

Let's go Markov

Markov Process

- Markov property: $P[s_{t+1}|s_t] = P[s_{t+1}|s_1, \dots, s_t]$
- State transition probabilities
 - $P_{ss'} = P[s_{t+1} = s' | s_t = s]$
 - Matrix rows have to sum to 1

Markov Reward Process

- Return: $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$
- Bellman Equation for MRPs
 - $v(s) = E[R_t | s_t = s] = E[r_{t+1} + \gamma v(s_{t+1}) | s_t = s]$
- Sum notation: $v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$
- Direct solution
 - $v = R + \gamma Pv$
 - $v = (1 - \gamma P)^{-1} R$

Policy: $\pi_t(a, s) = P[A_t = a | s_t = s]$

Markov Decision Process

Immediate reward: $R_{ss'}^a = r(s, a, s')$

Value function

- $V^\pi(s) = \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s]$
- $= \mathbb{E}[r_{t+1} | s_t = s] + \gamma \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s]$
- $V^\pi(s) = \sum_{a \in A} \pi(s, a) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s'))$

Iterative Policy Evaluation Algorithm

- Input π , the policy to be evaluated
- Initialize $V(s) = 0$, for all $s \in S$
- Repeat
 - $\Delta \leftarrow 0$
 - For each $s \in S$:
 - $v \leftarrow V(s)$
 - $V(s) \leftarrow \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')]$
 - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- until $\Delta < \theta$ (a small positive number)
- Output $V \approx V^\pi$

State-Action Value function

- $Q^\pi(s, a) = \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, A_t = a]$
- $V^\pi(s) = \sum_{a \in A} \pi(s, a) Q^\pi(s, a)$

Optimal value function

- Optimal State-Action Value function
 - $Q^*(s, a) = \max_{\pi} Q^\pi(s, a)$
 - $Q^*(s, a) = E[r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, A_t = a]$
- Bellman optimality equation for Q^*
 - $Q^*(s, a) = \sum_{s'} P_{ss'}^a (R_{ss'}^a + \gamma \max_{a'} Q^*(s', a'))$

Dynamic Programming

Policy Improvement

- Policy Improvement Theorem
 - $Q^\pi(s, \pi'(s)) \geq V^\pi(s) \rightarrow V^{\pi'}(s) \geq V^\pi(s)$
- Bellman Optimality Equation (BOE)
 - $V^\pi(s) = \max_{a \in A} Q^\pi(s, a)$

Bellman's Principle of Optimality Theorem: A policy achieves the optimal value iff for any state s' reachable from s , π achieves the optimal value from state s'

Value Iteration Algorithm

- Initialize V arbitrarily, e.g., $V(s) = 0$, for all $s \in S$
- Repeat
 - $\Delta \leftarrow 0$
 - For each $s \in S$:
 - $v \leftarrow V(s)$
 - $V(s) \leftarrow \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$
 - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- until $\Delta < \theta$ (a small positive number)
- Output a deterministic policy π such that
 - $\pi(s) = \arg \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$

Policy Iteration Algorithm

- 1 Initialization
 - $V(s) \in \mathbb{R}$ and $\pi(s) \in A(s)$ arbitrarily for all $s \in S$
- 2 Policy evaluation
 - Repeat
 - $\Delta \leftarrow 0$
 - For each $s \in S$:
 - $v \leftarrow V(s)$
 - $V(s) \leftarrow \sum_{s'} P_{ss'}^{\pi(s)} [R_{ss'}^{\pi(s)} + \gamma V(s')]$
 - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
 - until $\Delta < \theta$ (a small positive number)
- 3 Policy improvement
 - Policy-stable \leftarrow true
 - For each $s \in S$:
 - $b \leftarrow \pi(s)$
 - $\pi(s) \leftarrow \arg \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$
 - if $b \neq \pi(s)$, then policy-stable \leftarrow false
 - If policy-stable, then stop; else go to 2

Generalized Policy Iteration Algorithm

- Evaluation: $\pi \rightarrow V$
 - $V \rightarrow V^\pi$
- Improvement: $V \rightarrow \pi$
 - $\pi \rightarrow \text{greedy}(V)$

Model-Free Control

MC Policy Evaluation

- procedure MonteCarloEstimation(π)
 - Init
 - $V\text{-hat}(s) \leftarrow$ arbitrary value, for all $s \in S$
 - $\text{Returns}(s) \leftarrow$ an empty list, for all $s \in S$
 - EndInit
 - repeat
 - Get trace, τ , using π .
 - for all s appearing in τ do
 - $R \leftarrow$ return from first appearance of s in τ
 - Append R to $\text{Returns}(s)$
 - $V\text{-hat}(s) \leftarrow \text{average}(\text{Returns}(s))$
 - until forever

