

# Monte-Carlo Learning

Mattia Pellegrino, Ph.D Fellow  
[mattia.pellegrino@unipr.it](mailto:mattia.pellegrino@unipr.it)



# Summary

- Monte-Carlo Learning (Model free prediction)
- Monte-Carlo Learning (Model free control)

# MONTÉ-CARLO RF

- Monte Carlo methods learn directly from episodes of experience
- Monte Carlo is classified as model-free (no knowledge of MDP transitions/rewards)
- MC learns from complete episodes
- MC uses the mean return as value  $\rightarrow$  value = mean return
- MC can only be applied to episodic MDPs
  - All episodes must terminate

# MC POLICY EVALUATION

- We want to learn  $V_\pi$  from episodes of experience under policy  $\pi$

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

- Recap: the return is the total discounted reward:

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

- Recap: value function is the expected return:

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

- Monte-Carlo policy evaluation used empirical mean return instead of expected return

# POLICY EVALUATION: FIRST ITERATION

- To evaluate state  $s$
- The FIRST time-step  $t$  that state  $s$  is visited in an episode:
  - Counter incrementation:  $N(s) \rightarrow N(s) + 1$
  - Total Return incrementation:  $S(s) \leftarrow S(s) + G_t$
  - Value is the mean return:  $V(s) = S(s) / N(s)$
  - By the law of large numbers  $\rightarrow V(s) \rightarrow V_{\pi}(s)$  as  $N(s) \rightarrow \infty$

# POLICY EVALUATION: EVERY ITERATION

- To evaluate state  $s$
- EVERY time-step  $t$  that state  $s$  is visited in an episode:
  - Counter incrementation:  $N(s) \rightarrow N(s) + 1$
  - Total Return incrementation:  $S(s) \leftarrow S(s) + G_t$
  - Value is the mean return:  $V(s) = S(s) / N(s)$
  - By the law of large numbers  $\rightarrow V(s) \rightarrow V_{\pi}(s)$  as  $N(s) \rightarrow \infty$

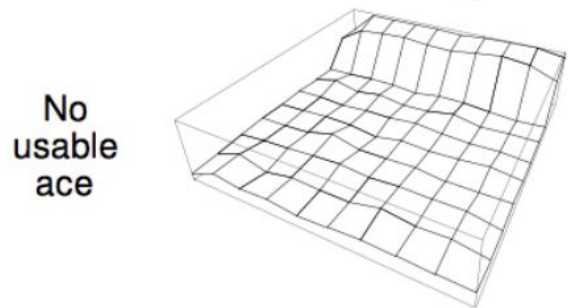
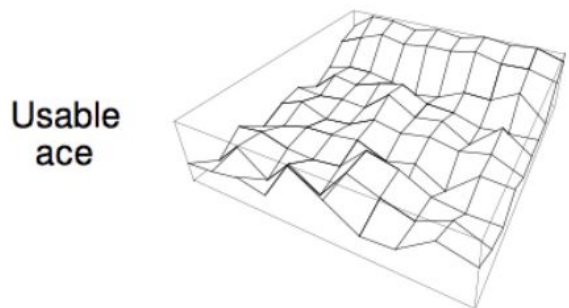
# EXAMPLE - BLACK JACK

- States (200 of them):
  - Current sum (12-21)
  - Dealer's showing card (ace-10)
  - Do I have a "useable" ace? (yes-no)
- Action stick: Stop receiving cards (and terminate)
- Action twist: Take another card (no replacement)
- Reward for stick:
  - +1 if sum of cards > sum of dealer cards
  - 0 if sum of cards = sum of dealer cards
  - -1 if sum of cards < sum of dealer cards
- Reward for twist:
  - -1 if sum of cards > 21 (and terminate)
  - 0 otherwise
- Transitions: automatically twist if sum of cards < 12

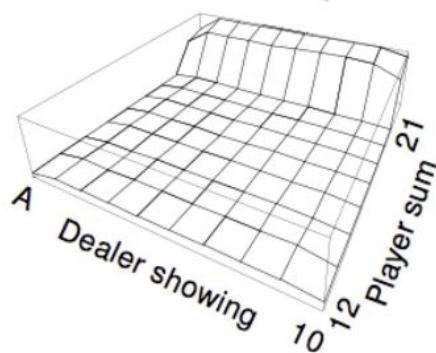
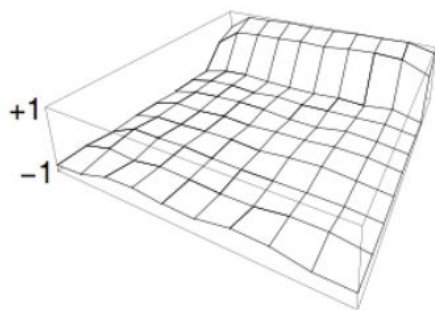


