

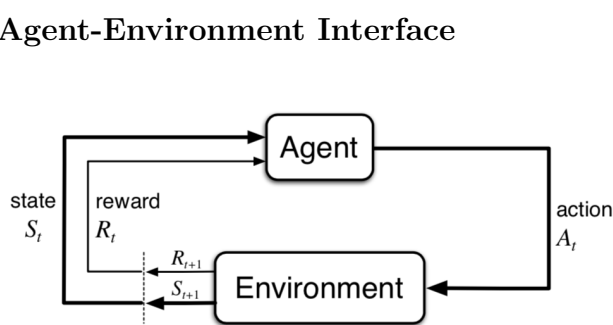
Reinforcement Learning (RL)

Concepts (no state)

Symbols

a action	λ trace decay ($0 \sim 1$)
A advantage	Pr transition distribution
α learning rate, step size	π policy ($state \rightarrow action$)
b bandit, baseline estimate	q quality
c cost, context (= state)	r reward
d difference error	R return
Δ difference	s state
D replay buffer	t time
E eligibility trace	τ trajectory
ϵ exploration rate ($0 \sim 1$)	θ parameter weight
Φ state transition function	v value
g gain	w weight
γ discount factor	y target / prediction
∇ gradient	$*$ optimal
H horizon	$'$ next / derivative
J expected return	$\hat{}$ estimation
L loss / regret	$\ \cdot \ $ vector norm

Agent-Environment Interface



General model

1. the *agent* at each time step t receives a representation of the environment's *state* $S_t \in S$
2. it selects an action $A_t \in A(s)$
3. from its action the agent receives a *reward* $R_{t+1} \in R \in \mathbb{R}$

General challenges

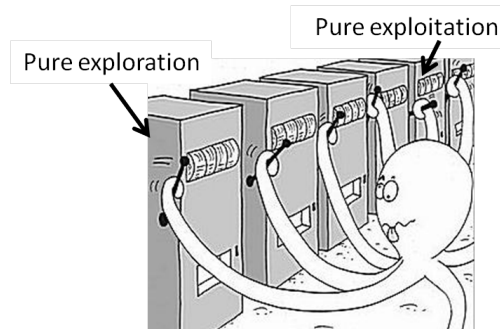
- choosing the right *Representation* of the problem / state
- *Generalization* beyond the cases trained on
- *Temporal Credit Assignment*: what caused this outcome?
- balance *Exploitation* vs *Exploration* (short vs long term)

Greedy

- 100% exploitation
- wants to pick best arm but won't know enough
- if 1st arm ok won't try others (stuck in local optimum)

Bandit algorithms

- *multi-armed bandit problem*: which slot arm gives most?
- no notion of state!
- this leaves just *exploration* vs *exploitation*



Bandits = heuristics. Image from research.microsoft.com

$$\begin{aligned}
 \underbrace{n_a}_{\text{action's use count}} &= \sum_{t=0}^T \underbrace{\begin{cases} 1 & a_t = a \\ 0 & \text{otherwise} \end{cases}}_{\text{tries for action so far}} \\
 \underbrace{\hat{r}_a}_{\text{action's expected reward}} &= \underbrace{\sum_t \frac{r_t}{n_a}}_{\text{average action reward so far}}
 \end{aligned}$$

(1)

Optimistic-Greedy

- like greedy but large initial \hat{r}_a
- + tries all
- no 2nd chance for arms unlucky on first try

ϵ -Greedy

- at probability ϵ pick randomly
- + not fooled by unlucky first try
- linear regret from constant exploration

Upper Confidence Bound (UCB)

try all for k rounds, then

$$a_t = \operatorname{argmax}_{a \in A} \left[\underbrace{\hat{r}_a}_{\text{reward}} + \underbrace{\sqrt{\frac{2 \log t}{n_a}}}_{\text{under-explored bonus}} \right] \quad (4)$$

- explore then gradually exploit
- + logarithmic regret

Contextual Bandit

Linear UCB (LinUCB)

$$\begin{aligned}
 x_{t,a} &= \Phi(s_{t,a}) \\
 \mathbb{E}[r_{t,a} | x_{t,a}] &= \theta_a \cdot x_{t,a} \\
 \theta_a &= A^{-1}b
 \end{aligned} \quad (5)$$

(2)

