# Reinforcement Learning (Tabular Methods) Reference Guide

Jian W Dong, based on Sutton & Barto: Reinforcement Learning - An Introduction 2E

Updated 2025.08.18

| Algorithm | On/Off Policy | Model-Free/ Model-Based | Control/ Predic-tion | Policy Update | Objective | Bootstrap | Target | Update Rule | When to Use |
|---|---|---|---|---|---|---|---|---|---|
| DYNAMIC PROGRAMMING (Chapter 4) | | | | | | | | | |
| Policy Eval | On | Model | Pred | None | $V^\pi$ | Yes | $\sum_a \pi(a\|s) \sum_{s',r} p(s',r\|s,a)[r + \gamma V(s')]$ | $V(s) \leftarrow$ Target | Known environment, need policy evaluation |
| Policy Iter | On | Model | Ctrl | Greedy | $V^*$ | Yes | Same as Policy Eval | Policy Eval + $\pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r\|s,a)[r + \gamma V(s')]$ | Known environment, guaranteed optimal policy |
| State-Value Iter | On | Model | Ctrl | Greedy | $V^*$ | Yes | $\max_a \sum_{s',r} p(s',r\|s,a)[r + \gamma V(s')]$ | $V(s) \leftarrow$ Target | Known environment, faster than policy iteration |
| Action-Value Iter | On | Model | Ctrl | Greedy | $Q^*$ | Yes | $\sum_{s',r} p(s',r\|s,a)[r + \gamma \max_{a'} Q(s',a')]$ | $Q(s,a) \leftarrow$ Target | Known environment, directly learns action values |
| MONTE CARLO METHODS (Chapter 5) | | | | | | | | | |
| First-Visit MC | On | Free | Pred | None | $V^\pi$ | No | $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ (from first visit) | $V(S_t) \leftarrow V(S_t) + \alpha[\text{Target} - V(S_t)]$ | Episodic tasks, unbiased estimates |
| Every-Visit MC | On | Free | Pred | None | $V^\pi$ | No | $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ (from each visit) | $V(S_t) \leftarrow V(S_t) + \alpha[\text{Target} - V(S_t)]$ | Episodic tasks, more data per episode |
| MC Exploring | On | Free | Ctrl | Greedy | $Q^*$ | No | $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ | $Q(s,a) \leftarrow$ average, $\pi(s) \leftarrow \arg\max_a Q(s,a)$ | Episodic, can ensure exploring starts |
| On-policy MC | On | Free | Ctrl | $\epsilon$-gr | $Q^\pi$ | No | $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ | $Q(s,a) \leftarrow$ average, $\epsilon$-greedy policy | Episodic, practical exploration |
| Off-policy MC Pred | Off | Free | Pred | None | $V^\pi$ | No | $\rho_{t:T-1} G_t$ where $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ | $V(s) \leftarrow$ weighted average | Learn about different policy than behavior |
| Off-policy MC Ctrl | Off | Free | Ctrl | Greedy | $Q^*$ | No | $\rho_{t+1:T-1} G_t$ where $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ | Weighted average + greedy policy | Learn optimal policy from suboptimal data |
| TEMPORAL-DIFFERENCE LEARNING (Chapter 6) | | | | | | | | | |
| TD(0) | On | Free | Pred | None | $V^\pi$ | Yes | $R_{t+1} + \gamma V(s_{t+1})$ | $V(s_t) \leftarrow V(s_t) + \alpha[\text{Target} - V(s_t)]$ | Online learning, fast updates |
| SARSA | On | Free | Ctrl | $\epsilon$-gr | $Q^\pi$ | Yes | $R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Safe online control, conservative |
| Q-learning | Off | Free | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | $R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a')$ | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Learn optimal policy, explo-ration/exploitation |
| Expected SARSA | On/Off | Free | Ctrl | $\epsilon$-gr | $Q^\pi$ | Yes | $R_{t+1} + \gamma \sum_{a'} \pi(a'\|s_{t+1}) Q(s_{t+1}, a')$ | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Lower variance than SARSA |
| Double Q-learning | Off | Free | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | Alternate: $R_{t+1} + \gamma Q_B(s_{t+1}, \arg\max_{a'} Q_A(s_{t+1}, a'))$ or $R_{t+1} + \gamma Q_A(s_{t+1}, \arg\max_{a'} Q_B(s_{t+1}, a'))$ | Randomly select: $Q_A \leftarrow Q_A + \alpha[\text{Target}_A - Q_A]$ or $Q_B \leftarrow Q_B + \alpha[\text{Target}_B - Q_B]$ | Avoid overestimation bias |
| n-STEP BOOTSTRAPPING (Chapter 7) | | | | | | | | | |
| n-step TD | On | Free | Pred | None | $V^\pi$ | Yes | $G_{t:t+n} = \sum_{i=0}^{n-1} \gamma^i R_{t+i+1} + \gamma^n V(s_{t+n})$ | $V(s_t) \leftarrow V(s_t) + \alpha[\text{Target} - V(s_t)]$ | Bridge MC and TD, tune bias/variance |
| n-step SARSA | On | Free | Ctrl | $\epsilon$-gr | $Q^\pi$ | Yes | $G_{t:t+n} = \sum_{i=0}^{n-1} \gamma^i R_{t+i+1} + \gamma^n Q(s_{t+n}, a_{t+n})$ | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Multi-step lookahead, on-policy |
| n-step Tree Backup | Off | Free | Ctrl | Any | $Q^\pi$ | Yes | $G_{t:t+n}^{tree} = R_{t+1} + \gamma[\sum_{a \neq A_{t+1}} \pi(a\|S_{t+1}) Q(S_{t+1}, a) + \pi(A_{t+1}\|S_{t+1}) G_{t+1:t+n}^{tree}]$ | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Off-policy without importance sampling |
| n-step $Q(\sigma)$ | On/Off | Free | Ctrl | Any | $Q^\pi$ | Yes | $\sigma$-weighted combination of SARSA and Tree Backup | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Unify on/off-policy methods |
| Off-policy n-step | Off | Free | Ctrl | Any | $Q^\pi$ | Yes | $\rho_{t+1:t+n-1} G_{t:t+n}$ (importance-weighted $n$-step return) | $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[\text{Target} - Q(s_t, a_t)]$ | Off-policy with multi-step returns |
| PLANNING AND LEARNING WITH TABULAR METHODS (Chapter 8) | | | | | | | | | |
| Dyna-Q | Off | Model | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | $r + \gamma \max_{a'} Q(s', a')$ (same as Q-learning) | Q-learning update + model learning + planning | Combine learning and planning |

