

# MARKOV DECISION PROBLEMS

MDP  $M = (S, A, T, R, \gamma)$  Discount Factor

States      Actions      Transition Function      Reward Function

\*  $R(s, a, s') = \text{Reward for } s \xrightarrow{a} s' \quad R \in [-R_{\max}, R_{\max}] \quad R_{\max} \geq 0$

$T(s, a, s') = \text{Probability of reaching } s' \text{ by starting at } s \text{ and taking action } a.$

\* Thus  $T(s, a, \cdot) = \text{Probability distribution over } S.$

## AGENT ENVIRONMENT INTERACTION:

At  $t=0,$

$s' \sim T(s^0, a^0, \cdot)$

$r^0 = R(s^0, a^0, s')$

Agent  $\xrightleftharpoons[r^t, s^{t+1}]{a^t}$  Environment

Resulting Trajectory =  $s^0 a^0 r^0 s^1 a^1 r^1 s^2 \dots$

Assume  $a^t$  is picked based on  $s^t$  alone

Policy  $\pi: S \rightarrow A$

$\pi \rightarrow$  Markovian, deterministic and stationary.

## POLICY:

$\Pi \rightarrow$  Denote the set of all policies.

$|\Pi| = k^n$        $n \rightarrow \text{states}$   
                          $k \rightarrow \text{actions}$

Q. Which  $\pi \in \Pi$  is a good policy?

Value function "Larger is better"

## $V^\pi(s)$ : VALUE OF STATE $s$ UNDER POLICY $\pi$

For  $s \in S, V^\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_\pi [r^0 + \gamma r^1 + \gamma^2 r^2 + \gamma^3 r^3 + \dots \mid s^0 = s]$

$\gamma = [0, 1) \Rightarrow$  Discount Factor

Large  $\gamma \Rightarrow$  Further lookahead

★★ Every MDP is guaranteed to have an optimal policy  $\pi^*$  such that  
 $\forall \pi \in \Pi, \forall s \in S: V^{\pi^*}(s) \geq V^\pi(s)$

## MDP PLANNING PROBLEM:

Given  $M = (S, A, T, R, \gamma)$  find a policy  $\pi^*$  from the set of all policies  $\Pi$  such that  $\forall s \in S, \forall \pi \in \Pi: V^{\pi^*}(s) \geq V^\pi(s)$

\* Every MDP is guaranteed to have a deterministic, markovian, stationary optimal policy

\* An MDP can have more than one optimal policy.

$\Rightarrow$  Value function of every optimal policy is the same unique  $V^*$

## ALTERNATIVE FORMULATIONS:

$\rightarrow R(s, a, s') \Rightarrow$  Random variable bounded in  $[-R_{\max}, R_{\max}]$

$\rightarrow R(s, a, s') \begin{cases} \rightarrow R(s, a) \\ \rightarrow R(s') \end{cases}$

$\rightarrow$  Combined  $T$  and  $R$ .

$\rightarrow$  Minimum "COST" instead of maximising "Reward"

\* **FINITE HORIZON REWARD**

$\mathbb{E}_\pi [r^0 + r^1 + r^2 + \dots + r^{T-1} \mid s^0 = s]$

\* **AVERAGE REWARD**

$\mathbb{E}_\pi [\lim_{M \rightarrow \infty} (r^0 + r^1 + \dots + r^{M-1}) / M \mid s^0 = s]$

**Total reward**

$V^\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_\pi [r^0 + r^1 + r^2 + \dots \mid s^0 = s]$   $\rightarrow$  Only be used on episodic tasks

## EPISODIC TASKS:

$\rightarrow$  Have special sink / terminal state  $s_T$  from which there are no outgoing transitions on rewards

\* From every non-terminal state and for every policy there is a non-zero probability of reaching terminal state in finite number of steps

