

# Executive Summary

---

- Markov Decision Process
  - Markov Chain  $\langle \mathcal{S}, \mathcal{P} \rangle$
  - Markov Reward Process  $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
  - MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ 
    - MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle + \text{policy } \pi = \text{MRP } \langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$
    - $\mathcal{P}_{s,s'}^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{P}_{ss'}^a$
    - $\mathcal{R}_s^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{R}_s^a$
  - Bellman Expectation Equations
  - Bellman Optimality Equations
- Planning (Lec 3)
  - Synchronous DP
    - Policy Iteration
      - Policy Evaluation: estimate  $v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$
      - Iterative PE:  $\mathbf{v}^{k+1} = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi \mathbf{v}^k$
      - Policy Improvement: generate  $\pi' \geq \pi$ 
        - Greedy PI:  $\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_\pi(s, a)$
      - Value Iteration:  $v_*(s) \leftarrow \max_{a \in \mathcal{A}} \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$
    - Asynchronous DP
      - In-place DP
      - Prioritised Sweeping
      - Real-time DP
    - Approximate DP (using function approximation)
      - Fitted Value Iteration:  $\hat{v}_k(s) = \max_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \hat{v}(s', \mathbf{w}_k))$  for  $s \in \tilde{\mathcal{S}} \subseteq \mathcal{S}$ 
        - train  $\hat{v}(\cdot, \mathbf{w}_{k+1})$  using targets  $\{\langle s, \hat{v}_k(s) \rangle\}$
    - Contraction Mapping
  - Reinforcement Learning
    - Model Free
      - Prediction (Lec 4)
        - Monte Carlo Policy Evaluation:  $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$
        - TD Learning
          - TD(0):  $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$
          - Forward TD( $\lambda$ ):  $V(S_t) \leftarrow V(S_t) + \alpha(G_t^\lambda - V(S_t))$
          - Backward TD( $\lambda$ ):  $V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$ 
            - $E_t(s) = \gamma \lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$
        - Control (Lec 5)
          - On-Policy
            - Monte-Carlo Policy Iteration:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(G_t - Q(S_t, A_t))$ 
              - GLIE Monte-Carlo Control
            - SARSA: On-Policy TD Learning based Policy Iteration:
 
$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A))$$
              - Forward Sarsa( $\lambda$ ):  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(q_t^\lambda - Q(S_t, A_t))$
              - Backward Sarsa( $\lambda$ ):  $Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t E_t(s, a)$ 
                - $E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + \mathbf{1}(S_t = s, A_t = a)$
            - Off-Policy
              - Off-Policy Monte Carlo (*impractical* as too high variance):
 
$$V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} \dots \frac{\pi(A_T | S_T)}{\mu(A_T | S_T)} G_t - V(S_t) \right)$$
              - Off-Policy TD Learning:  $V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$
              - Q-Learning:  $Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma \max_{a'} Q(S', a') - Q(S, A))$

- Value Function Approximation (*Lec 6*)
  - Incremental Methods
    - Incremental Prediction  $\hat{v}(S, \mathbf{w}) \approx v_\pi(S)$
    - MC:  $\Delta \mathbf{w} = \alpha(G_t - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
    - TD(0):  $\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
    - Forward TD( $\lambda$ ):  $\Delta \mathbf{w} = \alpha(G_t^\lambda - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
    - Backward TD( $\lambda$ ):  $\Delta \mathbf{w} = \alpha \delta_t E_t$ 
      - $E_t = \gamma \lambda E_{t-1} + \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
  - Incremental Control
    - MC:  $\Delta \mathbf{w} = \alpha(G_t - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
    - TD(0):  $\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
    - Forward-view TD( $\lambda$ ):  $\Delta \mathbf{w} = \alpha(q_t^\lambda - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
    - Backward-view TD( $\lambda$ ):  $\Delta \mathbf{w} = \alpha \delta_t E_t$ 
      - $E_t = \gamma \lambda E_{t-1} + \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
  - Batch Methods
- Policy Gradient Methods (*Lec 7*)
  - Policy Gradient Theorem
  - Monte-Carlo PG
    - REINFORCE:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) v_t]$
  - Compatible Function Approximation Theorem
    - Compatible Value Function:  $\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$
  - Actor-Critic PG
    - Q AC:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^w(s, a)]$
    - Advantage AC:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^w(s, a)]$
    - TD(0) AC:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) (r + \gamma V^w(s_{t+1}) - V^w(s_t))]$
    - Forward TD( $\lambda$ ) AC:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) (v_t^\lambda - V^w(s_t)) e]$
    - Backward TD( $\lambda$ ) AC:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \delta_t e_t]$ 
      - $e_t = \gamma \lambda e_{t-1} + \nabla_\theta \log \pi_\theta(s_t, a_t)$
  - Natural PG
    - Natural AC:  $\nabla_\theta^{nat} J(\theta) = G_\theta^{-1} \nabla_\theta J(\theta) = w$
- Model Based (*Lec 8*)
  - Model  $\mathcal{M} = \langle \mathcal{P}_\eta, \mathcal{R}_\eta \rangle$  where  $\mathcal{P}_\eta \approx \mathcal{P}$  and  $\mathcal{R}_\eta \approx \mathcal{R}$
  - Model Learning
    - Supervised learning problem  $\{S_t, A_t \rightarrow R_{t+1}, S_{t+1}\}_{t=1}^{T-1}$ 
      - $s, a \rightarrow r$ : regression
      - $s, a \rightarrow s'$ : density estimation
  - Table Lookup Model
  - Sample-based Planning
  - Dyna - Integrated Planning and Learning

## Part I: Elementary Reinforcement Learning

### Lecture 1: Introduction to Reinforcement Learning

#### The RL Problem

- A reward  $R_t$  is a single scalar feedback signal
- Reward Hypothesis:** All goals can be described by the maximisation of expected cumulative reward.
- Sequential Decision Making**
  - Actions may have long term consequences
  - Reward may be delayed
  - It may be better to sacrifice immediate reward to gain more long-term reward
- Formalism:
  - observation ( $O_t$ ), action ( $A_t$ ), reward ( $R_t$ )

- **Agent** interacts with the **Environment**
- **History:** sequence of observations, actions, rewards
  - $H_t = O_1, R_1, A_1, \dots, A_{t-1}, O_t, R_t$
  - Decides what happens next:
    - agent selects action
    - environment selects observations/rewards.
- **State:**  $S_t = f(H_t)$ 
  - **Environment State** ( $S_T^e$ ) - environment's private representation
  - **Agent State** ( $S_T^a$ ) - agent's internal representation
  - **Information State or Markov State** - contains all the useful information from the history
    - Once Markov state  $S_t$  is known, history may be thrown away
      - $H_{1:t} \rightarrow S_t \rightarrow H_{t+1:\infty}$
      - State is sufficient statistic of future
    - Eg.  $S_T^e$  and  $H_t$
- **Fully Observable Environments**
  - Agent directly observes environment state  $O_t = S_a^t = S_e^t$
  - Formally: *MDP*.
- **Partially Observable Environments**
  - Agent indirectly observes environment.
  - Formally: *POMDP*
  - Agent must construct its own state representation  $S_t^a$ 
    - Complete history:  $S_t^a = H_t$
    - Beliefs of environment state:  $S_t^a = (\mathbb{P}[S_t^e = s^1], \dots, \mathbb{P}[S_t^e = s^n])$
    - Recurrent neural network:  $S_t^a = \sigma(S_{t-1}^a W_s + O_t W_o)$
- **Components of an RL Agent**
  - Policy ( $\pi$ ): agent's behaviour function
    - Deterministic or Stochastic
  - Value function: how good is each state and/or action
    - $v_\pi(s) = \mathbb{E}_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$
  - Model: agent's representation of the environment
    - $\mathcal{P}$ : state transition matrix:  $\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$
    - $\mathcal{R}$ : rewards:  $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
    - Not required
  - Categorisation 1:
    - **Value Based:** Only value Function
    - **Policy Based:** Only Policy
    - **Actor Critic:** Both Policy and Value Function
  - Categorisation 2:
    - **Model Free**
    - **Model Based**
- Fundamental problems in sequential decision making
  - **Reinforcement Learning:**
    - The environment is initially unknown
    - The agent interacts with the environment
    - The agent improves its policy
  - **Planning:**
    - The model/dynamics of the environment is known
    - The agent performs computations with its model (without any external interaction)
    - The agent improves its policy
- Exploration and Exploitation
  - **Exploration** finds more information about the environment
  - **Exploitation** exploits known information to maximise reward
  - Both are important in RL
- Prediction and Control
  - **Prediction:** evaluate the future given a policy
  - **Control:** find the best policy to optimise the future
  - typically need to solve *prediction* to solve *control*

