

Optimal λ to use:

- tune up λ as a fⁿ of time?
- objective fⁿ to optimize (fⁿ of λ)

10/Nov/2018

Recap:

$$\underbrace{G_t - V(S_t)}_{\text{monte carlo error}} = \sum_{k=1}^{T-1} \gamma^{k-t} \delta_k \quad \text{--- (1)}$$

monte carlo error.

$$\text{Where } \delta_k = R_{k+1} + \gamma V(S_{k+1}) - V(S_k)$$

TD error.

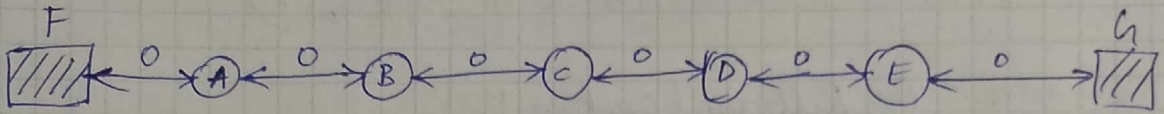
(1) holds provided V does not change with time.

DP-methods - value iteration / policy iteration.

TD	DP	Monte Carlo
- No System model	- Requires System model	- No system model
- Allows incremental updates	- Allows incremental updates	- Updates after completion of episode.

Although there is no theoretical proof, studies have shown TD is better than MC.

Eg: Random Walk.



F & G are terminating states. Once we reach terminating state, reward = 0

Markov Reward process: A MRP is an MDP w/o ^{actions} ~~rewards~~

We assume that all episodes start in the center state, C & proceed left or right by one step with equal prob

Episodes may terminate in either state F or state G.

When an episode terminates in G, a reward of 1 is obtained

When an episode terminates in F, a reward of 0 is obtained

Instance of a episode

C, 0, D, 0, C, 0, B, 0, C, 0, D, 0, E, 0, 1, G

We are in a undiscounted setting. True value of a state is the prob of terminating in state G when starting from that state.

Let π be our policy $\pi(\text{left} | \text{non-term state}) = \pi(\text{right} | \text{non-term state}) = 1/2$

Note $V_{\pi}(A) = E[R_1 + R_2 + \dots + R_T | S_0 = A]$

$= E[R_T | S_0 = A]$ $R_1 + R_2 + \dots + R_{T-1} = 0$ nonterminating state

If MRP terminates in F, $R_T = 0$

If MRP terminates in G, $R_T = 1$

$$E[R_T | S_0 = A] = E[I_{\{\text{termination in } G\}} | S_0 = A]$$

$$V_{\pi}(A) = P(\text{termination in } G \text{ given } S_0 = A)$$

We solve the Bellman Equation now

$$\begin{aligned} V_{\pi}(A) &= E[R_1 + V_{\pi}(\text{next state}) | S_0 = A] \\ &= 0 + \frac{1}{2} V_{\pi}(A) + \frac{1}{2} V_{\pi}(B) \\ &= 0 + 0 + \frac{1}{2} V_{\pi}(B) \end{aligned}$$

$$V_{\pi}(B) = 2 V_{\pi}(A) \quad \text{--- (1)}$$

$$\begin{aligned} V_{\pi}(B) &= E[R_1 + V_{\pi}(\text{next state}) | S_0 = B] \\ &= 0 + \frac{1}{2} V_{\pi}(A) + \frac{1}{2} V_{\pi}(C) \end{aligned}$$

$$\cancel{V_{\pi}(B)} = \left(2 - \frac{1}{2}\right) V_{\pi}(A) = \frac{1}{2} V_{\pi}(C)$$

$$3 \cancel{V_{\pi}(A)} = V_{\pi}(C) \quad \text{--- (2)}$$

$$\begin{aligned} V_{\pi}(C) &= E[R_1 + V_{\pi}(\text{next state}) | S_0 = C] \\ &= 0 + \frac{1}{2} V_{\pi}(B) + \frac{1}{2} V_{\pi}(D) \end{aligned}$$

$$3 V_{\pi}(A) = V_{\pi}(A) + \frac{1}{2} V_{\pi}(B)$$

$$4 V_{\pi}(A) = V_{\pi}(B) \quad \text{--- (3)}$$

$$\begin{aligned} V_{\pi}(D) &= E[R_1 + V_{\pi}(\text{next state}) | S_0 = D] \\ &= 0 + \frac{1}{2} V_{\pi}(C) + \frac{1}{2} V_{\pi}(E) \end{aligned}$$

