

Intro:

A RL agent can optionally consist 3 things:

- 1) Model { transitional Probability - deterministic or stochastic $P(S'|S, a)$
rewards - deterministic or stochastic $R(S)$ or $R(S, a)$

- The agent may not have access to model (how world works). This is the distinction between model-based and model-free agent.

- 2) Policy $\pi: S \rightarrow A$ - deterministic or stochastic

3) Value function $V_t(S) : E[G_t(S)] = E\left[\sum_{i=t}^{H-1} \gamma^{i-t} r_i\right]$

Horizon (H): length of the episode { $H = \infty$ infinite, no subscript t
 $H < \infty$ finite, have subscript t .

- For $H < \infty$, P, R, π, V can be stationary or non-stationary (t -dependent)

- For $H = \infty$, P, R, π, V are all stationary.

If V is stationary, $V(S) = V_0(S) = V_1(S) = \dots = V_t(S)$ for any t .

Model-Based Agents:

Markov process: (S, P) $P = P(S'|S)$ the key is memoryless

Markov Reward process: (S, P, R, γ) $P = P(S'|S)$

$$H = \infty: V(S) = R(S) + \gamma \sum_{S' \in S} P(S'|S) V(S')$$

- Analytical soln: $V = R + \gamma PV \rightarrow V = (I - \gamma P)^{-1} R \rightarrow O(|S|^3)$

- Iterative soln: $V_{\text{new}} = R + \gamma PV$ until $\|V_{\text{new}} - V\|_\infty < \epsilon$

$$H < \infty: V_t(S) = R(S) + \gamma \sum_{S' \in S} P(S'|S) V_{t+1}(S') \quad \text{w/ stationary } P, R$$

- DP soln: $V_H(S) = 0$ for all S , compute $V_t(S)$ backwards.

Markov Decision Process: (S, A, P, R, γ) $P = P(S'|S, a)$

$$V_t^{\pi}(s) = E_{\pi} [G_t | S_t = s] \text{ for } t < \infty \quad V^{\pi}(s) = V_0^{\pi}(s) \text{ for } H = \infty.$$

$$Q_t^{\pi}(s, a) = E_{\pi} [G_t | S_t = s, A_t = a] \text{ for } H < \infty \quad Q^{\pi}(s, a) = Q_0^{\pi}(s, a) \text{ for } H = \infty$$

$V_t^{\pi}(s)$ is $E[G_t]$ when taking actions according to π at $i \geq t$

$Q_t^{\pi}(s, a)$ is $E[G_t]$ when $A_t = a$ and taking actions according to π at $i \geq t+1$.

$$Q^{\pi}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^{\pi}(s')$$

Two algorithms to find π^* for MDP: (for both deterministic and stochastic)

↳ Policy Iteration : Compute optimal value and policy iteratively

- based on monotonic improvement in π $V^{\text{init}}(s) \geq V^{\pi}(s)$

def policy_iteration(P, R, γ, ε):

while ($\pi^* \neq \pi$). any():

$$\pi = \pi^*$$

$$V^{\pi} = \text{policy_eval}(P, R, \pi, \gamma, \epsilon)$$

$$\pi^* = \text{policy_improve}(P, R, V^{\pi}, \gamma)$$

$$\text{return } V^{\pi}, \pi$$

O(1S1) def policy_eval(P, R, π, γ, ε): * similar to iterative soln for MRP

$$R^{\pi}(s) = \sum_{a \in A} \pi(a|s) R(s, a)$$

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

while $\|V - V_{\text{old}}\|_{\infty} \geq \epsilon$:

$$V_{\text{old}} = V$$

$$V(s) = R^{\pi}(s) + \gamma \sum_{s' \in S} P^{\pi}(s'|s) V(s') \quad \forall s \in S$$

$$\text{return } V$$

O(1S1) def policy_improve(P, R, V^π, γ):

$$\hat{\pi}(s) = \underset{a \in A}{\text{argmax}} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^{\pi}(s')] \quad \forall s \in S$$

$$\text{return } \hat{\pi}(s)$$

- `policy_improve()` is based on $\pi^*(s) = \arg\max_{a \in A} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^*(s')]$

