# RL Lecture 4 — Model free Learning

What we did earlier using DP, required the agent to know the entire model dynamics ie. the state transition probabilities, reward functions, etc. But since in most real life problems, the agent wouldn't know what the entire dynamics of the environment are. What can we do?

That is discussed now!    ↳ Model-free Learning

Just like earlier, we break it into 2 components:

① Policy evaluation → model free prediction

② Finding the optimal policy → model free control
                                    (Lecture 5).

estimate the value function of an unknown MDP.

## Model free prediction

→ Monte Carlo Learning:
- Learns directly from episodes of experience
- Look only at complete episodes
- Idea→ value = mean of the sample returns across many episodes.

Problem: Applicable to only episodic MDPs, ie. those where episodes which terminate.

# Monte Carlo Policy Evaluation:

→ Policy already given ⇒ Π

→ Learn $V_\Pi$ from episodes of experience following Π.

※ $S_1, A_1, R_2, \ldots S_k$ ~ Π.

Earlier: we had __expected__ return:

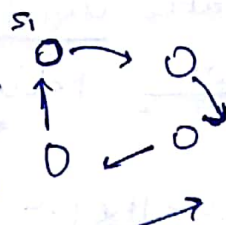$$V_\Pi(s) = E_\Pi(G_T \mid S_t = s).$$

But here, we'll have __empirical mean__ instead of expected.

Challenge: How to reset to a particular ~~step~~ state everytime and start our evaluation? And how to do this for every state?

2 methods:

① First visit: Imagine: circular MDP ie. Loop in an MDP.

Now, we consider only the time when a particular state, say $S_1$, is __1st visited__ in __an episode__.

After that, update the count of # visits in that state: $N(s) \pm = 1$

$$S(S) = S(3) + \boxed{G_t}$$

add reward from that time step onwards to the total reward for the state

estimated value for the state: $V(s) = \dfrac{S(s)}{N(s)}$

According to the law of large numbers, the mean of a large # samples actually converges to the true mean.

$$P\left(\left|\frac{X_1 + X_2 + \ldots X_n}{n} - \mu\right| \leq \varepsilon\right) \to 1$$

$X_1, X_2 \ldots X_n \to$ n samples of the dists

M → actual mean of the distribution

as $n \to \infty$, $\forall \varepsilon > 0$.

$\varepsilon \to$ any arbitrary small +ve value.

# Clarifications :

**1) Does # states decide the # episodes required?**

**No!**

From Central Limit Theorem, the variance decrease as a factor of $\frac{1}{n}$. Hence, it doesn't depend on the # states, just how many times a particular state is visited.

The trajectory in each episode should be such that it covers most of the states that we care about.

~~copy~~

**2) How to make sure that we reach all the states we care about?**

→ Problem of control. (next lecture).  make sure to

→ Here, we care about states under policy $\pi$ and we visit those states by following policy $\pi$.

**② Every visit** → instead of just the 1st visit, add on for every visit → i.e. multiple visits possible for each episode.

(Choice depends on the domain).

Incremental mean :

$$\mu_k = \frac{1}{k}\left(\sum_{i=1}^{k} x_i\right) = \frac{1}{k}\left[x_k + \sum_{i=1}^{k-1} x_k\right]$$

$$= \frac{1}{k}\left(x_k + (k-1)\mu_{k-1}\right) = \mu_{k-1} + \frac{1}{k}\underbrace{(x_k - \mu_{k-1})}_{\downarrow}$$

error term

Thus, we take a small step from our previous expectation following a error function.
↳ as difference bet 2 what we expected and what we actually got

# Incremental Monte Carlo updates:

$$N(s) = N(s) + 1.$$

$$V(s_t) = V(s_t) + \frac{1}{N(s)}(G_T - V(s_t)).$$

Better → add a weightage to discard old terms / values.

$$V(s_t) = V(s_t) + \alpha(G_T - V_{s_t})$$

$$= (1-\alpha) V(s_t) + \alpha \cdot G_T$$

We can weigh older terms accordingly.

# Temporal Difference Learning: (TD)

* Learn from incomplete episodes → bootstrapping
  → don't need to get to the end and figure out the reward, rather, go to an intermediate step and estimate the reward from there to the end.

So, we have to make an initial guess and after making / taking some steps, we make another guess of the final reward from that point. Then, we try to bring our initial guess closer to that guess.