MC Batch averaging

- $\mu = \frac{1}{k} \sum_j x_j$

MC Online averaging

- $\mu_k = \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$

Incremental Monte-Carlo Updates

- $V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (R_t - V(s_t))$
- non-stationary: $V(s_t) \leftarrow V(s_t) + \alpha (R_t - V(s_t))$

Monte-Carlo update $V(s_t) \leftarrow V(s_t) + \alpha (R_t - V(s_t))$

Temporal Difference update every time-step

- $V(s_t) \leftarrow V(s_t) + \alpha (r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$

Temporal Difference Error: $r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$

Temporal Difference Target: $r_{t+1} + \gamma V(s_{t+1})$

TD value function estimation Algorithm

- procedure TD-Estimation(π)
 - Init
 - $V\text{-hat}(s) \leftarrow$ arbitrary value, for all $s \in S$
 - EndInit
 - repeat(For each episode)
 - Initialize s
 - repeat(For each step of episode)
 - a action chosen from t at s
 - Take action a ; observe r , and next state, s'
 - $\delta \leftarrow r + \gamma V\text{-hat}(s') - V\text{-hat}(s)$
 - $V\text{-hat}(s) \leftarrow V\text{-hat}(s) + \alpha \delta$
 - $s \leftarrow s'$
 - until s is absorbing state
 - until Done
- end procedure

Greedy policy improvement

- $\pi'(s) = \arg \max_{a \in A} R_{sa}^{s'} + P_{ss'}^a V(s')$
- $\pi'(s) = \arg \max_a Q(s, a)$

On-Policy

- Soft control, ϵ -greedy, SARSA

Off-Policy

- Importance sampling, Q-Learning

ϵ -greedy policy with $\epsilon \in [0, 1]$

- $\pi(s, a) =$
 - $1 - \epsilon + \epsilon / |A(s)|$, if $a = a^* = \arg \max_a Q(s, a)$
 - $\epsilon / |A(s)|$, if $a \neq a^*$

On-policy ϵ -greedy first-visit MC control algorithm

- Code
 - Initialize, for all $s \in S, a \in A(s)$:
 - $Q(s, a) \leftarrow$ arbitrary
 - $\text{Returns}(s, a) \leftarrow$ empty list
 - $\pi(a|s) \leftarrow$ an arbitrary ϵ -soft policy
 - Repeat forever:
 - (a) Generate an episode using π
 - (b) For each pair s, a appearing in the episode:
 - $G \leftarrow$ return following the first occurrence of s, a
 - Append G to $\text{Returns}(s, a)$
 - $Q(s, a) \leftarrow \text{average}(\text{Returns}(s, a))$
 - (c) For each s in the episode:
 - $a^* \leftarrow \arg \max_a Q(s, a)$
 - For all $a \in A(s)$:
 - ♦ $\pi(a|s) \leftarrow$
 - ◊ $1 - \epsilon + \epsilon / |A(s)|$, if $a = a^*$
 - ◊ $\epsilon / |A(s)|$, if $a \neq a^*$

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,
 - $\lim_{k \rightarrow \infty} \pi_k(a, s)$
 - $= (a == \arg \max_{a'} Q_k(s, a'))$
 - '==' evaluates to 1 if true and 0 else
- e.g. if ϵ reduces to zero with $\epsilon_k = 1/k$

MC Batch Learning to Control

- procedure MonteCarloBatchOptimization(n)
 - Init
 - $Q\text{-hat}(s, a) \leftarrow$ arbitrary value, for all $s \in S, a \in A$.
 - $\pi \leftarrow \epsilon\text{-greedy}(Q\text{-hat})$
 - EndInit
 - repeat(For each batch)
 - $\text{Returns}(s, a) \leftarrow$ an empty list, for all $s \in S$
 - for $i = 1$ to n do
 - Get trace, τ , using π
 - for all (s, a) appearing in τ do
 - ♦ $R \leftarrow$ return from first appearance of (s, a) in τ .
 - ♦ Append R to $\text{Returns}(s, a)$
 - for all $s \in S, a \in A$ do
 - $Q\text{-hat}(s, a) \leftarrow \text{average}(\text{Returns}(s, a))$
 - $\pi \leftarrow \epsilon\text{-greedy}(Q\text{-hat})$

