

Reinforcement Learning:

MDP: Tuple $\langle S, A, P, R, \gamma \rangle$ | S : {States} A : {Actions}, P : transition prob $P_{s's'}^a = P(s_{t+1} = s' | s_t = s, A_t = a)$, R : Reward function $R_s = E(R_{t+1} | s_t = s, A_t = a)$, γ : discount factor

Return: Accumulated Reward for #steps \rightarrow total reward over trajectory, regardless if flip/stuck | G_t : discounted Reward, $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k}$ | $\gamma = 0$: near sighted, $\gamma = 1$: far sighted, $\gamma = 0.99$: typical

Policy Function: $\pi: S \rightarrow A$ w/ obj. to max G_t w/ $\pi^* = \arg \max_{\pi} E[\sum_{t=0}^{\infty} \gamma^t R_t | \pi] \rightarrow E(x)$ b/c prob policy take action & prob enter state

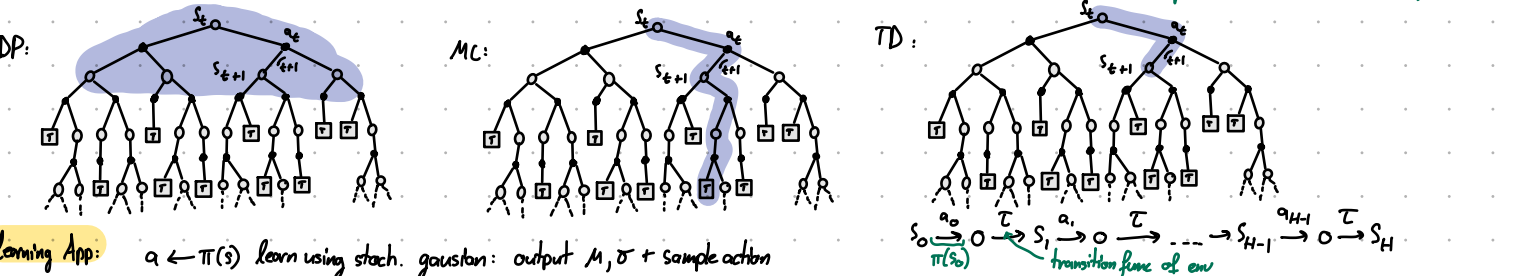
Value Functions: How good is a state s , w/ $V(s)$ = Expected Cumulative reward from following π from s | $V^*(s) = E[\sum_{t=0}^{\infty} \gamma^t r_t | s_t = s, \pi]$ | $\gamma = 0$: near sighted, $\gamma = 1$: far sighted, $\gamma = 0.99$: typical

Bellman: How to learn Value Function for whole env. $V^*(s) = E_x[r_t + \sum_{t=1}^{\infty} \gamma^t r_t | s_0 = s, \pi] = E_x[r_t + \gamma V(s_{t+1}) | s_0 = s, \pi]$ | $V_{k+1}(s) = \sum_{a \in A} \pi(a|s) [R_s^a + \gamma \sum_{s'} P_{ss'}^a V_k(s')]$ | $P_{ss'}^a$: Prob trans $s \rightarrow s'$ using action a , R_s^a : immediate reward taking action a in state s , $V_k(s')$: future reward state s' , $\pi(a|s)$: prob take a in state s

Monte Carlo: $V_0 \rightarrow V_1 \rightarrow \dots \rightarrow V_{\infty}$. Check all actions, following 1 policy until terminates. $V_{k+1} = V_k(s_t) + \alpha (G_t - V_k(s_t)) \rightarrow$ Consider avg return for all traj. start at s based on rolling avg

$$V_{k+1}(s) = \frac{1}{T} \sum_{i=1}^T G_i = \frac{1}{T} G_T + \frac{1}{T} \sum_{i=1}^{T-1} G_i = \frac{1}{T} G_T + \frac{T-1}{T} \frac{1}{T-1} \sum_{i=1}^{T-1} G_i = \frac{1}{T} G_T + \frac{T-1}{T} \frac{1}{T-1} \sum_{i=1}^{T-1} (G_i - V_k(s_t) + V_k(s_t)) \Rightarrow V_{k+1}(s) = V_k(s) + \alpha (G_T - V_k(s))$$

Temporal Distance: Execute 1 action, know value next state, then backup: $V_k(s_t) \leftarrow V_k(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$ | Equivalent to $G(t)$: im reward + disc. factor future



Learning App: $a \leftarrow \pi(s)$ learn using stoch. gaussian: output μ, σ + sample action

VPGL: Initialize π_{θ} w/ θ for $j=1$ to H :
 rollouts \leftarrow Execute π_{θ} m times, get $\{s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \dots, s_T^i\}$
 rollouts \leftarrow Replace rew. w/ ret, get $\{s_0^i, a_0^i, R_0^i, s_1^i, a_1^i, R_1^i, \dots, s_T^i\}$
 $b \leftarrow$ (old Baseline w/ $\frac{1}{m} \sum_{i=1}^m R(t)$)
 rollouts \leftarrow Replace R w/ $A = R - b$, get: $\{s_0^i, a_0^i, A_0^i, s_1^i, a_1^i, A_1^i, \dots, s_T^i\}$
 $\theta \leftarrow \theta - \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i, s_t^i) A_t^i$ A : advantage

RL Gradient: $a \leftarrow \pi(s)$ learn using stoch. gaussian: output μ, σ + sample action

$$\max_{\theta} \text{reward} = E[\sum_{t=0}^H R(s_t, a_t) | \pi_{\theta}] = \max_{\theta} \sum_{\tau} P(\tau, \theta) R(\tau) = \sum_{\tau} \nabla_{\theta} P(\tau, \theta) R(\tau) = \sum_{\tau} P(\tau, \theta) \nabla_{\theta} \ln(P(\tau, \theta)) R(\tau) = E_{P(\tau, \theta)} [\sum_{t=0}^{H-1} \log \pi(a_t | s_t, \theta)] R(\tau)$$

$$P(\tau, \theta) = \prod_{i=1}^{H-1} \tau(s_{t+1} | a_t, s_t) \pi(a_t | s_t, \theta) \Rightarrow \nabla_{\theta} \log(P(\tau, \theta)) = \sum_{i=1}^{H-1} \nabla_{\theta} \log(\tau(s_{t+1} | a_t, s_t)) + \sum_{i=1}^{H-1} \nabla_{\theta} \log(\pi(a_t | s_t, \theta))$$