## Lecture 2: Markov Decision Processes

---

- Almost all RL problems can be formalised as MDPs,
- **Markov Property** - state captures all relevant information from history
- **Markov Process or Chain:** memoryless random process:  $\langle \mathcal{S}, \mathcal{P} \rangle$ 
  - $\mathcal{S}$  - finite set of states
  - $\mathcal{P}$ : State Transition Matrix:  $\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$
- **Markov Reward Process:**  $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ 
  - $\mathcal{R}$ : reward function:  $\mathcal{R}_s = \mathbb{E}[R_{t+1} | S_t = s]$ 
    - reward from exiting a state
  - $\gamma \in [0, 1]$ : discount factor
- **Return:** total discounted reward from time-step  $t$ :  $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
- Why discount?
  - Mathematically convenient to discount rewards
  - Avoids **infinite returns** in cyclic Markov processes
  - Preference for more immediate reward
  - Uncertainty about the future may not be fully represented
  - May use  $\gamma = 1$  (no discounting)
    - If all sequences terminate
    - Average reward formulation
- **Value Function:**  $v(s) = \mathbb{E}[G_t | S_t = s]$
- **Bellman Equation for MRPs**
  - $v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$
  - $v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s')$
  - in matrix form:  $v = \mathcal{R} + \gamma \mathcal{P}v$
  - Solving for  $v$ 
    - Analytically -  $O(n^3)$
    - Iterative Methods
      - Dynamic programming
      - Monte-Carlo evaluation
      - Temporal-Difference learning
- **Markov Decision Processes:**  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ 
  - $\mathcal{A}$  - finite set of actions
  - $\mathcal{P}$ : state transition matrix:  $\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$
  - $\mathcal{R}$ : reward function:  $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
- $\pi$ : Policy:  $\pi(a | s) = \mathbb{P}[A_t = a | S_t = s]$ 
  - Policies are stationary (time-independent).
- Given an MDP  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and a policy  $\pi$ 
  - The state sequence  $S_1, S_2, \dots$  is a Markov process  $\langle \mathcal{S}, \mathcal{P}^\pi \rangle$
  - The state and reward sequence  $S_1, R_2, S_2, \dots$  is a Markov reward process  $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$ 
    - $\mathcal{P}_{ss'}^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{P}_{ss'}^a$
    - $\mathcal{R}_s^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{R}_s^a$
- Value functions:
  - **state-value function:**  $v^\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$
  - **action-value function:**  $q^\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$
- **Bellman Expectation Equation for MDPs**
  - for state-value function
    - $v_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$
    - $v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) q_\pi(s, a) = \sum_{a \in \mathcal{A}} \pi(a | s) (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s'))$
  - for action-value function
    - $q_\pi(s, a) = \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$
    - $q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a' | s') q_\pi(s', a')$
  - in matrix form:  $v_\pi = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v_\pi$
- **Optimal Value Functions**
  - $v_*(s) = \max_\pi v_\pi(s)$

- $q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$ 
  - This can be used to make decision on which action to take
- **Optimal Policy:**  $\pi_*$ 
  - *Partial Ordering* over policies:
  - Define a partial ordering over policies:  $\pi \geq \pi'$  if  $v_{\pi}(s) \geq v_{\pi'}(s), \forall s$
  - For any MDP,
    - There exists an optimal policy  $\pi_*$  that is better than or equal to all other policies,  $\pi_* \geq \pi, \forall \pi$
    - $v_{\pi_*}(s) = v_*(s)$
    - $q_{\pi_*}(s, a) = q_*(s, a)$
  - **Deterministic Optimal Policy**
    - Always exists for any MDP i.e. in fully observed setting
      - One such  $\pi_*$  can be obtained from  $q_*$
      - May not exist in partially observed setting.
        - Eg. when there is state-aliasing or when value function approximation is used.
- **Bellman Optimality Equation for  $v_*$** 
  - $v_*(s) = \max_a q_*(s, a) = \max_a \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$
  - $q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a'} q_*(s', a')$
  - max instead of expectation/summation over  $a \sim \pi(a | s)$
  - Only Iterative methods (No analytical solution)
    - Dynamic Programming
      - Value Iteration
      - Policy Iteration
    - Q-learning
    - Sarsa
- **Extensions to MDPs**
  - Infinite and continuous MDPs
  - **Partially observable MDPs (POMDPs):**  $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{P}, \mathcal{R}, \mathcal{Z}, \gamma \rangle$ 
    - $\mathcal{O}$ : finite set of observations
    - $\mathcal{Z}$ : observation function:  $\mathcal{Z}_{s' o}^a = \mathbb{P}[O_{t+1} = o | S_{t+1} = s', A_t = a]$
    - Belief State
  - **Undiscounted Average Reward MDPs**

## Lecture 3: Planning by Dynamic Programming

---

### Introduction

- **Planning:**
  - The MDP/model/dynamics of the environment is known
  - The agent performs computations with its model to improve its policy.
- MDP satisfy both requirements of DP
  - Optimal substructure - bellman equation
  - Overlapping subproblems properties - value function
- DP assumes full knowledge of MDP and is used for planning
  - For prediction:
    - *Input:* MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and policy  $\pi$  or MRP  $\langle \mathcal{S}, \mathcal{P}^{\pi}, \mathcal{R}^{\pi}, \gamma \rangle$
    - *Output:* value function  $v_{\pi}$
  - For control:
    - *Input:* MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
    - *Output:* optimal value function  $v_*$  and optimal policy  $\pi_*$

### Policy Evaluation

- **Iterative Policy Evaluation**
  - *Solution:* iterative application of Bellman expectation backup
  - $v_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a | s) (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s'))$
  - in matrix form:  $\mathbf{v}^{k+1} = \mathcal{R}^{\pi} + \gamma \mathcal{P}^{\pi} \mathbf{v}^k$
  - synchronous or asynchronous backups

## Policy Iteration

- **Policy Iteration**
  - **Policy Evaluation:** Any policy evaluation algorithm to estimate  $v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$ 
    - Eg. *Iterative Policy Evaluation*
    - Don't need to converge to  $v_\pi$  -> Other algorithms
    - different stopping criterion in iterative policy evaluation:
      - $\epsilon$ -convergence of value function
      - stop after  $k$  iterations
      - $k = 1$  gives *value iteration*
  - **Policy Improvement:** Any policy improvement algorithm to generate  $\pi' \geq \pi$ 
    - Eg. *Greedy Policy Improvement*
- **Greedy Policy Improvement:**
  - \* Acting greedily with respect to  $v_\pi: \pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_\pi(s, a)$ 
    - \* Improves policy
    - \* taking greedy action in next step only
      - \*  $q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_\pi(s, a) \geq q_\pi(s, \pi(s)) = v_\pi(s)$
    - \* improves value function if acting greedily in all steps
      - \*  $v_\pi(s) \leq q_\pi(s, \pi'(s)) = \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] = v_{\pi'}(s)$
    - \* If improvement stops:  $V_\pi(s) = V_*(s)$  and  $\pi$  is optimal.
    - \*  $q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_\pi(s, a) = q_\pi(s, \pi(s)) = v_\pi(s)$
    - \* Bellman Optimality Equations are satisfied
    - \* *Policy Iteration with Greedy Policy Improvement always converges*

