

Multi-Agent Systems

Introduction to Reinforcement Learning

Part 3: Model-free Methods:
Monte Carlo, SARSA and Q-learning

Eric Pauwels (CWI & VU)

December 12, 2023

Reading

- Sutton & Barto: chapters 5 & 6

Outline

Model-based versus Model-free

Model-free RL problems

Monte Carlo for Model-free Policy Evaluation

Temporal Difference Methods for Model-free Prediction

SARSA and Q-Learning for Model-free Control

Model-based versus Model-free

	Prediction <i>Estimation:</i> <i>Given π, what is v?</i>	(Optimal) Control <i>Optimisation:</i> <i>What is optimal π?</i>
model-based (MDP given)	Policy evaluation using Dyn. Programming (DP)	Policy improvement (+ Policy evaluation) = Policy iteration
model-free (MDP unknown)	Monte Carlo (MC) Temporal Diff ^{ing} (TD) = "impatient MC" <i>bootstrapping!</i>	Q-learning and Generalized Policy Iteration <i>"simultaneous"</i>

Outline

Model-based versus Model-free

Model-free RL problems

Monte Carlo for Model-free Policy Evaluation

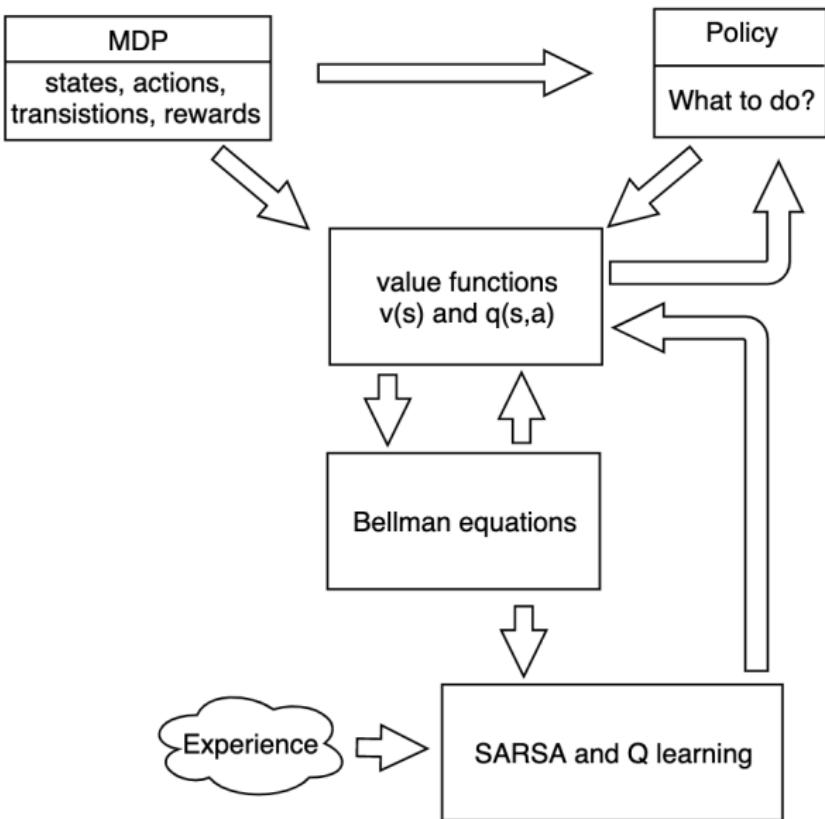
Temporal Difference Methods for Model-free Prediction

SARSA and Q-Learning for Model-free Control

Solving model-free RL problems

- The agent has **no prior knowledge** about states, actions, rewards, transitions!
- By **acting** in the world, the agent **gains experience**. An **experience** can be expressed as a 4-tuple: (s, a, r, s')
- Over time the agent collects a **list of experiences** that he uses to find better policies ... HOW?
 - **Directly:** policy improvement
 - **Indirectly:** estimate **value functions**, improve **policy using greedification**;

Overview



Definition of Greedification

- For a given action-value function $q(s, a)$, a corresponding **greedy policy** is a *deterministic* policy that picks (one of the) actions that **maximise the action value**:

$$\pi_g(s) = a^* := \arg \max_a q(s, a).$$

Greedification results in policy improvement

- To **improve the policy**, apply successive **iteration steps**:

$$\pi_k \xrightarrow{\text{eval}} v_k, q_k \xrightarrow{\text{greedify}} \pi_{k+1}$$

- Greedification implies **deterministic action choice** at each s :

$$\pi_{k+1}(s) = a^* \quad \text{iff} \quad q_k(s, a^*) = \max_a q_k(s, a)$$

- Hence **value function increases**, i.e. **policy has been improved!**

$$\begin{aligned} v_{k+1}(s) &= q_k(s, a^*) \quad (\pi_{k+1} \text{ is deterministic at } s) \\ &= \max_a q_k(s, a) \\ &\geq v_k(s) \quad (= \sum_{a'} \pi_k(a' | s) q_k(s, a')) \end{aligned}$$

Greedification requires $q(s, a)!$

model-free vs. model-based

Greedification: $\pi(s) := \arg \max_a q(s, a)$

Model-based: (a.k.a. planning, search)

- Suffices to estimate value function $v(s)$:
- $q(s, a)$ computed with Bellman's one-step look-ahead (i.e. use back-up diagram):

$$q(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v(s')]$$

Model-free: (a.k.a. (optimal) control)

- $q(s, a)$ needs to be estimated from collected experiences;
- algorithms shift focus from $v(s)$ to $q(s, a)$;

Conclusion: for **model-free** RL

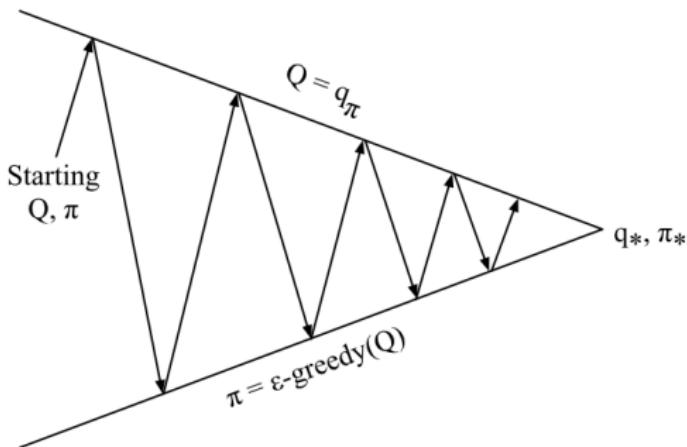
In contrast with to model-based, **model-free** RL is distinct:

- Focus on $q(s, a)$ rather than $v(s)$
- Make sure **all state-action pairs (s,a) are sampled**: importance of **exploration!**

