

# RL Lecture: Notation Cheat Sheet

M. Lindauer – Leibniz University Hannover

Symbol	Meaning
$S$	space of states
$s \in S$	one specific state
$A$	space of Actions ( <b>Warning:</b> overloaded notation with advantage function)
$a \in A$	one specific action
$P : S \times A \rightarrow S$	dynamics of environment
$R : S \times A \rightarrow \mathbb{R}$	reward function
$r \in \mathbb{R}$	concrete reward value
$\gamma \in [0, 1]$	discount factor
$T$	maximal time horizon ( <b>Warning:</b> overloaded notation with terminal states)
$t$	concrete time step $t \leq T$
$\rho_0 : S \rightarrow \mathbb{R}^+$	a distribution of start states
$\pi : S \rightarrow A$	policy (Note: could also be defined as $S \times A \rightarrow \mathbb{R}$ to highlight non-deterministic behavior)
$\pi^* : S \rightarrow A$	optimal policy
$G_t : S \rightarrow \mathbb{R}$	discounted sum of rewards from time step $t$ to horizon (also with $S \times A$ possible)
$G_t^{(n)} : S \rightarrow \mathbb{R}$	general $n$ -step return
$V^\pi : S \rightarrow \mathbb{R}$	state-value function: Expected return starting from a given state following policy $\pi$
$V^* : S \rightarrow \mathbb{R}$	expected return starting from a given state following the optimal policy
$N : S \rightarrow \mathbb{R}$	number of times $s$ was visited (also with $S \times A$ possible)
$\delta_t$	TD (temporal difference) error
$\lambda$	weight for TD( $\lambda$ )
$Q^\pi : S \times A \rightarrow \mathbb{R}$	state-action value function; follow $\pi$ after taking the given action
$Q^* : S \times A \rightarrow \mathbb{R}$	state-action value function; follow the optimal $\pi$ after taking the given action
$\epsilon$	probability to do a random exploration step
$\alpha$	step size for updating (e.g.) the $Q$ -function
$\mathbb{E}$	expectation
$\nabla$	gradient
$\Delta$	difference (e.g., update)
$\partial$	partial derivative
$\hat{\cdot}$	approximation of $\cdot$ (e.g., $\hat{V}$ or $\hat{Q}$ )
$\mathbf{x}$	feature vector
$\mathbf{w}$ or $\theta$	weight vector or tensor of function approximator
$\pi_\theta$	policy with policy network parameterized by $\theta$
$d : S \rightarrow \mathbb{R}$	stationary distribution over states
$\tau$ or $h$	Trajectory or history (state, action, reward, state, action, ...)
$b : S \rightarrow \mathbb{R}$	baseline (estimator) for a given state
$\mu$	mean
$\sigma$	standard deviation
$H$	entropy