2, Value iteration: compute value function iteratively, then find π^*

- based on $V^*(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^*(s')] = V^{\pi^*}$

def `value_iteration(P, R, γ, ε)`:

while $\|V - V_{old}\| > \epsilon$:

$$V_{old} = V$$

$$V(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V(s')] \quad \forall s \in S$$

$$\pi^* = \text{policy_improve}(P, R, V, \gamma)$$

return V, π^*

Above function are all for $H = \infty$, thus, V is stationary. For $H < \infty$ and non-stationary V_t , we could use DP instead of iterative case.

e.g. def `value_iteration(P, R, γ, H)`:

$$V_H^*(s) = 0 \quad \forall s \in S$$

for t in range($H, -1, -1$):

$$V_t^*(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V_{t+1}^*(s)] \quad \forall s \in S$$

$$\pi_t^*(s) = \arg\max_{a \in A} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V_{t+1}^*(s)] \quad \forall s \in S$$

return V_t^* and π_t^* $\forall t = 0..H$ and $s \in S$

Model-free Agents:

Policy Evaluation: 2 methods : Monte-Carlo (MC) or Temporal Difference (TD(0))

1, MC : first-visit (update only for 1st visit to s in each episode) vs. every-visit

$$V^{\pi}(s_{j,t}) = V^{\pi}(s_{j,t}) + \alpha (G_{j,t} - V^{\pi}(s_{j,t})) \quad s_{j,t} = j^{\text{th}} \text{ episode}, t^{\text{th}} \text{ time step.}$$

$$\alpha = \frac{1}{N(s_{j,t})} \quad \text{normally, } \alpha > \frac{1}{N(s_{j,t})} \quad \text{if more weight for recent occurrence.}$$

Properties: episodic only, unbiased (first-visit) and biased (every-visit), data efficient
high variance, need to wait for the end of episode to update.

\Rightarrow TD(0): Combination of MC and DP (bootstrapping)

Instead of G_t , we use $r_t + \gamma V^\pi(S_{t+1})$ as the approximation of G_t

Properties: works for non-episodic environment too, biased, Markov assumption, low variance, can update as soon as a state is visited, obviously smaller update compared to MC since it performs many updates in 1 episodes.

- Batch MC and TD(0): we can think of this as mini-batch whereas normal MC and TD(0) as SGD. We accumulate several episode of history, find the avg then update.

$$\hat{P}(S' | S, a) = \frac{1}{N(S, a)} \sum_{j=1}^n \sum_{t=1}^{T_j-1} \mathbb{1}(S_{j,t} = S, a_{j,t} = a, S_{j,t+1} = S')$$

$$\hat{r}(S, a) = \frac{1}{N(S, a)} \sum_{j=1}^n \sum_{t=1}^{T_j-1} \mathbb{1}(S_{j,t} = S, a_{j,t} = a) r_{j,t}$$

Control: (Policy optimization)

- In control, we work w/ $Q(s, a)$ instead of $V(s)$, we can update $Q(s, a)$ in the exact same ways as updating $V(s)$ above (in policy evaluation).

- MC: $Q^\pi(s_t, a_t) = \text{avg of } G_t (s = s_t, a = a_t)$

$$\left. \begin{array}{l} \text{- SARSA: } Q^\pi(s_t, a_t) = Q^\pi(s_t, a_t) + \alpha_t [r_t + \gamma Q^\pi(s_{t+1}, a_{t+1}) - Q^\pi(s_t, a_t)] \\ \text{- Q-Learning: } Q^\pi(s_t, a_t) = Q^\pi(s_t, a_t) + \alpha_t [r_t + \gamma \max_{a'} Q^\pi(s_{t+1}, a') - Q^\pi(s_t, a_t)] \end{array} \right\} \text{TD(0)}$$

- Since we don't know the model, we can't simply do exploitation every time, we need to do exploration $\rightarrow \epsilon$ -greedy algorithm: $\pi_\epsilon(s) = \begin{cases} \text{random } a \text{ w/ prob } \frac{\epsilon}{|A|} \\ \underset{a \in A}{\operatorname{argmax}} Q^\pi(s, a) \text{ w/ prob } 1 - \epsilon. \end{cases}$

- ϵ -greedy algorithm is used as the policy improvement.