Policy Iteration (PI) for Model-free RL

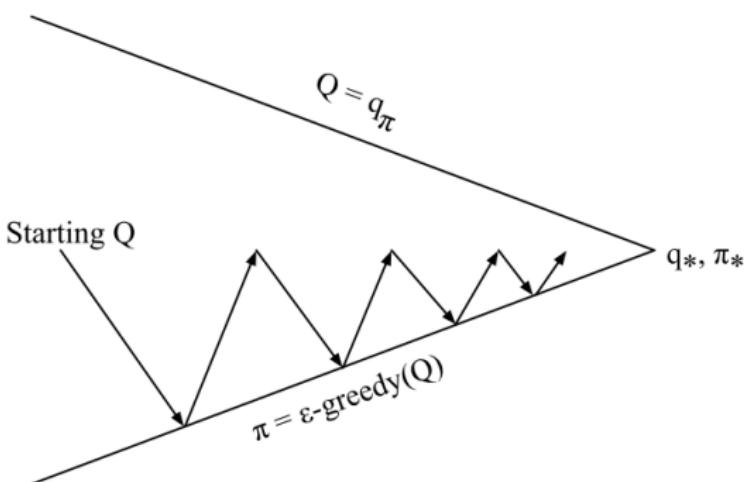
Goal: Find optimal policy π^* (aka optimal control)

- PI: Combine **policy evaluation** with **policy improvement**
 - Make Q -function **consistent with policy π** , i.e. $Q = q_\pi$;
 - **Greedify policy**: $\pi = \text{greedy}(Q)$;
- Introduce **exploration** in policy (e.g. $\pi = \varepsilon\text{-greedy}(Q)$);

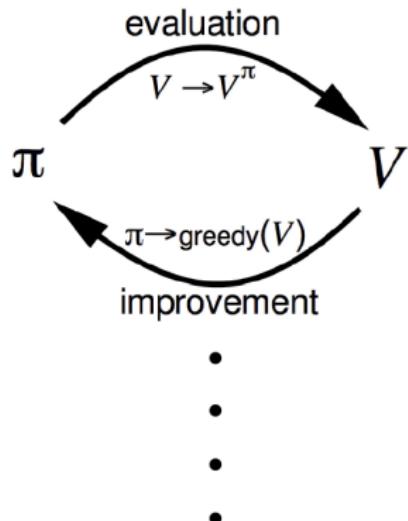
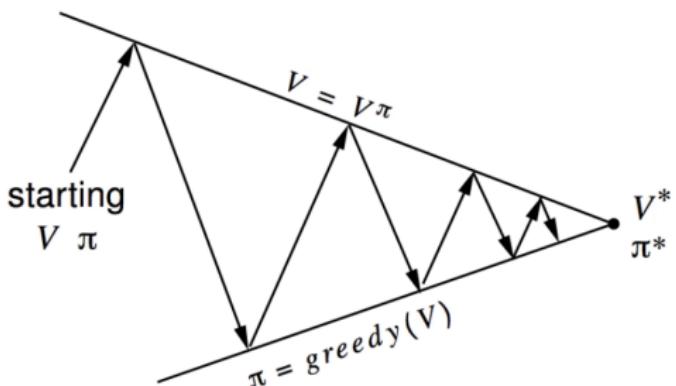


Generalized Policy Iteration (GPI)

- **Impatient greedification:** No insistence on making Q-function fully consistent with current policy π ;
Hence $Q \approx q_\pi$, but not necessarily $Q = q_\pi$;
- **Ensure exploration:** use $\pi = \varepsilon\text{-greedy}(Q)$, not $\pi = \text{greedy}(Q)$;
- Works if both processes continue updating all states;



Generalized Policy Iteration (GPI) schematically



Policy evaluation Estimate v_π

Any policy evaluation algorithm

Policy improvement Generate $\pi' \geq \pi$

Any policy improvement algorithm

π^* \longleftrightarrow V^*

http://blog.csdn.net/coffe_cream

Model-based versus Model-free

oo

Model-free RL problems

oooooooooooo

Monte Carlo for Model-free Policy Evaluation

●oooooooooooo

Temporal Difference Method

oooooooooooooooooooo

Outline

Model-based versus Model-free

Model-free RL problems

Monte Carlo for Model-free Policy Evaluation

Temporal Difference Methods for Model-free Prediction

SARSA and Q-Learning for Model-free Control

Monte Carlo methods

- **Monte Carlo methods** are versatile statistical techniques for estimating properties of complex systems via **random sampling**;
- Example 1: if $X \sim p(x)$ and φ some function, then for any (sufficiently large) i.i.d. sample: X_1, X_2, \dots, X_n :

$$E(\varphi(X)) = \int \varphi(x) p(x) dx \approx \frac{1}{n} \sum_{X_i \sim p} \varphi(X_i)$$

- Example 2: **Kullback-Leibler**

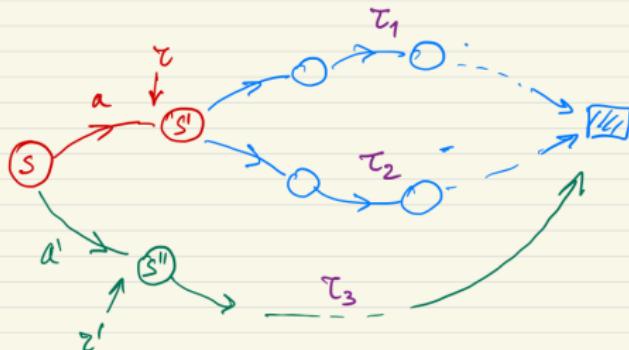
$$D_{KL}(f; g) = \int f(x) \log \frac{f(x)}{g(x)} dx \approx \frac{1}{n} \sum_{X_i \sim f} \log \frac{f(X_i)}{g(X_i)}$$

Model-free policy evaluation: MC-based estimation of v_π and $q_\pi(s, a)$

- Pick a **starting state s** (at **random** or systematic **sweep**)
- **Estimation of $v_\pi(s)$**
 - Use **policy π** to generate the **initial and all subsequent actions** (till terminal state)
 - Compute total return $r_{tot} = r_1 + r_2 + \dots + r_T$ along path:
 - r_{tot} yields **one sample value** for $v_\pi(s)$
- **Estimation of $q_\pi(s, a)$**
 - Apply **action a** in state s , observe reward r_1 and new state s'
 - Use **policy π** to generate **all subsequent actions** (from s' till terminal state)
 - Compute total return $r_{tot} = r_1 + r_2 + \dots + r_T$ along path:
 - r_{tot} yields **one sample value** for $q_\pi(s, a)$

Monte-Carlo for model-free policy evaluation

Monte Carlo: using **sampled** episodes to estimate $q(s, a)$!



