

- Feb 19, 2021 Notes: Temporal Difference

- Temporal Difference : Constantly adjust value function to decrease error

- ↳ Updates after every atomic experience

- ↳ You do not have to update in the sequence you get data

- ↳ You can replay experiences

- TD w/ Function Approx.

- ↳ Parameterize value function with weights u

- ↳ Semi-gradient

- In code: Input is not iterable of episodes, it is iterable on atomic experiences.

- ↳ In RL we used "iterate-updates" but here we cannot use it.

- ↳ Instead use "iterate-accumulate"

- ↳ Inputs:

- ① Data structure you iterate over

- ② Function that takes

- i. current sum (or other function)
 - ii. element in list we are looking at

- returns updated sum (or other function)

- ③ Initial value of accumulator

- ↳ You can accumulate any quantity over any data structure.

- ↳ See slide 17 for code related to TD.

- Functional Programming: Code looks clean/simple
 - ↳ Almost all bugs come from "imperative" style of programming with lots of loops.
- TD is the main way of implementing reinforcement learning
 - ↳ Returns can be expressed in recursive manner (like DP)
 - ↳ You can learn from experiences, you don't need Pa (like MC)
 - * Can overcome curse of dimensionality and curse of modeling
 - * TD reflects the way humans learn in real life
 - ↳ We make immediate updates (like TD), we don't wait and learn in bulk (as in MC)
- MC is unbiased estimate of VF, TD is biased
 - ↳ But no free lunch:
 - MC has high variance, TD has low variance
- Because MC is unbiased, it converges well even with function approx
- Theorem: Tabular TD prediction converges to true VF in the mean for constant α
 - ↳ Requires Robbins - Monro learning rate schedule
 - $$\sum_{n=1}^{\infty} \alpha_n = \infty \quad \text{and} \quad \sum_{n=1}^{\infty} \alpha_n^2 < \infty$$
 - ↳ TD target can move rapidly and can converge quickly

- Speed of convergence of TD vs. MC:
 - Compare speed of convergence, efficiency in use of limited data
 - With constant α , TD performs better than MC
 - Compare TD and MC for different Markov Decision Processes using infinite horizon
 - MC is sensitive to initial value, TD is not.
 - Try generating convergence graph
- So far we've focused on the mechanics of how TD and MC work.
 - Now we will talk about what it is learning?
 - Only have 16 episodes, but can replay experience = "Experience Replay." Randomly pull out traces for MC or reuse experience in TD
 - What is learned by TD and MC is different.
 - TD pulls atomic experiences
 - MC is simple Supervised learning
 - Gives you average return at each state
 - TD drives towards VF of MRP that is implied by Experience
 - Empirical Markov Reward Process
 - Get empirical transition probabilities
 - Takes advantage of Markov Properties
 - If process is not truly Markovian, MC is better
 - TD is not building MRP, put empirical MRP VF is the same as that constructed by TD
 - MC is a Naive Method, TD calculates VF of an implied MRP

- We summarize MC, DP, TD in terms of whether they:
 - ① Bootstrap: Update in VF utilize current or prior estimate of VF
 - ② Experiencing: Interaction with actual or simulated environment
 - ↳ Learning vs. planning

* See slide 26 for characterizations

Not wide \rightarrow • MC Backup is Deep - Backup VF through entire episode

Not wide \rightarrow • TD backup is Shallow - Backup VF through one atomic experience

Very wide \rightarrow • DP is Shallow

- Axial: You always want to use DP whenever possible

\hookrightarrow Approximate DP is a middle-ground between TD and RL

- Tabular n-step bootstrap:

$$V(s_t) \leftarrow V(s_t) + \alpha (G_{t,n} - V(s_t; v))$$

$\hookrightarrow \lambda$ serves as a tuning knob between TD and MC

- λ -Return Prediction Algorithm \longleftarrow This is an offline (batch update) algorithm

$$u_n = (1-\lambda) \lambda^{n-1} \text{ for all } n = 1, \dots, T-t+1$$

\rightarrow Instead of $G_{t,n}$ a valid target is a weight-averaged target

\hookrightarrow λ -return algorithm used as a link between TD and MC

- Online algorithms update after every step

\hookrightarrow We want an online algorithm that enables learning from TD to MC

- TD(λ) Algorithm:

- Slight detour: Eligibility Trace.
- ↳ Consider events occurring in time.
- ↳ Maintain count which decays with time
- Functions when you care about recent events more than older events
- Eligibility Trace is a memory function for each state
- ↳ At each time step, update memory function for all states.

$$E_t(s) = \gamma \lambda E_{t-1}(s) + I[s_t = s]$$

- See algebra proof on slide 39

 ↳ Telescoping algorithm.

 ↳ Try to write this out.

$$\underbrace{b_t^{(1)}}_{\text{lambda return algo update}} = v(s_t) = b_t + \dots \text{ See slide 40.}$$

- When generalizing to function approx, eligibility trace for parameters, not states.
- Try out RL assignments from this week to build foundation.