# EXAMPLE - BLACK JACK

After 10,000 episodes



After 500,000 episodes



- Policy:
  - Stick
    - Sum of cards  $\geq 20$
  - Twist
    - Otherwise



# INCREMENTAL MEAN

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

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

# INCREMENTAL UPDATES

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

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

$$V(S_t) \leftarrow V(S_t) + 1 / N(S_t) (G_t - V(S_t))$$

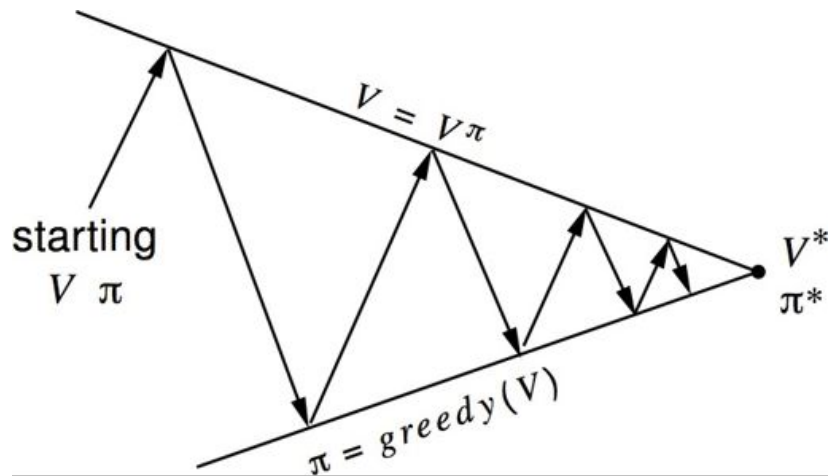
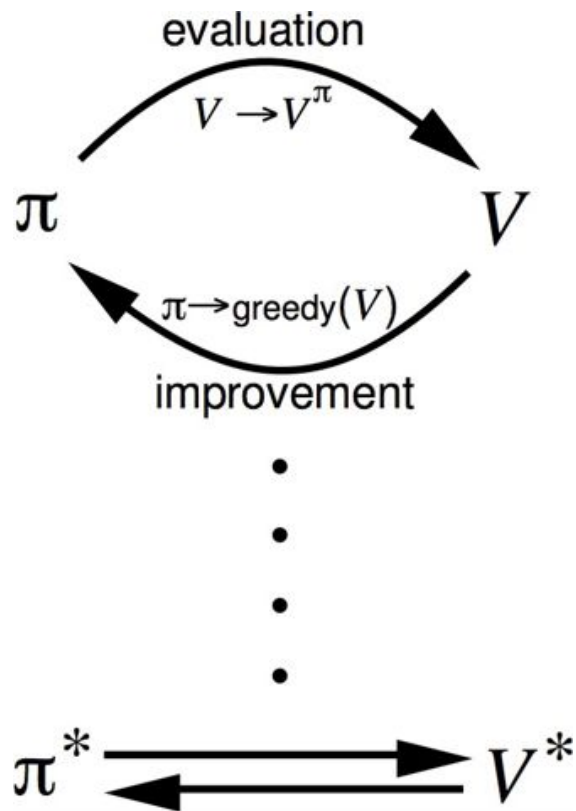
- In non-stationary problems, it can be useful to track a running mean, i.e. forget old episodes

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

# MONTECARLO – MODEL CONTROL

- Monte Carlo is an On-policy learning algorithm
  - “Learn on the job”
  - Learn about policy  $\pi$  from experience sampled from  $\pi$

