# Lecture 4: Model-free RL (Part I)

Max(Chong) Li

# Outline

- Monte Carlo Learning

- Monte Carlo Control

- Extensions: on- and off-policy evaluation

# Outline

- <span style="color:red">Monte Carlo Learning</span>

- Monte Carlo Control

- Extensions: on- and off-policy evaluation

# Monte Carlo RL

- MC methods learn directly from episodes of experience
- MC is *model-free*: no knowledge of MDP transitions / rewards
- MC learns from *complete* episodes: no bootstrapping
- MC uses the simplest possible idea: value = mean return
- Caveat: can only apply MC to *episodic* MDPs
  - All episodes must terminate

# MC Policy Evaluation

- Goal: learn $v_\pi$ from episodes of experience under policy $\pi$

$$S_1, A_1, R_2, ..., S_k \sim \pi$$

- Recall that the *return* is the total discounted reward:

$$G_t = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T$$

- Recall that the value function is the expected return:

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s]$$

- Monte-Carlo policy evaluation uses *empirical mean* return instead of *expected* return

# First Visit MC Policy Evaluation

- To evaluate state $s$
- The first time-step $t$ that state $s$ is visited in an episode,
- Increment counter $N(s) \leftarrow N(s) + 1$
- Increment total return $S(s) \leftarrow S(s) + G_t$
- Value is estimated by mean return $V(s) = S(s)/N(s)$
  By law of large numbers, $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

# Every Visit MC Policy Evaluation

- To evaluate state $s$
- Every time-step $t$ that state $s$ is visited in an episode,
- Increment counter $N(s) \leftarrow N(s) + 1$
- Increment total return $S(s) \leftarrow S(s) + G_t$
- Value is estimated by mean return $V(s) = S(s)/N(s)$
  Again, $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

# Recursive Mean Computation

The mean $\mu_1, \mu_2, \ldots$ of a sequence $x_1, x_2, \ldots$ can be computed incrementally,

$$
\begin{aligned}
\mu_k &= \frac{1}{k} \sum_{j=1}^{k} x_j \\
&= \frac{1}{k} \left( x_k + \sum_{j=1}^{k-1} x_j \right) \\
&= \frac{1}{k} \left( x_k + (k-1)\mu_{k-1} \right) \\
&= \mu_{k-1} + \frac{1}{k} \left( x_k - \mu_{k-1} \right)
\end{aligned}
$$

# Recursive MC Updates

- Update $V(s)$ incrementally after episode $S_1, A_1, R_2, ..., S_T$
  For each state $S_t$ with return $G_t$

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

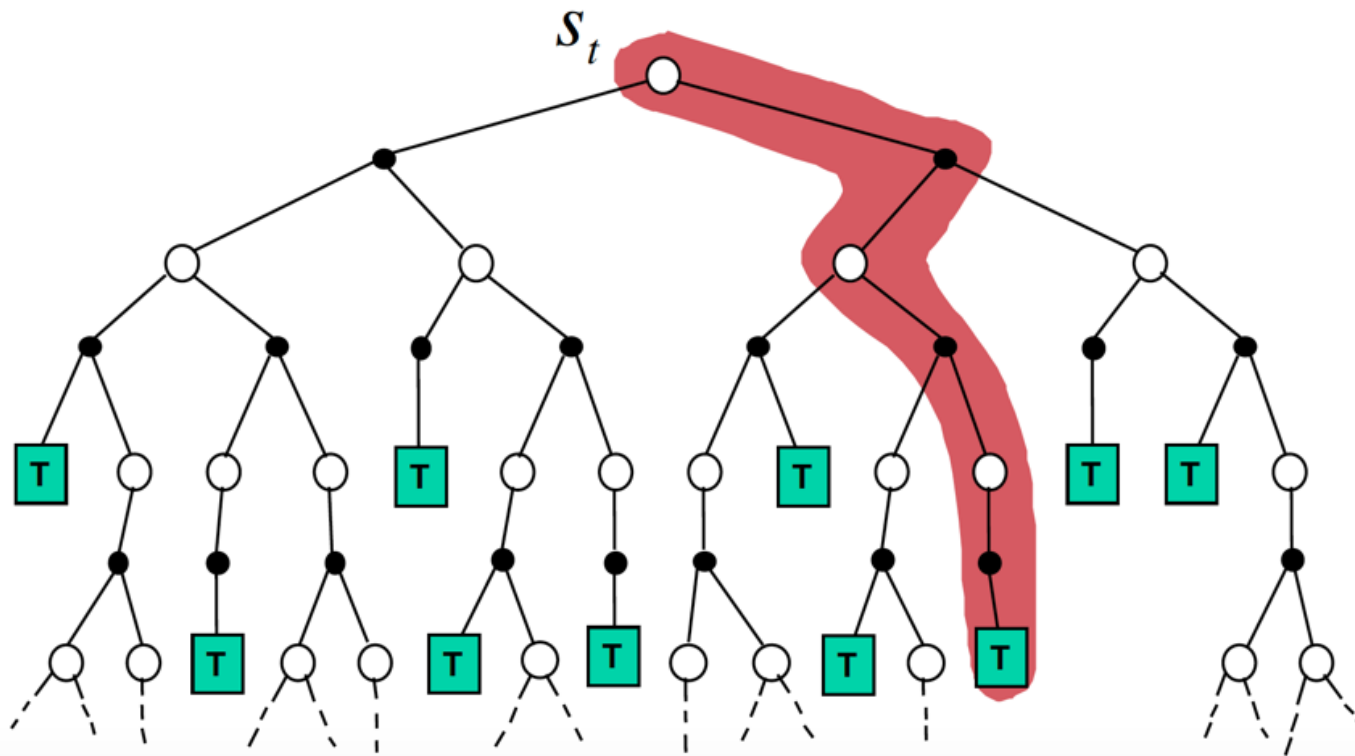$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$$

