# Reinforcement Learning Course by David Silver

Yaryl Gassimov Sutton's Notes

— 2023 —

## 4.1: Intro to Reinforcement Learning.

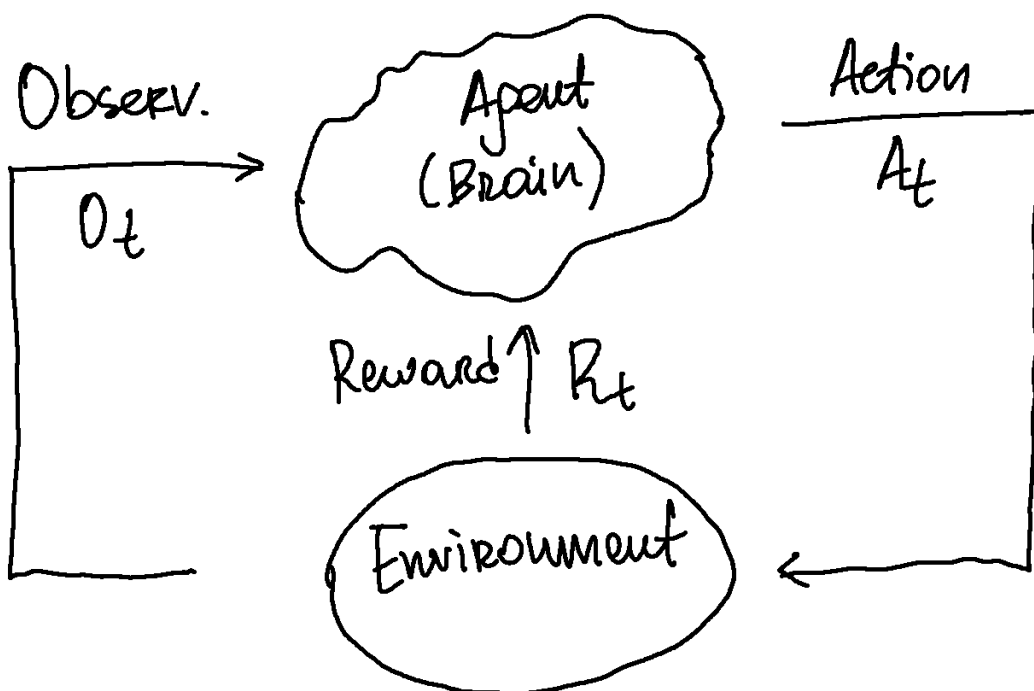Books: Sutton & Barto. An introduction to reinforcement learning. 1998.

Szepesvari. Algorithms for RL. 2010.
↳ more mathematical.

Example: • manage an investment portfolio

Reward $R_t$ — scalar feedback signal.

Def. (Reward hypothesis) All goals can be described by the maximisation of expected cumulative reward.

Goal: select actions to max. total future reward
• You cannot use greedy algos.

Observ.
$O_t$

Agent
(Brain)

Action
$A_t$

Reward ↑ $R_t$

Environment

The _history_ is the sequence:
$$H_t = A_1, O_1, R_1, \ldots, A_t, O_t, R_t$$

_State_ is the information used to determine what happens next.
$$S_t = f(H_t)$$

_Environment state_ $S_t^e$ is the environment's private representation.

not usually visible to the agent.

_Agent state_ $S_t^a$ is the agent's internal representation. $S_t^a = f(H_t)$

_Def._ A state $S_t$ is _Markov_ if and only if
$$\mathbb{P}[S_{t+1} | S_t] = \mathbb{P}[S_{t+1} | S_1, \ldots, S_t]$$

- "The future is independent of the past, given the present": $H_{1:t} \longrightarrow S_t \longrightarrow H_{t+1:\infty}$
- i.e. the state is sufficient statistic of the future.
- The environment state is Markov.
- The history is Markov.

$\rightarrow$ nice case.

Full observability : agent directly observes environment state. $O_t = S_t^a = S_t^e$.

Formally, this is Markov decision process (MDP).

Partial observability: agent indirectly observes environment. $S_t^a \neq S_t^e$

Formally, this is partially observable MDP (POMDP).

Agent must creates its own $S_t^a$, e.g.:
- Whole history
- Beliefs. $S_t^a = (P[S_t^e = s^1], \ldots, P[S_t^e = s^n])$.
- RNN. $S_t^a = \delta(S_{t-1}^a W_s + O_t W_o)$.

Inside an RL agent

- Policy : agent's behaviour. $a = \pi(s)$ ← policy
- Value function : prediction of future reward
- Model : predicts what env. will do next.

Maybe less, maybe more.

Policy may be stochastic: $\pi(a|s) = P[A=a|S=s]$

$V_\pi(s) = E_\pi[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \ldots | S_t = s]$

$\gamma$ — discount reward (e.g. 0.99).

About model:

Transitions: $P$ predict the next state (dynamics)

Rewards: $R$ predicts the next (immediate) Reward, e.g.

$$P_{ss'}^a = \mathbb{P}[S'=s' \mid S=s, A=a]$$

$$R_s^a = \mathbb{E}[R \mid S=s, A=a]$$

Categorizing RL agents: $\swarrow$ implicit

- Value based: No policy, value function.
- Policy based: policy, no value function.
- Actor critic: policy & value function.
- Model free: policy and/or value function, no model.
- Model based: policy and/or val. funct, model.

Exploitation — exploration.

## L2. Markov Decision Processes

Almost all RL problems can be formalized as MDP

**Def.** A state $S_t$ is Markov $\iff$

$$\mathbb{P}[S_{t+1} \mid S_t] = \mathbb{P}[S_{t+1} \mid S_1, \ldots, S_t]$$

State transition prob.: $\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' \mid S_t = s]$

State transition matrix:

$$P = \begin{bmatrix} P_{11} & \cdots & P_{1n} \\ \vdots & & \\ P_{n1} & \cdots & P_{nn} \end{bmatrix}$$

each row sums to 1.