Posterior / Thompson sampling

Pick actions by probability they maximize expected reward

Greedy in the Limit with Infinite Exploration (GLIE)

- infinitely explores
- converges on greedy

Markov Decision Process (MDP)

a 5-tuple (S, A, P, R, γ)

$$\begin{aligned}
 & \underbrace{s}_{\text{state}} \\
 & \underbrace{\pi_t(a|s)}_{\text{policy}} \\
 & p(s'|s, a) \\
 & r(s', s, a) \\
 & \underbrace{v_\pi(s)}_{\text{value function}} \\
 & \underbrace{q_\pi(s, a)}_{\text{action-value (Q) function}} \\
 & \underbrace{q_*(s, a)}_{\text{optimal Q function}} \\
 & \underbrace{v_*(s)}_{\text{state value under optimal policy}}
 \end{aligned}
 \in
 \begin{aligned}
 & \underbrace{S}_{\text{finite set of states}} \\
 & : \text{a mapping from a state to an action} \\
 & = \underbrace{Pr\{S_{t+1} = s' | S_t = s, A_t = a\}}_{\text{state transition probabilities}} \\
 & = \underbrace{\mathbb{E}[R_{t+1} | S_{t+1} = s', S_t = s, A_t = a]}_{\text{expected return for state-action-nextstate}} \\
 & = \underbrace{\mathbb{E}_\pi[G_t | S_t = s]}_{\text{expected gain under policy } \pi \text{ in state } s} \\
 & = \underbrace{\sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a)[r + \gamma v_\pi(s')]}_{\text{recursive definition (Bellman equation)}} \\
 & = \underbrace{\mathbb{E}_\pi[G_t | S_t = s, A_t = a]}_{\text{expected reward under policy } \pi \text{ for a state-action pair}} \\
 & = \underbrace{\sum_{s', r} p(s', r|s, a)[r + \gamma V_\pi(s')]}_{\text{recursive definition (Bellman equation)}} \\
 & = \max_a q_\pi(s, a) \\
 & = \max_a v_\pi(s) \\
 & = \max_{a \in A(s)} q_{\pi_*}(s, a) \\
 & \quad \text{expected return from its best action}
 \end{aligned}
 \tag{6}$$

Contraction Mapping

For metric space (X, d) and $f : X \rightarrow X$, f is a *contraction* given a *Lipschitz coefficient* $k \in [0, 1)$ where for all x / y in X :

$$d(f(x), f(y)) \leq kd(x, y) \tag{7}$$

Contraction Mapping theorem

For complete metric space (X, d) and contraction $f : X \rightarrow X$, there is only 1 fixed point x^* where $f(x^*) = x^*$.

For point x in X , and $f^n(x)$ inductively defined by $f^n(x) = f(f^{n-1}(x))$, $f^n(x) \rightarrow x^*$ as $n \rightarrow \infty$, yielding a unique optimal solution for DP.

Model-based Methods (known MDP)

Exhaustive Search

brute force, usually computationally unviable.

Dynamic Programming (DP)

- + bootstrap (learn mid-episode)

Find π_* for V / Q :

Policy Iteration

Initialize $V(s) \in \mathbb{R}$.e.g. 0, $\Delta \leftarrow 0$, $\pi(s) \in A$ for all $s \in S$

```

1. Policy Evaluation
while  $\Delta < \theta$  (e.g. 0.001) do
  foreach  $s \in S$  do
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  end
end

2. Policy Improvement
policy-stable  $\leftarrow$  true
while not policy-stable do
  foreach  $s \in S$  do
    old-action  $\leftarrow \pi(s)$ 
     $\pi(s) \leftarrow \arg\max_a \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$ 
    policy-stable  $\leftarrow$  old-action  $\neq \pi(s)$ 
  end
end

```

Policy iteration methods:

- gradient-based (policy gradient methods): gradient ascent
- gradient-free: simulated annealing, cross-entropy search or methods of evolutionary computation
- value search/iteration: stop after 1 state sweep
- async DP: update iteratively, no full sweeps
- generated policy iteration (GPI)

