

RL Cheat Sheet		Value Function $V(s)$ <ul style="list-style-type: none"> Long term value of state S State value function $V(s)$ of a MRP is expected reward from state s $V(s) = E(G_t S_t = s)$
Definitions <ul style="list-style-type: none"> State(S): Current Condition Reward(R): Instant Return from environment to appraise the last action Value(V): Expected Long-term reward with discount, as opposed to short-term reward R Action-Value(Q): This is similar to Value, except, it takes extra parameter, action A Policy(π): Approach of agent to determine next action based on current state Exploitation is about using already known info to maximize rewards Exploration is about exploring and capturing more information 		Action Value Function $q(s, a)$ <ul style="list-style-type: none"> $q_\pi(s, a) = E_\pi(G_t S_t = s, A_t = a)$
Discount Factor (γ) <ul style="list-style-type: none"> Varies between 0 to 1 Closer to 0 \rightarrow Agent tend to consider immediate reward Closer to 1 \rightarrow Agent tend to consider future reward with greater 		Bellmen Equation <ul style="list-style-type: none"> $V(s) = R(s) + \gamma E_{s' \in S} [V(s')]$ $V(s) = R(s) + \gamma \sum_{s' \in S} P_{ss'}(V(s'))$ $V = R + \gamma PV \rightarrow V = (I - \gamma P)^{-1} R$
Q-Learning <ol style="list-style-type: none"> Create reward matrix R where R_{sa} = reward for taking action a in state s and set γ parameter. Initialize Q matrix to 0 Set initial random state and assign this to current state Select One among all possible actions of current state <ol style="list-style-type: none"> Use this action to get new state Get maximum Q value for this state based on all previous actions Compute Q matrix using $Q_{sa} = R_{sa} + \gamma \times \text{Max}[Q_{s'a'}]$ $\forall s' \text{ accessible from } s$ $\forall a' \text{ available in } s'$ Repeat 4 until current state=goal state 		Markov Process <ul style="list-style-type: none"> Consists of $\langle s, p \rangle$ tuple where s are states and p is state transition matrix $P_{ss'} = P(s_{t+1} = s' s_t = s)$ $\mu_{t+1} = p^T \mu_t$ where $\mu_t = [\mu_{t,1} \dots \mu_{t,n}]^T$ Markov Reward Process <ul style="list-style-type: none"> Consists of $\langle s, p, R, \gamma \rangle$ tuple where R is reward γ is discount $R = E[R_{t+1} S_t = S] = R(s)$ $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0} \gamma^k R_{t+k+1}$ is Total discounted reward <p>Discounted Reward cons: Uncertainty may not be fully represented. Immediate rewards values > delayed, Avoid ∞ rewards in cycle</p> Markov Decision Process <ul style="list-style-type: none"> Consists of $\langle s, A, p, R, \gamma \rangle$ tuple where A is action $P_{ss'} = P(S_{t+1} = S' S_t = S, A_t = a)$ <p>Discounted Reward cons: Uncertainty may not be fully represented. Immediate rewards values > delayed, Avoid ∞ rewards in cycle</p>
Monte Carlo Policy Evaluation <ol style="list-style-type: none"> To evaluate Value(s) $V_\pi(S)$ At any time step t when state s is visited in an episode <ol style="list-style-type: none"> Increment Counter N(s) <- N(s)+1 Increment total return S(s) <- S(s)+G_t Value estimated is mean $V(s) = S(s)/N(s)$ 		Policy <ul style="list-style-type: none"> $\pi(a s) = P(a_t = a s_t = s)$ Either deterministic or stochastics. In deterministic P=1 for one a_t $P_\pi(s' s) = \sum \pi(a s) \times P(s' s, a)$ for stochastic process. One step expected reward $r_\pi = \sum_a \pi(a s) r(s, a)$ For rewards as function of transition states $r_\pi = \sum_a \pi(a s) \sum_s \pi(a s) \times \sum_{s'} P(s' s, a) \times r(s, a, s')$
Policy Gradients <ul style="list-style-type: none"> $p_\theta(s_1, a_1, s_2, a_2 \dots) = p(s_1) \times \prod_{t=1}^T p(s_{t+1} s_t, a_t) \pi_\theta(s_t, a_t)$ Goal is to $\theta^* = \arg \max_{\theta} E_{\tau \sim p_\theta(\tau)} [\sum_t r(s_t, a_t)] = \arg \max J(\theta)$ $J(\theta) = \frac{1}{N} \sum_i \sum_t r(s_i, a_i)$ $\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \times \nabla_\theta \log p_\theta(\tau)$ $\nabla_\theta J_\theta = \frac{1}{N} \sum_{i=1}^N \{ \sum_{t=1}^N \nabla_\theta \log \pi_\theta(a_{i,t} s_{i,t}) \sum_{t=1}^N r(s_{i,t}, a_{i,t}) \}$ $\theta \leftarrow \theta + \nabla_\theta J(\theta)$ $\log \pi_\theta(a_{i,t} s_{i,t})$ is the log probability of action, defines how likely are we going to see $a_{i,t}$ as action 		Relation Between V_π and q_π <ul style="list-style-type: none"> $V_\pi(s) = \sum_{a \in A} \pi(a s_t = s) q_\pi(s, a)$ $V_\pi(s) = \sum_{a \in A} \pi(a s) \times \{ r(s, a) + \gamma \times \sum_{s' \in S} P(s' s, a) V_\pi(s') \}$ $V_\pi(s) = r(s) + \gamma \sum_{a \in A} \pi(a s) \sum p(s' s, a) \times V_\pi(s')$ $q_\pi(s, a) = r(s, a) + \gamma \times \sum_{s' \in S} P(s' s, a) V_\pi(s)$ $q_\pi(s, a) = r(s, a) + \gamma \times \sum P(s' s, a) \{ \sum_{a' \in A} \pi(a' s') q_\pi(s', a') \}$ $q_\pi(s, a) = r(s, a) + \sum_{s \in S'} p(s' s, a) \sum_{a' \in A} q_\pi(s', a') \times P(a' s')$
Actor Critic <ul style="list-style-type: none"> Q-V = Advantage $Q^\pi(s_t, a_t) = \sum_{t'=t}^T E_{\pi_0} [r(s_t, a_t) s_t, a_t]$ Reward of action a_t in s_t $V^\pi(s_t) = E_{a_t \sim \pi_\theta(a_t s_t)} [Q^\pi(s_t, a_t)]$ total reward from st $A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$ how much better A_t is $\nabla_\theta J(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t} s_{i,t}) A^\pi(s_t, a_t)$ $Q^\pi(s_t, a_t) = r(s_t, a_t) + \sum_{t'=t+1}^T E_{\pi_\theta} [r(s'_t, a'_t) s_t, a_t]$ $= r(s_t, a_t) + E_{\pi_{t+1} \sim p(s_{t+1} s_t, a_t)} [V^\pi(s_{t+1})]$ $A^\pi(s_t, a_t) = r(s_t, a_t) + V^\pi(s_{t+1}) - V^\pi(s_t)$ 		Optimality Condition <ul style="list-style-type: none"> $V_\pi^*(s) = \max_{\pi} V_\pi(s) \forall s \in S$ similarly for $q_\pi^*(s, a)$