**Def.** A Markov Process (or Markov Chain) is a tuple $\langle S, P \rangle$:

- $S$ is a (finite) set of sets
- $P$ is a prob. matrix

# Markov Reward Process

**Def.** <u>MRP</u> is a tuple $\langle S, P, R, \gamma \rangle$:

     • $R$ — reward function, $R_s = \mathbb{E}[R_{t+1} | S_t = s]$

     • $\gamma$ — discount factor, $\gamma \in [0,1]$

\* give value to vertices

**Def.** The <u>return</u> $G_t$ is the ~~total~~ discounted $\rightarrow$ goal

reward from time-step $t$:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

**Def.** The <u>value function</u> $v(s)$ of an MRP is the

expected return starting from state $s$:

$$v(s) = \mathbb{E}[G_t | S_t = s]$$

Bellman Equation for MRPs:

     • immediate reward $R_{t+1}$

     • discounted value of successor state $\gamma v(S_{t+1})$

$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

$$v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$$

$$V = R + \gamma P v \quad \text{(in matrices)}$$

$$\begin{pmatrix} v(1) \\ --- \\ v(n) \end{pmatrix} = \begin{pmatrix} R_1 \\ --- \\ R_n \end{pmatrix} + \gamma \begin{pmatrix} P_{11} \cdots P_{1n} \\ \vdots \quad \ddots \quad --- \\ P_{n1} \cdots P_{nn} \end{pmatrix} \begin{pmatrix} v(1) \\ --- \\ v(n) \end{pmatrix}$$

\* $n$ - number of states.

It can be solved directly:

$$V = (I - \gamma P)^{-1} R$$

$O(n^3)$ — comp. complexity.

## Markov Decision Processes

<u>Def.</u> A MDP is a tuple $\langle S, A, P, R, \gamma \rangle$:

- $A$ – a finite set of actions
- $P$ is a state transition prob. matrix:
$$P_{ss'}^a = \mathbb{P}[S_{t+1} = s' \mid S_t = s, A_t = a]$$

- $R$ is a reward function
$$R_s^a = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$$

**Def.** A _policy_ $\pi$ is a distribution over action, given states:

$$\pi(a|s) = \mathbb{P}[A_t = a \mid S_t = s]$$

A policy defines the behaviour of an agent. There, policy depends only on $s$, not time i.e. policy are stationary: $\forall t > 0, A_t \sim \pi(\cdot | S_t)$

**Def.** A _state-value function_ $v_\pi(s)$ of an MDP is the expected return starting from state $s$, and then following $\pi$

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

**Def.** The action-value function $q_\pi(s,a)$

$$q_\pi(s,a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$$

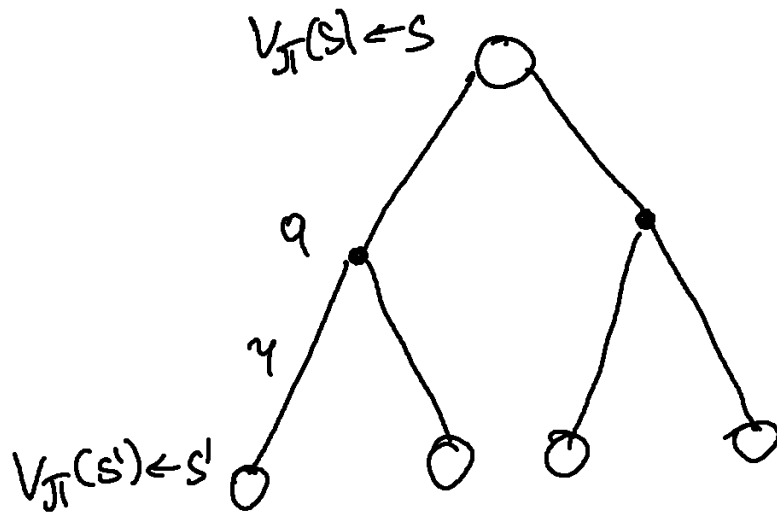expected return taking action $a$ and following $\pi$ from state $s$.

Bellman Expectation Equation:

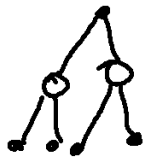$$v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s]$$

$$q_\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

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

$$q_\pi(s,a) = \mathcal{R}_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a \, V_\pi(s')$$



$$V_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left( \mathcal{R}_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a \, V_\pi(s') \right)$$

$$q_\pi(s,a) = \mathcal{R}_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a \sum_{a \in \mathcal{A}} \pi(a|s') \, q_\pi(s',a)$$

mean by $\pi$

In MRP: $\qquad V_\pi = R^\pi + \gamma \mathcal{P}^\pi V_\pi$

$$V_\pi = (I - \gamma \mathcal{P}^\pi)^{-1} R^\pi$$

**Def.** The optimal state-value function $V_*(s)$ is the max. value function over all policies:

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

The optimal action-value function $q_*(s,a)$ is the max. action-value function over all $\pi$:

$$q_*(s,a) = \max_{\pi} q_{\pi}(s,a)$$

\* it's max reward, not the best path.
↳ if it's known then we can find the best path.
↳ we should find $q_*$ then go by increasing $q_*$ path.

