

## Summary of Notation

Capital letters are used for random variables, whereas 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). Matrices are bold capitals.

$\backslash\defeq$	$\doteq$	equality relationship that is true by definition
$\backslashapprox$	$\approx$	approximately equal
$\backslashpropto$	$\propto$	proportional to
$\backslashPr\{X\!=\!x\}$	$\Pr\{X=x\}$	probability that a random variable $X$ takes on the value $x$
$X\sim p$	$X \sim p$	random variable $X$ selected from distribution $p(x) \doteq \Pr\{X=x\}$
$\mathbb{E}\{X\}$	$\mathbb{E}[X]$	expectation of a random variable $X$ , i.e., $\mathbb{E}[X] \doteq \sum_x p(x)x$
$\arg\max_a f(a)$	$\arg\max_a f(a)$	a value of $a$ at which $f(a)$ takes its maximal value
$\ln x$	$\ln x$	natural logarithm of $x$
$e^x$	$e^x$	the base of the natural logarithm, $e \approx 2.71828$ , carried to power $x$
$\mathbb{R}$	$\mathbb{R}$	set of real numbers
$f:\mathcal{X}\rightarrow\mathcal{Y}$	$f:\mathcal{X}\rightarrow\mathcal{Y}$	function $f$ from elements of set $\mathcal{X}$ to elements of set $\mathcal{Y}$
$\leftarrow$	$\leftarrow$	assignment
$(a,b]$	$(a,b]$	the real interval between $a$ and $b$ including $b$ but not including $a$
$\varepsilon$	$\varepsilon$	probability of taking a random action in an $\varepsilon$ -greedy policy
$\alpha, \beta$	$\alpha, \beta$	step-size parameters
$\gamma$	$\gamma$	discount-rate parameter
$\lambda$	$\lambda$	decay-rate parameter for eligibility traces
$\mathbb{I}_{\text{predicate}}$	$\mathbb{I}_{\text{predicate}}$	indicator function ( $\mathbb{I}_{\text{predicate}} \doteq 1$ if the <i>predicate</i> is true, else 0)
In a multi-arm bandit problem:		
$k$	$k$	number of actions (arms)
$t$	$t$	discrete time step or play number
$q_*(a)$	$q_*(a)$	true value (expected reward) of action $a$
$Q_t(a)$	$Q_t(a)$	estimate at time $t$ of $q_*(a)$
$N_t(a)$	$N_t(a)$	number of times action $a$ has been selected up prior to time $t$
$H_t(a)$	$H_t(a)$	learned preference for selecting action $a$ at time $t$
$\pi_t(a)$	$\pi_t(a)$	probability of selecting action $a$ at time $t$
$\bar{R}_t$	$\bar{R}_t$	estimate at time $t$ of the expected reward given $\pi_t$
In a Markov Decision Process:		
$s, s'$	$s, s'$	states
$a$	$a$	an action
$r$	$r$	a reward
$S$	$S$	set of all nonterminal states
$S^+$	$S^+$	set of all states, including the terminal state
$\mathcal{A}(s)$	$\mathcal{A}(s)$	set of all actions available in state $s$
$\mathcal{R}$	$\mathcal{R}$	set of all possible rewards, a finite subset of $\mathbb{R}$
$\subset$	$\subset$	subset of; e.g., $\mathcal{R} \subset \mathbb{R}$

$\$ \backslash \text{in} \$$	$\in$	is an element of; e.g., $s \in \mathcal{S}$ , $r \in \mathcal{R}$
$\$ \backslash \mathcal{S}   \$$	$ \mathcal{S} $	number of elements in set $\mathcal{S}$
$\$ t \$$	$t$	discrete time step