Off-policy Evaluation & Control: When it's expensive to apply a policy.

- On-policy = apply π_1 , update Q_1^π - off-policy = apply π_2 , update Q_2^π ,

- Importance Sampling : say we want $E_{\pi^*}[f(x)]$ but we can only get $E_{\pi^{\text{exp}}}[f(x)]$

in RL setting, $\pi = \pi_t$, $P = \pi_t$, $x = \text{state}$ $f(x) = V^{\pi_t}(s)$ or $Q^{\pi_t}(s)$

$$E_{\pi^*}[f(x)] = E_{\pi^{\text{exp}}}\left[\frac{\pi(x)}{P(x)} f(x)\right]$$

$$\left\{ \begin{array}{l} \text{MC: } V^{\pi_t}(s) \approx \frac{1}{n} \sum_{j=1}^n \frac{P(h_j | \pi_t, s)}{P(h_j | \pi^*, s)} G(h_j) = \frac{1}{n} \sum_{j=1}^n G(h_j) \frac{\prod_{t=1}^{T_j-1} \pi_t(a_{j,t} | s_{j,t})}{\prod_{t=1}^{T_j} \pi^*(a_{j,t} | s_{j,t})} \\ \text{where } P(h_j | \pi_t, s) = \prod_{t=1}^{T_j-1} \pi(a_{j,t} | s_{j,t}) P(r_{j,t} | s_{j,t}, a_{j,t}) P(s_{j,t+1} | s_{j,t}, a_{j,t}) \\ \text{TD(0): } V^{\pi_t}(s) \approx V^{\pi_t}(s) + \alpha \left[\frac{\pi_t(a|s)}{\pi^*(a|s)} (r + \gamma V^{\pi_t}(s') - V^{\pi_t}(s)) \right] \end{array} \right.$$

- Q-learning is an off-policy control because $a_{t+1} = \arg\max_{a'} Q^{\pi_t}(s_{t+1}, a')$ instead of $a_{t+1} = \pi_t(s_{t+1})$

Maximization Bias: when policy is not biased but we are systematically biased (+ve) wrt $V^*(s)$

- Use Double Q-learning (significant speed up during training) $\hat{V}(s) \geq V^*(s)$

Have 2 Q functions, during training update one of Q using the $\arg\max_{a'} Q^*$ of the other Q.
take ϵ -greedy wrt $Q_1 + Q_2$.

Value Function Approximation (VFA): $\hat{V}(s, w) \approx V(s)$ $\hat{Q}(s, a, w) \approx Q(s, a, w)$

- For large/infinite state space, impossible to use tabular approach.

- Typical functions: linear, decision tree, nearest neighbors, neural network ...

- Typically a regression problem, loss function $J(w) = E_{\pi} [(Q_{\pi}(s, a) - \hat{Q}(s, a, w))^2]$

optimize through gradient descent, $w = w + \alpha \Delta w$

- MC: $\Delta w = (G_t(s) - \hat{Q}_t(s, w)) \nabla_w \hat{Q}(s, a, w)$

- SARSA: $\Delta w = (r + \gamma \hat{Q}(s', a', w) - \hat{Q}(s, a, w)) \nabla_w \hat{Q}(s, a, w)$

- Q-learning: $\Delta w = (r + \gamma \max_{a'} \hat{Q}(s', a', w) - \hat{Q}(s, a, w)) \nabla_w \hat{Q}(s, a, w)$

Linear Feature Representation: $\hat{Q}(s, a, w) = w^T x(s) = \sum_{j=1}^{\infty} w_j x_j(s) \rightarrow \nabla_w \hat{Q}(s, a, w) = x(s)$

Neural Network: $\nabla_w \hat{Q}(s, a, w)$ can come from backpropagation

For linear VFA, MC and SARSA are guaranteed to converge to the min MSE.

For non-linear VFA, we don't have guarantee for any of the 3 methods.