**Def.** $\pi \geq \pi'$ if $V_{\pi}(s) \geq V_{\pi'}(s), \forall s$

**Theor.** For any MDP
 • There exists an optimal policy $\pi_*$ that is better than or equal to all other policies, $\pi_* \geq \pi, \forall \pi$
 • All optimal policies achieve the optimal value function, $V_{\pi_*}(s) = V_*(s)$
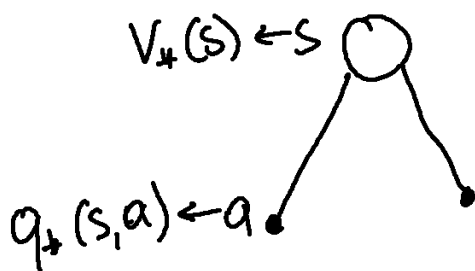 • All optimal policies achieve the optimal action-value function, $q_{\pi_*}(s,a) = q_*(s,a)$

$$\pi_*(a|s) = \begin{cases} 1, & \text{if } a = \underset{a \in \mathcal{A}}{\arg\max} \ q_*(a,s) \\ 0, & \text{otherwise} \end{cases}$$
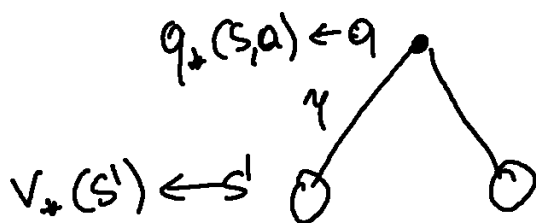
There's always optimal deterministic solution of MDP.

How do we find $q_*$? By looking backwards.

Bellman Optimality Equation



$V_*(s) \leftarrow s$

$q_*(s,a) \leftarrow a$

$$V_*(s) = \max_a q_*(s,a)$$

$q_*(s,a) \leftarrow a$

$\gamma$
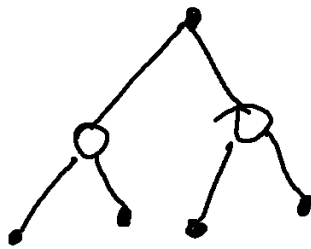
$V_*(s') \leftarrow s'$

$$q_*(s,a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$$

$$V_*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$$

$$q_*(s,a) = R_s^a + \sum_{s' \in S} P_{ss'}^a \max_{a'} q_*(s', a')$$

Bellman Optimality Equation is non-linear.
Many iterative solutions.

Extensions to MDPs

## L.3. Planning by Dynamic Programming

### Introduction

DP assumes *full knowledge* of the MDP.
It's used for _planning_ in an MDP.

$\ast$ Viterbi algorithm — what's

### Policy Evaluation

Problem: evaluate policy $\pi$

Solution: iterative application of Bellman expectation backup

$$V_1 \longrightarrow V_2 \longrightarrow \dots \longrightarrow V_\pi$$

Using synchronous backups

$$V_{k+1}(s) = \sum_{a \in d} \pi(a|s) \left( \mathcal{R}_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a V_k(s') \right)$$

$$V^{k+1} = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi V^k$$

### Policy Iteration

- Given a policy $\pi$
  - Evaluate $\pi$: $V_\pi(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = S]$
  - Improve $\pi$: $\pi' = greedy(V_\pi)$.

In general, need more iterations of improvement/evaluation

Always converges to $\Pi_*$.

Modified policy iteration:
- Does policy evaluation need to converge to $V_\pi$?
- Introduce stopping condition:
    - $\epsilon$ - convergence
    - $k$ iterations.

## Value iteration

**Theor.** (Principle of Optimality)

A policy $\pi(a|s)$ achieves the optimal value from state $s$, $V_\pi(s) = V_*(s) \iff$ for $\forall$ state $s'$ reachable from $s$, $\pi$ achieves the optimal value from state $s'$, $V_\pi(s') = V_*(s')$

$$V_*(s) \leftarrow \max_{a \in d} R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$$

Intuition: start w/ final rewards and work backwards.

Problem: find optimal policy $\pi$

Solution: iterative application of Bellman optimality backup

$$V_1 \longrightarrow V_2 \longrightarrow \ldots \longrightarrow V_*$$

$$V_{k+1}(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right)$$

in matrices: $\quad V_{k+1} = \max_{a \in A} \left( R^a + \gamma P^a V_k \right)$

| Problem | Bellman Equation | Algorithm |
|---|---|---|
| Prediction | Bellman Expectation Equation | Iterative Policy Eval. |
| Control | Bellman Exp. Eq. + Greedy Policy Improvement | Policy Iteration |
| Control | Bellman Opt. Eq. | Value Iter. |

Algos based on state-value function: $O(mn^2)$ per iter. $\rightarrow V_\pi(s), v_*(s)$

Algos based on action-value function: $O(m^2 n^2)$ per iter. $\hookrightarrow q_\pi(s,a), q_*(s,a)$

# Extensions to Dynamic Programming

Three simple ideas for asynchronous DP:

1) In-place dp
2) Prioritised sweeping
3) Real-time dp

1) In synchr. you save 2 versions of v.
In-place stores only one.

2) Use magnitude of Bellman error to guide state selection:

$$\left| \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v(s') \right) - v(s) \right|$$

Backup the state with largest remaining Bellman error

3) Only states that are relevant to agent
Use agent's XP to guide the selection of states
After each time step $S_t, A_t, R_{t+1}$
Backup $S_t$

$$v(S_t) \leftarrow \max_{a \in A} \left( R_{S_t}^a + \gamma \sum_{s' \in S} P_{S_t s'}^a v(s') \right)$$

Large DP suffers curse of dimensionality,
so we'll use sampling.

## Contraction mapping

$\hookrightarrow$ Math is in lecture notes.

## L4. Model-free prediction

### Introduction

↳ Estimate the value function of unknown MDP.
  ↳ policy evaluation

## Monte-Carlo learning

Can only be applied to episodic MDPs:

- All episodes must terminate

Goal: learn $V_\pi$ from episodes of XP under $\pi$

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

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

Value function: $V_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$ ←

Monte-Carlo uses empirical mean return.

