

Temporal difference learning

Learn directly from episodes of experience

Model free like MC

Learns from incomplete episodes.

No requirement of termination

Bootstrap \rightarrow update guess towards guess

TD updates the estimates as it progresses.

TD is one step lookahead

TD requires that reward is calculated at every transition

TD is better than MC.

TD exploits Markov property

MC does not exploit Markov property

TD is a prediction algorithm, which means it only predicts V given π . It does not optimize π

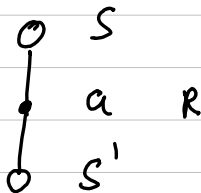
The goal here is only to find V values for all states given a policy π .

Incremental every visit Monte Carlo Prediction

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

In T.D., we replace G_t by estimated
reward

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$



$$\text{Estimated return at } S = \text{Reward at } S' + \text{Goodness of } S'$$

So, we need only info about 2 states and action for TD. We don't need complete episode

Bias & Variance in TD & MC

TD can learn before knowing the final outcome unlike MC which must wait until the end of the episode before return is known

Return $G_t = R_{t+1} + \gamma R_{t+2} \dots$ of MC is unbiased estimate.

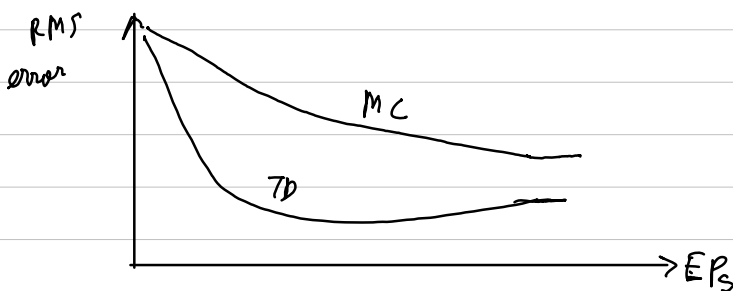
But $R_{t+1} + \gamma V(S_{t+1})$ is biased estimate

G_t depends on many random actions, transitions and rewards

TD depends on one action, transition and reward

MC has high variance, zero bias

TD has low variance, some bias



In Monte carlo, G can drastically differ between episodes. Also due to delayed rewards, a lot of variance is introduced

TD updates are performed at each time step, allowing the value function to converge smoothly. Each small update contributes only a portion of the total change needed which smooths the learning process and reduces the variance of each update.

Frequent updates lead to incremental improvements that average out fluctuations caused by noisy rewards or random events.

TD relies on immediate next reward, that is local information rather than the entire episode. This helps avoid large swings in estimates that might result from unusual rewards further along an episode, thus keeping the update variance lower.

Batch updating TD & MC

Updates are postponed till a collection of experiences has been gathered. Then values are updated in one go.

Accumulate error in various transitions and update in bulk

Reduce variance and increase stability in learning.

faster convergence as updates are in bulk

Makes model more robust

Useful for non stationary and complex environments
Average out the variability in episodic returns.

Monte Carlo

Complete episodes

Actual return G_t

Low Bias

High Variance

Need episode to terminate

No Bootstrapping

Computationally expensive

doesn't exploit Markov property

Temporal difference learning

Incomplete episodes also

Estimated return $R_t + V(s_{t+1})$

High bias

Low Variance

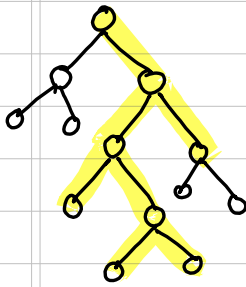
Can work for continuing tasks

Bootstrapping

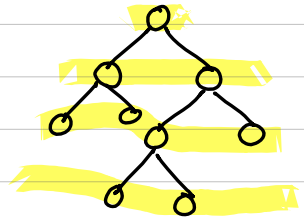
Cheaper computations

Exploits Markov property

Comparison of Methods

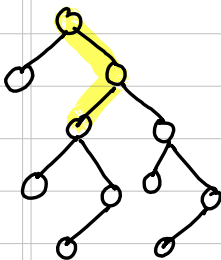
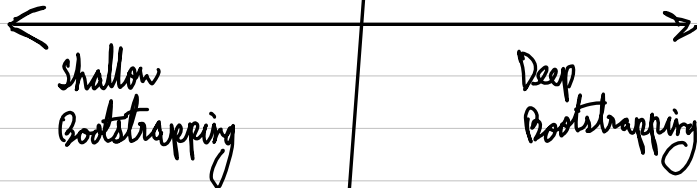


D.P.



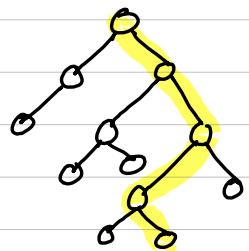
Exhaustive Search

full Backups



T.D.L.

$$V(s) \leftarrow V(s) + \alpha (R_t + \gamma V(s') - V(s))$$



M.C.