Goal: Learn $V_\pi$ online from experience under $\pi$.

Here, we update $V(s_t)$ towards the ~~expected~~ estimated return:

Simple
TD learning: $V(s_t) = V(s_t) + (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$
TD(0)

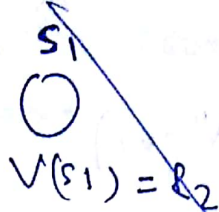$\underbrace{R_{t+1} + \gamma V(s_{t+1})}_{\text{TD target}}$

↳ TD error

## Advantages / Disadvantages
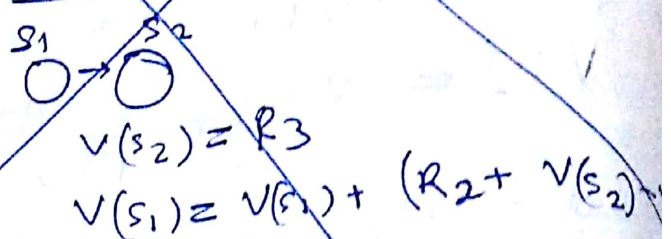
→ TD can learn before knowing the final outcome. ie. online after every step.

→ TD can learn without the final outcome.
  le. it can be applied to processes that do not end.
  hence, learning from      (episodes)
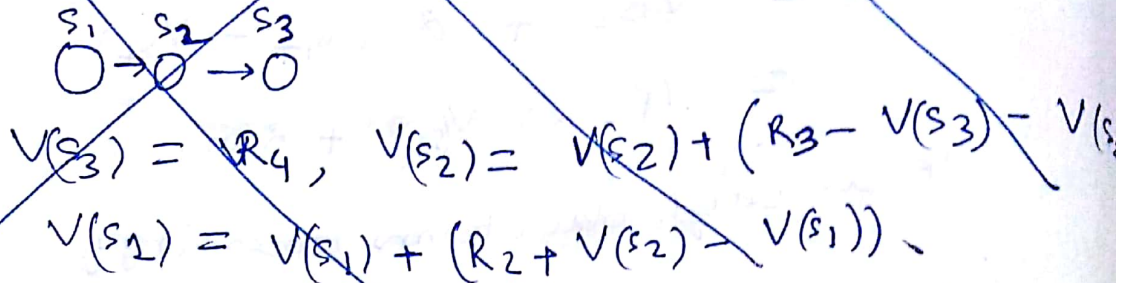           incomplete sequences.

## TD update:

Step 1 →

$S_1$

$\bigcirc$

$V(s_1) = R_2$

**Step 2:**

$S_1 \quad S_2$

$\bigcirc \rightarrow \bigcirc$

$V(s_2) = R_3$

$V(s_1) = V(s_1) + (R_2 + V(s_2))$

**Step 3:**

$S_1 \quad S_2 \quad S_3$

$\bigcirc \rightarrow \bigcirc \rightarrow \bigcirc$

$V(s_3) = R_4, \quad V(s_2) = V(s_2) + (R_3 - V(s_3)) - V(s$

$V(s_1) = V(s_1) + (R_2 + V(s_2) - V(s_1))$.

i.e. update after each step.

TD → less noise

$G_T = R_{t+1} + \gamma R_{t+2} + \ldots \quad \gamma^{T-1} R_T \rightarrow$ unbiased estimate of $V_\pi(s_t)$.

(bias of an estimate is given by:

$$Bias(\hat{\Theta}_m) = E(\hat{\Theta}_m) - \Theta, \quad \hat{\Theta} \rightarrow \text{unbiased estimator when}$$

$\hookrightarrow$ expectation

$E(\hat{\Theta}_m) = \Theta$ i.e.

$Bias(\hat{\Theta}_m) = 0$

$R_{t+1} + \gamma V_\pi(s_{t+1})$ unbiased estimate of $V_\pi(s_t)$.

But, we use: $R_{t+1} + \gamma V(s_{t+1}) \rightarrow$ biased estimate as.

we don't know the true future $V_\pi(s_{t+1})$.

But $G_T = R_{t+1} + R_{t+2}\gamma + \ldots R_T \gamma^{T-1}$

$\hookrightarrow$ adds noise at each stage. (term)

$R_{t+1} + \gamma V(s_{t+1}) \rightarrow$ noise just in one term.