$V(s) \longrightarrow V_\pi(s)$ , $V(s) = S(s)/N(s)$

as $N(s) \to \infty$

↑ sum of total return where s was visited

↳ number of times s is visited first time

- First-visit-MC evaluation

## Every-visit MC policy Estimation:

$N(s)$ increments every time $s$ encountered

$S(s) \leftarrow S(s) + G_t$ every-time $s$ encountered

$V(s) = S(s)/N(s)$

The mean $\mu_1, \mu_2,...$ of sequence $x_1, x_2,...$ can be computed incrementally:

$$\mu_k = \frac{1}{k} \sum_{j=1}^{k} x_j = \mu_{k-1} + \frac{1}{k}(x_k - \mu_{k-1})$$

## Incremental MC Updates

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 running mean, i.e. forget old episodes:

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

## Temporal-Difference Learning

TD learns from incomplete episodes, by bootstraping.
TD updates a guess towards a guess.

# Simple TD learning algo TD(0)

- Update $V(S_t)$ toward estimated return

$$R_{t+1} + \gamma V(S_{t+1}) \quad \text{(TD target)}$$

$$V(S_t) \leftarrow V(S_t) + \alpha \underbrace{\left(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)\right)}_{\delta_t - \text{TD error}}$$

TD ~~target~~ is biased estimate of $V_\pi(S_t)$
It has lower variance than the return

MC has high variance, zero bias
TD has low variance, some bias

MC converges to solution w/ minimum MSE
TD converges to the solution of max likelihood
Markov Model.
$\rightarrow$ exploits Markov property

$\rightarrow$ doesn't exploit Markov property

n-Step return

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$G_t^{(\infty)} - MC$$

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

Averaging n-Step Returns

e.g. $\frac{1}{2} G^{(2)} + \frac{1}{2} G^{(4)}$

## TD($\lambda$)

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^\lambda - V(S_t) \right)$$

if terminal state, $\lambda^{T-t-1}$

Can only be computed for complete episodes.

$\lambda = 1 - MC$
$\lambda = 0 - TD(0)$

Eligibility traces
- Frequency heuristic : assign credit to most freq. state
- Recency heuristic : — / — / — / — / recent state

$$E_0(s) = 0$$

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

Backward view TD($\lambda$)

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

$$V(S) \leftarrow V(S) + \alpha \delta_t E_t(s)$$

When $\lambda = 0 \implies E_t(s) = \mathbb{1}(S_t = S) = 1 \implies TD(0)$

$$\lambda = 1 \implies MC$$

__Theor.__ The sum of offline updates is identical for forward-view and backward-view TD($\lambda$):

$$\sum_{t=1}^{T} \alpha \delta_t E_t(s) = \sum_{t=1}^{T} \alpha \left( G_t^\lambda - V(S_t) \right) \mathbb{1}(S_t = S)$$

## L5. Model-free control

## Introduction

Model-free prediction: evaluates policy $\pi$
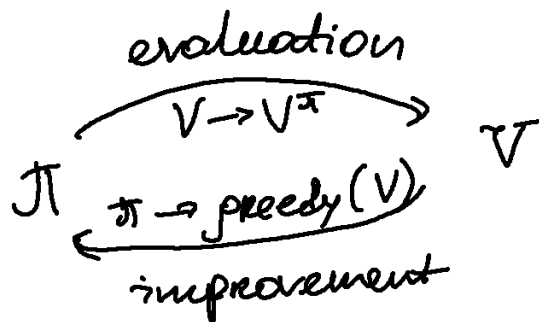find best possible policy

┌ Note: Try model easy football and try RL ┐

On-policy : "learn on the job"
learn about $\pi$ from XP using $\pi$

Off-policy: "Look over someone's sholder"
learn about policy $\pi$ from XP sampled
from $\mu \neq \pi$

evaluation
$$V \to V^\pi$$

$\pi$   $\pi \to$ greedy(V)   $V$

improvement

$\longrightarrow$ Exploration problem

Greedy Action Selection $\longrightarrow$ Bandit problem

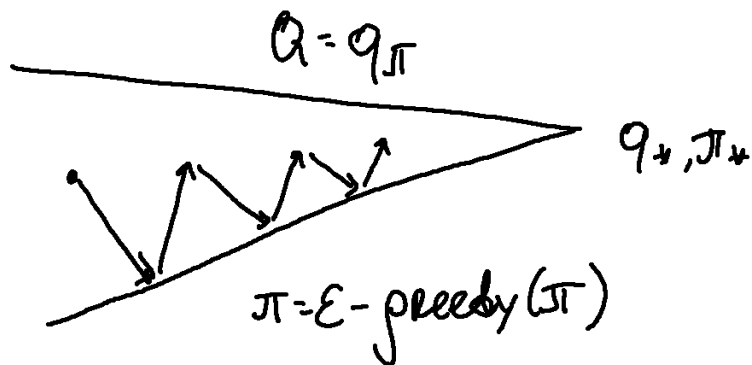$\varepsilon$-greedy exploration:

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

**Theorem.** $\forall$ $\varepsilon$-greedy policy $\pi$, the $\varepsilon$-greedy $\pi'$ withe respect to $q_\pi$ is an improvement,

$$V_{\pi'}(s) \geq V_\pi(s).$$

policy improvement theorem

Idea: update $V$ every episode:



$$Q = q_\pi$$

$$q_*, \pi_*$$

$$\pi = \varepsilon - \text{greedy}(\pi)$$

**Def.** Greedy in the Limit with Infinite Exploration (GLIE):

• all state-action pairs are explored infinitely many times, $\lim\limits_{k \to \infty} N_k(s,a) = \infty$

• the policy converges on a greedy policy,

$$\lim\limits_{k \to \infty} \pi_k(a|s) = \mathbb{1}(a = \underset{a' \in \mathcal{A}}{\text{argmax}} \, Q_k(s,a'))$$
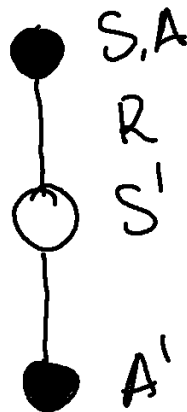
e.g. $\varepsilon$ reduces to zero at $a_k = 1/k$.