## POLICY EVALUATION :

$$V^\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_\pi [r^0 + \gamma r^1 + \gamma^2 r^2 \dots | s^0 = s]$$

$$= \sum_{s' \in S} T(s, \pi(s), s') \mathbb{E}_\pi [r^0 + \gamma r^1 + \dots | s^0 = s, s^1 = s']$$

$$= \sum_{s' \in S} T(s, \pi(s), s') R(s, \pi(s), s') + \gamma \sum_{s' \in S} T(s, \pi(s), s') \mathbb{E}_\pi [r^1 + \gamma r^2 + \dots | s^1 = s']$$

$$= \sum_{s' \in S} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma V^\pi(s') \} \Rightarrow \text{BELLMAN'S EQUATIONS}$$

$n$ -equations,  $n$ -unknowns -  $V^\pi(s_1), V^\pi(s_2) \dots V^\pi(s_n)$

① If  $\gamma < 1 \Rightarrow$  UNIQUE SOLUTION

POLICY EVALUATION : Computing  $V^\pi$  for a policy  $\pi$

## ACTION VALUE FUNCTION :

$$Q^\pi(s, a) \stackrel{\text{def}}{=} \mathbb{E} [r^0 + \gamma r^1 + \gamma^2 r^2 \dots | s^0 = s; a^0 = a; a^t = \pi(s^t) \text{ for } t \geq 1]$$

$$Q^\pi(s, a) = \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma V^\pi(s') \}$$

\* All optimal policies have the same action value function  $Q^*$

ALGORITHMS : ① Bellman Optimality ② Value Iteration ③ Linear Programming Formulation

## CONTRACTION MAPPING :

A mapping  $T: X \rightarrow X$  is called a contraction mapping with a contraction factor  $L$  if  $\forall u, v \in X$ ,

$$\|Tu - Tv\| \leq L \|u - v\|$$

Fixed point:  $Tx^* = x^*$

## BANACH'S FIXED POINT THEOREM :

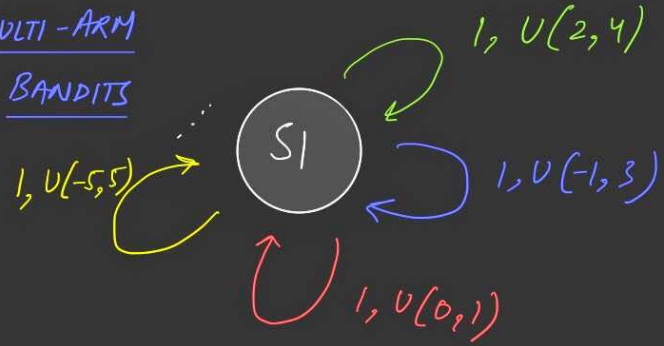
Let  $(X, \|\cdot\|)$  be a Banach space and let  $T: X \rightarrow X$  be a contraction mapping with  $L \in [0, 1)$  then:

①  $T$  has a unique fixed point  $x^* \in X$

② For  $x \in X$ ,  $m \geq 0$ :  $\|T^m x - x^*\| \leq L^m \|x - x^*\|$

MULTI-ARM

BANDITS



# BELLMAN OPTIMALITY OPERATOR

$B^* : (S \rightarrow \mathbb{R}) \rightarrow (S \rightarrow \mathbb{R})$  for an MDP  $(S, A, T, R, \gamma)$  is defined as:

For  $F: S \rightarrow \mathbb{R}$  and  $s \in S$ :

$$(B^*(F))(s) \stackrel{\text{def}}{=} \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma F(s')\}$$

Since  $S = \{s_1, s_2, \dots, s_n\}$ , we may equivalently view  $B^*$  as a mapping from  $\mathbb{R}^n \rightarrow \mathbb{R}^n$

② max norm  $\|\cdot\|_\infty$  of  $F = (f_1, f_2, \dots, f_n) \in \mathbb{R}^n \Rightarrow \|F\|_\infty = \max \{|f_1|, |f_2|, \dots, |f_n|\}$

Already established  $\rightarrow (\mathbb{R}^n, \|\cdot\|_\infty)$  is a Banach Space.