- In non-stationary problems, it can be useful to track a running mean, i.e. forget old episodes.

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

$$v_{t+1}(s) = (1 - \alpha)^t v_1(s) + \sum_{i=1}^{t} \alpha(1 - \alpha)^{t-i} G_i,$$
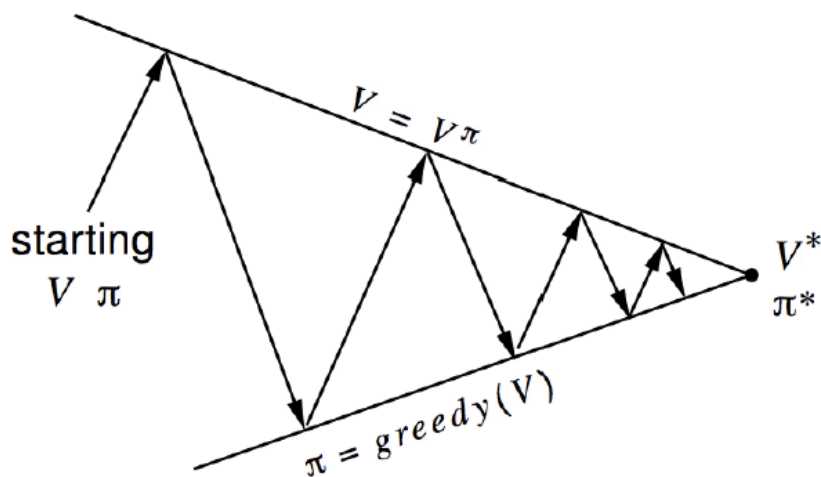
# MC Backup

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t - V(S_t) \right)$$
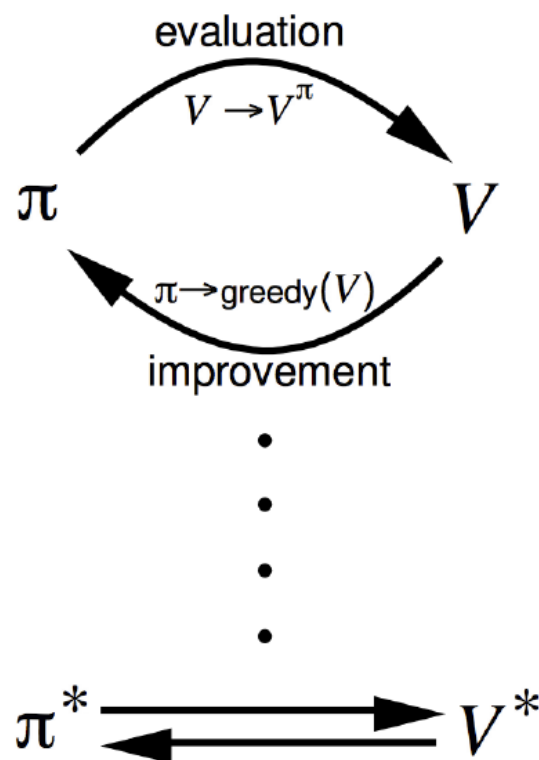
# Outline