## Value Iteration

- **Principle of Optimality:** For a policy  $\pi(a | s)$  and any state  $s$ ,  $v_\pi(s) = v_*(s)$  iff for any state  $s'$  reachable from  $s$ ,  $v_\pi(s') = v_*(s')$ .
- **Algorithm:** Directly solve for  $v_*$  by iteratively applying this update to all states:
  - **Bellman Optimality Backup:**  $v_*(s) \leftarrow \max_{a \in \mathcal{A}} \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$
  - May be synchronous or asynchronous.
- **Intuition:** start with final rewards and work backwards.
- No explicit policy unlike policy iteration.
  - Intermediate value functions may not correspond to any policy.
- Still works with loopy, stochastic MDPs.
- Equivalent to Policy Iteration with single iteration of Iterative Policy Evaluation and Greedy Policy Improvement.

## Summary of Synchronous DP Algorithms

| Problem    | Bellman Equation   | Algorithm                   |
|------------|--|-----------------------------|
| Prediction | Bellman Expectation Equation                             | Iterative Policy Evaluation |
| Control    | Bellman Expectation Equation + Greedy Policy Improvement | Policy Iteration            |
| Control    | Bellman Optimality Equation                              | Value Iteration             |

- Algorithms are based on state-value function  $v_\pi(s)$  or  $v_*(s)$
- Complexity  $O(mn^2)$  per iteration, for  $m$  actions and  $n$  states
- Could also apply to action-value function  $q_\pi(s, a)$  or  $q_*(s, a)$
- Complexity  $O(m^2n^2)$  per iteration

## Extensions to Dynamic Programming

- **Asynchronous Dynamic Programming Algorithms**
  - Synchronous DP uses synchronous backups i.e. all states are backed up in parallel
  - Ideas for Asynchronous DP:
    - **In-place dynamic programming**
      - Synchronous VI keeps two copies of value function:  $v_{new}(s)$  and  $v_{old}(s)$
      - for all  $s$ :
      - $v_{new}(s) \leftarrow \max_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{old}(s'))$

- $v_{new}(s) \leftarrow v_{old}(s)$
- In-place VI only keeps one copy:
  - for all  $s$ :
  - $v(s) \leftarrow \max_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v(s'))$
- **Prioritised sweeping**
  - Bellman Error:  $|\max_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v(s')) - v(s)|$
  - Backup state with largest remaining Bellman Error
- **Real-time dynamic programming**
  - Idea: Only backup states that are relevant to agent
  - Use agent's experience to guide the selection of states
  - At each time-step sample  $S_t, A_t, R_{t+1}$  and backup state  $S_t$
- Backup type
  - **Full Width Backup**
    - For each backup (sync or async)
      - Every successor state and action is considered
      - Requires knowledge of the MDP transitions and reward function
    - Used by DP
    - Problem: *Curse of Dimensionality*- even one full-width backup may be too expensive
  - **Sample Backup**
    - Sample rewards and transitions  $\langle S, A, R, S' \rangle$
    - Pros:
      - Model-free: no advance knowledge of MDP required
      - Breaks the curse of dimensionality through sampling
      - Cost of backup is constant, independent of  $n = |\mathcal{S}|$
- **Approximate DP**
  - Apply DP to a function approximator  $\hat{v}(s; w)$  for the value function
  - Example: **Fitted Value Iteration**
    - Repeat at iteration  $k$ :
    - Sample states  $\tilde{\mathcal{S}} \subseteq \mathcal{S}$
    - For each state  $s \in \tilde{\mathcal{S}}$ , estimate target value using Bellman optimality equation,
      - $\tilde{v}_k(s) = \max_{a \in \mathcal{A}} (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \hat{v}(s', \mathbf{w}_k))$
      - Train next value function  $\hat{v}(\cdot, \mathbf{w}_{k+1})$  using targets  $\{(s, \tilde{v}_k(s))\}$

## Contraction Mapping

- $l_\infty$  norm distance between value functions:  $\|u - v\|_\infty = \max_{s \in \mathcal{S}} |u(s) - v(s)|$
- *Bellman Expectation Backup is a Contraction*
  - Bellman expectation backup operator  $T^\pi: T^\pi(v) = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v$
  - $T^\pi$  is a  $\gamma$ -contraction:  $\|T^\pi(u) - T^\pi(v)\|_\infty \leq \gamma \|u - v\|_\infty$
- *Bellman Optimality Backup is a Contraction*
  - Bellman optimality backup operator  $T^*: T^*(v) = \max_{a \in \mathcal{A}} \mathcal{R}^a + \gamma \mathcal{P}^a v$
  - $T^*$  is a  $\gamma$ -contraction:  $\|T^*(u) - T^*(v)\|_\infty \leq \gamma \|u - v\|_\infty$
- **Contraction Mapping Theorem**
  - For any metric space  $\mathcal{V}$  that is complete and hence closed under an operator  $T(v)$ , where  $T$  is a  $\gamma$ -contraction,
    - $T$  converges to a *unique fixed point* at a *linear convergence rate* of  $\gamma$
    - $T^\pi$  converges to  $v_\pi$  (a fixed point of  $T^\pi$  by Bellman expectation equation)
      - Iterative policy evaluation converges on  $v_\pi$
      - Policy iteration converges on  $v_*$
    - $T^*$  converges to  $v_*$  (a fixed point of  $T^*$  By Bellman Optimality equation)
      - Value iteration converges on  $v_*$

## Lecture 4: Model-Free Prediction

### Monte-Carlo Learning

- **Monte Carlo Methods**
  - MC methods learn directly from complete episodes of experience

- uses **complete episodes**: no bootstrapping
  - Model-free: no knowledge of MDP transitions / rewards
  - MC uses the simplest possible idea: value = mean return
  - **Can only apply MC to episodic MDPs**
- **Monte-Carlo Policy Evaluation**
  - Goal: learn  $v_\pi$  from episodes of experience under policy  $\pi$ 
    - $S_1, A_1, R_2, \dots, S_k \sim \pi$
    - return = total discounted reward:  $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$
    - value function = expected return:  $v_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$
    - Monte-Carlo policy evaluation uses empirical mean return instead of expected return
      - $V(s) = S(s)/N(s)$  where
        - $N(s)$  - counter for number of times  $s$  is visited
        - $S(s)$  - total return from  $s$
      - By law of large numbers,  $V(s) \rightarrow v_\pi(s)$  as  $N(s) \rightarrow \infty$
  - **First-Visit** Monte-Carlo Policy Evaluation
    - For *first time-step*  $t$  that  $s$  is visited in an episode:
      - Increment counter  $N(s) \leftarrow N(s) + 1$
      - Increment total return  $S(s) \leftarrow S(s) + G_t$
  - **Every-Visit** Monte-Carlo Policy Evaluation
    - For *every time-step*  $t$  that  $s$  is visited in an episode:
      - Increment counter  $N(s) \leftarrow N(s) + 1$
      - Increment total return  $S(s) \leftarrow S(s) + G_t$
  - **Incremental Monte-Carlo Updates**
    - After each episode, for each state  $S_t$  with return  $G_t$ :
      - $N(S_t) \leftarrow N(S_t) + 1$
      - $V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)}(G_t - V(S_t))$
    - In non-stationary problems where it can be useful to forget old episodes.
      - $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$