CLAIM:  $B^*$  is a contraction mapping in the  $(\mathbb{R}^n, \|\cdot\|_\infty)$  Banach space with contraction factor  $\gamma$ .

PRE REQUISITE:  $|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$

PROOF:  $\|B^*(F) - B^*(G)\|_\infty \leq \max_{s \in S} |(B^*(F))(s) - (B^*(G))(s)|$

$$= \max_{s \in S} \left| \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma F(s')\} - \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma G(s')\} \right|$$

$$= \gamma \max_{s \in S} \max_{a \in A} \left| \sum_{s' \in S} T(s, a, s') \{F(s') - G(s')\} \right|$$

$$= \gamma \max_{s \in S} \max_{a \in A} \sum_{s' \in S} T(s, a, s') \|F - G\|_\infty = \gamma \|F - G\|_\infty \Rightarrow \text{Fixed point for } B^*$$

BELLMAN OPTIMALITY EQUATIONS  $\leftarrow$

$$V^*(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma V^*(s')\}$$

\* Not linear like

Bellman Equations

$\Rightarrow V^*$  is the value function for every policy  $\pi^* : S \rightarrow A$  that satisfies for all  $s \in S$ :

$$\pi^*(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma V^*(s')\}$$



## ① VALUE ITERATION

$$V_0 \xrightarrow{B^*} V_1 \xrightarrow{B^*} V_2 \xrightarrow{B^*} \dots$$

\* stop when  $V_t \approx V_{t+1}$  (upto machine precision)

\*  $V^* \rightarrow Q^*$   $Q^*(s, a) = \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma V^*(s') \}$

\*  $Q^* \rightarrow \pi^*$   $\pi^*(s) = \operatorname{argmax}_{a \in A} Q^*(s, a)$

\*  $\pi^* \rightarrow V^*$  Solve Bellman Equations for  $\pi^*$

## ② LINEAR PROGRAMMING FORMULATION

\* Optimise linear function of variables

\* Subject to linear constraints

Create  $n$  variables  $V(s_1), V(s_2) \dots V(s_n) \rightarrow V^*$  will be the solution

Linear constraint for max  $\rightarrow V(s) \geq \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma V(s') \}$  for  $s \in S$   
 $a \in A$

### Vector Comparison

$$X \succeq Y \Leftrightarrow \forall s \in S : X(s) \geq Y(s)$$

$$X \succ Y \Leftrightarrow X \succeq Y \text{ and } \exists s \in S : X(s) > Y(s)$$

\* Sometimes vectors are incomparable.

$$* V \preceq B^*(V)$$

$$V \preceq \lim_{l \rightarrow \infty} (B^*)^l(V) = V^*$$

$$\Rightarrow \sum_{s \in S} V(s) \geq \sum_{s \in S} V^*(s)$$

### $B^*$ Preserves $\succeq$

TO PROVE:  $(B^*(X))(s) - (B^*(Y))(s) \geq 0$

$$(B^*(X))(s) - (B^*(Y))(s)$$

$$= \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma X(s') \}$$

$$- \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma Y(s') \}$$

$$\geq \gamma \min_{a \in A} \sum_{s' \in S} T(s, a, s') \{ X(s') - Y(s') \} \geq 0$$

$$* \max_a f(a) - \max_a g(a) \geq \min_a (f(a) - g(a))$$

Goal: Maximise  $\left(-\sum_{s \in S} V(s)\right)$

Subject to  $V(s) \geq \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma V(s')\} \quad \forall s \in S, a \in A$

LP  $\rightarrow$   $n$  variables,  $nk$  constraints

## POLICY IMPROVEMENT:

Given  $\pi$ , Pick one or more improvable states, and in them, switch to an arbitrary improvement action. Resulting state  $\rightarrow \pi'$

① For  $\pi \in \Pi, s \in S$

IMPROVING ACTIONS  $\leftarrow IA(\pi, s) \stackrel{\text{def}}{=} \{a \in A : Q^\pi(s, a) > V^\pi(s)\}$

