# Model Free Control

## Model Free Reinforcement Learning

- Last lecture:
- Model-Free Prediction
- Estimate the value function of an unknown MDP
- This lecture:
- Model-Free Control
- Optimise the value function of an unknown MDP

## Uses of Model-Free Control

Some example problems that can be modelled as MDPs:

- ✓ Elevator                    ✓ Quake
- ✓ Parallel parking           ✓ Robocup Soccer
- ✓ Ship Steering              ✓ Portfolio management
- ✓ Bioreactor                 ✓ Protein Folding
- ✓ Helicopter                 ✓ Robot walking
- ✓ Aeroplane logistics       ✓ Game of Go

- For most of this problems, either:
- MDP model is unknown, but experience can be sampled
- MDP model is known, but is too big to use, except by sampling
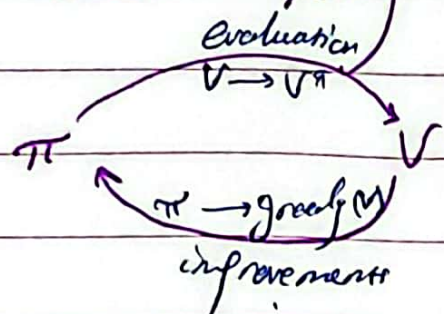- Model Free Control can solve these problems.

# On and off-policy Learning

- **On-policy Learning:** 1) "Learn on the job"
  2) Learn about policy $\pi$ from experience sampled from $\pi$
- **Off-policy Learning:** 1) "Look over someone's shoulder"
  2) Learn about policy $\pi$ from experience sampled from $\mu$

## Generalized Policy Iteration (Refresher)

evaluation
$V \rightarrow V^\pi$

$\pi \rightarrow$ greedy(V)
improvement

$\pi^* \longrightarrow V^*$

Policy Evaluation Estimate $V_\pi$
eg Iterative Policy evaluation

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

## Generalized Policy Iteration with Monte-Carlo Evaluation

Policy Evaluation Monte-Carlo Policy evaluation, $V = V_\pi$ ?

Policy improvement $\varepsilon$-Greedy policy improvement?
do not have exploration

needs dynamics MDP $\rightarrow Q = V_\pi$

## Model Free Policy Iteration Using Action-Value Function

- Greedy policy improvement over $V(s)$ requires model of MDP
$$\pi'(s) = \arg\max_{a \in A} R_s^a + P_{ss'}^a V(s')$$

- Greedy policy improvement over $Q(s,a)$ is model-free
$$\pi'(s) = \arg\max_{a \in A} Q(s,a)$$

# ε-Greedy Exploration

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

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

# ε-Greedy Policy improvement

## Theorem:

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

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

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

$$\geq \frac{\varepsilon}{m} \sum_{a \in A} q_\pi(s,a) + (1-\varepsilon) \sum_{a \in A} \frac{\pi(a|s) - \varepsilon/m}{1-\varepsilon} q_\pi(s,a) \qquad \pi(a|s) \geq \frac{\varepsilon}{m} \; \forall's$$

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

therefore from policy improvement theorem $V_{\pi'}(s) \geq V_\pi(s)$

# GLIE

## Definition

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

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

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

## GLIE Monte-Carlo Control

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,

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

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

Improve policy based on new action-value function

$$\varepsilon \leftarrow \frac{1}{k}$$

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

### Theorem:

GLIE Monte-Carlo Control converges to the optimal action-value function

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

# Off Policy Learning

- Evaluate target policy $\pi(a|s)$ to compute $V_\pi(s)$ or $q_\pi(s,a)$
- while following behaviour policy $\mu(a|s)$

$$\{S_1, A_1, R_2, \ldots, S_T\} \sim \mu$$

- why is this important?
- learn from observing humans or other agents
- Re use experience generated from old policies $\pi_1, \pi_2, \ldots, \pi_{t-1}$
- learn about optimal policy while following exploratory policy
- learn about multiple policies while following one policy

## Importance Sampling

Estimate the expectation of a different distribution

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

## Off Policy Monte-Carlo

use returns generated from $\mu$ to evaluate $\pi$

weight return $G_t$ according to similarity between policies

Multiply importance sampling corrections along whole episode

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

Update value towards corrected return

$$V(S_t) \leftarrow V(S_t) + \alpha \left(G_t^{\pi/\mu} - V(S_t)\right)$$

Cannot use if $\mu$ is zero when $\pi$ is non-zero

Importance sampling $\underset{\text{Can}}{\text{and}}$ dramatically increase $\underset{\text{variance}}{\text{variance}}$

$\hookrightarrow$ extremely high variance and for off policy you have do use TD