$$G(\tau_1) = \sum_{\tau_1} r_i \quad \left. \right\} \rightarrow q_{\pi}(s, a)$$

$$G(\tau_2) = \sum_{\tau_2} r_i \quad \left. \right\} \rightarrow q_{\pi}(s, a')$$

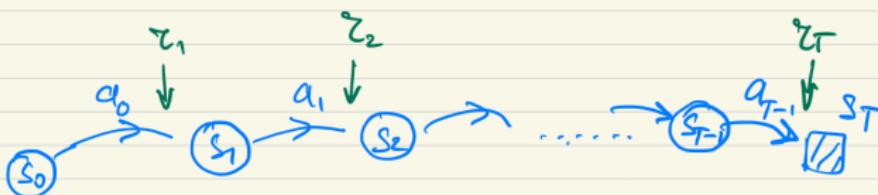
$$G(\tau_3) = \sum_{\tau_3} r_i \quad \rightarrow q_{\pi}(s, a')$$

If $a, a' \sim \pi$:

$$G(\tau_1), G(\tau_2), G(\tau_3) \rightarrow V_{\pi}(s)$$

MC estimation of $q(s, a)$ along trajectory

$$\mathcal{T} = \{ s_0, a_0, r_1, \underbrace{s_1, a_1, r_2, \underbrace{s_2, a_2, \dots, \underbrace{s_{T-1}, a_{T-1}, r_T, s_T}_{\text{trajectory}}}_{r_{T-1}}} \}$$



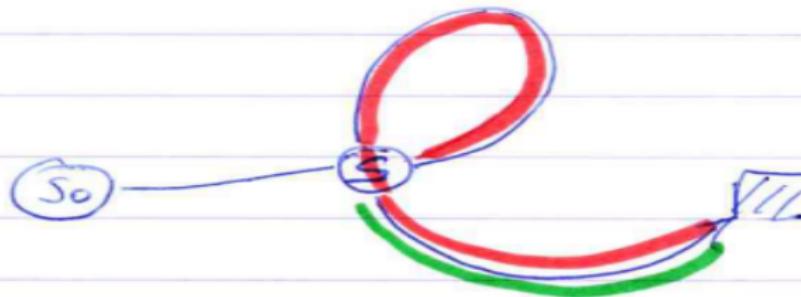
$$G_0 \rightarrow G_1 \rightarrow G_2 \rightarrow \dots$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$q_{\pi}(s_0, a_0) \quad q_{\pi}(s_1, a_1) \quad q_{\pi}(s_2, a_2) \quad \dots$$

MC policy evaluation: **First** versus **every visit** estimate

- **First-visit MC:** average returns only for the first time s is visited in an episode
- **Every-visit MC:** average returns for every time s is visited in an episode:
 - More sample efficient: more samples per episode;
 - Samples no longer independent
- Both converge asymptotically

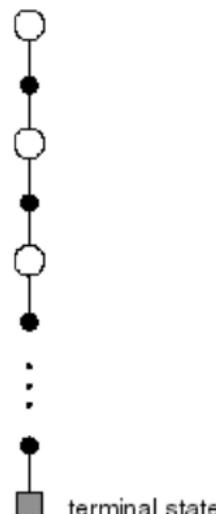


Monte-Carlo estimation of q -values

- Can learn q_π by averaging returns obtained when following π after taking action a in state s
- Converges asymptotically if every (s, a) visited infinitely often
 - Requires **explicit exploration** of actions not favored by π
 - **Possible solutions:**
 - **Exploring starts:** every (s, a) has a non-zero probability of being the starting pair
 - **Soft policies:** $\pi(a|s) > 0$ for all (s, a) . E.g. ϵ -greedy

Monte-Carlo backup diagram

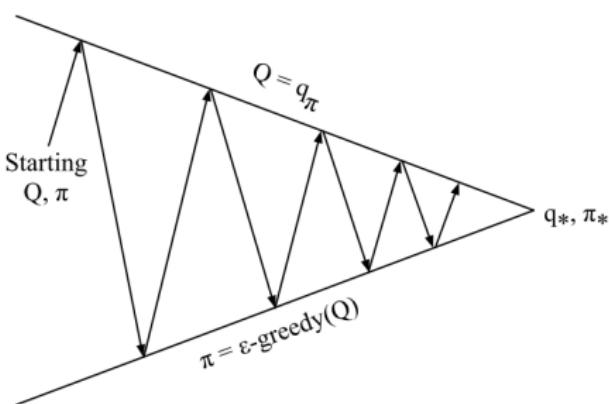
- Unlike dynamic programming, MC has only one choice at each state (i.e. the one actually taken!)
- Unlike dynamic programming, entire episode included: **MC does not bootstrap**
- Bellman equations are **NOT** used!



Model-free MC: From Evaluation to Control

Control: Find optimal policy π^* :

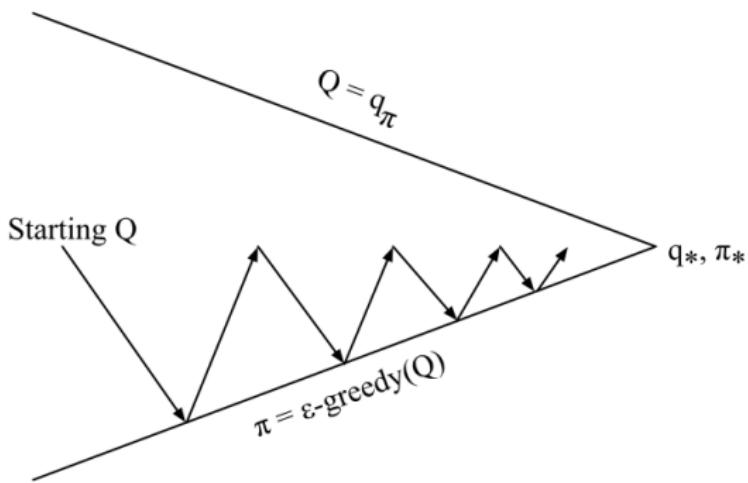
- Combine **MC-based policy evaluation** with **improvement** (e.g. greedification)
- Introduce **exploration** in policy
 - (e.g. ϵ -greedy)



Policy evaluation Monte-Carlo policy evaluation, $Q = q_\pi$

Policy improvement ϵ -greedy policy improvement

Monte Carlo for Model-free Control (2)



Every episode:

Policy evaluation Monte-Carlo policy evaluation, $Q \approx q_\pi$

Policy improvement ϵ -greedy policy improvement

Monte-Carlo for **model-free estimation and control**

- MC provides one way to perform **model-free reinforcement learning** (RL): finding optimal policies without an explicit model for the MDP;
- MC for RL learns from **complete sample returns** in episodic tasks:
- Computes **value functions** using **direct sampling rather than Bellman equations**;

Convergence condition: Greedy in the Limit with Infinite Exploration (GLIE)

Def (GLIE)

- All state-action pairs are explored infinitely often:

$$\lim_{t \rightarrow \infty} N_t(s, a) = +\infty.$$

- The policy converges to a greedy policy:

$$\lim_{t \rightarrow \infty} \pi(a | s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q(s, a') \\ 0 & \text{otherwise} \end{cases}$$

Example: ϵ -greedy is GLIE iff $\epsilon \downarrow 0$.