② For  $\pi \in \Pi$

IMPROVABLE STATES  $\leftarrow IS(\pi) \stackrel{\text{def}}{=} \{s \in S : |IA(\pi, s)| \geq 1\}$

$\pi' \in \Pi$  is obtained by Policy Improvement  $\rightarrow \begin{cases} \forall s \in S : \pi'(s) = \pi(s) \text{ or } \pi'(s) \in IA(\pi, s) \\ \exists s \in S : \pi'(s) \in IA(\pi, s) \end{cases}$

## POLICY IMPROVEMENT THEOREM

(1) If  $IS(\pi) = \emptyset$ , then  $\pi$  is optimal

(2) If  $\pi'$  is obtained by policy improvement on  $\pi$ , then  $\pi' \succ \pi$

\* If  $IS(\pi) \neq \emptyset$  then  $\exists \pi' \in \Pi$  such that  $\pi' \succ \pi$

\* But  $\Pi$  has finite policies ( $k^n$ )

\*  $\Rightarrow \exists \pi^*$  such that  $IS(\pi^*) = \emptyset$

\*  $IS(\pi^*) = \emptyset \Leftrightarrow B^*(V^{\pi^*}) = V^{\pi^*}$

Convention  $\rightarrow V^*$  for  $V^{\pi^*}$

\* SWITCHING

STRATEGY

Path taken and iterations depend on which improved state is chosen.

- ① Random policy iteration
- ② Simple Policy iteration

## POLICY ITERATION ALGORITHM

$\pi \leftarrow$  Arbitrary Policy

While  $\pi$  has improvable states

$\pi' \leftarrow$  Policy improvement( $\pi$ )

$\pi \leftarrow \pi'$

Return  $\pi$

## \* HOWARD'S POLICY ITERATION:

$\Rightarrow$  Greedy  $\rightarrow$  Switch all improvable states

## \* RANDOM POLICY ITERATION:

$\Rightarrow$  Switch a non-empty subset of improvable states chosen uniformly at random.  $\rightarrow 2^n - 1$  subsets

# PROOF OF POLICY IMPROVEMENT THEOREM

BELLMAN OPERATOR  $B^\pi: (S \rightarrow R) \rightarrow (S \rightarrow R)$  \* It is a contraction mapping.

$$(B^\pi(x))(s) \stackrel{\text{def}}{=} \sum_{s' \in S} T(s, \pi(s), s') (R(s, \pi(s), s') + \gamma x(s'))$$

$$\text{For } X: S \rightarrow R: \lim_{l \rightarrow \infty} (B^\pi)^l(x) = V^\pi$$

$$\text{For } X: S \rightarrow R, Y: S \rightarrow R: X \geq Y \Rightarrow B^\pi(X) \geq B^\pi(Y)$$

$$\text{Observe that: } B^{\pi'}(V^\pi)(s) = Q^\pi(s, \pi'(s))$$

PROOF:  $IS(\pi) = \emptyset$

$$\Rightarrow \forall \pi' \in \Pi: V^\pi \geq B^{\pi'}(V^\pi)$$

$$\Rightarrow \forall \pi' \in \Pi: V^\pi \geq B^{\pi'}(V^\pi) \geq (B^{\pi'})^2(V^\pi) \dots \geq \lim_{l \rightarrow \infty} (B^{\pi'})^l(V^\pi) \xrightarrow{\text{arrow}} V^{\pi'}$$

$$IS(\pi) \neq \emptyset$$

$$\Rightarrow B^{\pi'}(V^\pi) \geq V^\pi$$

$$\Rightarrow (B^{\pi'})^2(V^\pi) \geq B^{\pi'}(V^\pi) \geq V^\pi$$

$$\Rightarrow V^{\pi'} \geq V^\pi$$

## COMPUTATIONAL COMPLEXITY:

RUNNING TIME BOUNDS  $\rightarrow$  Value Iteration  $\Rightarrow$  Upper bound:  $\text{poly}(n, k, B, \frac{1}{1-\gamma})$