| Algorithm | On/Off Policy | Model-Free/ Model-Based | Control/ Predic-tion | Policy Update | Objective | Bootstrap | Target | Update Rule | When to Use |
|---|---|---|---|---|---|---|---|---|---|
| Dyna-Q+ | Off | Model | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | $r + \kappa\sqrt{\tau} + \gamma \max_{a'} Q(s', a')$ | Same as Dyna-Q with exploration bonus | Handle changing environments |
| Prioritized Sweeping | Off | Model | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | $r + \gamma \max_{a'} Q(s', a')$ | Updates prioritized by $|\text{Target} - Q(s, a)| > \theta$ | Efficient planning, focus important updates |
| Trajectory Sampling | Off | Model | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | $r + \gamma \max_{a'} Q(s', a')$ | Sample long trajectories vs. uniform sweeping | Better state distribution for planning |
| Real-time DP | On | Model | Ctrl | Greedy | $V^*$ | Yes | $\max_a \sum_{s', r} p(s', r|s_t, a)[r + \gamma V(s')]$ | $V(s_t) \leftarrow$ Target only for visited states | Online DP, focus on relevant states |
| **ELIGIBILITY TRACES (Chapter 12)** | | | | | | | | | |
| Offline $\lambda$-return | On | Free | Pred | None | $V^\pi$ | Yes | $G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_t$ | $V(s_t) \leftarrow V(s_t) + \alpha[G_t^\lambda - V(s_t)]$ (offline at episode end) | Theoretical foundation for TD($\lambda$), offline learning |
| TD($\lambda$) | On | Free | Pred | None | $V^\pi$ | Yes | $R_{t+1} + \gamma V(s_{t+1})$ (TD error: $\delta_t$) | $V(s) \leftarrow V(s) + \alpha\delta_t e_t(s)$ where $e_t(s) = \gamma\lambda e_{t-1}(s) + \mathbf{1}_{s_t=s}$ | Credit assignment, faster learning |
| SARSA($\lambda$) | On | Free | Ctrl | $\epsilon$-gr | $Q^\pi$ | Yes | $R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ (TD error: $\delta_t$) | $Q(s, a) \leftarrow Q(s, a) + \alpha\delta_t e_t(s, a)$ | On-policy with eligibility traces |
| Q($\lambda$) | Off | Free | Ctrl | $\epsilon$-gr | $Q^*$ | Yes | $R_{t+1} + \gamma \max_a Q(s_{t+1}, a)$ (TD error: $\delta_t$) | Watkins's Q($\lambda$): traces reset if non-greedy action | Off-policy with traces (limited) |
| True Online TD($\lambda$) | On | Free | Pred | None | $V^\pi$ | Yes | $R_{t+1} + \gamma V(s_{t+1})$ (TD error: $\delta_t$) | Modified update with trace correction term | More accurate trace implementation |

**Notation:** $V(s)$: State value, $Q(s, a)$: Action-value, $\pi(a|s)$: Policy, $\alpha$: Learning rate, $\gamma$: Discount, $\epsilon$: Exploration, $G_t$: Return, $\rho$: Importance ratio, $e_t(s)$: Eligibility trace, $\lambda$: Trace decay, $\delta_t$: TD error
**Abbreviations:** Model = Model-based, Free = Model-free, Ctrl = Control, Pred = Prediction, $\epsilon$-gr = $\epsilon$-greedy