- Monte Carlo Learning

- <span style="color:red">Monte Carlo Control</span>

- Extensions: on- and off-policy evaluation
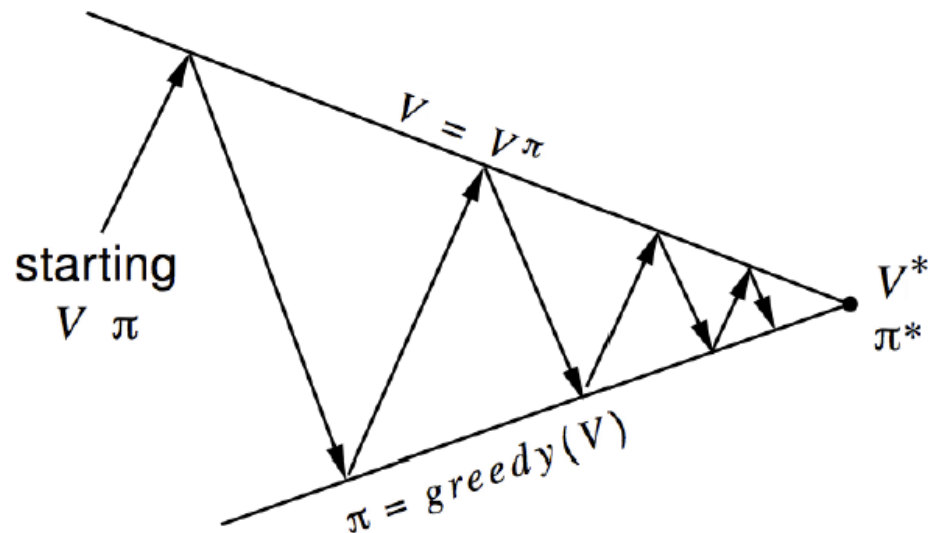
# Generalized Policy Iteration (GPI) Revisited



Policy evaluation  Estimate $v_\pi$
   e.g. Iterative policy evaluation

Policy improvement  Generate $\pi' \geq \pi$
   e.g. Greedy policy improvement

# GPI with MC Evaluation



Policy evaluation Monte-Carlo policy evaluation, $V = v_\pi$?

Policy improvement Greedy policy improvement?

# Model Free GPI with MC Evaluation

- Greedy policy improvement over $V(s)$ requires model of MDP

$$\pi'(s) = \underset{a \in \mathcal{A}}{\text{argmax}} \; \mathcal{R}_s^a + \mathcal{P}_{ss'}^a V(s')$$

- Greedy policy improvement over $Q(s, a)$ is model-free

$$\pi'(s) = \underset{a \in \mathcal{A}}{\text{argmax}} \; Q(s, a)$$

# Ɛ – Greedy Policy

- Simplest idea for ensuring continual exploration
- All $m$ actions are tried with non-zero probability
- With probability $1 - \epsilon$ choose the greedy action
- With probability $\epsilon$ choose an action at random

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{if } a^* = \underset{a \in \mathcal{A}}{\text{argmax }} Q(s, a) \\ \epsilon/m & \text{otherwise} \end{cases}$$

# Ɛ- Greedy Policy Improvement

Theorem: For any Ɛ-greedy policy $\pi$, the Ɛ-greedy policy $\pi'$ with respect to $q_\pi$ is an improvement, that is, $v_{\pi'}(s) \geq v_\pi(s)$

# Proof

$$v_\pi(s) = \sum_{a\in\mathcal{A}} \pi(a|s)q_\pi(s,a)$$

- $q_\pi(s,\pi'(s)) = \sum_{a\in\mathcal{A}} \pi'(a|s)q_\pi(s,a)$

$$\boxed{m = |\mathcal{A}(s)|}$$

$$= \epsilon/m \sum_{a\in\mathcal{A}} q_\pi(s,a) + (1-\epsilon)\max_{a\in\mathcal{A}} q_\pi(s,a)$$