$\$ T, T(t) \$$	$T, T(t)$	final time step of an episode, or of the episode including time $t$
$\$ A\_t \$$	$A_t$	action at time $t$
$\$ S\_t \$$	$S_t$	state at time $t$ , typically due, stochastically, to $S_{t-1}$ and $A_t$
$\$ R\_t \$$	$R_t$	reward at time $t$ , typically due, stochastically, to $S_{t-1}$ and $A_t$
$\$ \backslash \pi \$$	$\pi$	policy (decision-making rule)
$\$ \backslash \pi(s) \$$	$\pi(s)$	action taken in state $s$ under <i>deterministic</i> policy $\pi$
$\$ \backslash \pi(a s) \$$	$\pi(a s)$	probability of taking action $a$ in state $s$ under <i>stochastic</i> policy $\pi$
$\$ G\_t \$$	$G_t$	return following time $t$
$\$ h \$$	$h$	horizon, the time step one looks up to in a forward view
$\$ G_{\{t:t+n\}}, G_{\{t:h\}} \$$	$G_{t:t+n}, G_{t:h}$	$n$ -step return from $t+1$ to $t+n$ , or to $h$ (discounted and corrected)
$\$ \bar{G}_{\{t:h\}} \$$	$\bar{G}_{t:h}$	flat return (undiscounted and uncorrected) from $t+1$ to $h$
$\$ G^{\backslash \lambda}_{\text{bda}_t} \$$	$G_t^\lambda$	$\lambda$ -return
$\$ G^{\backslash \lambda}_{\text{bda}_{\{t:h\}}} \$$	$G_{t:h}^\lambda$	truncated, corrected $\lambda$ -return
$\$ G^{\backslash \lambda}_{\text{bda}_t^s, a} \$$	$G_t^{s,a}$	$\lambda$ -return, corrected by estimated state, or action, values
$\$ \backslash p(s', r   s, a) \$$	$p(s', r   s, a)$	probability of transition to state $s'$ with reward $r$ , from state $s$ taking action $a$
$\$ \backslash p(s'   s, a) \$$	$p(s'   s, a)$	probability of transition to state $s'$ , from state $s$ taking action $a$
$\$ r(s, a) \$$	$r(s, a)$	expected immediate reward from state $s$ after action $a$
$\$ r(s, a, s') \$$	$r(s, a, s')$	expected immediate reward on transition from $s$ to $s'$ under policy $\pi$
$\$ \backslash v_\pi(s) \$$	$v_\pi(s)$	value of state $s$ under policy $\pi$ (expected return)
$\$ \backslash v_\star(s) \$$	$v_\star(s)$	value of state $s$ under the optimal policy
$\$ \backslash q_\pi(s, a) \$$	$q_\pi(s, a)$	value of taking action $a$ in state $s$ under policy $\pi$
$\$ \backslash q_\star(s, a) \$$	$q_\star(s, a)$	value of taking action $a$ in state $s$ under the optimal policy
$\$ V, V\_t \$$	$V, V_t$	array estimates of state-value function $v_\pi$ or $v_\star$
$\$ Q, Q\_t \$$	$Q, Q_t$	array estimates of action-value function $q_\pi$ or $q_\star$
$\$ \bar{V}_t(s) \$$	$\bar{V}_t(s)$	expected approximate action value, e.g., $\bar{V}_t(s) \doteq \sum_a \pi(a s) Q_t(s, a)$
$\$ U\_t \$$	$U_t$	target for estimate at time $t$
$\$ \delta\_t \$$	$\delta_t$	temporal-difference (TD) error at $t$ (a random variable)
$\$ \delta^s\_t, \delta^a\_t \$$	$\delta_t^s, \delta_t^a$	state- and action-specific forms of the TD error
$\$ n \$$	$n$	in $n$ -step methods, $n$ is the number of steps of bootstrapping
$\$ d \$$	$d$	dimensionality—the number of components of $\mathbf{w}$
