}$$

Summary of Forward and Backward $TD(\lambda)$

|                  | $\lambda = 0$ | $\lambda \in (0,1)$ | $\lambda = 1$ |
|------------------|---------------|---------------------|---------------|
| Forward View     | $TD(0)$       | $TD(\lambda)$       | $TD(1)$       |
| Backward View    | $TD(0)$       | Forward TD $(\lambda)$ | $MC$       |
| Exact Online     | $TD(0)$       | Exact $TD(\lambda)$ | Exact $TD(1)$ |

$TD(\lambda)$ is mix between Monte-Car
and $TD(0)$