# POLICY ITERATION (REFRESH)



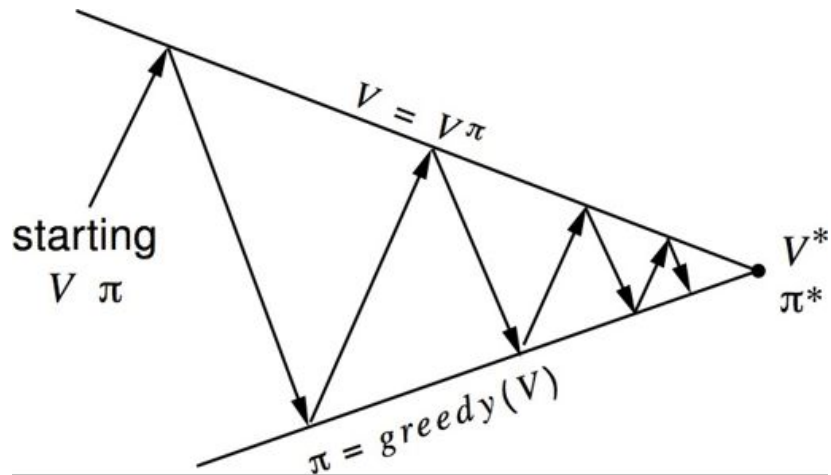
Policy evaluation

→ Estimate  $v_\pi$

Policy improvement

→ Generate  $\pi' \geq \pi$

# POLICY ITERATION (MONTECARLO)



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

Policy improvement: Greedy?

# MONTE-CARLO – POLICY ITERATION WITH ACTION-VALUE FUNCTION

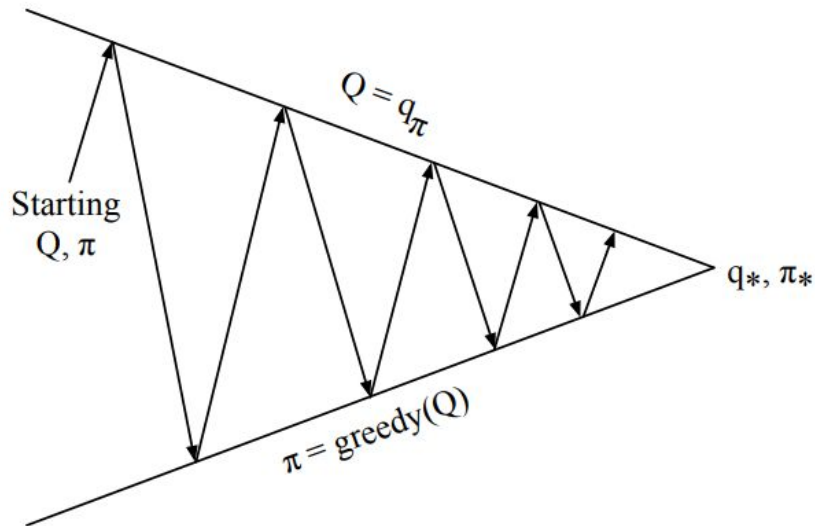
- Greedy policy improvement over  $V(s)$  requires a model of MDP

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

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

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$$

# POLICY ITERATION (MONTECARLO)

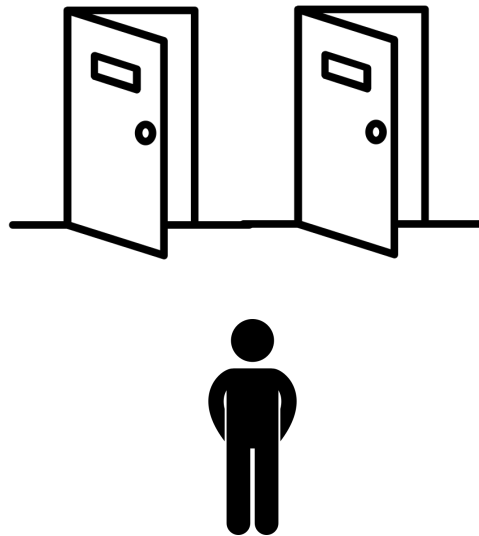


Policy evaluation Monte-Carlo policy evaluation  $Q = q_\pi$

Policy improvement: Greedy?

# EXAMPLE OF GREEDY ACTION SELECTION

- There are two doors in front of you
- Open left door (reward 0)  $\rightarrow V(\text{left}) = 0$
- Open left door (reward +1)  $\rightarrow V(\text{right}) = +1$
- Open left door (reward +3)  $\rightarrow V(\text{right}) = +2$
- Open left door (reward +2)  $\rightarrow V(\text{right}) = +2$
- We've chosen the best door?