$\$ d' \$$	$d'$	alternate dimensionality—the number of components of $\boldsymbol{\theta}$
$\$ \mathbf{w}, \mathbf{w}_t \$$	$\mathbf{w}, \mathbf{w}_t$	$d$ -vector of weights underlying an approximate value function
$\$ \backslash > \$ w_i, w_{\{t,i\}} \$$	$w_i, w_{t,i}$	$i$ th component of learnable weight vector
$\$ \hat{v}(s, \mathbf{w}) \$$	$\hat{v}(s, \mathbf{w})$	approximate value of state $s$ given weight vector $\mathbf{w}$
$\$ v_{\mathbf{w}}(s) \$$	$v_{\mathbf{w}}(s)$	alternate notation for $\hat{v}(s, \mathbf{w})$

$\hat{q}(s, a, \mathbf{w})$	$\hat{q}(s, a, \mathbf{w})$	approximate value of state–action pair $s, a$ given weight vector $\mathbf{w}$
$\nabla \hat{v}(s, \mathbf{w})$	$\nabla \hat{v}(s, \mathbf{w})$	column vector of partial derivatives of $\hat{v}(s, \mathbf{w})$ with respect to $\mathbf{w}$
$\nabla \hat{q}(s, a, \mathbf{w})$	$\nabla \hat{q}(s, a, \mathbf{w})$	column vector of partial derivatives of $\hat{q}(s, a, \mathbf{w})$ with respect to $\mathbf{w}$
$\mathbf{x}(s)$	$\mathbf{x}(s)$	vector of features visible when in state $s$
$\mathbf{x}(s, a)$	$\mathbf{x}(s, a)$	vector of features visible when in state $s$ taking action $a$
$x_i(s), x_i(s, a)$	$x_i(s), x_i(s, a)$	$i$ th component of vector $\mathbf{x}(s)$ or $\mathbf{x}(s, a)$
$\mathbf{x}_t$	$\mathbf{x}_t$	shorthand for $\mathbf{x}(S_t)$ or $\mathbf{x}(S_t, A_t)$
$\mathbf{w}^\top \mathbf{x}$	$\mathbf{w}^\top \mathbf{x}$	inner product of vectors, $\mathbf{w}^\top \mathbf{x} \doteq \sum_i w_i x_i$ ; e.g., $\hat{v}(s, \mathbf{w}) \doteq \mathbf{w}^\top \mathbf{x}(s)$
$\mathbf{v}, \mathbf{v}_t$	$\mathbf{v}, \mathbf{v}_t$	secondary $d$ -vector of weights, used to learn $\mathbf{w}$
$\mathbf{z}_t$	$\mathbf{z}_t$	$d$ -vector of eligibility traces at time $t$
$\boldsymbol{\theta}, \boldsymbol{\theta}_t$	$\boldsymbol{\theta}, \boldsymbol{\theta}_t$	parameter vector of target policy
$\pi(a s, \boldsymbol{\theta})$	$\pi(a s, \boldsymbol{\theta})$	probability of taking action $a$ in state $s$ given parameter vector $\boldsymbol{\theta}$
$\pi_{\boldsymbol{\theta}}$	$\pi_{\boldsymbol{\theta}}$	policy corresponding to parameter $\boldsymbol{\theta}$
$\nabla \pi(a s, \boldsymbol{\theta})$	$\nabla \pi(a s, \boldsymbol{\theta})$	column vector of partial derivatives of $\pi(a s, \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$
$J(\boldsymbol{\theta})$	$J(\boldsymbol{\theta})$	performance measure for the policy $\pi_{\boldsymbol{\theta}}$
$\nabla J(\boldsymbol{\theta})$	$\nabla J(\boldsymbol{\theta})$	column vector of partial derivatives of $J(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$
$h(s, a, \boldsymbol{\theta})$	$h(s, a, \boldsymbol{\theta})$	preference for selecting action $a$ in state $s$ based on $\boldsymbol{\theta}$
$b(a s)$	$b(a s)$	behavior policy used to select actions while learning about $\pi$