$$4 V_{\pi}(A) = \frac{3}{2} V_{\pi}(A) + \frac{1}{2} V_{\pi}(E)$$

$$V_{\pi}(E) = 5 V_{\pi}(A)$$

$$V_{\pi}(E) = E[\pi_1 + V_{\pi}(\text{next state}) | S_0 = E]$$

$$= \frac{1}{2} + \frac{1}{2} V_{\pi}(A) + \frac{1}{2} V_{\pi}(B)$$

$$5 V_{\pi}(A) = \frac{1}{2} + 0 + \frac{1}{2} V_{\pi}(A)$$

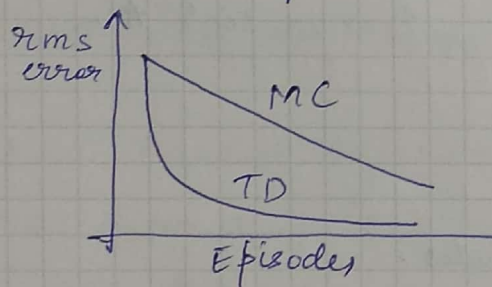
$$3 V_{\pi}(A) = \frac{1}{2} \Rightarrow V_{\pi}(A) = \frac{1}{6}$$

$$\therefore V_{\pi}(B) = \frac{2}{6}, V_{\pi}(C) = \frac{3}{6}, V_{\pi}(D) = \frac{4}{6}, V_{\pi}(E) = \frac{5}{6}$$

Indicate the prob of ending in a terminaly state

Apply TD & MC methods:

TD(0) = 100 episodes

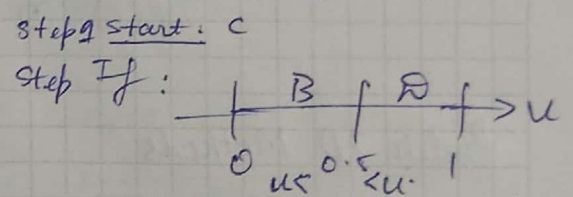


To Do : program.

$U \sim U[0, 1]$

step 1 start : c

step If :



TD has more oscillations,

MC has less oscillation.
since there aggregation

Example: Designing a prediction from data

Consider we have 8 different episodes

Episode 1 : A, 0, B, 0

5 : B, 1

2 : B, 1,

6 : B, 1

3 : B, 1

7 : B, 1

4 : B, 1

8 : B, 0

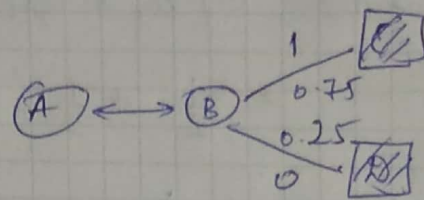
Predict ~~the~~ value of A?

Predict value of B?

$$V_{\pi}(B) = \frac{6}{8} = \frac{3}{4}$$

$$V_{\pi}(A) = \frac{3}{4}$$

$$\begin{aligned} V_{\pi}(A) &= E[r_t + V_{\pi}(\text{next, state})] = r_t + V_{\pi}(B) \\ &= 0 + \frac{3}{4} \\ &= \frac{3}{4} \end{aligned}$$



TD would predict:

$$V_{\pi}(A) = V_{\pi}(B) = \frac{3}{4}$$

MC would predict:

$$\begin{aligned} V_{\pi}(A) &= 0 \quad \text{Since it sees at } \text{step} = 1 = 0 \\ V_{\pi}(B) &= 0 \quad \text{It never sees A.} \end{aligned}$$

Batch TD Methods

Certainty Equivalence estimates:

Suppose we estimate transition prob $i \rightarrow j$ as

$$\left(\frac{\text{no of transitions } i \rightarrow j}{\text{Total no of transitions out of } i} \right)$$

* Expected reward = sample average of rewards on these transitions, then value function can be estimated assuming the above estimated model is the true model.

N_{ij} = no of transitions from $i \rightarrow j$

N_i = — — — — — out of i

$$P_{ij}^{(n)} = \frac{N_{ij}}{N_i}$$

R_i^n = n^{th} sample reward in state i

$$R_i = \frac{1}{N_i} \sum_{n=1}^{N_i} R_i^n$$

$$V(i) = R_i + \sum_j P_{ij} (V(j))$$

$$V(i) = \frac{1}{N_i} \sum_{n=1}^{N_i} R_i + \sum_j \frac{N_{ij}}{N_i} V(j)$$

Maximum Likelihood Parameter Estimation

Consider a random experiment whose we know the outcomes from a parameterized $P_\theta(\cdot)$ of distributions where θ is the parameter.