# GLIE MC Control

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

- $\forall\, S_t$ and $A_t$ in the episode

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

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

- Improve policy based on new action-value function

$$\varepsilon \leftarrow 1/k$$

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

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

## Updating Action-Value Functions w/ Sarsa



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

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

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

$$\sum_{n=1}^{\infty} \alpha_t = \infty ; \quad \sum_{n=1}^{\infty} \alpha_t^2 < \infty. \qquad \alpha_t:$$

### n-Step Sarsa

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

$$Q(S_t, A_t) \longleftarrow Q(S_t, A_t) + \alpha \left( q_{(t)}^n - Q(S_t, A_t) \right)$$

### Sarsa $(\lambda)$

$$q_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$$

$$Q(S_t, A_t) \longleftarrow Q(S_t, A_t) + \alpha \left( q_t^\lambda - Q(S_t, A_t) \right)$$

↰ forward-view (not online because wait till the end of an episode)

### Backward View Sarsa $(\lambda)$

Use **eligibility traces**

$$E_0(s,a) = 0$$
$$E_t(s,a) = \gamma \lambda E_{t-1}(s,a) + \mathbb{1}(S_t = s, A_t = a)$$

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$
$$Q(s,a) \leftarrow Q(s,a) + \alpha \delta_t E_t(s,a)$$

## Off-policy learning

Evaluate $\pi(a|s)$ to compute $V_\pi(s)$ or $q_\pi(s,a)$ while following $\mu(a|s)$.

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

Important because e.g. see what people did and do better by robot.

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

## Importance Sampling for Off-policy MC

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

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

Very high variance, so don't work

## Importance Sampling for Off-policy TD

$$V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

Much lower variance

## Q-learning

$$A_{t+1} \sim \mu(\cdot|S_t)$$
$$A' \sim \pi(\cdot|S_t)$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t) \right)$$

$\exists \; \pi$ is greedy with respect to $Q(s,a)$

$\exists \; \mu$ is $\varepsilon$-greedy w.r.t. $Q(s,a)$

Then, $R_{t+1} + \gamma Q(S_{t+1}, A') = R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a')$

# Q-learning Control Algorithm (SARSAMAX)

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

**Theor.** Q-learning control converges to the optimal action-value function, $Q(S,a) \rightarrow q_*(S,a)$

## L6. Value Function Approximation.

### Introduction

$$\hat{v}(s,w) \approx v_\pi(s)$$

$$\hat{q}(s,a,w) \approx q_\pi(s,a)$$

$w$ - parameters (of NN, for example)

we'll consider approximators as • linear comb. of features

• NN

we require that methods are suitable for non-iid and non-stationary data.

### Incremental Methods

$\rightarrow$ Value Function Approx. using SGD

VFA Represent state by feature vector:

$$x(s) = \begin{pmatrix} x_1(s) \\ --- \\ x_n(s) \end{pmatrix}$$

Linear VFA: $\hat{v}(s,w) = x(s)^T w = \sum_{j=1}^{n} x_j(s) w_j$

Update $= \alpha \left( v_\pi(s) - \hat{v}(s,w) \right) x(s)$

$=$ step-size $\times$ error $\times$ feature value

Table lookup is special case of linear VFA:

$$x^{table}(S) = \begin{pmatrix} \mathbb{1}(S=S_1) \\ --- \\ \mathbb{1}(S=S_n) \end{pmatrix}$$

Incremental Prediction Algorithms:

- For MC:

$$\Delta w = \alpha \left( G_t - \hat{v}(S_t, w) \right) \nabla_w \hat{v}(S_t, w)$$

- For TD(0):

$$\Delta w = \alpha \left( R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w) \right) \nabla_w \hat{v}(S_t, w)$$

- For TD($\lambda$):

Forward: $\Delta w = \alpha \left( G_t^\lambda - \hat{v}(S_t, w) \right) \nabla_w \hat{v}(S_t, w)$
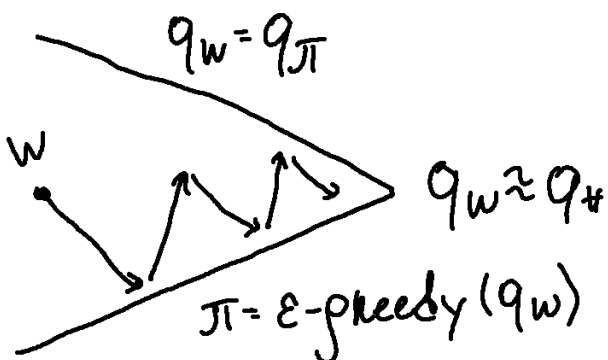
Backward:

$$\delta_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)$$

$$E_t = \gamma \lambda E_{t-1} + x(S_t)$$

$$\Delta w = \alpha \delta_t E_t$$

Control with VFA



$q_w = q_\pi$

$q_w \approx q_*$

$\pi = \varepsilon\text{-greedy}(q_w)$

Approximate policy evaluation

$$\hat{q}(\cdot, \cdot, w) \approx q_\pi$$

# Action-value Function Approximation

AVFA
$$\hat{q}(S,A,w) \approx q_\pi(S,A)$$

$$\Delta w = \alpha \left( q_\pi(S,A) - \hat{q}(S,A,w) \right) \nabla_w \hat{q}(S,A,w)$$

feature: $X(S,A) = \begin{pmatrix} x_1(S,A) \\ \cdots \\ x_n(S,A) \end{pmatrix}$

linear AVFA: $\hat{q}(S,A,w) = X(S,A)^T w$

Increremental Control Algorithms:
Same as prediction, but use
$\hat{q}(S_t, A_t, w)$ instead of $\hat{v}(S_t, w)$

( Mountain Car problem )

Convergence of Prediction Algos

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

Gradient TD   all   ✓

# Convergence of Control Algos

| Algo | Table | Linear | Non-Linear |
|---|---|---|---|
| MC | V | (V) | X |
| Sarsa | V | (V) | X |
| Q-learning | V | X | X |
| Gradient Q-learning | V | V | X |