$$\geq \epsilon/m \sum_{a\in\mathcal{A}} q_\pi(s,a) + (1-\epsilon)\sum_{a\in\mathcal{A}} \frac{\pi(a|s) - \epsilon/m}{1-\epsilon} q_\pi(s,a)$$
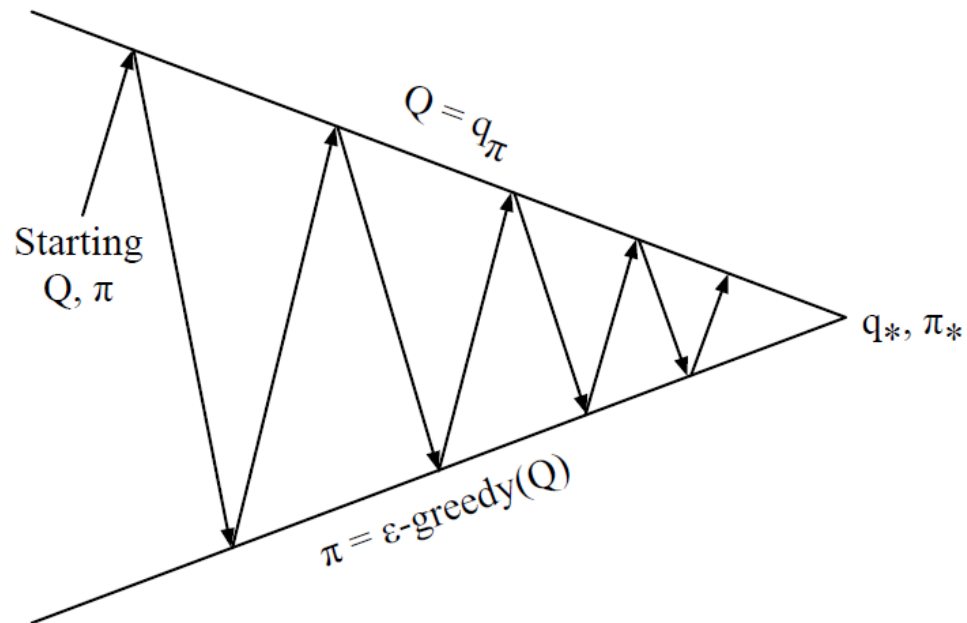
$$= \sum_{a\in\mathcal{A}} \pi(a|s)q_\pi(s,a) = v_\pi(s)$$

- Policy improvement theorem (Sutton, Section 4.2),

$$v_{\pi'}(s) \geq v_\pi(s)$$

# MC Policy Iteration



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

Policy improvement  $\epsilon$-greedy policy improvement

# GLIE Condition for Policies

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

All state-action pairs are explored infinitely many times,

$$\lim_{k \to \infty} N_k(s, a) = \infty$$

The policy converges on a greedy policy,

$$\lim_{k \to \infty} \pi_k(a|s) = \mathbf{1}(a = \underset{a' \in \mathcal{A}}{\operatorname{argmax}} \, Q_k(s, a'))$$

For example, $\epsilon$-greedy is GLIE if $\epsilon$ reduces to zero at $\epsilon_k = \frac{1}{k}$

# GLIE MC Policy Iteration

- Sample $k$th episode using $\pi$: $\{S_1, A_1, R_2, ..., S_T\} \sim \pi$
- 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

$$\epsilon \leftarrow 1/k$$

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

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

# Outline

- Monte Carlo Learning

- Monte Carlo Control

- <span style="color:red">Extensions: on- and off-policy evaluation</span>

# On- and Off-Policy

- **On-policy** learning

  "Learn on the job"
  Learn about policy $\pi$ from experience sampled from $\pi$

- **Off-policy** learning

  "Look over someone's shoulder"
  Learn about policy $\pi$ from experience sampled from $\mu$

# Important Sampling

- Estimate the expectation of a different distribution

$$\mathbb{E}_{X \sim P}[f(X)] = \sum P(X) f(X)$$

$$= \sum Q(X) \frac{P(X)}{Q(X)} f(X)$$

$$= \mathbb{E}_{X \sim Q} \left[ \frac{P(X)}{Q(X)} f(X) \right]$$

# Importance Sampling in MC Policy Iteration

- Use returns generated from $\mu$ to evaluate $\pi$
- Weight return $G_t$ according to similarity between policies
- Multiply importance sampling corrections along whole episode

weight

$$G_t^{\pi/\mu} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t$$

- Update value towards *corrected* return

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

- Cannot use if $\mu$ is zero when $\pi$ is non-zero
- Importance sampling can dramatically increase variance

# Ordinary & Weighted Importance Sampling

- Ordinary importance sampling: when importance sampling is done as a simple average of returns; the estimator is unbiased but the variance could be in general unbounded.

- Weighted importance sampling: when importance sampling is done as a weighted average of returns ; the estimator is biased (the bias converges asymptotically to zero) but the variance is dramatically low.

- See Chapter 5.5 in Sutton's book (2nd Edition)

# Example: Blackjack State Value

- See Example 5.4 in Sutton's book (2nd Edition)