## Temporal-Difference Learning

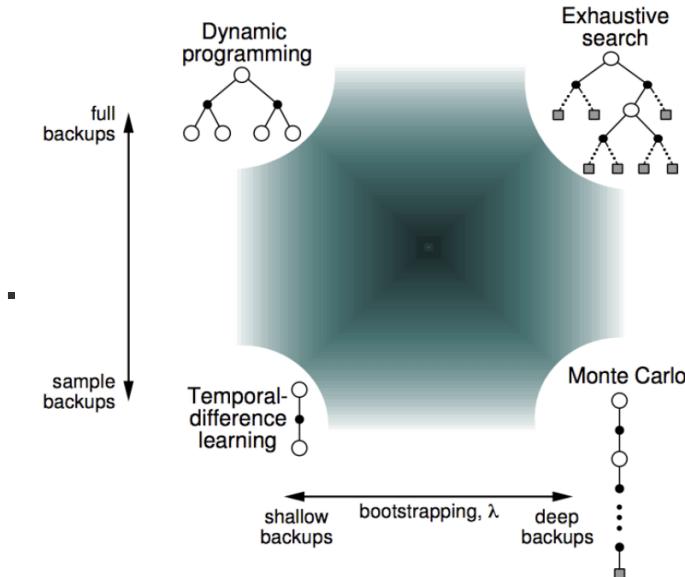
- **TD Methods**
  - TD methods learn directly from episodes of experience
    - works with **incomplete episodes** by bootstrapping
  - Model-free: no knowledge of MDP transitions / rewards
  - TD updates a guess towards a guess
- **TD(0):**
  - Update  $V(S_t)$  toward estimated return:  $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$ 
    - **TD target:**  $R_{t+1} + \gamma V(S_{t+1})$
    - **TD Error:**  $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$
- **MC v/s TD(0)**
  - Update:
    - MC: Update toward actual return:  $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$
    - TD(0): Update toward estimated return:  $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$
  - **Bias/Variance Tradeoff**
    - Return  $G_t$  is an unbiased estimate of  $v_\pi(S_t)$ 
      - But depends on many random actions, transitions, rewards => higher variance
    - TD target  $R_{t+1} + \gamma V(S_{t+1})$  is a biased estimate of  $v_\pi(S_t)$ 
      - But depends on one random actions, transitions, rewards => lower variance
  - **Certainty Equivalence of Batch MC and TD**
    - MC converges to solution with minimum mean-squared error
      - Finds the best fit to the observed returns
        - $\sum_{k=1}^K \sum_{t=1}^{T_k} (G_t^k - V(s_t^k))^2$
    - TD (0) converges to solution of max likelihood Markov model
      - Solves the MDP  $\langle \mathcal{S}, \mathcal{A}, \hat{\mathcal{P}}, \hat{\mathcal{R}}, \gamma \rangle$  that best fits the data

$$\hat{\mathcal{P}}_{s,s'}^a = \frac{1}{N(s,a)} \sum_{k=1}^K \sum_{t=1}^{T_k} \mathbf{1}(s_t^k, a_t^k, s_{t+1}^k = s, a, s')$$

$$\hat{\mathcal{R}}_s^a = \frac{1}{N(s,a)} \sum_{k=1}^K \sum_{t=1}^{T_k} \mathbf{1}(s_t^k, a_t^k = s, a) r_t^k$$

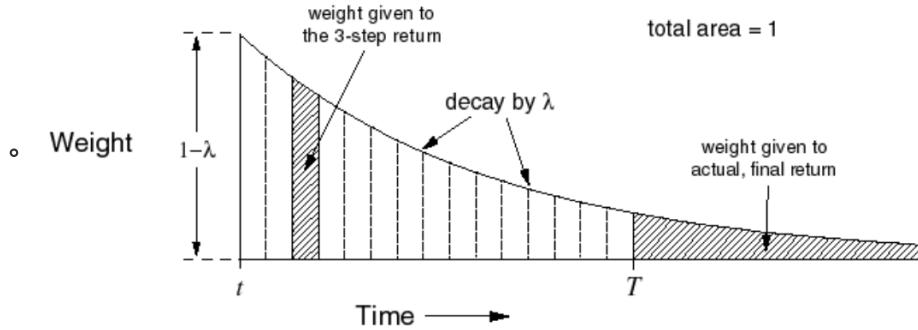
- Advantages and Disadvantages of MC vs. TD

- Shallow vs deep Backup:
  - TD can learn online after every step before knowing the final outcome
  - MC must wait until end of episode when return is known
- Episodic v/s Continuing environments:
  - TD can learn from incomplete episodes and works with non-episodic environments.
  - MC can only learn from complete episodes and only works for episodic (terminating) environments
- sample efficiency: TD is more sample-efficient
- bias-variance tradeoff: MC has higher variance but zero bias
- convergence:
  - \* MC has good convergence properties even with function approximation
  - \* TD may not converge with function approximation
  - \* TD is sensitive to initial value
- Markov Property:
  - TD exploits Markov property and is more efficient in Markov environments
  - MC does not and is more effective in non-Markov environments.
- Summary of Backups:**
  - MC:  $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$
  - TD(0):  $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$
  - DP:  $V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$
- Bootstrapping and Sampling**
  - Bootstrapping: update involves an estimate: eg. DP, TD
  - Sampling: update samples an expectation: eg. MC, TD



## TD( $\lambda$ ) and Multi-Step Returns

- n-Step Return:**  $G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$ 
  - TD( $n$ ):  $V(S_t) \leftarrow V(S_t) + \alpha(G_t^{(n)} - V(S_t))$ 
    - $n = 0$  - TD(0)
    - $n = \infty$  - MC
- $\lambda$ -return:**  $G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$ 
  - geometric weighting are *memoryless* => doesn't require storing anything extra for all  $G_t^{(n)}$



- **Forward-view TD( $\lambda$ ):**  $V(S_t) \leftarrow V(S_t) + \alpha (G_t^\lambda - V(S_t))$

- **Eligibility Traces**

- **Credit Assignment problem:**

- **Frequency heuristic:** assign credit to most frequent states
    - **Recency heuristic:** assign credit to most recent states

- Eligibility Traces combine both heuristics

- $E_0(s) = 0$
    - $E_t(s) = \gamma \lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$

- **Backward-view TD( $\lambda$ ):**  $V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$

- Idea:

- Keep an eligibility trace for every state  $s$
    - Update  $V(s)$  in proportion to TD-error  $\delta_t$  and eligibility trace  $E_t(s)$
- Backward-view provides mechanism for using  $\lambda$ -return that allows
  - \* Updating online, every step, from incomplete sequences
- Eligibility trace allows earlier states and actions in the episode to get updates even if the reward comes at the end of the episode.
  - $\lambda$  controls how far back in the episode the updates go and hence the bias-variance tradeoff as further back you go more random variables will the updates pass through.