GLIE Monte Carlo Model-free Control

- Sample k th episode using π : $\{S_1, A_1, R_2, \dots, S_T\} \sim \pi$
- For each state S_t and action A_t in the episode,

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$$

- Improve policy based on new action-value function

$$\epsilon \leftarrow 1/k$$

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

Theorem

GLIE Monte-Carlo control converges to the optimal action-value function, $Q(s, a) \rightarrow q_(s, a)$*

Control based on MC+exploring starts: Algorithm

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

Notice: exploration based on exploring starts, not ϵ -greedy!

Model-based versus Model-free

oo

Model-free RL problems

oooooooooooo

Monte Carlo for Model-free Policy Evaluation

oooooooooooooooooooo

Temporal Difference Method

●oooooooooooooooooooo

Outline

Model-based versus Model-free

Model-free RL problems

Monte Carlo for Model-free Policy Evaluation

Temporal Difference Methods for Model-free Prediction

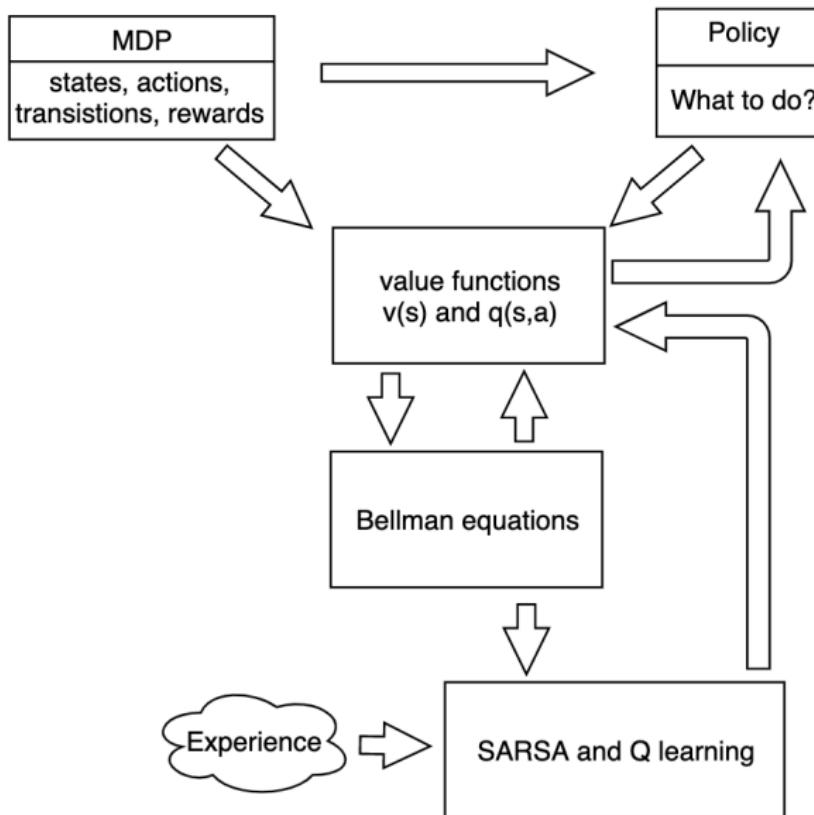
SARSA and Q-Learning for Model-free Control

Solving model-free RL using **Temporal Differencing**

- By **acting in the world**, the agent gains experience.
An **experience** can be expressed as a 4-tuple: (s, a, r, s')
- Over time the agent collects a list of experiences that he uses to find value functions (and corresponding policy) ... HOW?
- **Key observations:**
 - Bellman eqs. link values in neighbouring states along paths!
 - This backing-up improves learning efficiency!
 - Basis for RL algo's (**SARSA**, Q-learning, etc)

Temporal Differencing (TD): Use Bellman eqs to propagate values!

Exploiting the Bellman equations



Temporal-differencing: Exploiting Bellman eqs.

- For reasons of simplicity: let's start with v rather q
- Recall Bellman equation for v ($\gamma = 1$)

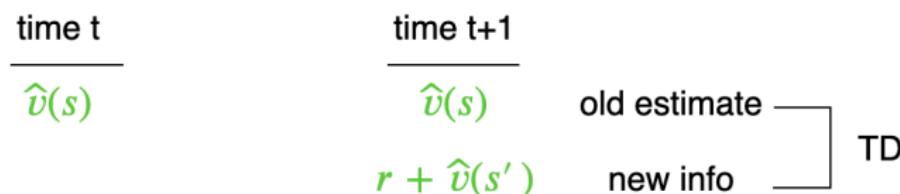
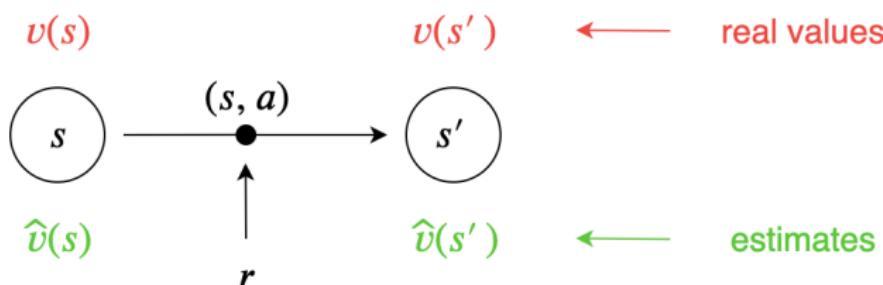
$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a, s') + v_{\pi}(s')]$$

- Sampling "*takes care*" of probabilities:

$$s \xrightarrow{\pi} a \quad (s, a) \xrightarrow{p} s'$$

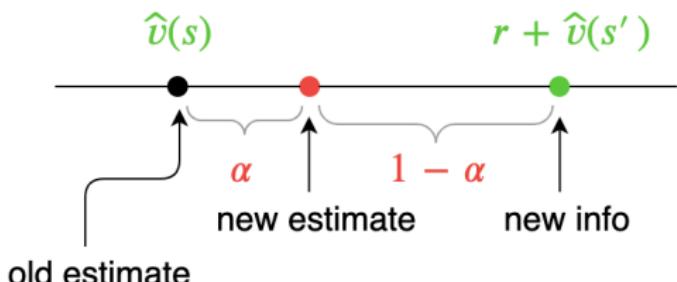
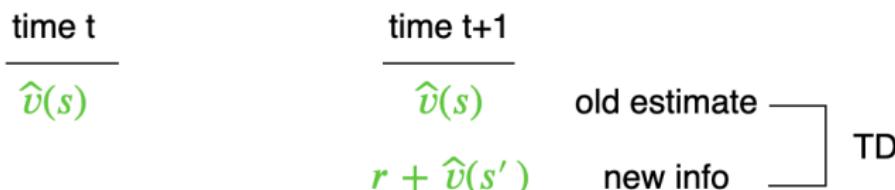
$$v_{\pi}(s) = r(s, a, s') + v_{\pi}(s')$$

Temporal-differencing: Exploiting Bellman eqs.