$b(s)$	$b(s)$	a baseline function $b : \mathcal{S} \mapsto \mathbb{R}$ for policy-gradient methods
$b$	$b$	branching factor for an MDP or search tree
$\rho_{t:h}$	$\rho_{t:h}$	importance sampling ratio for time $t$ through time $h$
$\rho_t$	$\rho_t$	importance sampling ratio for time $t$ alone, $\rho_t \doteq \rho_{t:t}$
$r(\pi)$	$r(\pi)$	average reward (reward rate) for policy $\pi$
$\bar{R}_t$	$\bar{R}_t$	estimate of $r(\pi)$ at time $t$
$\mu(s)$	$\mu(s)$	on-policy distribution over states
$\boldsymbol{\mu}$	$\boldsymbol{\mu}$	$ \mathcal{S} $ -vector of the $\mu(s)$ for all $s \in \mathcal{S}$
$\ v\ _\mu^2$	$\ v\ _\mu^2$	$\mu$ -weighted squared norm of value function $v$ , i.e., $\ v\ _\mu^2 \doteq \sum_s \mu(s) v(s)^2$
$\eta(s)$	$\eta(s)$	expected number of visits to state $s$ per episode
$\Pi$	$\Pi$	projection operator for value functions
$B_\pi$	$B_\pi$	Bellman operator for value functions
$\mathbf{A}$	$\mathbf{A}$	$d \times d$ matrix $\mathbf{A} \doteq \mathbb{E}[\mathbf{x}_t(\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top]$
$\mathbf{b}$	$\mathbf{b}$	$d$ -dimensional vector $\mathbf{b} \doteq \mathbb{E}[R_{t+1} \mathbf{x}_t]$
$\mathbf{w}_{\text{TD}}$	$\mathbf{w}_{\text{TD}}$	TD fixed point $\mathbf{w}_{\text{TD}} \doteq \mathbf{A}^{-1} \mathbf{b}$ (a $d$ -vector)
$\mathbf{I}$	$\mathbf{I}$	identity matrix
$\mathbf{P}$	$\mathbf{P}$	$ \mathcal{S}  \times  \mathcal{S} $ matrix of state-transition probabilities under $\pi$
$\mathbf{D}$	$\mathbf{D}$	$ \mathcal{S}  \times  \mathcal{S} $ diagonal matrix with $\boldsymbol{\mu}$ on its diagonal
$\mathbf{X}$	$\mathbf{X}$	$ \mathcal{S}  \times d$ matrix with the $\mathbf{x}(s)$ as its rows
$\bar{\delta}_{\mathbf{w}}(s)$	$\bar{\delta}_{\mathbf{w}}(s)$	Bellman error (expected TD error) for $v_{\mathbf{w}}$ at state $s$
$\bar{\delta}_{\mathbf{w}}, \text{BE}$	$\bar{\delta}_{\mathbf{w}}, \text{BE}$	Bellman error vector, with components $\bar{\delta}_{\mathbf{w}}(s)$

$\$ \backslash \text{MSVE}(\backslash \mathbf{w}) \$$	$\overline{\text{VE}}(\mathbf{w})$	mean square value error $\overline{\text{VE}}(\mathbf{w}) \doteq \ v_{\mathbf{w}} - v_{\pi}\ _{\mu}^2$
$\$ \backslash \text{MSBE}(\backslash \mathbf{w}) \$$	$\overline{\text{BE}}(\mathbf{w})$	mean square Bellman error $\overline{\text{BE}}(\mathbf{w}) \doteq \ \bar{\delta}_{\mathbf{w}}\ _{\mu}^2$
$\$ \backslash \text{MSPBE}(\backslash \mathbf{w}) \$$	$\overline{\text{PBE}}(\mathbf{w})$	mean square projected Bellman error $\overline{\text{PBE}}(\mathbf{w}) \doteq \ \Pi \bar{\delta}_{\mathbf{w}}\ _{\mu}^2$
$\$ \backslash \text{MSTDE}(\backslash \mathbf{w}) \$$	$\overline{\text{TDE}}(\mathbf{w})$	mean square temporal-difference error $\overline{\text{TDE}}(\mathbf{w}) \doteq \mathbb{E}_b[\rho_t \delta_t^2]$
$\$ \backslash \text{MSRE}(\backslash \mathbf{w}) \$$	$\overline{\text{RE}}(\mathbf{w})$	mean square return error