- When  $\lambda = 0$ , only current state is updated

- $E_t(s) = \mathbf{1}(S_t = s)$
    - $V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$ 
      - Equivalent to TD(0) update:  $V(S_t) \leftarrow V(S_t) + \alpha \delta_t$

- When  $\lambda = 1$ , credit is deferred until end of episode

- In an episodic environment, it has the same update as TD( $\lambda$ )
    - The sum of offline updates is identical for forward-view and backward-view TD( $\lambda$ )
      - $\sum_{t=1}^T \alpha \delta_t E_t(s) = \sum_{t=1}^T \alpha (G_t^\lambda - V(S_t)) \mathbf{1}(S_t = s)$

- **TODO:** Finish This

- 

| Offline updates | $\lambda = 0$ | $\lambda \in (0, 1)$         | $\lambda = 1$      |
|-----------------|---------------|------------------------------|--------------------|
| Backward view   | TD(0)<br>     | TD( $\lambda$ )<br>          | TD(1)<br>          |
| Forward view    | TD(0)         | Forward TD( $\lambda$ )      | MC                 |
| Online updates  | $\lambda = 0$ | $\lambda \in (0, 1)$         | $\lambda = 1$      |
| Backward view   | TD(0)<br>     | TD( $\lambda$ )<br>*<br>     | TD(1)<br>*<br>     |
| Forward view    | TD(0)<br>     | Forward TD( $\lambda$ )<br>  | MC<br>             |
| Exact Online    | TD(0)         | Exact Online TD( $\lambda$ ) | Exact Online TD(1) |

## Lecture 5: Model-Free Control

### Introduction

- **On-policy learning:** Learn about policy  $\pi$  from experience sampled from  $\pi$
- **Off-policy learning:** Learn about policy  $\pi$  from experience sampled from  $\mu$

### On-Policy Monte-Carlo Control

- **Model-Free Policy Iteration**

- First attempt:
  - *Policy Evaluation:* Monte-Carlo Policy Evaluation
  - *Policy improvement:* Greedy policy improvement over  $V(s)$
- *Problem 1:* Greedy policy improvement over  $V(s)$  requires model of MDP
  - $\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} \mathcal{R}_s^a + \mathcal{P}_{ss'}^a V(s')$
  - *Solution:* Use  $Q(s, a)$  instead as Greedy policy improvement over it is model-free
    - $\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$
- *Problem 2:* If acting greedily with Monte-Carlo policy evaluation will not explore all states and end up with incorrect estimates of the value function.
  - *Solution:* Ensure Exploration
    - **$\epsilon$ -Greedy Exploration**
      - With probability  $1 - \epsilon$  choose the greedy action
      - With probability  $\epsilon$  choose an action at random
    - **$\epsilon$ -Greedy Policy Improvement:** For any  $\epsilon$ -greedy policy  $\pi$ , the  $\epsilon$ -greedy policy  $\pi'$  with respect to  $q_\pi$  is an improvement,  $v_{\pi'}(s) \geq v_\pi(s)$

- **Monte-Carlo Policy Iteration**

- *Algorithm:*
  - *Policy Evaluation:* Monte-Carlo Policy Evaluation,  $Q = q_\pi$
  - *Policy improvement:*  $\epsilon$ -Greedy policy improvement over  $Q(s, a)$
- *Problem:* Monte-Carlo is inefficient
  - *Solution:* find  $Q \approx q_\pi$  using a few or even just one episode

- **Greedy in the Limit with Infinite Exploration (GLIE)**

- All state-action pairs are explored infinitely many times,
  - $\lim_{k \rightarrow \infty} N_k(s, a) = \infty$
- The policy converges on a greedy policy,
  - $\lim_{k \rightarrow \infty} \pi_k(a | s) = \mathbf{1} \left( a = \operatorname{argmax}_{a' \in \mathcal{A}} Q_k(s, a') \right)$
- Example:  $\epsilon$ -greedy is GLIE if  $\epsilon$  decays to zero as  $\epsilon_k = 1/k$

- **GLIE Monte-Carlo Control**

- *Algorithm:*
  - Sample  $k$ th episode using  $\pi : \{S_1, A_1, R_2, \dots, S_T\} \sim \pi$
  - 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
    - $\epsilon \leftarrow 1/k$
    - $\pi \leftarrow \epsilon$ -greedy( $Q$ )
- Just generates a single episode in each control (outer) loop iteration
- The statistics collected in the policy evaluation step will be for changing (improving) policies.
- Converges to  $q_*(s, a)$

## Sarsa: On-Policy Temporal-Difference Learning

- Pros of TD over MC:
  - Lower variance
  - Online
  - Incomplete sequences
- Idea: Use TD instead of MC in control loop
- **SARSA:**
  - $Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma Q(S', A') - Q(S, A))$ 
    - $A$  and  $A'$  are both chosen using the *same* policy derived from the  $Q(S, A)$ 
      - eg.  $\epsilon$ -greedy
  - Algorithm:
    - **Every time-step:**
      - *Policy evaluation:* Sarsa,  $Q \approx q_\pi$

- Policy improvement:  $\epsilon$ -greedy policy improvement
  - Converges to  $q_*(s, a)$
- **n-Step Sarsa**
  - n-step Q-return:  $q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$
  - n-step Sarsa updates  $Q(s, a)$  towards the  $n$ -step Q-return
  - $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^{(n)} - Q(S_t, A_t))$ 
    - $n = 1$ : Sarsa
    - $n = \infty$ : MC
- **Sarsa( $\lambda$ )**
  - $q_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$
  - **Forward-view Sarsa( $\lambda$ )**:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^\lambda - Q(S_t, A_t))$
  - **Backward-view Sarsa( $\lambda$ )**:  $Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t E_t(s, a)$ 
    - $\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$
    - Use an Eligibility Trace for each state-action pair
      - $E_0(s, a) = 0$
      - $E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + \mathbf{1}(S_t = s, A_t = a)$
    - $Q(s, a)$  is updated **for every state  $s$  and action  $a$**  in proportion to TD-error  $\delta_t$  and eligibility trace  $E_t(s, a)$

## Off-Policy Learning

- *Idea:*
  - Evaluate target policy  $\pi(a | s)$  to compute  $v_\pi(s)$  or  $q_\pi(s, a)$
  - While following behaviour policy  $\mu(a | s)$ 
    - $\{S_1, A_1, R_2, \dots, S_T\} \sim \mu$
- *Advantages:*
  - Learn from observing humans or other agents
  - Re-use experience generated from old policies
  - Learn about optimal policy while following exploratory policy
  - Learn about multiple policies while following one policy
- **Importance Sampling**:  $E_{x \sim p(x)}[f(x)] = \int p(x)f(x)dx = E_{x \sim q(x)} \left[ \frac{p(x)}{q(x)} f(x) \right]$
- **Importance Sampling for Off-Policy Monte-Carlo**
  - $G_t^{\pi/\mu} = \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} \frac{\pi(A_{t+1} | S_{t+1})}{\mu(A_{t+1} | S_{t+1})} \dots \frac{\pi(A_T | S_T)}{\mu(A_T | S_T)} G_t$ 
    - Multiply importance sampling corrections along whole episode
  - $V(S_t) \leftarrow V(S_t) + \alpha (G_t^{\pi/\mu} - V(S_t))$
  - Cannot use if  $\mu$  is zero when  $\pi$  is non-zero
  - *Problem*: Importance sampling can dramatically increase variance => **Cannot off-policy Monte-Carlo in practice**
- **Importance Sampling for Off-Policy TD**
  - $V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$ 
    - Only need a single importance sampling corrections

## Q-Learning

- **No importance sampling**
- *Idea*:
  - Next action is chosen using behaviour policy  $A_{t+1} \sim \mu(\cdot | S_t)$
  - But successor action taken from target policy  $A' \sim \pi(\cdot | S_t)$
  - $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t))$
- Q-learning allows both behaviour and target policies to improve
- The target policy  $\pi$  is greedy w.r.t.  $Q(s, a)$ :  $\pi(S_{t+1}) = \operatorname{argmax}_{a'} Q(S_{t+1}, a')$
- The behaviour policy  $\mu$  can be  $\epsilon$ -greedy w.r.t.  $Q(s, a)$
- **Q-learning target**:  $R_{t+1} + \gamma Q(S_{t+1}, A') = R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a')$
- $Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma \max_{a'} Q(S', a') - Q(S, A))$