# $\epsilon$ -Greedy Exploration

- 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

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

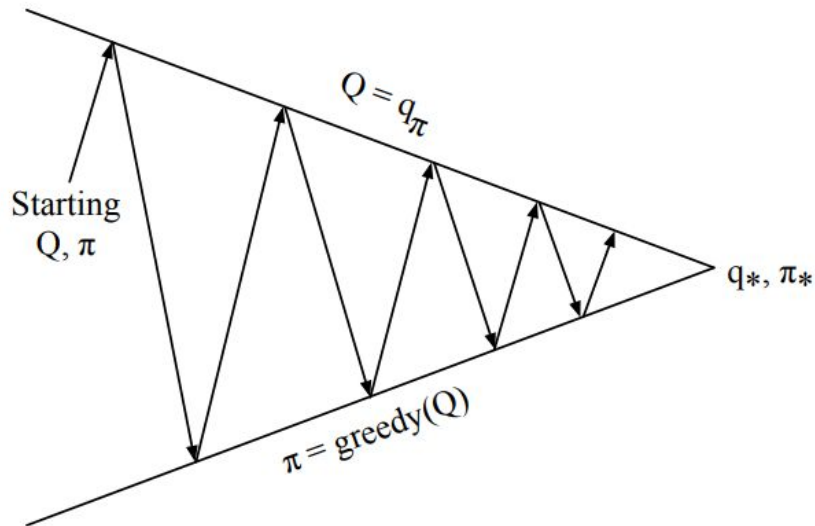
# $\epsilon$ -Greedy Improvement

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

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

- From policy improvement theorem  $v_{\pi'}(s) \geq v_\pi(s)$

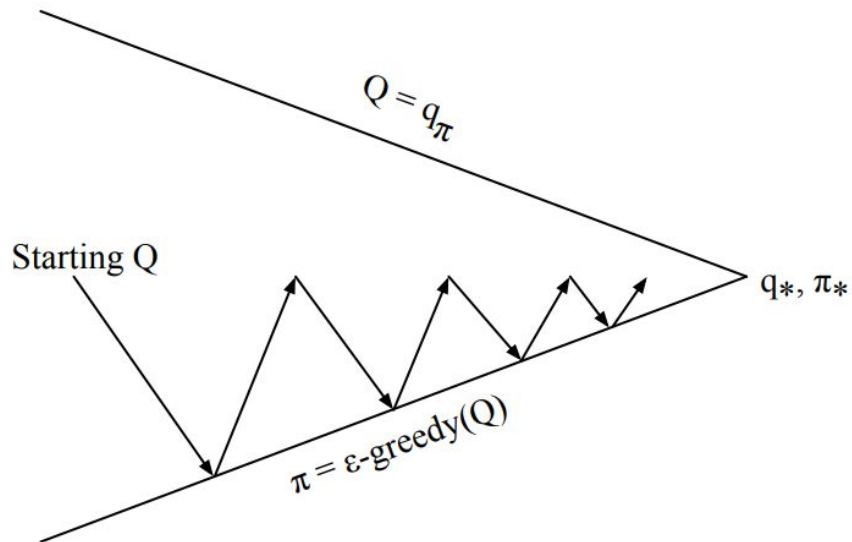
# POLICY ITERATION (MONTE-CARLO)



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

Policy improvement:  $\epsilon$ -Greedy Policy Improvement

# MONTÉ-CARLO CONTROL



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

Policy improvement:  $\epsilon$ -Greedy Policy Improvement

# GLIE

- 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

$$\lim_{k \rightarrow +\infty} \pi_k(a|s) = 1(a = \operatorname{argmax}_{a' \in A} Q_k(s, a'))$$

# GLIE - MONTE-CARLO

- Sample kth 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

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

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + 1/N(S_t, A_t) (G_t - Q(S_t, A_t))$$

- Improve policy based on new action-value function

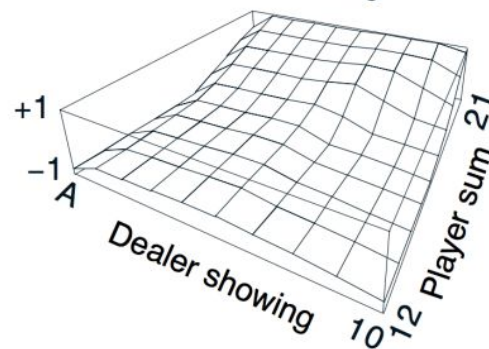
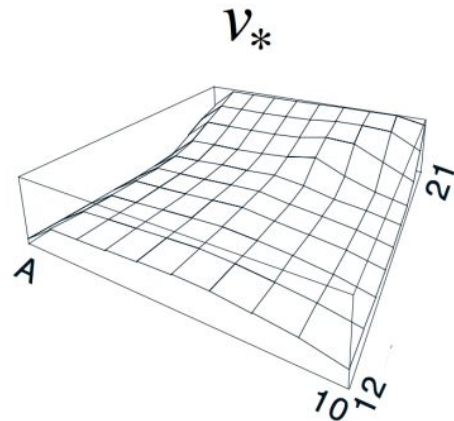
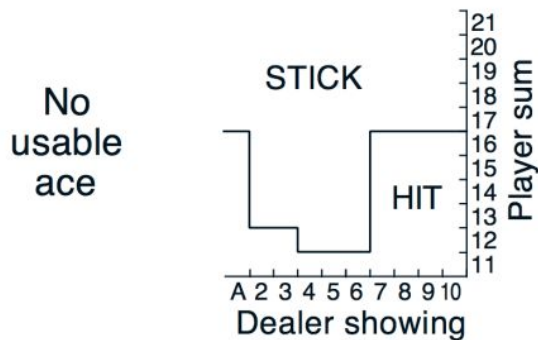
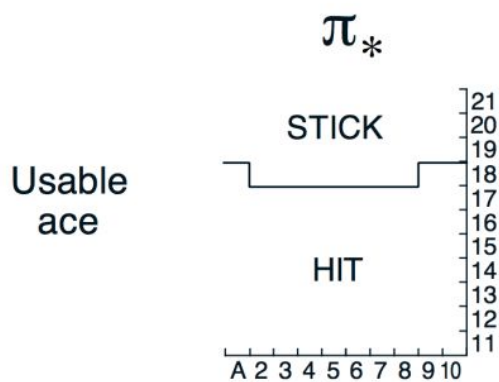
$$\varepsilon \leftarrow 1/k$$