Value Iteration

ditch $V(s)$ convergence for policy improvement and truncated policy eval step in 1 operation:

```

Initialize  $V(s) \in \mathbb{R}$ .e.g. 0,  $\Delta \leftarrow 0$ 
while  $\Delta < \theta$  (e.g. 0.001) do
  foreach  $s \in S$  do
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \max_a \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  end
end

output: deterministic policy  $\pi \approx \pi_*$  where
 $\pi(s) = \arg\max_a \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$ 

```

Model-free methods

Monte Carlo (MC) Methods

- uses **averaging sample returns** per state-action pair
- episodic
- + sampling

Initialize for all $s \in S, a \in A(s)$:

$Q(s, a) \leftarrow$ arbitrary
 $\pi(s) \leftarrow$ arbitrary
 $Returns(s, a) \leftarrow$ empty list

```

while forever do
  Pick  $S_0 \in S$  and  $A_0 \in A(S_0)$ , all  $p(s, a) > 0$ 
  Generate an episode starting at  $S_0, A_0$  following  $\pi$ 
  foreach pair  $s, a$  in the episode do
     $G \leftarrow$  return for first occurrence of  $s, a$ 
    Append  $G$  to  $Returns(s, a)$ 
     $Q(s, a) \leftarrow average(Returns(s, a))$ 
  end
  foreach  $s$  in the episode do
     $\pi(s) \leftarrow \arg\max_a Q(s, a)$ 
  end
end

```

estimate for non-stationary problems:

$$V(S_t) + = \alpha[G_t - V(S_t)] \tag{8}$$

for learning rate α , how much we want to forget about past experiences.

Temporal Difference (TD)

- + DP's bootstrap
- + MC's sampling
- substitutes expected discounted reward G_t from the episode with an estimation:

$$V(S_t) + = \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \tag{9}$$

State-action-reward-state-action (SARSA)

- on-policy TD control
- can use priors

Initialize $Q(s, a)$ arbitrarily and
 $Q(\text{terminal} - \text{state},) = 0$

```

foreach episode  $\in$  episodes do
  Pick  $a$  from  $s$  by policy from  $Q$  (e.g.  $\epsilon$ -greedy)
  while  $s$  is not terminal do
    Take action  $a$ , observe  $r, s'$ 
    Pick  $a'$  from  $s'$  by policy from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $y \leftarrow r + \gamma Q(s', a')$ 
     $Q(s, a) + = \alpha[y - Q(s, a)]$ 
     $s \leftarrow s'$ 
     $a \leftarrow a'$ 
  end
end

```

n -step Sarsa (n -step TD) for n -step Q-Return:

$$\begin{aligned}
 q_t^{(n)} &= \gamma^n Q(S_{t+n}) + \sum_{i=1}^n \gamma^{i-1} R_{t+i} \\
 Q(s_t, a_t) &+ = \alpha [q_t^{(n)} - Q(s_t, a_t)]
 \end{aligned}
 \tag{10}$$

Forward View Sarsa(λ) (/ TD(λ))

$$\begin{aligned}
 q_t^\lambda &= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)} \\
 Q(s_t, a_t) &+ = \alpha [q_t^\lambda - Q(s_t, a_t)]
 \end{aligned}
 \tag{11}$$

Backward View Sarsa(λ) (/ TD(λ))

- + more efficient
- + can update at every time-step
- eligibility traces : find cause in frequency vs recency

$$\begin{aligned} E_0(s, a) &= 0 \\ E_t(s, a) &= \gamma \lambda E_{t-1}(s, a) + \mathbf{1}(S_t = s, A_t = a) \\ \delta_t &= R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \\ Q(s, a) &+ = \alpha \delta_t E_t(s, a) \end{aligned} \quad (12)$$

Linear Function Approximation

- + efficient
- + generalize

update temporal difference error, minimize squared loss:

$$\begin{aligned} \delta &\leftarrow r_t + \gamma \theta^T \Phi(s_{t+1}) - \theta^T \Phi(s_t) \\ \theta &+ = \alpha \delta \Phi(s_t) \\ J(\theta) &= ||\delta||^2 \end{aligned} \quad (13)$$

Q Learning

$$\begin{aligned} \delta &\leftarrow r_t + \gamma \underset{a \in A}{\operatorname{argmax}} Q(\Phi(s_{t+1}, a); \theta_i^-) - Q(\Phi(s_t, a); \theta_i) \\ Q(s, a) &+ = \alpha [r + \gamma \underset{a' \in A}{\operatorname{max}} Q(s', a') - Q(s, a)] \end{aligned} \quad (14)$$