## Summary

Table 1: Relationship Between DP and TD

|  | Full Backup (DP)   | Sample Backup (TD)   |
|--|--|--|
| Bellman Expectation Equation for $v_\pi(s)$    | <p>A tree diagram representing the Bellman expectation equation for <math>v_\pi(s)</math>. The root node is a white circle labeled <math>v_\pi(s) \leftrightarrow s</math>. It has two children, both black circles labeled <math>a</math>. Each child has two children, both white circles labeled <math>r</math>. Each <math>r</math> node has two children, both white circles labeled <math>v_\pi(s') \leftrightarrow s'</math>.</p> <p><b>Iterative Policy Evaluation</b><br/> <math display="block">V(s) \leftarrow \mathbb{E}[R + \gamma V(S')   s]</math></p>                      | <p>A vertical chain of nodes. The top node is a white circle. It has a black circle labeled <math>\bullet</math> below it, which in turn has a white circle below it.</p> <p><b>TD Learning</b><br/> <math display="block">V(S) \xleftarrow{\alpha} R + \gamma V(S')</math></p>  |
| Bellman Expectation Equation for $q_\pi(s, a)$ | <p>A tree diagram representing the Bellman expectation equation for <math>q_\pi(s, a)</math>. The root node is a black circle labeled <math>q_\pi(s, a) \leftrightarrow s, a</math>. It has two children, both white circles labeled <math>r</math>. Each <math>r</math> node has two children, both black circles labeled <math>s'</math>. Each <math>s'</math> node has two children, both black circles labeled <math>a'</math>.</p> <p><b>Q-Policy Iteration</b><br/> <math display="block">Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A')   s, a]</math></p>                      | <p>A vertical chain of nodes. The top node is a black circle labeled <math>\bullet</math>, with a label "S,A" above it. It has a white circle labeled "R" below it, which in turn has a white circle labeled "S'" below it, and finally a black circle labeled "A'" at the bottom.</p> <p><b>Sarsa</b><br/> <math display="block">Q(S, A) \xleftarrow{\alpha} R + \gamma Q(S', A')</math></p>  |
| Bellman Expectation Equation for $q_*(s, a)$   | <p>A tree diagram representing the Bellman expectation equation for <math>q_*(s, a)</math>. The root node is a black circle labeled <math>q_*(s, a) \leftrightarrow s, a</math>. It has two children, both white circles labeled <math>r</math>. Each <math>r</math> node has two children, both white circles labeled <math>s'</math>. Each <math>s'</math> node has two children, both black circles labeled <math>a'</math>.</p> <p><b>Q-Value Iteration</b><br/> <math display="block">Q(s, a) \leftarrow \mathbb{E}[R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')   s, a]</math></p> | <p>A vertical chain of nodes. The top node is a black circle labeled <math>\bullet</math>, with a label "S,A" above it. It has a white circle labeled "R" below it, which in turn has a white circle labeled "S'" below it. Below "S'" are three black circles labeled "A'" at the bottom.</p> <p><b>Q-Learning</b><br/> <math display="block">Q(S, A) \xleftarrow{\alpha} R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')</math></p> |

where  $x \xleftarrow{\alpha} y \equiv x \leftarrow x + \alpha(y - x)$

## Part II: Reinforcement Learning in Practice

### Lecture 6: Value Function Approximation

#### Introduction

- So far we have represented value function by a *lookup table*
- *Problem:* with large MDPs,
  - Too many states and/or actions to store in memory
  - Too slow to learn the value of each state individually
- Solution for large MDPs:
  - Estimate value function with function approximation
    - $\hat{v}(s; w) \approx v_\pi(s)$  or  $\hat{q}(s, a; w) \approx q_\pi(s, a)$
  - Generalise from seen states to unseen states
  - Update parameter  $w$  using MC or TD learning
- Will consider **differentiable** function approximators,
  - Linear Combination of Features
  - Neural Network
- The training method should be suitable for data that is
  - **non-stationary**- if estimating  $v_\pi$ , the policy  $\pi$  will keep changing
  - **non-iid**- because it's not supervised learning, data arrives in a trajectory

#### Incremental Methods

- Let  $J(w)$  be the differentiable function of parameter vector  $w$ .

#### Incremental Prediction Algorithms

- Let  $x(S) = [x_1(S), \dots, x_n(S)]^T$  be the state feature vector

- **Linear Value Function Approximation:**

- $\hat{v}(S, \mathbf{w}) = \mathbf{x}(S)^\top \mathbf{w} = \sum_{j=1}^n \mathbf{x}_j(S) \mathbf{w}_j$
- $J(\mathbf{w}) = \mathbb{E}_\pi \left[ (v_\pi(S) - \mathbf{x}(S)^\top \mathbf{w})^2 \right]$ 
  - Quadratic in parameters  $\mathbf{w}$  => SGD converges to global optimum
- Update rule:
  - $\nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w}) = \mathbf{x}(S)$
  - $\Delta \mathbf{w} = \alpha (v_\pi(S) - \hat{v}(S, \mathbf{w})) \mathbf{x}(S)$
- Update Complexity = step-size  $\times$  prediction error  $\times$  feature value

- **Incremental Prediction Algorithms**

- Train  $\hat{v}(S, \mathbf{w}) \approx v_\pi(S)$
- supervised learning on "training data":  $\langle S_1, Y_1 \rangle, \langle S_2, Y_2 \rangle, \dots, \langle S_T, Y_T \rangle$ 
  - where  $Y_t$  is an sample/estimate of true value  $v_\pi(S_t)$
  - different algorithms use different estimates:
    - MC: Uses return  $G_t$  as an unbiased, noisy sample
    - TD(0): Uses TD-target  $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$  for a biased sample
    - TD( $\lambda$ ): Uses  $\lambda$ -return  $G_t^\lambda$  for a biased sample
- **MC:**  $\Delta \mathbf{w} = \alpha (G_t - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$ 
  - Monte-Carlo evaluation converges to a local optimum
  - Even with non-linear value function approximation
- **TD(0):**  $\Delta \mathbf{w} = \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$ 
  - Converges (close) to global optimum for linear value function approximation
- **Forward-view TD( $\lambda$ ):**  $\Delta \mathbf{w} = \alpha (G_t^\lambda - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
- **Backward-view TD( $\lambda$ ):**  $\Delta \mathbf{w} = \alpha \delta_t E_t$ 
  - $\delta_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})$
  - $E_t = \gamma \lambda E_{t-1} + \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
  - This update can be applied online, to incomplete sequences.
- Forward view and backward view linear TD( $\lambda$ ) are equivalent

## Incremental Control Algorithms

- Let  $x(S, A) = [x_1(S, A), \dots, x_n(S, A)]^T$  be the state-action feature vector

- **Linear Action-Value Function Approximation:**

- $\hat{q}(S, A, \mathbf{w}) = \mathbf{x}(S, A)^\top \mathbf{w} = \sum_{j=1}^n \mathbf{x}_j(S, A) \mathbf{w}_j$
- Update rule:
  - $\nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w}) = \mathbf{x}(S, A)$
  - $\Delta \mathbf{w} = \alpha (q_\pi(S, A) - \hat{q}(S, A, \mathbf{w})) \mathbf{x}(S, A)$
- **Incremental Control Algorithms**
- different algorithms use different estimates for the true value  $q_\pi(s, a)$ :
  - MC: Uses return  $G_t$  as an unbiased, noisy sample
  - TD(0): Uses TD-target  $R_{t+1} + \gamma \hat{q}(S_{t+1}, \mathbf{w})$  for a biased sample
  - TD( $\lambda$ ): Uses  $\lambda$ -return  $G_t^\lambda$  for a biased sample
- **MC:**  $\Delta \mathbf{w} = \alpha (G_t - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
- **TD(0):**  $\Delta \mathbf{w} = \alpha (R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
- **Forward-view TD( $\lambda$ ):**  $\Delta \mathbf{w} = \alpha (q_t^\lambda - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
- **Backward-view TD( $\lambda$ ):**  $\Delta \mathbf{w} = \alpha \delta_t E_t$ 
  - $\delta_t = R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})$
  - $E_t = \gamma \lambda E_{t-1} + \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$
  - This update can be applied online, to incomplete sequences.