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

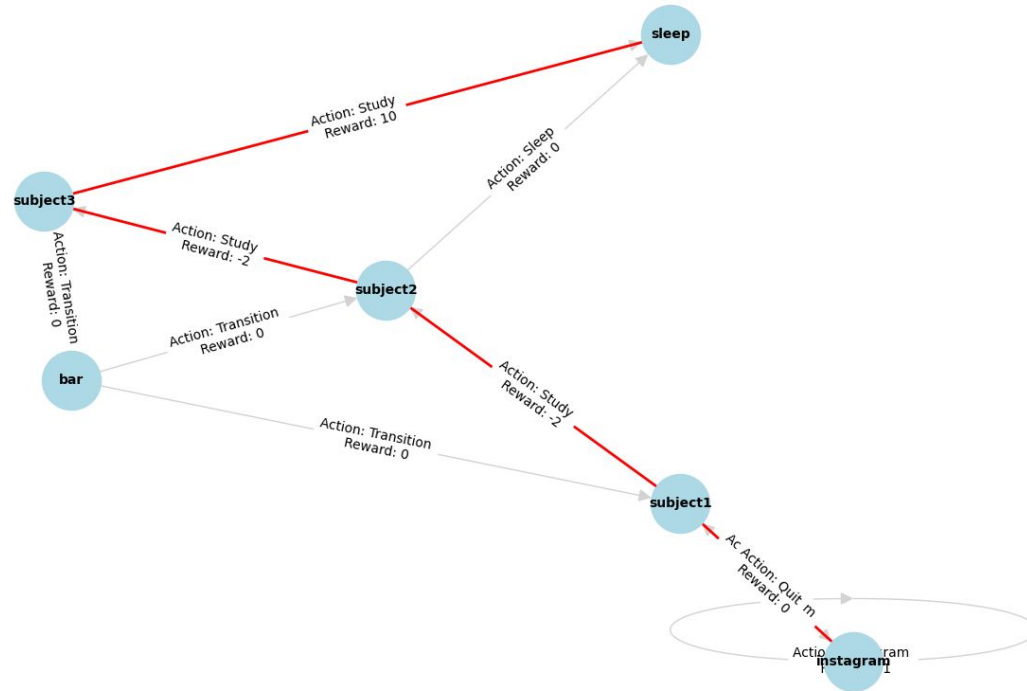
- Theorem: GLIE Monte-Carlo converges to the optimal action-value function

$$Q(s,a) \rightarrow q_*(s,a)$$

# MONTÉ-CARLO CONTROL (BLACKJACK)



# ASSIGNMENT: MONTE CARLO LEARNING – A STUDENT DECISION-MAKING PROBLEM





# ASSIGNMENT: BUILDING A STUDENT DECISION-MAKING PROBLEM

- **Understand Monte Carlo Learning:**
  - Read and comprehend the provided code.
  - Add comments to the code to explain each part.
- **Visualize the Student Chain Problem:**
  - Use NetworkX to build a graph of the states.
  - Highlight the best policy **in red** on the graph.

# ASSIGNMENT: MONTE CARLO LEARNING – A STUDENT DECISION-MAKING PROBLEM

- **Submission:**

- Submit the Python Jupyter script containing the code and comments
- <https://shorturl.at/JMn8r> - Expiration date: **1st June 2024**

- **Note:** Ensure the code is well-commented and understandable. Avoid using external guidance or assistance, as the purpose of the assignment is to demonstrate understanding and implementation independently.