Deep Q-learning :

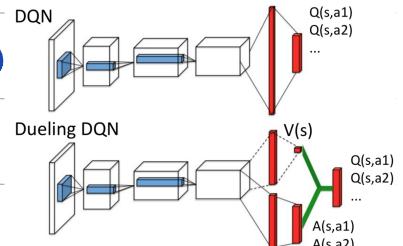
- Deep Q-Network (DQN) : similar to \hat{q} -learning except w/ the added CNN feature extractor
- In Atari paper, besides the preprocessing steps, 2 additional tricks are proposed.
 - Experience Replay : randomly pick a batch of (s_t, a_t, r_t, s_{t+1}) from memory
 - Greater data efficiency , less correlation , less oscillation / divergence.
 - Can use prioritized replay to more oftenly update problematic (s, a, s')
 - Target Network : keep 2 sets of weights . w^- is the target , update very C steps to w . w is update at each step .
 - Avoid "chasing a moving target" \rightarrow hyperparameter C .

- Algorithm 1 deep Q-learning

- 1: Initialize replay memory D with a fixed capacity
- 2: Initialize action-value function \hat{q} with random weights w
- 3: Initialize target action-value function \hat{q} with weights $w^- = w$
- 4: **for** episode $m = 1, \dots, M$ **do**
- 5: Observe initial frame x_1 and preprocess frame to get state s_1
- 6: **for** time step $t = 1, \dots, T$ **do**
- 7: Select action $a_t = \begin{cases} \text{random action} & \text{with probability } \epsilon \\ \arg \max_a \hat{q}(s_t, a, w) & \text{otherwise} \end{cases}$
- 8: Execute action a_t in emulator and observe reward r_t and image x_{t+1}
- 9: Preprocess s_t, x_{t+1} to get s_{t+1} , and store transition (s_t, a_t, r_t, s_{t+1}) in D
- 10: Sample uniformly a random minibatch of N transitions (s_j, a_j, r_j, s_{j+1}) from D
- 11: Set $y_j = r_j$ if episode ends at step $j + 1$;
- 12: otherwise set $y_j = r_j + \gamma \max_{a'} \hat{q}(s_{j+1}, a', w^-)$ \leftarrow wrt w^-
- 13: Perform a stochastic gradient descent step on $J(w) = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{q}(s_j, a_j, w))^2$ w.r.t. parameters w
- 14: Every C steps reset $w^- = w$

- Double DQN: similar to double \hat{q} -learning

- DQN target : $r + \gamma \max_{a'} \hat{q}(s', a', w^-)$ \leftarrow for picking policy
- DDQN target : $r + \gamma \hat{q}(s', \arg \max_{a'} \hat{q}(s', a', w), w^-)$ \leftarrow for evaluating policy .
- Dual DQN : $\hat{q}(s, a, w, w_v, w_A) = \hat{V}(s, w, w_v) + \hat{A}(s, a, w, w_A)$
- Advantage function : $A^{\pi}(s, a) = V^{\pi}(s) - Q^{\pi}(s, a)$
- Measure each action's effect on the value function .



$$- A(s, \alpha^*) = 0 \quad A(s, \alpha) \geq 0 \text{ for } \alpha \neq \alpha^*$$

- To make Q-function identifiable:

$$\hat{q} = \hat{V}(s, w, w_v) + (\hat{A}(s, \alpha, w, w_A) - \max_{\alpha' \in A} A(s, \alpha', w, w_A)) \quad \text{or}$$

$$\hat{q} = \hat{V}(s, w, w_v) + (\hat{A}(s, \alpha, w, w_A) - \frac{1}{|A|} \sum_{\alpha'} A(s, \alpha', w, w_A))$$

Imitation Learning:

- we have an human to give a few demonstrations to the RL agent $\rightarrow \pi^*$ for limited states.
- $P(s'|s, a)$ may / may not be known, reward function is not known.

3 types of Imitation Learning:

1) Behavioural Cloning : like supervised learning \rightarrow Given s , predict a .

- drawbacks**
- In supervised learning, we assume data is i.i.d. In RL, data is highly correlated.
 - Agent doesn't know what to do when entered in an unseen state
 - Error is compounding in RL, unlike SL.