MC Iterative Learning to Control algorithm

- procedure MonteCarloIterativeOptimization(n)
 - Init
 - $Q\text{-hat}(s, a) \leftarrow$ arbitrary value, for all $s \in S, a \in A$.
 - $\pi \leftarrow \epsilon\text{-greedy}(Q\text{-hat})$
 - EndInit
 - for $i = 1$ to n do
 - Get trace, τ , using π
 - for all (s, a) appearing in τ do
 - $R \leftarrow$ return from first appearance of (s, a) in τ .
 - $Q\text{-hat}(s, a) \leftarrow Q\text{-hat}(s, a) + \alpha [R - Q\text{-hat}(s, a)]$
 - $\pi \leftarrow \epsilon\text{-greedy}(Q\text{-hat})$
 - return greedy($Q\text{-hat}$)

Sarsa: $Q(S,A) \leftarrow Q(S,A) + \alpha(r + \gamma Q(S',A') - Q(S,A))$

Sarsa algorithm

- Initialize $Q(s,a)$, $\forall s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
- Repeat (for each episode):
 - Initialize S
 - Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 - Repeat (for each step of episode):
 - Take action A , observe R, S'
 - Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)
 - $Q(S,A) \leftarrow Q(S,A) + \alpha(R + \gamma Q(S',A') - Q(S,A))$
 - $S \leftarrow S'; A \leftarrow A'$
 - until S is terminal

Robbins-Munro sequence of step-sizes α_t :

- $\sum_{t=1}^{\infty} \alpha_t = \infty \wedge \sum_{t=1}^{\infty} \alpha_t^2 < \infty$
- e.g. $\alpha_t = \alpha/t$ for $\alpha > 0$

Q-Learning:

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

Q-Learning algorithm

- Initialize $Q(s,a)$, $\forall s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
- Repeat (for each episode):
 - Initialize S
 - Repeat (for each step of episode):
 - Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 - Take action A , observe R, S'
 - $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma \max_a Q(S',a) - Q(S,A)]$
 - $S \leftarrow S'$
 - until S is terminal

	Full Backup (DP)	Sample Backup (TD)
Bellman Expectation Equation for $v_{\pi}(s)$	Iterative policy evaluation $V(s) \leftarrow E[R + \gamma V(S') s]$	TD Learning $V(s) \leftarrow R + \gamma V(S')$
Bellman Expectation Equation for $q_{\pi}(s,a)$	Q-Policy Iteration $Q(s,a) \leftarrow E[R + \gamma Q(S',A') s,a]$	Sarsa $Q(S,A) \leftarrow R + \gamma Q(S',A')$
Bellman Optimality Equation for $q^*(s,a)$	Q-Value Iteration $Q(s,a) \leftarrow E[R + \gamma \max_{a'} Q(S',a') s,a]$	Q-Learning $Q(S,A) \leftarrow R + \gamma \max_{a'} Q(S',a')$

Function Approximation

Estimate value function with function approximation

- $V^{\pi}(s) \approx \hat{V}(s, w)$
- $Q^{\pi}(s, a) \approx \hat{Q}(s, a, w)$

Stochastic Gradient Descent

- $J(w) = E[V^{\pi}(s) - \hat{V}(s, w)]^2$
- $\Delta w = -1/2 \alpha \nabla_w J(w)$
 - $= \alpha (V^{\pi}(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)$

Represent state by a feature vector

- $x(s) = (x_1(s) \dots x_n(s))^T$
- $\hat{V}(s, w) = x(s)^T w$

Radial Basis Functions: $\phi_s(i) = \exp\left(-\frac{\|s - c_i\|^2}{2\sigma_i^2}\right)$

Update rule: $\Delta w = \alpha (V^{\pi}(s) - \hat{V}(s, w)) x(s)$

- $\nabla_w \hat{V}(s, w) = x(s)$

Value function estimation using linear function approximation

- Initialize $w = 0, k = 1$
- loop
 - Sample tuple (s_k, a_k, s_{k+1}) given π
 - Update weights:
 - $w = w + \alpha(r + \gamma x(s_{k+1})^T w - x(s_k)^T w) x(s_k)$
 - $k = k + 1$
- end loop

Deep Q Learning

Q-Learning vs Lin FA Q-Learning

- Initialize $Q(s,a)$, $\forall s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
- Repeat (for each episode):
 - Initialize S
 - Repeat (for each step of episode):
 - Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 - Take action A , observe R, S'
 - $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma \max_{a'} Q(S',a') - Q(S,A)]$
 - $S \leftarrow S'$
 - until S is terminal

DQN TD error(w) = $r_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a'; w) - Q(S_t, a; w)$

$\Delta w = \alpha [r + \gamma \max_{a'} Q(S_{t+1}, a'; w) - Q(S_t, a; w)] \nabla_w Q(S_t, a; w)$

