

Summary of Notation

Capital letters are used for random variables and major algorithm variables. Lower case letters are used for the values of random variables and for scalar functions. Quantities that are required to be real-valued vectors are written in bold and in lower case (even if random variables).

s	state
a	action
\mathcal{S}	set of all nonterminal states
\mathcal{S}^+	set of all states, including the terminal state
$\mathcal{A}(s)$	set of actions possible in state s
\mathcal{R}	set of possible rewards
t	discrete time step
T	final time step of an episode
S_t	state at t
A_t	action at t
R_t	reward at t , dependent, like S_t , on A_{t-1} and S_{t-1}
G_t	return (cumulative discounted reward) following t
$G_t^{(n)}$	n -step return (Section 7.1)
G_t^λ	λ -return (Section 7.2)
π	policy, decision-making rule
$\pi(s)$	action taken in state s under <i>deterministic</i> policy π
$\pi(a s)$	probability of taking action a in state s under <i>stochastic</i> policy π
$p(s', r s, a)$	probability of transitioning to state s' , with reward r , from s, a
$v_\pi(s)$	value of state s under policy π (expected return)
$v_*(s)$	value of state s under the optimal policy
$q_\pi(s, a)$	value of taking action a in state s under policy π
$q_*(s, a)$	value of taking action a in state s under the optimal policy
$V_t(s)$	estimate (a random variable) of $v_\pi(s)$ or $v_*(s)$
$Q_t(s, a)$	estimate (a random variable) of $q_\pi(s, a)$ or $q_*(s, a)$
$\hat{v}(s, \mathbf{w})$	approximate value of state s given a vector of weights \mathbf{w}
$\hat{q}(s, a, \mathbf{w})$	approximate value of state-action pair s, a given weights \mathbf{w}
\mathbf{w}, \mathbf{w}_t	vector of (possibly learned) <i>weights</i> underlying an approximate value function
$\mathbf{x}(s)$	vector of features visible when in state s
$\mathbf{w}^\top \mathbf{x}$	inner product of vectors, $\mathbf{w}^\top \mathbf{x} = \sum_i w_i x_i$; e.g., $\hat{v}(s, \mathbf{w}) = \mathbf{w}^\top \mathbf{x}(s)$

δ_t	temporal-difference error at t (a random variable, even though not upper case)
$E_t(s)$	eligibility trace for state s at t
$E_t(s, a)$	eligibility trace for a state–action pair
\mathbf{e}_t	eligibility trace vector at t
γ	discount-rate parameter
ε	probability of random action in ε -greedy policy
α, β	step-size parameters
λ	decay-rate parameter for eligibility traces