2) Inverse RL: recover the reward function from human's demonstrations.

$$\begin{aligned} & \text{- use function approximation, say } R(s) = w^T x(s) \quad V^\pi(s) = E_\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right] \\ & \text{find } w \text{ s.t. } E_{\pi^*} \left[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid s_0 = s \right] \geq E_\pi \left[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid s_0 = s \right] \quad \forall \pi \\ & \text{a.k.a. } w^{*\top} M(\pi^* \mid s_0 = s) \geq w^{*\top} M(\pi \mid s_0 = s) \quad \forall \pi, \forall s \end{aligned}$$

3) Apprenticeship Learning: generate a good policy from the recovered reward

$$\| M(\pi \mid s_0 = s) - M(\pi^* \mid s_0 = s) \|_1 \leq \epsilon \quad \text{complicated, see notes}$$

↑ Value-based Approaches: find optimal value at each state, policy = picking a that produced max V

↓ Policy-based Approaches: find optimal policy directly by optimizing V wrt π .

- Policy-based approaches can be applied to stochastic policy. It's guaranteed to converge to local min.

Effective in high-dimensional / continuous state. Not sample-efficient because once policy is updated, we can't use the previous episodes \rightarrow can be solved by Importance Sampling.

Policy-Based Approaches:

Given a trajectory $T : (S_0, a_0, \dots, S_T, a_T)$, probability of it under policy π_θ is

$$\pi_\theta(T) = P(S_0) \sum_{t=1}^T \pi_\theta(a_t | S_t) P(S_{t+1} | S_t, a_t) \quad \text{where } \pi_\theta(a_t | S_t) = P(a_t | S_t; \theta)$$

objective function is the value function : $\theta^* = \underset{\theta}{\operatorname{argmax}} \left[E_{T \sim \pi_\theta(T)} [V^{\pi_\theta}(S_0)] \right]$

- Finite horizon: $J(\theta) = V^{\pi_\theta}(S_0) = E_{\pi_\theta}[S_0]$ (fixed initial state)

$$J(\theta) = \sum_S d(S) E_{\pi_\theta}[S] \quad (\text{stochastic initial state})$$

- Infinite horizon: $J(\theta) = \sum_S d(S) E_{\pi_\theta}[S]$

$$V(\theta) = \sum_I P(I; \theta) R(I) = \sum_I P(I; \theta) \sum_{t=1}^T \gamma^t r(S_t, a_t) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \gamma^t r(S_t^{(i)}, a_t^{(i)})$$

$$\nabla_\theta V(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left[\nabla_\theta \log \pi_\theta(a_t^{(i)} | S_t^{(i)}) \sum_{t'=t}^T \gamma^{t'} r(S_{t'}^{(i)}, a_{t'}^{(i)}) \right] \rightarrow \text{high variance}$$

Reinforce Algorithm: For $I^{(i)}$ in $i=1 \dots N$, For $t=1 \dots T$, $\theta = \theta + \alpha \nabla_\theta V(\theta)$

$$\text{Baseline: } \bar{V}_\theta V(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T [\bar{V}_\theta \log \pi_\theta(a_t^{(i)} | S_t^{(i)}) A_t^{(i)}] \rightarrow \text{lower variance.}$$

$$A_t^{(i)} = \sum_{t'=t}^T \gamma^{t'} r(S_{t'}^{(i)}, a_{t'}^{(i)}) - b(S_t) \leftarrow \text{doesn't affect result as long as } b = f(S_t) \text{ only.}$$

$$\text{A natural choice for } b(S_t) \text{ is } V(S_t) \approx \hat{V}(S_t, w)$$

Actor-Critic Algorithm: developed from baseline intuition \rightarrow learn θ and w simultaneously

$$-\text{Learn } w \text{ by minimizing } b(S_t, w) - G_t \rightarrow \min_w \sum_{i=1}^N \sum_{t=1}^T (b(S_t, w) - G_t)^2 \rightarrow \text{critic}$$

$$-\text{Learn } \theta \text{ by } \theta = \theta + \alpha \nabla_\theta V_\theta \text{ as shown above} \rightarrow \text{actor}$$