Hence, TD → less noisy.

footer_navigationScanned by CamScanner

MC → high variance, 0 bias
  → Good convergence
  → independent g initializat⁰.

TD → some bias, low variance.
  → more efficient
  → $TD(0) \rightarrow V_\pi(s)$
  → sensitive to init value.

Batch MC / TD

Both converge as experience $\rightarrow 0$
In batch case,
  → repeatedly sample
    $k \in [1, K]$.
Do they converge?
(eg. run 3 episodes &
repeatedly iterate & learn
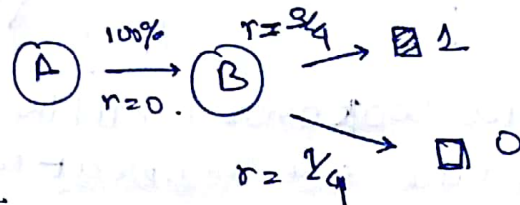from those 3).

eg. These are observed
  1  A 0     B 0
  2  B 1
    : :
  7  B 1
  8  B 0

Monte carlo for B = $\frac{6}{8}$

Monte carlo for A = $\frac{0}{1}$ •

TD estimate → tries to build an
MDP that best explains the data.



MC → reduces the mean
      squared error.
  • Bit fit to actual returns received → $\sum_{k=1}^{K} \sum_{t=1}^{T_k} \left( g_t^k - V(s_t^k) \right)$

TD → converges to the MDP that best
                        fits the data
  → sol² to max
      likelihood Markov model
  → jull by v counting
      normal

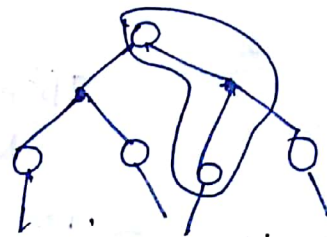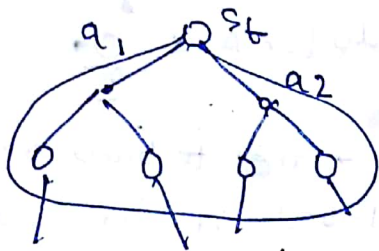$$P_{ss'}^{a} = \frac{1}{N(s,a)} \sum_{k=1}^{K} \sum_{t=1}^{T_k} I \left( s_t^k, a_t^k, s_{t+1}^k = s, a, s' \right)$$

$$R_{s}^{a} = \frac{1}{N(s,a)} \sum_{k=1}^{K} \sum_{t=1}^{T_k} I \left( s_t^k, a_t^k = s, a \right) r_t^k.$$

TD makes use of markov property → much more efficient
in those environments.
situations maybe non-markov.

ie. Partially observed $\overline{MDP}$ → under the hood, there
is an MDP, but what is observed is partial → so
TD(0) doesn't work well on such environments &
MC still finds a reasonable solution.

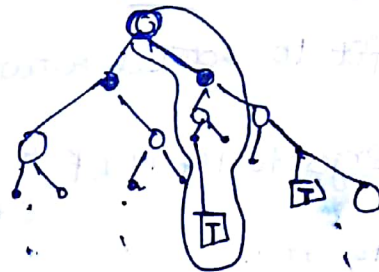In DP → we did a one-step lookahead over all possible
next states.                    whereas, foreg. in TD, we just
                                lookover one step, one
                                              sample



In MC, we lookover until the entire terminate
but only take one sample at each stage.

We could take all the way to
the bottom for DP too, which
would be called exhaustive
tree search



→ Bootstrap → don't use the real value but the estimated
value of the next step to update the current step.
   • Not by MC
   • By DP & TD

→ sampling → don't do fullwidth backup & only take
              samples.
   • MC & TD samples
   • DP doesn't.

Exhausti search → Bad!