DQN algorithm

- Initialize replay memory D to capacity N
- Initialize action-value function Q with random weights
- for episode = 1, M do
 - Initialize state s_t
 - for $t = 1, T$ do
 - With probability ϵ select a random action a_t
 - otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
 - Execute action a_t and observe reward r_t and state s_{t+1}
 - Store transition (s_t, a_t, r_t, s_{t+1}) in D
 - Set $s_{t+1} = s_t$
 - Sample random minibatch of transitions (s_t, a_t, r_t, s_{t+1}) from D
 - Set $y_t =$
 - r_t for terminal s_{t+1}
 - $r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta)$ for non-terminal s_{t+1}
 - Perform a gradient descent step on $(y_t - Q(s_t, a_t; \theta))^2$
 - end for

DDQN: $Q^*(s_t, a_t) \approx r_t + \gamma Q'(s_{t+1}, \arg \max_{a'} Q(s_{t+1}, a'))$

Tricks of the Trade

Formatting inputs

- Frame Stacking
- Input normalization
- Frames of reference/Coordinate Systems

Control of training

- Bootstrapping
- Reward normalization or Reward standardization
- Target networks
- Weight initialization
- Gradient clipping
- Considering (mini-)batch sizes
- Revisiting continuous actions

Policy Gradients

$p_{\theta}(\theta) = p_{\theta}(s_1, a_1, \dots, s_T, a_T)$

$p_{\tau}(\theta) = p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$

Infinite horizon: $\theta^* =$

$\arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\tau \sim p_{\theta}(\tau)}$

Finite horizon: $\theta^* =$

$\arg \max_{\theta} \sum_{t=0}^T \mathbb{E} [r(s_t, a_t)]_{(s_t, a_t) \sim p_{\theta}(\tau)}$

Finite Difference Policy Gradient update rule:

- $\theta_k \leftarrow \theta_k + \alpha (J(\theta + u_k \epsilon) - J(\theta)) / \epsilon$

Direct Policy gradients

- $\theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_t r(s_t, a_t) \right]_{\tau \sim p_{\theta}(\tau)}$

- $J(\theta) = \mathbb{E} \left[\sum_t r(s_t, a_t) \right]_{\tau \sim p_{\theta}(\tau)}$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_t r(s_{i,t}, a_{i,t})$$

- $J(\theta) = \sum_s p^{\pi}(s) V^{\pi}(s)$

$\nabla_{\theta} J(\theta)$

$$= E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]$$

REINFORCE algorithm:

- Sample $\{\tau^i\}$ from $\pi_{\theta}(a_t | s_t)$ (run the policy)
- $\nabla_{\theta} J(\theta) \approx$

$$\sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \right) \left(\sum_t r(s_t^i, a_t^i) \right)$$

- $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Gaussian Distribution

- $p(x | \mu, \Sigma) = (2\pi)^{-D/2} |\Sigma|^{-1/2} \exp(-1/2 (x - \mu)^T \Sigma^{-1} (x - \mu))$
- $\nabla_{\theta} J(\theta)$

$$\approx \frac{1}{N} \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \right) \left(\sum_t r(s_t^i, a_t^i) \right)$$

Example: $\pi_{\theta}(a_t | s_t) = N(\text{fneural network}(s_t); \Sigma)$

- $\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) = -\frac{1}{2} \Sigma^{-1} (f(s_t) - a_t) \frac{df}{d\theta}$

Maximum likelihood

- $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) r(\tau_i)$
- $\nabla_{\theta} J_{ML}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i)$

Actor-Critic Methods

Actor: policy-based learning

Critic: value-based learning

Advantage Actor-Critic (A2C) Algorithm

- Sample $\{s_i, a_i\}$ from $\pi_{\theta}(a | s)$ (run it on the robot)
- Fit $\hat{V}_{\phi}^{\pi}(s)$ to sampled reward sums
 - $y_{i,t} \approx r(s_{i,t}, a_{i,t}) + \hat{V}_{\phi}^{\pi}(s_{i,t+1})$
 - $\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_{\phi}^{\pi}(s_i) - y_i \right\|^2$

- Evaluate $\hat{A}^{\pi}(s_i, a_i) = r(s_i, a_i) + \hat{V}_{\phi}^{\pi}(s_i') - \hat{V}_{\phi}^{\pi}(s_i)$

- $\nabla_{\theta} J(\theta) \approx \sum_i \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \hat{A}^{\pi}(s_i, a_i)$

- $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Advantage element: $Q(s, a) = V(s) + A(s, a)$