Off-Policy Policy Gradient: as we update the policy, we are under a different distribution.

$$J(\theta) = E_{T \sim \pi_\theta(T)} [r(T)] \rightarrow J(\theta') = E_{T \sim \pi_\theta(T)} \left[\frac{\pi_{\theta'}(T)}{\pi_\theta(T)} r(T) \right]$$

$$\nabla_\theta J(\theta') = E_{T \sim \pi_\theta(T)} \left[\sum_{t=1}^T \left(\nabla_\theta \log \pi_{\theta'}(a_t | S_t) \frac{\pi_\theta(a_t | S_t)}{\pi_{\theta'}(a_t | S_t)} \right) \left(\sum_{t'=t}^T \gamma^{t'} r(S_{t'}, a_{t'}) - b(S_t) \right) \right]$$

This allows us to use the replay memories.

N-step Estimator: The $J(\theta)$ and $\nabla_\theta J(\theta)$ above are based on MC, we can use TD as well.

$$\hat{G}_t^{(1)} = r_t + \gamma V(S_{t+1}) \rightarrow \text{high bias, low variance} \quad (\text{TD}(0))$$

$$\hat{G}_t^{(2)} = r_t + \gamma r_{t+1} + \gamma^2 V(S_{t+2}) \quad n\text{-step is a hyperparameter}$$

$$\hat{G}_t^{(00)} = G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \rightarrow \text{low bias, high variance.} \quad (\text{MC})$$

$$\text{Discrete action space : } \pi_\theta(a | s) = \frac{\exp[\theta^\top \phi(s, a)]}{\sum_a \exp[\theta^\top \phi(s, a)]} \quad \nabla_\theta \log \pi_\theta(a | s) = \phi(s, a) - \sum_a \pi_\theta(a | s) \phi(s, a)$$

Continuous action space : $a \sim N(\mu(s), \sigma^2)$ where $\mu(s) = \theta^T \phi(s)$ $\nabla_\theta \log \pi_\theta(a|s) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$

Trust Region Policy Optimization : solves the problem that it's hard to pick appropriate learning rate.

- Complicated, see note for details.

Exploration vs. Exploitation :

- Pure exploitation : stuck in suboptimal policy e.g. Greedy
- Pure exploration : never use learned information to make good decision.
- Computationally Efficient : Compute a decision in a reasonable amount of time e.g. A.V.
- Sample Efficient : Compute a decision w/ a reasonable amount of data e.g. medical
- Greedy, ϵ -Greedy are not good. Decay ϵ -greedy, being optimistic are good algorithms

Optimistic : choose action that may have the best upper bound $P(Q(a) > \hat{Q}(a) + u) < \delta$

- actions w/ high uncertainty (haven't been visited much) are often selected.
- Intuition : if action is the best action \rightarrow perfect, else \rightarrow we gained info and adjust mean and uncertainty.
- Optimistic initialization : assign all $Q(s,a)$ to be $\frac{R_{max}}{1-\gamma}$ (best it could be)

Multi-Armed Bandit (MAB) : simpler than MDP, only 1 state so everything is i.i.d.

- Evaluation metric : regret

$$\text{- Action value} : Q(a) = E[r|a] \quad \text{- optimal} : V^* = Q(a^*) = \max_{a \in A} Q(a)$$

$$\text{- Gap} \Delta_a = V^* - Q(a) \quad \text{- Regret} : L_t = E[V^* - Q(a_t)]$$

$$\text{- Objective} : \text{minimize total regret} : L_t = E\left[\sum_{t=1}^T V^* - Q(a_t)\right] = TV^* - E\left[\sum_{t=1}^T Q(a_t)\right]$$

Upper Confidence Bound (UCB) :

$$\text{- pick } a_t = \arg\max_{a \in A} \left\{ \hat{Q}_t(a) + \hat{U}_t(a) \right\} \quad \text{where} \quad \hat{U}_t(a) = \sqrt{\frac{-\log P}{2N_t(a)}} = \sqrt{\frac{2\log t}{N_t(a)}} \quad P = t^{-4}$$

$$\hat{Q}_t(a) = \frac{1}{N_t(a)} \sum_{t'=1}^T r_{t'} \mathbb{1}(a_t = a)$$

- $\hat{U}_t(a) = \sqrt{\frac{2\log t}{N_t(a)}}$ is based on Hoeffding's Inequality $P(E[X] > \hat{X}_t + u) \leq \exp(-2Nu^2)$

- The only assumption is that reward is bounded.