Algoritims that are in the middle, full width & shallower
or
sample          deep
backups         backups.

$$TD(\lambda)$$

n-step predicts → look n steps into the future [instead of one in TD(0)]

$n=1$  → $TD(0)$:  $G_T^{(1)} = R_{t+1} + \gamma V(S_{t+1})$

$n=2$

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

$\vdots$

$n=\infty$  MC → $G_T^{(\infty)} = R_{t+1} + \cdots \cdots + R_T \gamma^{T-1}$

n-step return:

→  $G_T^{(n)} = R_{t+1} + \gamma R_{t+2} \cdots \quad R_{t+n} \gamma^{n-1} + \gamma^n V(S_t)$

n step TD

→  $V(S_t) = V(S_t) + \alpha.(G_t^{(n)} - V(S_t))$

what value of n to be used? → Hard to choose & depends on problem

online / offline updates
↓                        → Update the value only after reaching the end.
update the
value
immediately          Use average n-step returns.
                        ↙ over several n.

how to do this efficiently for all n.     eg  $\dfrac{G^{(2)} + G^{(4)}}{2}$ → get the best g all n.

$$TD(\lambda)$$

$\lambda \to 0 < \lambda < 1$   (λ) return → geometrically weighted average g all n

decaying for each specific λ.

(1-λ) → normalizing factor:  $(1-\lambda) + (1-\lambda)\lambda + \cdots \cdots$

sum  $= (1-\lambda)\dfrac{1}{(1-\lambda)} = 1$

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

$$TD(\lambda) \Rightarrow V(s_t) = V(s_t) + \alpha (G_T^\lambda - V(s_t))$$

Why geometric mean? → computational complexity

as
$$(1-\lambda), (1-\lambda)\lambda, (1-\lambda)\lambda^2 \leftarrow$$

→ memoryless → doesn't require storing separate values for each timestep!

(just store one value and keep multiplying $\lambda$).

This is the forward view $TD(\lambda)$ → needs to wait until completing the ~~future~~ episode → suffers from same disadvantages as M.C. $(\lambda = 1 \to MC)$

Backward view → $TD(\lambda)$ : compute the above stuff more efficiently without having to wait till the completion of the episode.
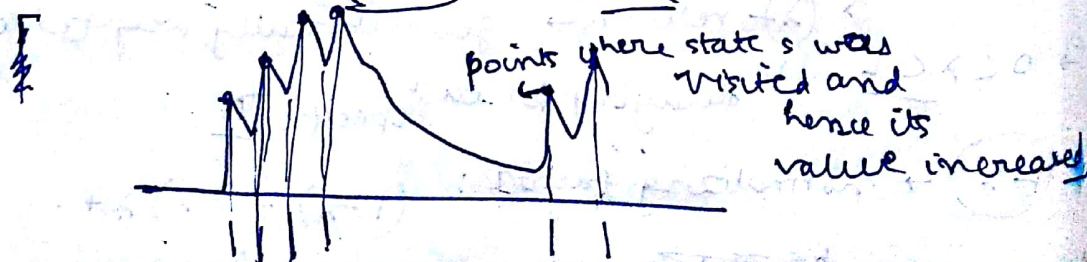
# Eligibility traces

→ Frequency heuristic → assign importance to most frequent states

→ Recency heuristic → importance to most recent states

eligibility trace → combines both

$$E_0(s) = 0$$

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

→ whenever state is visited, we increase the trace & if not visited, we (reduce) its value



points where state s was visited and hence its value increased

$I(S_t = s) \rightarrow$ Identity function $= \begin{cases} L & \text{if } S_t = s \\ o & \text{else.} \end{cases}$

For backward view:

→ Keep ET for each state

→ update $V(s)$ uin proportion to $E_t(s)$ & TD-error $S_t$

$S_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t).$

for each state.

$V(s) = V(s) + \alpha E_t(s) S_t$

→ $\lambda \rightarrow$ decaying factor.

→ $\lambda = 0 \rightarrow$ we squash the value to $0$ of $E(s)$ of $S_t \neq s$

$\therefore E_t(s) = I(S_t = s)$

only

$E$ TD(0) Thus, value function updated if We observe that $\downarrow$ visit
for a state in that
state time step.

$V(s) = V(s) + \alpha S_t E_t(s)$.

TD($\lambda$)

$\lambda = 1 \rightarrow$ same total update as MC.

Theorem → total updates same for backward as well as forward views.
offline

$\sum_{t=1}^{T} \alpha S_t E_t(s) = \sum_{t=1}^{T} \alpha (G_t^{\lambda} - V(S_t)) I(S_t = s)$

TD($\lambda$) → Spectrum. beth TD(0) & MC

$E_t(s) \rightarrow$ can be computed efficiently too.

Offline <u>online</u>
...step k.