- **Convergence of Prediction Algorithms**

-

| On/Off-Policy | Algorithm       | Table Lookup | Linear | Non-Linear |
|---------------|-----------------|--------------|--------|------------|
| On-Policy     | MC              | ✓            | ✓      | ✓          |
|               | TD(0)           | ✓            | ✓      | ✗          |
|               | TD( $\lambda$ ) | ✓            | ✓      | ✗          |
| Off-Policy    | MC              | ✓            | ✓      | ✓          |
|               | TD(0)           | ✓            | ✗      | ✗          |
|               | TD( $\lambda$ ) | ✓            | ✗      | ✗          |

- Gradient TD Learning

- TD diverges when off-policy or using non-linear function approximation because it doesn't follow the gradient of any objective function.
- Gradient TD follows true gradient of projected Bellman error and fixes convergence issues.

- Convergence of Control Algorithms

| Algorithm           | Table Lookup | Linear | Non-Linear |
|---------------------|--------------|--------|------------|
| Monte-Carlo Control | ✓            | (✓)    | ✗          |
| Sarsa               | ✓            | (✓)    | ✗          |
| Q-learning          | ✓            | ✗      | ✗          |
| Gradient Q-learning | ✓            | ✓      | ✗          |

- (✓) - chatters around near-optimal value function

## Batch Methods

- SGD with Experience Replay

- Given experience consisting of  $\langle \text{state}, \text{value} \rangle$  pairs
- $\mathcal{D} = \{\langle s_1, v_1^\pi \rangle, \langle s_2, v_2^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle\}$
- Repeat:
  - Sample state, value from experience  $\langle s, v^\pi \rangle \sim \mathcal{D}$
  - Apply SGD update:  $\Delta \mathbf{w} = \alpha (v^\pi - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$
- Converges to least squares solution:  $\mathbf{w}^\pi = \underset{\mathbf{w}}{\operatorname{argmin}} LS(\mathbf{w})$

- DQN

- Two ideas prevent DQN from blowing up unlike Sarsa and Q-Learning:
  - experience replay** - eliminates correlated samples
  - fixed Q-targets** - eliminates the moving target
- Algorithm:
  - Take action  $a_t$  according to  $\epsilon$ -greedy policy
  - Store transition  $(s_t, a_t, r_{t+1}, s_{t+1})$  in replay memory  $\mathcal{D}$
  - Sample random mini-batch of transitions  $(s, a, r, s')$  from  $\mathcal{D}$
  - Compute Q-learning targets w.r.t. old, fixed parameters  $w^-$
  - Optimise MSE between Q-network and Q-learning targets
    - $\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} [(r + \gamma \max_{a'} Q(s', a'; w^-) - Q(s, a; w_i))^2]$

## Lecture 7: Policy Gradient Methods

### Introduction

- Policy Based RL

- Parametrise the policy directly  $\pi_\theta(s, a) = P(a | s, \theta)$
- still **model-free**
- Pros:
  - Better convergence** properties
  - Effective in **high-dimensional or continuous action spaces** as no need to do a *max* over actions
  - Can learn **stochastic policies**

- If the environment is partially observed then it is possible that only stochastic policies are optimal and all deterministic policies are suboptimal.
  - Example: when state aliasing occurs
  - Example: if function approximation is used and the features used limit the view of the world
- Cons:
  - Typically converge to a **local optimum**
  - Evaluating a policy is typically **inefficient and high variance**
    - Resolved using value function in Actor-Critic methods.
- **Policy Objective Functions**

- For *episodic environments*
  - use the *start value*:  $J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}_{\pi_\theta}[v_1]$
- For *\_continuing environments*
  - use the *average value*:  $J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s)V^{\pi_\theta}(s)$
  - Or the *average reward per time-step*:  $J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a)\mathcal{R}_s^a$
- where  $d^{\pi_\theta}(s)$  is stationary distribution of Markov chain for  $\pi_\theta$
- Policy gradient is essentially the same for all three

## Monte-Carlo Policy Gradient

- **Likelihood-ratio trick**:  $\nabla_\theta \pi_\theta(s, a) = \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)$
- **Score Function**:  $\nabla_\theta \log \pi_\theta(s, a)$ 
  - Example: Softmax Policy:  $\pi_\theta(s, a) \propto e^{\phi(s, a)^T \theta}$ 
    - $\nabla_\theta \log \pi_\theta(s, a) = \phi(s, a) - \mathbb{E}_{\pi_\theta}[\phi(s, \cdot)]$
  - Example: Gaussian Policy:  $a \sim \mathcal{N}(\mu(s), \sigma^2)$  where  $\mu(s) = \phi(s)^T \theta$
  - $\nabla_\theta \log \pi_\theta(s, a) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$
- **Policy Gradient for One-step MDPs**
  - Contextual Bandits: state,  $s \sim d(s)$ , reward,  $r = \mathcal{R}_{s,a}$
  - $J(\theta) = \mathbb{E}_{\pi_\theta}[r] = \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s, a) \mathcal{R}_{s,a}$ 
    - All three objective functions will be equivalent
  - $\nabla_\theta J(\theta) = \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a) \mathcal{R}_{s,a} = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) r]$
- **Policy Gradient Theorem**
  - For any differentiable policy  $\pi_\theta(s, a)$ , for any of the policy objective functions  $J = J_1, J_{avR}$ , or  $\frac{1}{1-\gamma} J_{avV}$  the policy gradient is
    - $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)]$
- **REINFORCE: Monte-Carlo Policy Gradient**
  - $\nabla_\theta t = \alpha \nabla_\theta \log \pi_\theta(s, a) v_t$  where the return,  $v_t$ , is an unbiased estimator of  $Q^{\pi_\theta}(s, a)$

## Actor-Critic Policy Gradient

- Reduce Variance using Critic:  $Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$
- Actor-critic algorithms maintain two sets of parameters
  - **Critic**: Updates action-value function parameters  $w$
  - **Actor**: Updates policy parameters  $\theta$ , in direction suggested by critic
- Actor-critic algorithms follow an approximate policy gradient
  - $\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$
  - $\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$
- Critic is doing policy evaluation. Eg. MC PE, TD learning,  $TD(\lambda)$
- **Compatible Function Approximation Theorem**
  - If the following two conditions are satisfied:
    - Value function approximator is **compatible** to the policy:  $\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$
    - i.e. features for the value function approximator are the score function itself.
      - Eg. if  $Q_w(s, a) = \phi(s, a)^T w$  then,  $\phi(s, a) = \nabla_\theta \log \pi_\theta(s, a)$
    - Value function parameters  $w$  minimise the mean-squared error:  $\varepsilon = \mathbb{E}_{\pi_\theta}[(Q^{\pi_\theta}(s, a) - Q_w(s, a))^2]$
  - Then the policy gradient is exact:  $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$
- **Reducing Variance Using a Baseline**