The 1-step reward r is new information gained from experience.

Temporal differencing for v : Exploiting Bellman eqs.



$$\hat{v}_{\text{new}}(s) = (1 - \alpha)\hat{v}_{\text{old}}(s) + \alpha(r + \hat{v}(s'))$$

Temporal Differencing (TD) versus Monte Carlo (MC)

Given learning rate α

- MC update rule:

$$v_{\pi}^{new}(S_t) \leftarrow (1 - \alpha)v_{\pi}^{old}(S_t) + \alpha G_t(S_t)$$

$$v_{\pi}^{new}(S_t) \leftarrow v_{\pi}^{old}(S_t) + \alpha [G_t(S_t) - v_{\pi}^{old}(S_t)]$$

- TD update rule: uses a different update target:

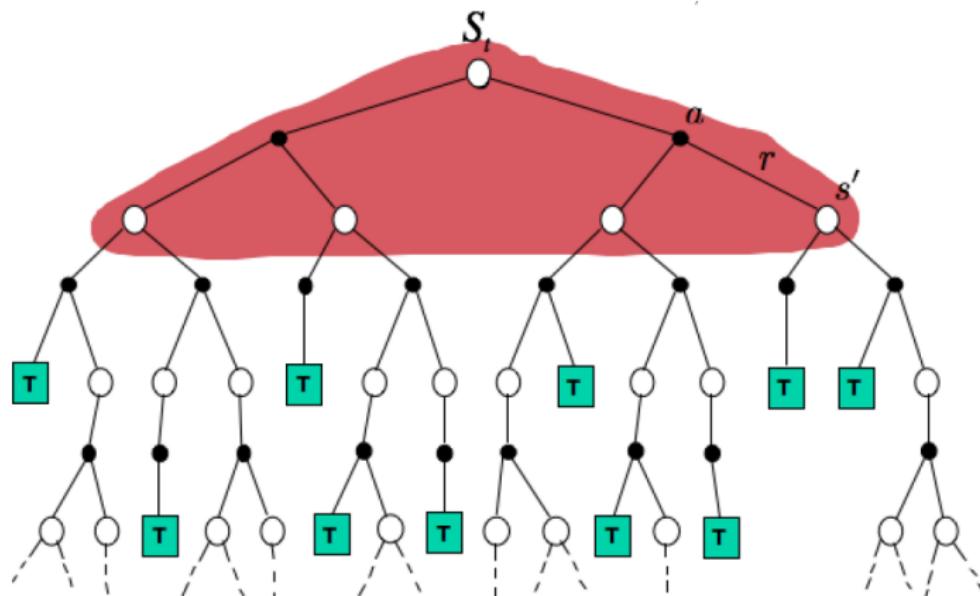
$$v_{\pi}^{new}(S_t) \leftarrow (1 - \alpha)v_{\pi}^{old}(S_t) + \alpha[R_{t+1} + \gamma v_{\pi}(S_{t+1})]$$

$$v_{\pi}^{new}(S_t) \leftarrow v_{\pi}^{old}(S_t) + \alpha[R_{t+1} + \gamma v_{\pi}(S_{t+1}) - v_{\pi}^{old}(S_t)]$$

- TD is a bootstrapping method:

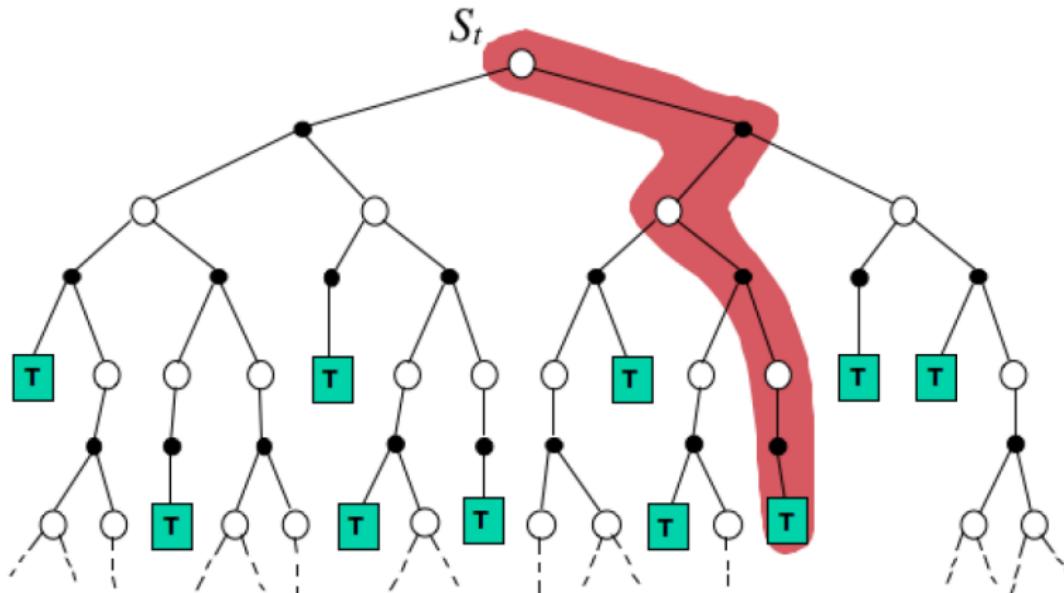
bases updates on existing estimates, like DP