(V) = chatters around

# Batch Methods

Gradient descent is not sample efficient

Batch try to find the best fitting value function

experience

$$D = \{ \langle S_1, V_1^{\pi} \rangle, \ldots, \langle S_T, V_T^{\pi} \rangle \}$$

## Least squares algo:

$$LS(w) = \sum_{t=1}^{T} \left( V_t^{\pi} - \hat{V}(S_t, w) \right)^2 =$$

$$= \mathbb{E}_D \left[ \left( V^{\pi} - \hat{V}(S, w) \right)^2 \right] \longrightarrow \min$$

# SGD with Experience Replay

Repeat: 1. Sample state, value from XP

$$\langle S, V^\pi \rangle \sim D$$

2. Apply SGD update:

$$\Delta w = \alpha (V^\pi - \hat{v}(S,w)) \nabla_w \hat{v}(S,w)$$

Converges to least squares solution:

$$w^\pi = \underset{w}{argmin}\ LS(w)$$

Linear LS-prediction: — condition (?)

$$\alpha \sum_{t=1}^{T} X(S_t)(V_t^\pi - X(S_t)^T w) = 0;$$

$$w = \left( \sum_{t=1}^{T} X(S_t) X(S_t)^T \right)^{-1} \sum_{t=1}^{T} X(S_t) V_t^\pi$$

LSMC: $V_t^\pi \approx G_t$

LSTD: $V_t^\pi \approx R_{t+1} + \gamma \hat{v}(S_{t+1}, w)$

LSTD($\lambda$): $V_t^\pi \approx G_t^\lambda$

Converges on- and off-policy both MC & TD.

LS Policy Iteration.

## L7. Policy Gradient Methods.

### Introduction

$$\pi_\theta(s,a) = \mathbb{P}[a|s,\theta]$$

$+$: policy methods can be stochastic

in partially observed env. (using features) we lose Markov Property $\implies$ no ~~deterministic~~ deterministic solution.

Goal: given $\pi_\theta(s,a)$ w/ $\theta$, find best $\theta$

To measure quality:

- start value: $J_1(\theta) = V^{\pi_\theta}(S_1) = \mathbb{E}_{\pi_\theta}(V_1)$

- average value: $J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$

- average reward per time-step:

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s,a) R_s^a$$

$d^{\pi_\theta}(s)$ — stationary distribution.

Find $\theta$ that maximizes $J$.

# Finite Difference Policy Gradient

$$\Delta \theta = \alpha \nabla_\theta J(\theta)$$

$$\nabla_\theta J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ --- \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix} \quad ; \quad \frac{\partial J(\theta)}{\partial \theta_k} \simeq \frac{J(\theta + \varepsilon u_k) - J(\theta)}{\varepsilon}$$

## MC Policy Gradient

likelihood ratios : $\nabla_\theta \pi_\theta(s,a) = \pi_\theta(s,a) \nabla_\theta \log \pi_\theta(s,a)$

score function : $\nabla_\theta \log \pi_\theta(s,a)$.

### softmax policy

$$\pi_\theta(s,a) \propto e^{\phi(s,a)^T \theta} \quad \longleftarrow \text{linear comb.}$$

$$\nabla_\theta \log \pi_\theta(s,a) = \phi(s,a) - \mathbb{E}_{\pi_\theta}[\phi(s,\cdot)]$$

### gaussian policy

$\phi(s)$

$\uparrow$ features

$\phi(s,a)$

$$\mu(s) = \phi(s)^T \theta \quad \longleftarrow \text{can be parametrized}$$

$$a \sim N(\mu(s), \sigma^2)$$

$$\nabla_\theta \log \pi_\theta(s,a) = \frac{(a - \mu(s)) \phi(s)}{\sigma^2}$$

**Theorem** For ∀ diff. policy $\pi_\theta(s, a)$, for any of the policy obj. functions $J \in \{J_1, J_{avR},$ $\frac{1}{1-\gamma} J_{avV}\}$, the policy gradient is

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\, Q^{\pi_\theta}(s, a)\right]$$

Algorithm idea: use return $v_t$ as an unbiased sample of $Q^{\pi_\theta}(s_t, a_t)$

$$\Delta\theta_t = \alpha \nabla_\theta \log \pi_\theta(s_t, a_t)\, v_t$$

50:27 ← full algo.

## Actor-Critic Policy Gradient

We use **critic** to estimate action-value function

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$$

critic params : $w$
actor params : $\theta$

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a)\, Q_w(s, a)\right]$$

$$\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a)\, Q_w(s, a)$$

The critic is solving policy evaluation problem. We can do this using MC, TD, TD($\lambda$), LS.

## Theorem (Compatible Function Approximation Theor)

If the following conditions are satisfied:

① $\nabla_w Q_w(s,a) = \nabla_\theta \log \pi_\theta(s,a)$

② w minimizes MSE
$$\varepsilon = \mathbb{E}_{\pi_\theta}\left[(Q^{\pi_\theta}(s,a) - Q_w(s,a))^2\right]$$

Then the policy gradient is exact,
$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s,a) \, Q_w(s,a)\right]$$

## Reducing Variance Using a Baseline

↳ Reduce variance \wo changing expectation

We subtract baseline function $B(S)$

A good baseline $B(S) = V^{\pi_\theta}(s)$

$$A^{\pi_\theta}(s,a) = Q^{\pi_\theta}(s,a) - V^{\pi_\theta}(s)$$

advantage function

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s,a) \, A^{\pi_\theta}(s,a)\right]$$

# Estimating the Advantage Function

We can use TD error as unbiased estimate of the $A^{\pi_\theta}(s,a)$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s,a)\, \delta^{\pi_\theta}\right]$$

$$\delta^v_v = r + \gamma V_v(s') - V_v(s)$$

in practice, can use approx. TD error.