Replay Memory

update θ by SGD

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim D} \delta \nabla_{\theta_i} Q(\Phi(s_t, a_t); \theta) \quad (15)$$

Deep Q Learning (DQL)

Made by *DeepMind*, uses a deep neural net (*Q-network*) for the *Q* function. Keeps *N* observations in a *memory* to train on.

$$\begin{aligned} y &= r_t + \gamma \underset{a \in A}{\operatorname{max}} Q(\Phi(s_{t+1}, a); \theta_i^-) \\ \underbrace{\nabla_i J(\theta_i)}_{\text{loss function}} &= \mathbb{E}_{(s, a, r, s') \sim U(D)} [\underbrace{r}_{\text{memory}} + \underbrace{\gamma Q(s', a')}_{\text{target}} - \underbrace{Q(s, a; \theta_i)}_{\text{prediction}}]^2 \end{aligned} \quad (16)$$

for network weights θ and experience replay history $U(D)$.

```
Initialize replay memory D with capacity N
Initialize Q(s, a) arbitrarily
foreach episode  $\in$  episodes do
    Pick a from s by policy from Q (e.g.  $\epsilon$ -greedy)
    while s is not terminal do
        Take action a, observer r, s'
        Store transition (s, a, r, s') in D
        Sample random transitions from D
         $y_i \leftarrow \begin{cases} r_j & \text{for terminal } s'_j \\ r_j + \gamma \max_a Q(s', a'; \theta) & \text{otherwise} \end{cases}$ 
        Perform gradient descent step on  $(y_j - Q(s_j, a_j; \theta))^2$ 
         $s \leftarrow s'$ 
    end
end
```

Prioritized Replay Memory learn esp. from high loss (traumas)

$$p(s_t, a_t, r_t, s_{t+1}) \propto r_t + \gamma \max_{a \in A} Q(\Phi(s_{t+1}, a); \theta_i^-) \quad (17)$$

Double DQN

solve bias from reused θ , faster

$$y = r_t + \gamma Q(\Phi(s_{t+1}, \underset{a \in A}{\operatorname{argmax}} Q(\Phi(s_{t+1}, a); \theta_i)); \theta_i^-) \quad (18)$$

Direct Policy Search

learn π .

Policy Gradient

- + simpler than *Q* or *V*
- + allows stochastic policy (rock-paper-scissors)
- local optimum
- less sample efficient

$$L = \mu(\log p(a|s) \cdot R(s))$$

Likelihood Ratio

$$\begin{aligned} R(\tau) &= \sum_{t=0}^T R(s_t, a_t) \\ \text{return state-action trajectory:} \\ \text{expected return:} \\ J(\theta) &= \mathbb{E}[\sum_{t=0}^T R(s_t, a_t); \pi_\theta] \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \\ \text{find } \theta \text{ to max:} \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \\ \text{yielding:} \\ \max_{\theta} J(\theta) &= \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ \text{gradient chasing reward:} \\ \nabla_{\theta} J(\theta) &= \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau) \\ \text{samled over m trajectories:} \\ &= \frac{1}{m} \sum_{i=1}^m P(\tau; \theta) \nabla_{\theta} \log P(\tau_i; \theta) R(\tau_i) \\ \text{combined:} \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau_i) \end{aligned} \quad (20)$$

REINFORCE

- check policy/episode's states/actions/rewards
- calc episode return for collected rewards
- update model params toward policy gradient

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R \quad (21)$$

desired loss function:

$$\frac{1}{m} \sum_{t=1}^m \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R \quad (22)$$

Baselined REINFORCE

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \underbrace{(R - V_{\Phi}(s_t))}_{\text{baselined reward}} \quad (23)$$

Actor-critic

- + less variance than Baselined REINFORCE
- actor (makes policy): policy gradient
- critic: policy iteration

$$\underbrace{A_{\pi}(s, a)}_{\text{advantage}} = Q(s, a) - V(s) \quad (24)$$

Asynchronous Advantage Actor Critic (A3C)

- 5-step Q-Value estimation
- shares params between actor/critic
- + run parallel, one policy
- + no more need for DQN's replay policy

$$(19) \quad \text{github.com/tycho01/Reinforcement-Learning-Cheat-Sheet}$$