- Subtract a baseline function  $B(s)$  from the policy gradient
  - Reduce variance, without changing expectation
    - $\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) B(s)] = \sum_{s \in \mathcal{S}} d^{\pi_\theta}(s) \sum_a \nabla_\theta \pi_\theta(s, a) B(s) = 0$
- Advantage function**  $A^{\pi_\theta}(s, a) = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s)$ 
  - Using  $B(s) = V^{\pi_\theta}(s)$
  - $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^{\pi_\theta}(s, a)]$
- Estimating Advantage function**
  - Option 1: Use two function approximators:  $V_v(s) \approx V^{\pi_\theta}(s)$ ,  $Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$
  - Option 2: Use a single function approximator:  $V_v(s) \approx V^{\pi_\theta}(s)$ 
    - TD error for the true value function  $\delta^{\pi_\theta} = r + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)$ , is an unbiased estimate of the advantage function.
      - \*  $\mathbb{E}_{\pi_\theta} [\delta^{\pi_\theta} | s, a] = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) = A^{\pi_\theta}(s, a)$
    - approximate TD error:  $\delta_v = r + \gamma V_v(s') - V_v(s)$
    - $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \delta_v]$
- Critics at Different Time-Scales**
  - MC:**  $\Delta w = \alpha (v_t - V_w(s)) \nabla_w V_w(s)$
  - TD(0):**  $\Delta w = \alpha (r + \gamma V(s') - V_w(s)) \nabla_w V_w(s)$
  - Forward-view TD( $\lambda$ ):**  $\Delta w = \alpha (v_t^\lambda - V_w(s)) \nabla_w V_w(s)$
  - Backward-view TD( $\lambda$ ):**  $\Delta w = \alpha \delta_t e_t$ 
    - $\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$
    - $e_t = \gamma \lambda e_{t-1} + \nabla_w V_w(s)$
    - This update can be applied online, to incomplete sequences.
  - Where  $\nabla_w V_w(s) = \phi(s)$  for linear  $V_w = w^T \phi(s)$ .
- Actors at Different Time-Scales**
  - Here  $\theta$  are the parameters for policy  $\pi_\theta$  and  $V_w$  is the function approximator for  $V^{\pi_\theta}$  which in turn is used to estimate  $A^{\pi_\theta}$ .
  - MC:**  $\Delta \theta = \alpha (v_t - V_w(s)) \nabla_\theta \log \pi_\theta(s_t, a_t)$
  - TD(0):**  $\Delta \theta = \alpha (r + \gamma V_w(s_{t+1}) - V_w(s_t)) \nabla_\theta \log \pi_\theta(s_t, a_t)$
  - Forward-view TD( $\lambda$ ):**  $\Delta \theta = \alpha (v_t^\lambda - V_w(s_t)) \nabla_\theta \log \pi_\theta(s_t, a_t)$
  - Backward-view TD( $\lambda$ ):**  $\Delta \theta = \alpha \delta_t e_t$ 
    - $\delta_t = r_{t+1} + \gamma V_w(s_{t+1}) - V_w(s_t)$
    - $e_t = \gamma \lambda e_{t-1} + \nabla_\theta \log \pi_\theta(s_t, a_t)$
    - This update can be applied online, to incomplete sequences unlike MC or forward-view TD( $\lambda$ ).

## Natural Policy Gradient

- A policy can often be reparametrised without changing action probabilities
  - Eg. increasing score of all actions in a softmax policy
- The vanilla gradient is sensitive to these reparametrisations
- Natural Policy Gradient**
  - parametrisation independent
  - It finds ascent direction that is closest to vanilla gradient, when changing policy by a small, fixed amount
    - $\nabla_\theta^{nat} \pi_\theta(s, a) = G_\theta^{-1} \nabla_\theta \pi_\theta(s, a)$
  - where  $G_\theta$  is the **Fisher information matrix**
    - $G_\theta = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)^T]$
  - With a compatible function approximator,  $A_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)^T w$ ,
    - $\nabla_\theta J(\theta) = G_\theta w$
  - Natural Policy Gradient:  $\nabla_\theta^{nat} J(\theta) = w$ 
    - i.e. update actor parameters in direction of critic parameters

## Summary

$$\begin{aligned}
\nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \textcolor{red}{v_t}] && \text{REINFORCE} \\
&= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \textcolor{red}{Q^w(s, a)]} && \text{Q Actor-Critic} \\
&= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \textcolor{red}{A^w(s, a)]} && \text{Advantage Actor-Critic} \\
&= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] && \text{TD Actor-Critic} \\
&= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta e] && \text{TD}(\lambda) \text{ Actor-Critic} \\
G_{\theta}^{-1} \nabla_{\theta} J(\theta) &= w && \text{Natural Actor-Critic}
\end{aligned}$$

## Lecture 8: Integrating Learning and Planning

---

### Introduction

- **Model-free RL**
  - Value-based Methods: learn value function directly from experience
  - Policy-based Methods: learn policy directly from experience
- **Model-Based RL:** learn model directly from experience
  - and use planning to construct a value function or policy

### Model-Based Reinforcement Learning

- **Pros:**
  - Easier to estimate model directly from experience than value fn/policy
  - Can efficiently learn model by supervised learning methods
  - Can reason about model uncertainty
- **Cons:**
  - First learn a model, then construct a value function => two sources of approximation error
- **Model**
  - A model  $\mathcal{M}$  is a representation of an MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$  parametrized by  $\eta$ 
    - Assume state space  $\mathcal{S}$  and action space  $\mathcal{A}$  are known
    - So  $\mathcal{M} = \langle \mathcal{P}_\eta, \mathcal{R}_\eta \rangle$  where  $\mathcal{P}_\eta \approx \mathcal{P}$  and  $\mathcal{R}_\eta \approx \mathcal{R}$
  - Typically assume conditional independence between state transitions and rewards
    - $\mathbb{P}[S_{t+1}, R_{t+1} | S_t, A_t] = \mathbb{P}[S_{t+1} | S_t, A_t] \mathbb{P}[R_{t+1} | S_t, A_t]$
- **Model Learning**
  - Goal: estimate model  $\mathcal{M}_\eta$  from experience  $\{S_1, A_1, R_2, \dots, S_T\}$
  - This is a supervised learning problem  $\{S_t, A_t \rightarrow R_{t+1}, S_{t+1}\}_{t=1}^{T-1}$ 
    - Learning  $s, a \rightarrow r$  is a regression problem
    - Learning  $s, a \rightarrow s'$  is a density estimation problem
  - Find parameters  $\eta$  that minimise empirical loss
- Examples of Models:
  - Table Lookup Model
  - Linear Expectation Model
  - Linear Gaussian Model
  - Gaussian Process Model
  - Deep Belief Network Model
- **Table Lookup Model**
  - Parametric:
    - Count visits  $N(s, a)$  to each state action pair
    - $\hat{\mathcal{P}}_{s,s'}^a = \frac{1}{N(s,a)} \sum_{t=1}^T \mathbf{1}(S_t, A_t, S_{t+1} = s, a, s')$
    - $\hat{\mathcal{R}}_s^a = \frac{1}{N(s,a)} \sum_{t=1}^T \mathbf{1}(S_t, A_t = s, a) R_t$
  - Non-parametric:
    - At each time-step t, record experience tuple  $\langle S_t, A_t, R_{t+1}, S_{t+1} \rangle$
    - To sample model, randomly pick tuple matching  $\langle s, a, \dots \rangle$
  - **Planning with Model**
    - Given a model  $\mathcal{M}_\eta = \langle \mathcal{P}_\eta, \mathcal{R}_\eta \rangle$
    - Solve the MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}_\eta, \mathcal{R}_\eta \rangle$ 
      - Algorithms: Value iteration, Policy iteration, Tree search etc.
  - **Sample-based Planning**
    - Use model to **only** to generate samples (*simulated experience*)

- Apply model-free RL to samples, e.g.: Monte-Carlo control, Sarsa, Q-learning
- Sample-based planning methods are often more efficient

## Integrated Architectures

- Experience
  - Real experience Sampled from environment (true MDP)
    - $S' \sim \mathcal{P}_{s,s'}^a$
    - $R = \mathcal{R}_s^a$
  - Simulated experience Sampled from model (approximate MDP)
    - $S' \sim \mathcal{P}_\eta(S' | S, A)$
    - $R = \mathcal{R}_\eta(R | S, A)$
- Model-Free RL
  - No model
  - **Learn** value function (and/or policy) from real experience
- Model-Based RL (using Sample-Based Planning)
  - Learn a model from real experience
  - **Plan** value function (and/or policy) from simulated experience
- Dyna
  - Learn a model from real experience
  - **Learn and plan** value

## Simulation-Based Search

### Lecture 9: Exploration and Exploitation

---

### Lecture 10: Case study - RL in games

---