Bayesian Bandit :

- Assumes rewards follows a known distribution, e.g. beta distribution for Bernoulli event
- we could use UCB approach. E.g. For Gaussian, $\hat{a}_t = \arg \max_{a \in A} \left\{ \hat{u}_a + C \frac{\sigma_a}{\sqrt{N(a)}} \right\}$
- we could use sampling instead of the upper bound.

- Thompson Sampling: sample the reward for each action, pick the max reward and update distribution. For $\text{beta}(d, B)$, $d++$ if reward = 1 or $B++$ if reward = 0

Provably Approximately Correct (PAC):

- Minimizing regret doesn't help us to determine whether we make infrequent large mistakes or frequent small mistakes. The former is worse.
- Defn: An algorithm is PAC if $P[Q(a) \geq Q(a^*) - \epsilon] \geq 1 - \delta$ on all but a polynomial number, usually in terms of ϵ, δ, N .
- UCB, TS are PAC algorithms, Greed is not.
- Generalization from MAB to MDP involves model approximation, but most is similar.
 - There are different states, so $P(s'|s, a)$ needs to be approximated
 - Instead of approximating $Q(a)$, we need to approximate $Q(s, a)$.

Safe RL: Combination of off-policy Evaluation/Control and sample-efficient algorithms

- an application of Safe RL is the Medical Industry, we can't cheaply change treatment on patients (off-policy methods) and we want figure out optimal solution quickly (sample-efficient algorithms).
- Defn: Given historical data $D = \{s_i, a_i, r_i, \dots\}_{i=1}^{ID}$ and current policy π_b , can we find new policy π_e s.t. $P(V^{\pi_e} \geq V^{\pi_b}) \geq 1 - \delta$ where $V^{\pi_e} = \mathbb{E} \left[\sum_{t=1}^T \gamma^t R_t | \pi_e \right]$

Importance Sampling:

$$\text{Recall: } V^{\pi_e}(s) \approx \frac{1}{n} \sum_{i=1}^n \frac{p(h_i | \pi_e, s)}{p(h_i | \pi_b, s)} G(h_i) = \frac{1}{n} \sum_{i=1}^n \left(\frac{p(h_i | \pi_e, s)}{p(h_i | \pi_b, s)} \frac{\pi_e(a_i | s_i)}{\pi_b(a_i | s_i)} \right) \left(\sum_{t=1}^{l_i} \gamma^t R_t^{(i)} \right)$$

- Variations that achieves better results:

$$\text{- Per-Decision IS: } V^{\pi_e}(s) = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{l_i} \left(\gamma^t R_t^{(i)} \frac{\pi_e(a_t^{(i)} | s_t^{(i)})}{\pi_b(a_t^{(i)} | s_t^{(i)})} \right) \quad \begin{array}{l} \text{intuition: future action can't affect} \\ \text{past rewards.} \end{array}$$

- weighted IS: let $w_i = \frac{1}{\prod_{t=1}^L \pi_b(a_t^{(i)} | s_t^{(i)})}$, $V^{\text{re}}(s) = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i \left(\sum_{t=1}^L y_t^{(i)} R_t^{(i)} \right)$

- Doubly Robust IS: Combined raw return $R_t^{(i)}$ w/ $\hat{Q}^{\text{re}}(s_t^{(i)}, a_t^{(i)})$

- IS is not biased but high variance. Model-based estimate ($\hat{Q}(s, a)$) is biased but low var