Dynamic Programming (DP): Backup diagram



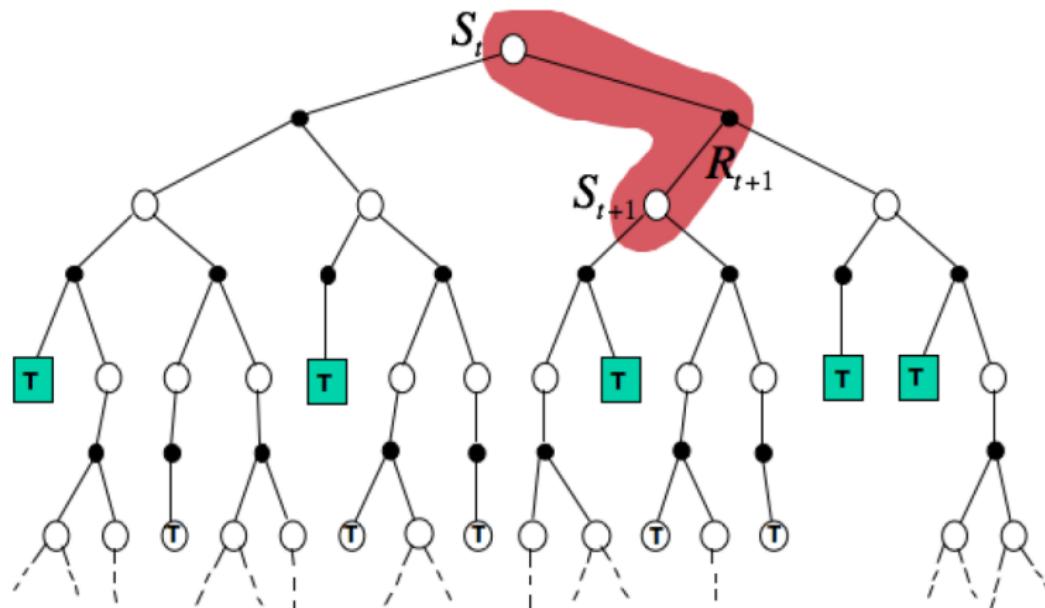
Monte-Carlo (MC): Backup diagram

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



Temporal Difference (TD): Backup diagram

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



Temporal-difference (TD) methods

- **DP** exploits Bellman equation but **requires model**
- **MC** doesn't require model but **doesn't exploit Bellman equation**
- **TD methods** can get the best of both worlds: **exploit Bellman equation without requiring a model**
- TD is therefore **core algorithm** for **model-free RL**

Policy evaluation: Estimating v_π using TD(0)

Initialize $V(s)$ arbitrarily, π to the policy to be evaluated

Repeat (for each episode):

 Initialize s

 Repeat (for each step of episode):

$a \leftarrow$ action given by π for s

 Take action a ; observe reward, r , and next state, s'

$V(s) \leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)]$

$s \leftarrow s'$

 until s is terminal

TD(0) backup diagram

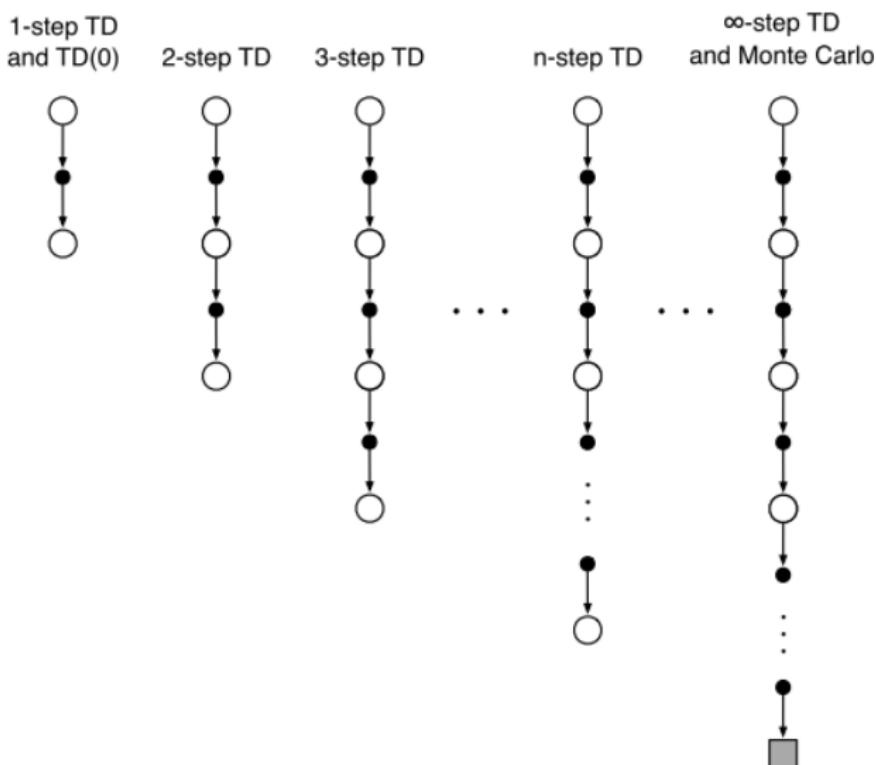
- Sampling: unlike DP but like MC, only one choice at each state
- Bootstrapping: like DP but unlike MC, use estimate from next state



Advantages of TD prediction methods

- TD methods require **only experience**, not a model
- TD, but not MC, methods can be **fully incremental**
- Learn **before knowing the final outcome**:
 - **Efficient**: less memory and peak computation
 - Learn from **incomplete sequences**
- Both MC and TD converge but TD tends to be faster

n-step TD: Spectrum between TD(0) and Monte Carlo



Issues when addressing model-free RL problems

1. We need to compute $q(s, a)$ rather than $v(s)$

- To **improve policy**, for every s , need to know $\max_a q(s, a)$:

$$\text{construct } \pi : s \mapsto a_{\max} \quad \text{where} \quad q(s, a_{\max}) = \max_{a'} q(s, a')$$

- Model-based:** $q(s, a)$ can be computed from $v(s)$:

$$q(s, a) = r(s, a, s') + v(s')$$

- Model-free:** $q(s, a)$ needs to be estimated explicitly!

Hence: SARSA and Q-learning focus on $q(s, a)$

2. Keep exploring!

- Balance exploration versus exploitation
- E.g. ϵ -greedy, soft-max, etc.

Model-based versus Model-free

oo

Model-free RL problems

oooooooooooo

Monte Carlo for Model-free Policy Evaluation

oooooooooooooooooooo

Temporal Difference Methods

oooooooooooooooooooo

Outline

Model-based versus Model-free

Model-free RL problems

Monte Carlo for Model-free Policy Evaluation

Temporal Difference Methods for Model-free Prediction

SARSA and Q-Learning for Model-free Control

Apply TD methods for estimating $q(s, a)$

TD methods Exploit correlations between **successor** states:

- **State value function** $v_\pi(s)$

$$v_\pi(S_t) \leftarrow v_\pi(S_t) + \alpha \underbrace{(R_{t+1} + \gamma v_\pi(S_{t+1}) - v_\pi(S_t))}_{\text{new info}}$$

- **State value function** $q_\pi(s, a)$:

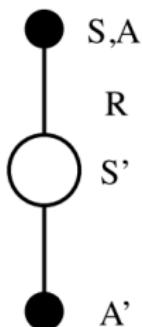
- **SARSA**: policy evaluation

$$q_\pi(s, a) \leftarrow ??$$

- **Q-learning**: optimal state-action value

$$q^*(s, a) \leftarrow ??$$

SARSA: TD estimation of Q-values



- **SARSA** In state s , take action a , transition to state s' , use policy π to select next action a'
- **Bellman:** Two estimate for state-action value:
 $q_\pi(s, a)$ and $r(s, a, s') + q_\pi(s', a')$

SARSA: TD for action values

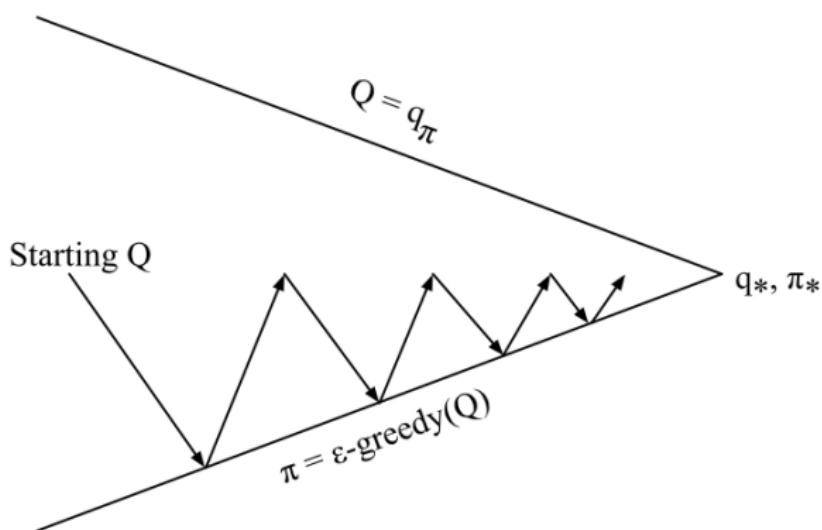
- **TD:** bootstrap update for **value function** v_π

$$v_\pi(s) \leftarrow v_\pi(s) + \alpha \underbrace{(r(s, a, s') + \gamma v_\pi(s') - v_\pi(s))}_{\text{new info}}$$

- **SARSA:** bootstrap update for **action value function** q_π

$$q_\pi(s, a) \leftarrow q_\pi(s, a) + \alpha \underbrace{(r(s, a, s') + q_\pi(s', a') - q_\pi(s, a))}_{\text{new info}}$$

SARSA model-free control: Schematically



Every **time-step**:

Policy evaluation **Sarsa**, $Q \approx q_\pi$

Policy improvement ϵ -greedy policy improvement

SARSA: Pseudo-code for model-free control

Initialize $Q(s, a)$ arbitrarily

Repeat (for each episode):

 Initialize s

 Choose a from s using policy derived from Q (e.g., ε -greedy)

 Repeat (for each step of episode):

 Take action a , observe r, s'

 Choose a' from s' using policy derived from Q (e.g., ε -greedy)

$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$

$s \leftarrow s'; a \leftarrow a'$;

 until s is terminal

SARSA convergence

Theorem

*Sarsa converges to the optimal action-value function,
 $Q(s, a) \rightarrow q_*(s, a)$, under the following conditions:*

- GLIE sequence of policies $\pi_t(a|s)$
- Robbins-Monro sequence of step-sizes α_t

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

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

Applying TD to model-free RL problems

Recall: Model-free, hence focus on $q(s, a)$ and use Bellman!

1. SARSA: Evaluating a given policy π

- $q_\pi(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a, s') + \sum_{a'} \pi(a' | s') q_\pi(s', a') \right]$
- **Sample:** $q_\pi(s, a) = r(s, a, s') + q_\pi(s', a')$
- **SARSA update rule** (sample-based):

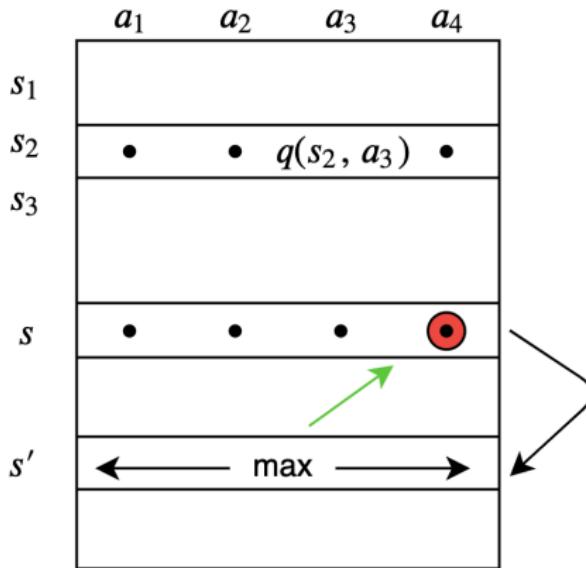
$$q_\pi(s, a) \leftarrow (1 - \alpha)q_\pi(s, a) + \alpha(r(s, a, s') + q_\pi(s', a'))$$

2. Q-learning: Estimating the **optimal** value function $q^*(s, a)$

- $q^*(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a, s') + \max_{a'} q^*(s', a') \right]$
- **Sample:** $q^*(s, a) = r(s, a, s') + \max_{a'} q^*(s', a')$
- **Q-learning update rule** (sample-based):

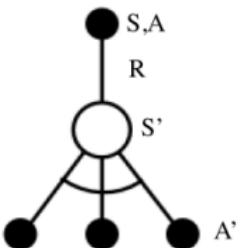
$$q^*(s, a) \leftarrow (1 - \alpha)q^*(s, a) + \alpha(r(s, a, s') + \max_{a'} q^*(s', a'))$$

Q-learning



Q-learning:

Bootstrap with **best** action (greedy!), not actual action



$$Q(S, A) \leftarrow Q(S, A) + \alpha \left(R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$

Convergence when TD-error vanishes: Recall optimality equation:

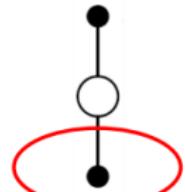
$$q^*(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a, s') + \gamma \max_{a'} q^*(s', a') \right]$$

Sampling version: $q^*(s, a) \approx r(s, a, s') + \gamma \max_{a'} q^*(s', a')$

SARSA vs. Q-learning: On-Policy vs. Off-Policy

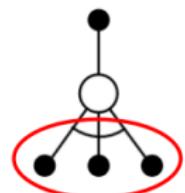
Sarsa:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$