$\rightarrow$  Linear Programming  $\Rightarrow \text{poly}(n, k, B)$

$\hookrightarrow \text{poly}(n, k) \cdot \exp(O(\sqrt{n \log(n)}))$  Expected

### POLICY ITERATION

	$k=2$	General $k$	
Deterministic	$O(\frac{2^n}{n})$	$O(\frac{k^n}{n})$	} $\rightarrow$ UPPER BOUND
Randomised	$1.7172^n$	$O(\frac{k}{2})^n$	

## REINFORCEMENT LEARNING PROBLEM:

Agent does not know Transition Probability and sometimes reward function

★ Can the agent eventually take optimal actions?

$$h^t = (s^0, a^0, r^0, s^1, \dots, s^t) \Rightarrow \text{History}$$

★ Learning Algorithm  $L$  is a mapping from set of all histories to set of probability distribution over arms.

### LEARNING PROBLEM

### ★ CONTROL PROBLEM

$$\lim_{T \rightarrow \infty} \frac{1}{T} \left( \sum_{t=0}^{T-1} \mathbb{P}\{a^t \sim L(h^t) \text{ is an optimal action for } s^t\} \right) = 1$$

↓  
we want to be sublinear.

### PREDICTION PROBLEM

★ We are give  $\pi$  that agent follows. AIM: To estimate  $V^\pi$ .

PROBLEM: Can we construct  $L$  such that  $\lim_{t \rightarrow \infty} \hat{V}^t = V^\pi$ ?

### ASSUMPTION 1: IRREDUCIBILITY

If there is a directed path from  $s$  to  $s'$  for every  $s, s' \in S$ ,  
↳ comes from non-zero-probability transition under  $\pi$ .

then  $M$  is irreducible under  $\pi$ .  $M$  is irreducible if it is irreducible under ALL  $\pi$ .

### ASSUMPTION 2: APERIODIC

$X(s, t) \rightarrow$  set of all states possible after time  $t$  following policy  $\pi$ .

$Y(s) \rightarrow$  set of all  $t$  such that  $s \in X(s, t)$

$p(s) \rightarrow \text{GCD}(Y(s))$

$M$  is aperiodic under  $\pi$  if for all  $s \in S$ :  $p(s) = 1$

If  $M$  is aperiodic under all  $\pi \in \Pi$ , then  $M$  is aperiodic.

ERGODICITY: Mdp which is APERIODIC and IRREDUCIBLE.

★ Every policy induces a unique steady state distribution  $\mu^\pi: S \rightarrow (0, 1)$ , subject to  $\sum_{s \in S} \mu^\pi(s) = 1$ ,  $\Rightarrow$  For Ergodic MDP  
which is independent of start.

$$\mu^\pi(s) = \lim_{t \rightarrow \infty} p(s, t)$$



## MODEL BASED APPROACH:

Estimate of MDP  $\hat{T} \rightarrow T$   
 $\hat{R} \rightarrow R$

\* At convergence, acting optimally for  $MDP(s, A, \hat{T}, \hat{R}, \gamma)$  must be optimal for the original  $MDP(s, A, T, R, \gamma)$

\* we must visit every state-action pair infinitely often  
Because reward is stochastic.

\* Algorithm discussed in slides  $\Rightarrow$  Takes sub-linear number of sub-optimal steps.

$\rightarrow$  Needs Irreducibility, not Aperiodicity.

### MODEL BASED

$\rightarrow$  uses  $\Theta(|S|^2 |A|)$  memory.

## PREDICTION

MONTÉ CARLO methods estimate based on sample averages.

NOTATION:  $s \rightarrow$  state  
 $i \rightarrow$  Episode Number  
 $j \rightarrow$  occurrence in an episode.