- Let $X = R$, $Y = \hat{Q}^{\text{re}}$ and $E[Y] = \hat{V}^{\text{re}}$, $E[X - Y + E[Y]] = E[X]$

$$V^{\text{re}}(s) = \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^L y_t^{(i)} w_t^{(i)} (R_t^{(i)} - \hat{Q}^{\text{re}}(s_t^{(i)}, a_t^{(i)})) + y_t^{(i)} \hat{V}^{\text{re}}(s_t^{(i)}) \quad \text{where } w_t^{(i)} = \frac{1}{\prod_{t=1}^L \pi_b(a_t^{(i)} | s_t^{(i)})}$$

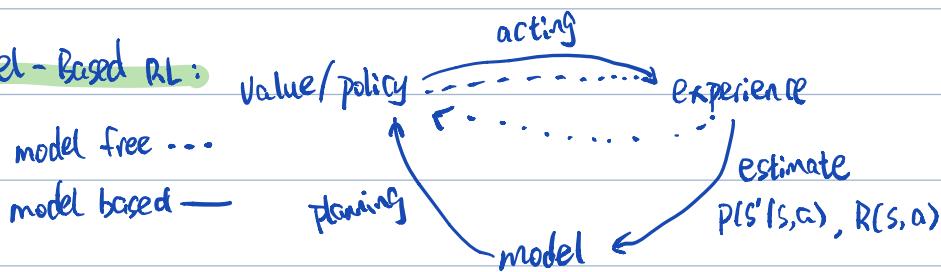
- Limitation: only works when $\pi_b(a|s) > 0$ when $\pi_b(a|s) > 0 \quad \forall s, a$

Hoeffding's Inequality:

$$E[X_i] \geq \frac{1}{n} \sum_{i=1}^n X_i - b \sqrt{\frac{\ln(1/\delta)}{2n}} \quad X_i = w_i \sum_{t=1}^L y_t^{(i)} R_t^{(i)} \quad w_i = \frac{1}{\prod_{t=1}^L \pi_b(a_t^{(i)} | s_t^{(i)})}$$

$b = \max_i (G_i w_i) \leftarrow \text{Could explode as } w_i \text{ could be very large if } \pi_b(a_t^{(i)} | s_t^{(i)}) \gg \pi_b(a_t^{(j)} | s_t^{(j)})$

Model-Based RL:



model-based is usually more sample efficient

- Just like estimate $\hat{Q}(s, a, w)$ and $\hat{\pi}_Q(s)$, we need to first define a function approximation (tabular, linear, NN...) and a loss function (MSE, KL Divergence...).

MC Tree Search: simulation-based approach for the planning stage.

- To avoid exponential growing # of nodes, we used simulation and we always use the current state as the root, so each time we are solving a sub MDP starting from the current state.

Algorithm 1 General MCTS algorithm

```

1: function MCTS( $s_0$ )
2:   Create root node  $v_0$  corresponding to  $s_0$ 
3:   while within computational budget do
4:      $v_k \leftarrow \text{TreePolicy}(v_0)$ 
5:      $\Delta \leftarrow \text{Simulation}(v_k)$ 
6:      $\text{Backprop}(v_k, \Delta)$ 
return  $\arg \max_a Q(s_0, a)$ 

```

Tree policy: pick the immediate leaf from current state

$$\text{- e.g. } a = \arg \max_a Q(v_b, a) + \sqrt{\frac{2 \log N(v)}{N(v, a)}}$$

Simulation / Rollout Policy: simulate until termination.

- e.g. random, follow the current policy, use a NN...

Backprop step updates $N(s)$, $N(s, a)$, $Q(s, a)$

- Advantage of MC TR: 1) parallelization 2) partial MDP 3) no exponential growth due to sampling

4) no domain-specific engineering as we are just sampling.

- For multiplayer games (e.g. GD), we can have self-play. Say $\pi = (\pi_B, \pi_W)$ (black, white) and $v^\pi(s) = P[\text{Black wins} | s]$, then $V^*(s) = \max_{\pi_B} \min_{\pi_W} V^\pi(s)$. When picking tree policy, we use UCB for black moves and LCB for white moves ($a = \arg \max_{a \in A} Q(V_0, a) - \sqrt{\frac{2 \log N(a)}{N(V_0, a)}}$)