# Policy Gradient w/ Eligibility Traces

$$\Delta\theta = \alpha(V_t^{\lambda} - V_v(s_t))\nabla_\theta \log \pi_\theta(s_t, a_t)$$

$$\delta^v = r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t)$$

$$e_{t+1} = \lambda e_t + \nabla_\theta \log \pi_\theta(s,a)$$

$$\Delta\theta = \alpha \delta^v e_t$$

Can be applied online, to incomplete sequences.

# Natural Policy Gradient     Fisher information matrix

$$\nabla_\theta^{nat} \pi_\theta(s,a) = G_\theta^{-1} \nabla_\theta \pi_\theta(s,a)$$

$$G_\theta = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s,a) \nabla_\theta \log \pi_\theta(s,a)^T\right]$$

## L8. Integrating Learning and Planning.

### Introduction

learn model from XP.

## Model-based Reinforcement Learning

- a model $\mathcal{M}$ is a representation of MDP $\langle S, A, P, R \rangle$, parametrized by $\eta$.
- assume that $S$ and $A$ are known
- $\mathcal{M} = \langle P_\eta, R_\eta \rangle$ represent state transition & reward.
- assume independence:

$$\mathbb{P}[S_{t+1}, R_{t+1} \mid S_t, A_t] = \mathbb{P}[S_{t+1} \mid S_t, A_t] \, \mathbb{P}[R_{t+1} \mid S_t, A_t]$$

Goal: Estimate $\mathcal{M}_\eta$ from XP $\{S_1, A_1, R_1, \ldots, S_T\}$

$$S_1, A_1 \longrightarrow R_2, S_2$$

$$\cdots \cdots$$

$$S_{T-1}, A_{T-1} \longrightarrow R_T, S_T$$

- $S, a \longrightarrow r$ — regression problem
- $S, a \longrightarrow S'$ — density estimation problem
- pick loss function
- find params $\eta$ that minimizes loss

# Table lookup model

$$P^a_{s,s'} = \frac{1}{N(s,a)} \sum_{t=1}^{T} \mathbb{1}(S_t, A_t, S_{t+1} = s, a, s')$$

$$R^a_s = \frac{1}{N(s,a)} \sum_{t=1}^{T} \mathbb{1}(S_t, A_t = s, a) R_t$$

Alternatively, record all $\langle S_t, A_t, R_{t+1}, S_{t+1} \rangle$ and then sample $\langle s, a, \cdot, \cdot \rangle$

Planning with a model:
- value iteration
- policy iteration
- tree search

## Sample-Based planning

Use model to generate samples

$$S_{t+1} \sim P_\eta(S_{t+1} | S_t, A_t)$$

$$R_{t+1} \sim R_\eta(R_{t+1} | S_t, A_t)$$

Apply model-free RL: MC control, Sarsa, Q-learning
→ Usually more efficient.

# Integrated Architectures

We consider two sources of XP:
- Real : sampled from env.
- Simulated : sampled from model $M_\eta$

Dyna - architecture :
- learn and plan value function (and/or policy) from real and simulated XP.

| Dyna-Q algo 54:11 |

The changed environment is <u>harder/easier</u>
↳ need exploration: Dyna-Q+
revisit states that haven't been visited for a while.

# Simulation-Based Search

Forward Search : build search tree with the curr. $s_t$
use model of MDP to look ahead
just sub-MDP from now.

Simulate from now → apply model-free RL

# Monte-Carlo Tree Search (Evaluation)

Given $\mathcal{M}_\vartheta$

Simulate: $\{ \underline{S_t}, A_t^k, R_{t+1}^k, S_{t+1}^k, \ldots, S_T^k \}_{k=1}^K$ ~

$\sim \mathcal{M}_{\vartheta, \pi}$

Build search tree

Evaluate states $Q(s, a)$:

$$Q(s,a) = \frac{1}{N(s,a)} \sum_{k=1}^{k} \sum_{u=t}^{T} \mathbb{1}(S_u, A_u = s, a)\, G_u \longrightarrow$$

$$\xrightarrow{P} q_\pi(s, a)$$

Select action: $a_t = \underset{a \in \mathcal{A}}{\arg\max}\, Q(S_t, a)$

- - - - - - - - - - - - - - - - - - - - - - -

MC control applied to simulated XP

$$Q(S, A) \longrightarrow q_*(S, A)$$

maximin — the best way to act is to think that
    opponent is acting the best way too

TD search applies Sarsa to sub-MDP.

# TD Search

For each step of simulation, update $Q(S,A)$ using Sarsa:

$$\Delta Q(S,A) = \alpha(R + \gamma Q(S',A') - Q(S,A))$$

select actions based action values $Q(Sa)$.

e.g. $\varepsilon$-greedy

may also use function approx. for $Q$.

# Dyna-2

stores: long-term memory
short-term (working) memory

$\rightarrow$ updated from real XP
updated from sim. XP

## L9. Exploration & Exploitation

### Introduction

Approaches to exploration:

1. Random expl. ($\varepsilon$-greedy)
2. Optimism in the face of uncertainity
   (prefer to explore uncertainity)
3. Information state space
   (agent's information as part of its state)

State-action vs. Parameter exploration
pick diff. A each time S is visited

$$\pi(A|S,\underline{u})$$

We'll focus on state-action exploration.

### Multi-Armed Bandits

tuple $\langle \mathcal{A}, \mathcal{R} \rangle$

$\mathcal{A}$-actions ("arms")

$\mathcal{R}^a(u) = \mathbb{P}[R = u | A = a]$ — unknown

$A_t \in \mathcal{A}$ generates $R_t \sim \mathcal{R}^{A_t}$

$$\sum_{\tau=1}^{t} R_\tau \longrightarrow max$$

Def. Action-value: $q(a) = \mathbb{E}[R|A=a]$

Def. Opt. value $V_*$: $V_* = q(a^*) = \max_{a \in \mathcal{A}} q(a)$

Def. Regret: $I_t = \mathbb{E}[V_* - q(A_t)]$

Def. Total regret: $L_t = \mathbb{E}\left[\sum_{J=1}^{t} V_* - q(A_J)\right]$

Count $N_t(a)$ ; Gap $\Delta_a = V_* - q(a)$

$$L_t = \sum_{a \in \mathcal{A}} \mathbb{E}[N_t(a)] \Delta_a$$

Good algos ensures small counts for large gaps.