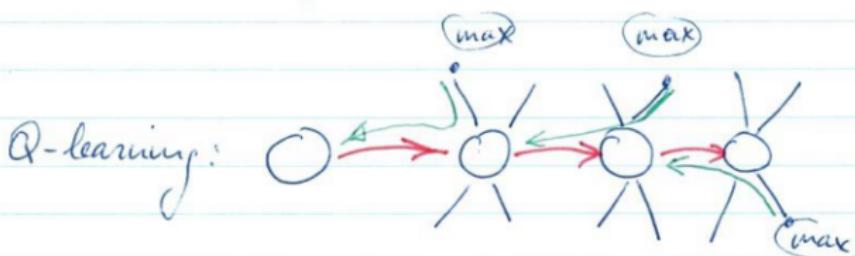
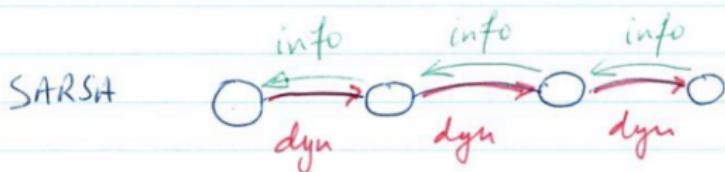
Q-learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$



SARSA (on-policy) vs. Q-Learning (off-policy)

SARSA vs. Q-learning
(ON- vs OFF-Policy)



Q-learning: behaviour policy different from backup policy

ON-policy vs. OFF-policy

ON-policy: behaviour policy = target policy

- Info gathered to improve policy, depends on that policy;

Feedback loop: policy \longleftrightarrow data

- Sample inefficient: experiences (s, a, r, s', a') cannot be re-used when policy changes since a' depends on actual policy!

OFF-policy: behaviour policy \neq target policy

- Improved sample efficiency: can re-use all samples for training (replay buffer);
 - E.g. in Q-learning: use data generated by behaviour policy π to learn value function for optimal policy π^* .
- Greedy updating (max!) prone to biased estimation! (e.g. Double Q-learning)

SARSA vs. Q-learning

Q-Learning addresses many of the limitations of SARSA

- Both based on TD learning
- SARSA updates depend on policy used to gather experiences;
- IN Q-learning the update is: $\max'_a Q(s', a')$
 - only depends s' , not on action a' which taken at that state,
 - hence, only depends on the environment ($s \xrightarrow{a} s'$) and **not on the policy** (that will change and updated as part of learning process)
- QL updates don't depend on actual policy, but on optimal policy.

Q-learning ALGO: off-policy TD control

Make TD off-policy: bootstrap with best action, not actual action:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

Initialize $Q(s, a)$ arbitrarily

Repeat (for each episode):

 Initialize s

 Repeat (for each step of episode):

 Choose a from s using policy derived from Q (e.g., ϵ -greedy)

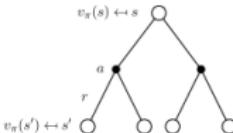
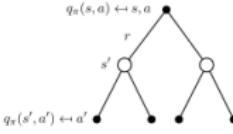
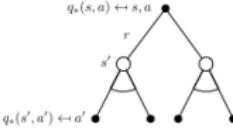
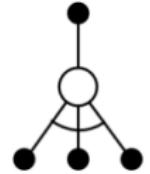
 Take action a , observe r, s'

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

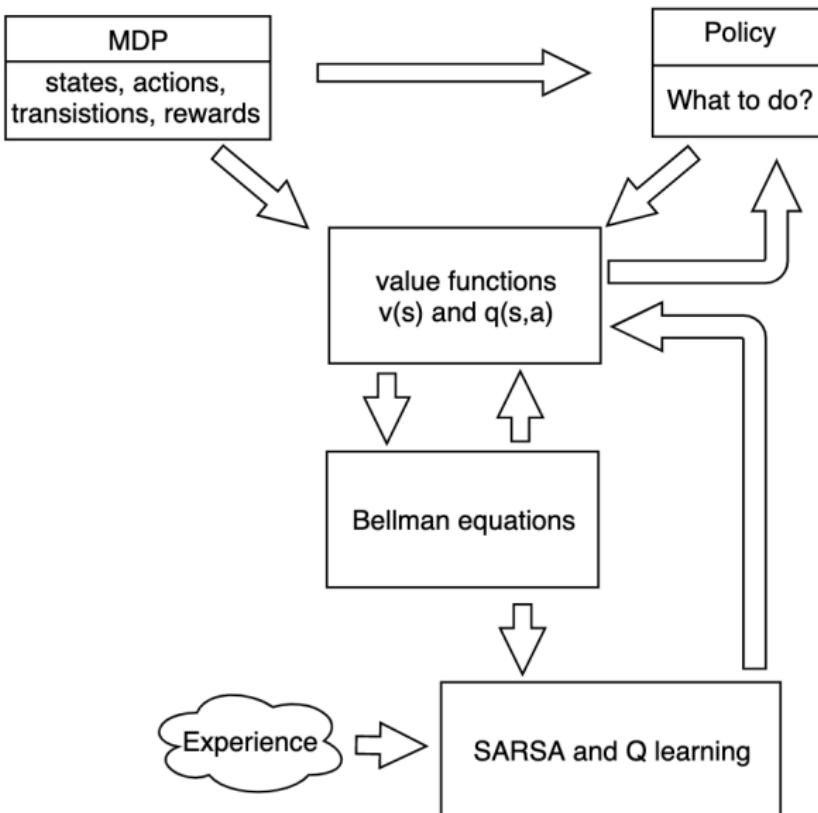
$s \leftarrow s'$;

 until s is terminal

Dynamic Program^{ing} (DP) vs. Temporal Diff^{ing} (TD)

	<i>Full Backup (DP)</i>	<i>Sample Backup (TD)</i>
Bellman Expectation Equation for $v_\pi(s)$	$v_\pi(s) \leftarrow s$  <p>$v_\pi(s') \leftarrow s'$</p> <p>Iterative Policy Evaluation</p>	 <p>TD Learning</p>
Bellman Expectation Equation for $q_\pi(s, a)$	$q_\pi(s, a) \leftarrow s, a$  <p>$q_\pi(s', a') \leftarrow a'$</p> <p>Q-Policy Iteration</p>	 <p>Sarsa</p>
Bellman Optimality Equation for $q_*(s, a)$	$q_*(s, a) \leftarrow s, a$  <p>$q_*(s', a') \leftarrow a'$</p> <p>Q-Value Iteration</p>	 <p>Q-Learning</p>

Summary



Summary: Model-based versus Model-free

	Prediction <i>Estimation:</i> <i>Given π, what is v?</i>	(Optimal) Control <i>Optimisation:</i> <i>What is optimal π?</i>
model-based (MDP given)	Policy evaluation using Dyn. Programming (DP)	Policy improvement (+ Policy evaluation) = Policy iteration
model-free (MDP unknown)	Monte Carlo (MC) Temporal Diff ^{ing} (TD) = "impatient MC" <i>bootstrapping!</i>	Q-Learning Generalized Policy Iteration <i>"simultaneous"</i>