Goal: Estimate θ given that a certain outcome (k) has occurred.

Thus, parameter θ will maximize the prob of occurrence of ' k '

Let $\hat{\theta}_{ML}(k) \equiv$ parameter θ that maximises the likelihood $P_\theta(k)$ wrt θ

2 maximizations possible in general

- max likelihood wrt θ given observation k
- max likelihood wrt k given θ

Example (max posterior)
Mercedes-Benz R & D India

Consider a biased coin showing heads w.p p
suppose the coin is tossed n -times

Let X = no of heads in n -tosses.

$$\text{For } 0 \leq k \leq n, \rightarrow P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

Let \hat{p}_{ML} = value of p that maximises $P(X=k)$

$$\begin{aligned} \frac{dP(X=k)}{dp} &= \frac{d}{dp} \left(\binom{n}{k} p^k (1-p)^{n-k} \right) \\ &= \binom{n}{k} \left(k p^{k-1} (1-p)^{n-k} - (n-k) p^k (1-p)^{n-k-1} \right) \\ &= \binom{n}{k} \left(\frac{k}{p} - \frac{(n-k)}{(1-p)} \right) p^k (1-p)^{n-k} \\ &= \binom{n}{k} (k - np) p^{k-1} (1-p)^{n-k-1} \\ &= 0 \quad \text{for } p = \frac{k}{n} \end{aligned}$$

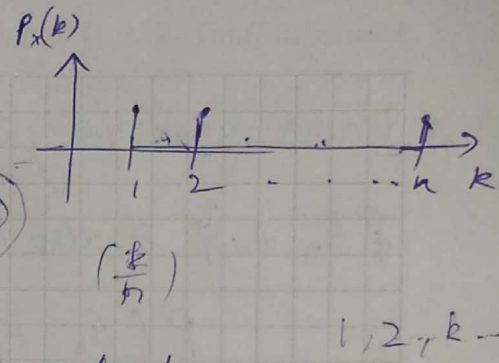
$$\begin{aligned} \left. \frac{d^2 P(X=k)}{dp^2} \right|_{p=k/n} &= -\binom{n}{k} n p^{k-1} (1-p)^{n-k-1} \\ &\quad + \binom{n}{k} (k - np) \frac{d}{dp} (p^{k-1} (1-p)^{n-k-1}) \\ &< 0 \end{aligned}$$

Example: (max likelihood)

Suppose X is drawn at random from numbers 1 to n with each possibility equally likely. Assume ' n ' is unknown but that $X=k$ is observed. Find the ML estimation of n given $X=k$ is observed.

Note: $p_x(k) = p(X=k)$

$$= \frac{1}{n} \mathbb{I}_{\{1 \leq k \leq n\}}$$



One can view

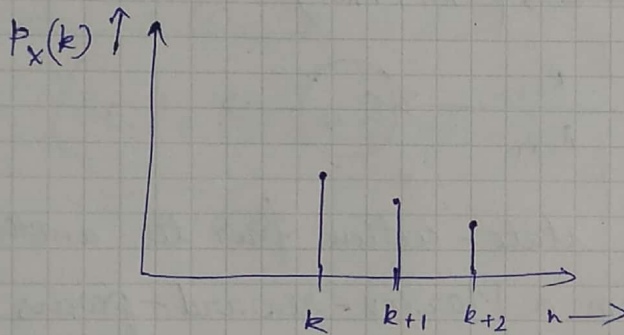
$p_x(k)$ as a function of n & not k

$$p_x(k) = 0 \quad \text{if} \quad n \leq k-1$$

$$= \frac{1}{k} \quad \text{if} \quad n = k$$

$$= \frac{1}{k+1} \quad \text{if} \quad n = k+1$$

$$\mathbb{I}_{\{1 \leq k \leq n\}}$$



The max value of n , $\hat{n}_{ml} = k$

$$(\theta = k)$$

TD - algorithm for prediction

SARSA - on-policy TD control.

Idea:

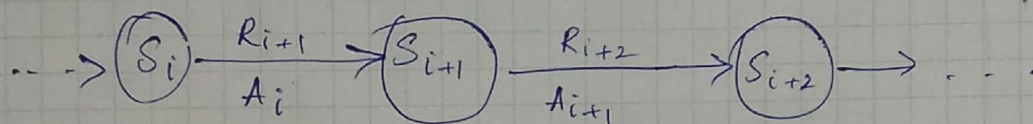
Use generalize policy iteration together with TD learning.