Greedy can lock onto a subopt. solution.

Optimistic initialization: $Q(a) = r_{max}$

Then act greedily: $A_t = \text{argmax}_{a \in \mathcal{A}} Q_t(a)$

Optimistic greedy has linear total regret.

$\varepsilon$-greedy explores forever
has linear total regret.

pick a decay schedule for $\varepsilon_1, \varepsilon_2, \ldots$

e.g.        $c > 0$

$d = \min_{a | \Delta_a > 0} \Delta_a$

$\varepsilon_t = \min\left\{1, \frac{c|\mathcal{A}|}{d^2 t}\right\}$

Decaying $\varepsilon_t$-greedy has logarithmic total regret.

## Theor. (Lai & Robbins)

Asymptotic total regret is at least logarithmic in numbers of steps:

$$\lim_{t \to \infty} L_t \geq \log t \sum_{a \mid \Delta_a > 0} \frac{\Delta_a}{KL(R^a \| R^{a^*})}$$

## Upper Confidence Bounds

Estimate upper conf. bound $U_t(a)$

$$A_t = \underset{a \in \mathcal{A}}{argmax}\ Q_t(a) + U_t(a)$$

## Theor. (Hoeffding's Inequality)

Let $X_1, \ldots, X_n$ be i.i.d random variables in $[0,1]$, and let $\bar{X} = \frac{1}{t} \sum_{j=1}^{t} X_j$. Then:

$$P\left[ E[X] > \bar{X}_t + u \right] \leq e^{-2tu^2}$$

Condition on selecting $a$:

$$P\left[ q(a) > Q_t(a) + U_t(a) \right] \leq e^{-2N_t(a)U_t(a)^2}$$

Pick $p$ :   $e^{-2N_t(a)U_t(a)^2} = p$

Solve $U$:   $U_t(a) = \sqrt{\dfrac{-\log p}{2N_t(a)}}$

Reduce $p$, e.g. $p = t^{-4}$

Ensures we select opt. act. $t \to \infty$: $U_t(a) = \sqrt{\dfrac{2\log t}{N_t(a)}}$

UCB1 algo: $A_t = \underset{a \in \mathcal{A}}{argmax} \; Q_t(a) + \sqrt{\dfrac{2 \log t}{N_t(a)}}$

Theor. $\underset{t \to \infty}{\lim} L_t \leq 8 \log t \underset{a \mid \Delta_a > 0}{\sum} \Delta_a$

# Bayesian Bandits

Exploits prior knowledge.

e.g. Gaussians

## Probability matching

$\pi(a) = \mathbb{P}\left[ Q(a) = \underset{a}{max} \; Q(a') \mid R_1, \dots, R_{t-1} \right]$

( select action $a$ according to $\pi(a)$ )

## Thompson sampling:

$\pi(a) = \mathbb{E}\left[ \mathbb{1}(Q(a) = \underset{a'}{max} \; Q(a')) \mid R_1, \dots, R_{t-1} \right]$

Sample from posterior and select max.

# Value of Information

Information State Space

$\tilde{S}$ — summary of all information

Action $A$ causes a transition to a new

$\tilde{S}'$ with prob. $\tilde{P}^A_{\tilde{S}\tilde{S}'}$

MDP $\tilde{\mathcal{M}} = \langle \tilde{S}, \mathcal{A}, \tilde{P}, \mathcal{R}, \gamma \rangle$

In Bernoulli case, MDP can be solved by DP. The solution is Gittins Index.
→ e.g. drug problem (work / don't work).

## Contextual Bandits.

tuple $\langle A, S, R \rangle$

## MDPs

- for unknown or poorly est. state, replace reward function with $r_{max}$.

- $\tilde{S} = \langle S, I \rangle$ accumulated information

## L10. Classic Games.

<u>best response</u> $\pi_*^i(\pi^{-i})$ is optimal policy against all other opponents' fixed policies

<u>Nash equilibrium</u>: $\pi^i = \pi_*^i(\pi^{-i})$

→ every player's policy is a best response
→ if found, game is solved.

Nash equilibrium is fixed-point of self-play RL

For general games, Nash equil. isn't unique, but we'll look at classic games where its unique.

Two-Player Zero-Sum Games

White & Black $R^1 + R^2 = 0$
↳ rewards

Methods of finding Nash eq.:
- Game tree search
- Self-play RL

Perfect information: all visible

Imperfect information: not all visible

# Minimax Search

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

$$V_*(s) = \max_{\pi_1} \min_{\pi_2} V_\pi(s)$$

minimax policy $\pi = \langle \pi_1, \pi_2 \rangle$ is a Nash eq.

instead, we use value function $v(s,w) \approx V_*(s)$

Chinook solved checkers in 2007

We can apply RL algorithms (MC, TD(0), TD($\lambda$)) to games by making them self-play.

Policy improvement w/ afterstates

$$q_*(s,a) = V_*(succ(s,a))$$

$$A_t = \arg\max_a V_*(succ(S_t, a)) \qquad \text{for white}$$

$$A_t = \arg\min_a -\;/\;-\;/\;-\;/\;- \qquad \text{for black}$$

TD performed poorly in chess, checkers because of tactical nature of games.

## TD Root

update value towards sucessor search value

$$V_+(S_t, w) = \min_{s \in leaves(S_t)} \max v(s, w)$$

## TD Leaf

update search value towards sucessor search value

## Tree Strap

— / — / — / deeper — / —

## Simulation-Based Search

UCT algo: MC + UCB algorithm

→ Smooth UCT Search
→ Imperfect information games