$1(s, i, j) \rightarrow$  If  $s$  is visited at least  $j$  times on episode  $i$   
 $G(s, i, j) \rightarrow$  Discounted long term reward starting from  $j^{\text{th}}$  visit of  $s$  on episode  $i$ .

### FIRST VISIT MONTÉ CARLO

$$\hat{V}_{\text{First-visit}}^T(s) = \frac{\sum_{i=1}^T G(s, i, 1)}{\sum_{i=1}^T 1(s, i, 1)}$$

### EVERY VISIT MONTÉ CARLO

$$\hat{V}_{\text{Every-visit}}^T(s) = \frac{\sum_{i=1}^T \sum_{j=1}^{\infty} G(s, i, j)}{\sum_{i=1}^T \sum_{j=1}^{\infty} 1(s, i, j)}$$

$$\lim_{T \rightarrow \infty} \hat{V}_{\text{First-visit}}^T = V^{\pi}$$

$$\lim_{T \rightarrow \infty} \hat{V}_{\text{Every-visit}}^T = V^{\pi}$$

$\Rightarrow$  More careful analysis

$$\lim_{T \rightarrow \infty} \hat{V}_{\text{second-visit}}^T = V^{\pi}$$

$$\lim_{T \rightarrow \infty} \hat{V}_{\text{last-visit}}^T \neq V^{\pi}$$



## FIRST VISIT MONTE CARLO

$$\hat{V}^t(s) = \frac{1}{t} \sum_{i=1}^t G(s, i, 1)$$

## STOCHASTIC APPROXIMATION

Let the sequence  $(\alpha_t)_{t \geq 1}$  satisfy:

- ①  $\sum_{t=1}^{\infty} \alpha_t = \infty$
- ②  $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$

For  $t \geq 1$ :

$$\hat{V}^t(s) \leftarrow (1 - \alpha_t) \hat{V}^{t-1}(s) + \alpha_t G(s, t, 1)$$

$\alpha_t \rightarrow$  learning rate.

## ONLINE IMPLEMENTATION

$$\hat{V}^t(s) = \frac{1}{t} \sum_{i=1}^t G(s, t, i)$$

$$= \frac{1}{t} \left( \sum_{i=1}^{t-1} G(s, t, i) + G(s, t, 1) \right)$$

$$= \frac{1}{t} \left( (t-1) \hat{V}^{t-1}(s) + G(s, t, 1) \right)$$

$$= (1 - \alpha_t) \hat{V}^{t-1}(s) + \alpha_t G(s, t, 1) \quad \alpha_t = \frac{1}{t}$$

Then  $\lim_{t \rightarrow \infty} \hat{V}^t(s) = V^*(s)$

## BOOTSTRAPPING

Replacing  $M$  by older estimates of the value function of state encountered in the episode.  
 ↓  
 Monte Carlo Estimate

## TEMPORAL DIFFERENCE LEARNING

$$\hat{V}^{t+1}(s^t) \leftarrow \hat{V}^t(s^t) + \alpha_{t+1} \{ r^t + \gamma \hat{V}^t(s^{t+1}) - \hat{V}^t(s^t) \}$$

$$r^t + \gamma \hat{V}^t(s^{t+1}) - \hat{V}^t(s^t) \Rightarrow \text{TEMPORAL DIFFERENCE PREDICTION ERROR.}$$

$$\hat{V}^t(s_T) = 0 \text{ For episodic Tasks}$$

↑  
Terminal States.

## ERROR

### FIRST VISIT MONTE CARLO

$$\text{Error}_{\text{First}}(V, s) \stackrel{\text{def}}{=} \sum_{i=1}^T 1(s, i, 1) (V(s) - G(s, i, 1))^2$$

### EVERY VISIT MONTE CARLO

$$\text{Error}_{\text{Every}}(V, s) \stackrel{\text{def}}{=} \sum_{i=1}^T \sum_{j=1}^{\infty} 1(s, i, j) (V(s) - G(s, i, j))^2$$

\* For TD estimate for a finite set of episodes,  $V^t$  depends on  $\hat{V}_0$ . Therefore to remove that dependence, convert finite to infinite episodes by repeating the episodes.