Under given policy π , the sequence of states $\{s_n\}$ is markov.

- Under given policy π , the sequence of state-action tuples $\{ (S_n, A_n) \}$ is also Markov.

$$\begin{aligned}
 & P(S_{n+1}=j, A_{n+1}=a \mid S_n=i, A_n=b) \\
 &= P(A_{n+1}=a \mid S_{n+1}=j, S_n=i, A_n=b) P(S_{n+1}=j \mid S_n=i, A_n=b) \\
 &= P(A_{n+1}=a \mid S_{n+1}=j) P(S_{n+1}=j \mid S_n=i, A_n=b) \\
 &= \pi(a \mid j) P_\pi(i, b, j)
 \end{aligned}$$

$$(S_n, A_n) \xrightarrow{\pi(A_{n+1} \mid S_{n+1}) P_\pi(S_n, A_n, S_{n+1})} (S_{n+1}, A_{n+1})$$



Transitions from one state-action pair to another state-action pair form a Markov-reward-process

- state to state ~~to~~ transitions (under a given policy) estimate value function $V^\pi(\cdot)$
- state to state-action functions (---) estimate action-value function $Q^\pi(\cdot, \cdot)$

TD(0) on the joint markov chain gives,

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

Assumption: All state action tuples are visited infinitely often.

$$Q(s, a) = R(s, a) + \gamma \sum_{(s', a')} p(s, a, s') \pi(a' | s') Q(s', a') - Q(s, a)$$

one stable fixed pt of DDE is $RHS = 0$

$$Q(s, a) = R(s, a) + \gamma \sum_{(s', a')} p(s, a, s') \pi(a' | s') Q(s', a')$$

Bellman eqn for the state-action markov chain.

The TD(0) update for the joint (state-action) Markov Chain is done after every transition from a non-terminal state s_t . If s_{t+1} is a terminal state, then $Q(s_{t+1}, a_{t+1}) \equiv 0$

$$\underbrace{(s_t, a_t, R_{t+1}, s_{t+1}, a_{t+1})}_{\text{SARSA algo}} \dots$$

- continually estimate q_π for the behavior policy π & then update π greedily w.r.t q_π [using GPI or some variant]

- Adopt and ϵ -greedy policy.

SARSA

Algo param: step size $\alpha \in (0, 1]$, small $\epsilon > 0$

Initialize $Q(s, a)$, for all $s \in S$, $a \in A(s)$ arbitrarily
expect that $Q(\text{terminal}, \cdot) = 0$

↓

loop for each episode

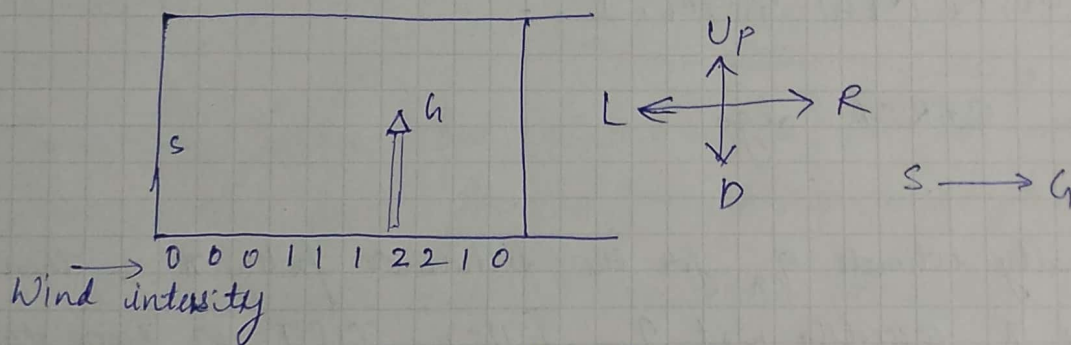
1. Initialize S
2. Choose A from S using policy derived from Q (eg: ^{greedy} policy)

loop for each step of episode:

- Take action A , observe R, S'
- Choose A' from S' using policy derived from Q (eg: ^{greedy} policy)
- $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
- $S \leftarrow S', A \leftarrow A'$

Until S is terminal state

Eg: Windy Grid World



the number indicates the offset

$\epsilon = 0.1$ $S \rightarrow G$ $\alpha = 0.5$, $Q(S, a) = 0$ $V(S, a)$ initialization

Each episode ~ 17 after 8000 steps