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

Sample Backup

SARSA, Q also are here

Model dependant	Model free
D.P. Exhaustive search	MC, TDL, SARSA Q-L
require model	No model required
Learn through equations	learn through episodes
Info of model, P , γ defined	learning through model interaction

Bootstrapping \rightarrow guess towards guess

D.P.	Exhaustive search
Only all nodes of a single path explored	All nodes of all the paths explored
Fast	Slow (infeasible)

TDL, SARSA, Q _L	MC
Estimate the reward	G_t direct reward calculated
High bias, low variance	Low Bias, High Variance
Useful when intermediate rewards known	Useful when model is black box

Bootstrapping makes the process faster

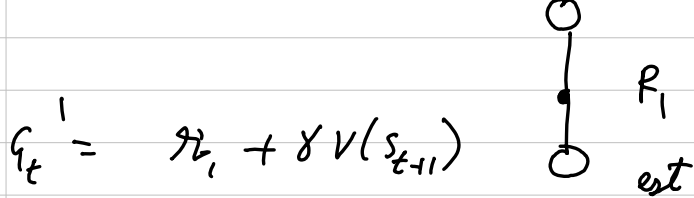
n - step TDL

MC uses Deep backups
TDL uses No backups

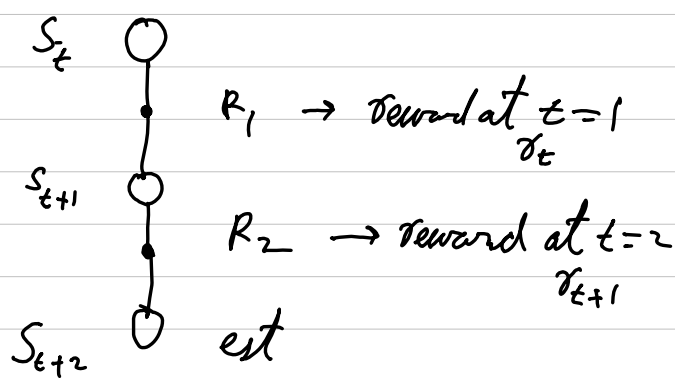
Middle ground is n-step bootstrapping

Instead of using only 1 reward & future estimates we can use n rewards.

$n=1$: Normal



$n=2$: 2 step look ahead



$$\text{goodness} = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2})$$

(Expected reward)

$$\text{goodness} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n V(s_{t+n})$$

That is, the observed return becomes

$$G_t^n = \sum_{i=0}^n r_{t+i} \gamma^i$$

If episode terminates beforehand \rightarrow

$$G_t^n = \sum_{i=0}^{\min(T, n)} r_{t+i} \gamma^i$$

$$\therefore V(s_t) \leftarrow V(s_t) + \alpha (G_t^n - \gamma^n V(s_{t+n}) - V(s_t))$$

Basic idea of n step, is that the V value is updated after n episodes.

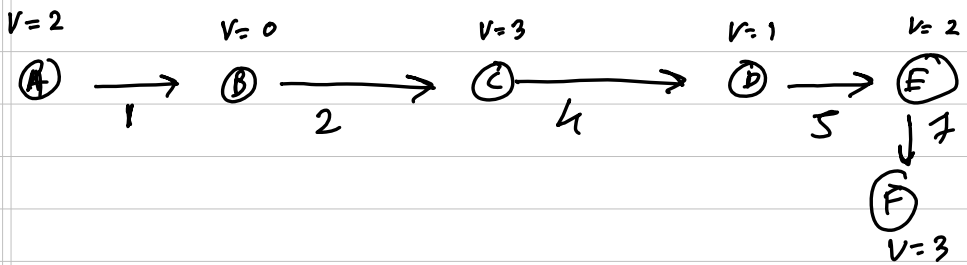
Update V

	TDL	n-step	Monte carlo
no of ep	1	n	all
		actually min(n, all)	

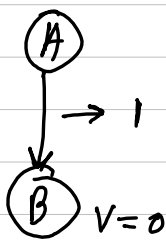
choosing n can be difficult

Book keeping (sliding window etc) required

Consider a path followed by an episode
the rewards obtained and previous V values
 $\alpha = 0.5, \gamma = 0.9$



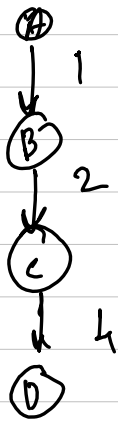
TDL



$$V_A = V_A - \alpha (\gamma V_B + r_t - V_A)$$

$$= 2 - 0.5 (1 + 0 - 2)$$

3 step TDL



$$V_A = V_A - \alpha (\gamma^3 V_D + r_t + \gamma V_B + \gamma^2 V_C - V_A)$$

$$= 2 - 0.5 (1 + 2 \times 0.9 + 4 \times 0.9^2 + 1 \times 0.9^3 - 2)$$

Monte Carlo

$$V_A = V_A - \alpha (G_t - V_A)$$

$$= 2 - 0.5 (1 + 2 \times 0.9 + 4 \times 0.9^2 + 5 \times 0.9^3 + 7 \times 0.9^4 - 2)$$

Parameters of T.D.L.

① Learning rate $\alpha \rightarrow$ Between 0-1

Large learning rate might lead to fluctuating training results

Small learning rate adjusts slowly, which will take more time to converge

$\alpha = 1$: fully trust most recent information discarding previous knowledge

$\alpha = 0$: Ignore new information, never learns or updates the knowledge

② Gamma $\gamma \rightarrow$ The discount rate. Represents how much we are valuing future rewards.

Bigger the discount rate, more we value future rewards.

Generally between 0-1

High γ (close to 1)	future rewards valued (farsighted)
Low γ (close to 0)	Intermediate rewards prioritized (Nearsighted)

③ Epsilon $\epsilon \rightarrow$ Exploration vs Exploitation

Prob $\epsilon \rightarrow$ Explore

Prob $1-\epsilon \rightarrow$ exploit

Larger ϵ means more exploration while training