### BATCH TD(0) ESTIMATE:

$$\hat{V}_{\text{Batch-TD}(0)}^T = V^* \text{ on MLE.}$$

## 3 CONTROL ALGORITHMS

From state  $s^t$ ,  $a^t \sim \pi^t(s^t)$

Update  $\hat{Q}_t$  after observing  $s^t, a^t, r^t, s^{t+1}$

$$\hat{Q}^{t+1}(s^t, a^t) \leftarrow \hat{Q}^t(s^t, a^t) + \alpha_{t+1} [\text{Target} - \hat{Q}^t(s^t, a^t)]$$

$$\lim_{t \rightarrow \infty} \hat{Q}^t = Q^* \leftarrow \pi^t \text{ is } \epsilon\text{-greedy w.r.t. } \hat{Q}^t$$

$$\text{Q-LEARNING} \Rightarrow r^t + \gamma \max_{a \in A} \hat{Q}^t(s^{t+1}, a) \rightarrow \text{OFF POLICY}$$

$$\text{SARSA} \Rightarrow r^t + \gamma \hat{Q}^t(s^{t+1}, a^{t+1})$$

$$\text{EXPECTED SARSA} \Rightarrow r^t + \gamma \sum_{a \in A} \pi^t(s^{t+1}, a) \hat{Q}^t(s^{t+1}, a) \rightarrow \text{ON POLICY}$$

For  $\pi^t = \pi \Rightarrow$  Time invariant

$$\text{Q-learning} \rightarrow Q^*$$

SARSA

$$\text{EXPECTED SARSA} \rightarrow Q^{\pi^*}$$

## MULTI-STEP RETURNS

$s_1 \rightarrow s_2 \rightarrow s_3$

2-step return

$$V^{\text{new}}(s_1) \leftarrow V^{\text{old}}(s_1) + \alpha \{ \gamma + \gamma^2 V^{\text{old}}(s_3) - V^{\text{old}}(s_1) \}$$

### n-STEP TD

$$V^{t+n}(s^t) \leftarrow V^{t+n-1}(s^t) + \alpha \{ G_{t:t+n} - V^{t+n-1}(s^t) \}$$

## n-STEP RETURN

$$G_{t:t+n} \stackrel{\text{def}}{=} r^t + \gamma r^{t+1} + \gamma^2 r^{t+2} + \dots + \gamma^n r^{t+n} + \gamma^n V^{t+n-1}(s^{t+n})$$

For episodic tasks  $\rightarrow G_{t:t+n} = G_{t:t+n'}$   
 $n' \rightarrow \text{END.}$

For  $n \geq 1$

$$\lim_{t \rightarrow \infty} V^t = V^*$$

## COMBINING RETURNS:

$\rightarrow$  Convex combinations of  $G_{t:t+k}$  where  $k \geq 1$

$\downarrow$   
weighted average  
where  $w_i \geq 0$ .

## LINEAR ARCHITECTURE

$$\hat{V}(w, s) = w \cdot \chi(s) \quad \chi: S \rightarrow \mathbb{R}^d$$

$$w^* = \underset{w \in \mathbb{R}^d}{\text{argmin}} \text{MSVE}(w) \quad \text{MSVE}(w) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{s \in S} \chi^\top(s) \{ V^*(s) - \hat{V}(w, s) \}^2$$

$$\text{MSVE}(w_\lambda^\infty) \leq \left( \frac{1-\gamma\lambda}{1-\gamma} \right) \text{MSVE}(w^*)$$

## TILE CODING

\* A tiling partitions  $x$  into equal width regions called tiles  
Multiple tiles are created with an offset from previous one.

## REPRESENTING $\hat{Q}$

$$\hat{Q}(s, a) = \sum_{j=1}^d T_{aj}(\chi_j(s))$$

## ★ COMBINATION OF FEATURES

To represent more complex functions