CS 4700: Foundations of Artificial Intelligence

> Fall 2019 Prof. Haym Hirsh

Lecture 12 October 28, 2019

Multi-Armed Bandit

What strategy do I use to pick a sequence of a_i ?



View Multi-Armed Bandit as a Single-State MDP



Pull each

Algorithm:

arm once

For $i \leftarrow 1$ to $n \{ Sum_i \leftarrow R(arm_i); N_i \leftarrow 1 \}; N \leftarrow n /* Initialization */$

Loop Forever

$$best \leftarrow \operatorname{argmax}_{1 \le i \le n} \left[\frac{\operatorname{Sum}_{i}}{\operatorname{N}_{i}} + \frac{g(\operatorname{N})}{\sqrt{\operatorname{N}_{i}}} \right] \qquad \begin{array}{c} g(\operatorname{N}) = \sqrt{2} \log \left(1 + \operatorname{N} \log^{2} \operatorname{N} \right) \\ g(\operatorname{N}) = c \sqrt{\ln(\operatorname{N})} \\ g(\operatorname{N}) = c \sqrt{\ln(\operatorname{N})} \\ \end{array}$$

 $Sum_{best} \leftarrow Sum_{best} + r; N_{best} \leftarrow N_{best} + 1; N \leftarrow N+1$

Pull each

Algorithm:

arm once

For $i \leftarrow 1$ to $n \{ Sum_i \leftarrow R(arm_i); N_i \leftarrow 1 \}; N \leftarrow n /* Initialization */$

Loop Forever

 $best \leftarrow \underset{1 \le i \le n}{\operatorname{argmax}} \begin{bmatrix} \frac{\operatorname{Sum}_{i}}{\operatorname{N}_{i}} + \sqrt{\frac{2 \ln(\operatorname{N})}{\operatorname{N}_{i}}} \end{bmatrix} g(\operatorname{N}) = \sqrt{2} \log(1 + \operatorname{N} \log^{2} \operatorname{N})$ $g(\operatorname{N}) = \operatorname{c} \sqrt{\ln(\operatorname{N})}$ $[\operatorname{c} = \sqrt{2}]$ pull arm a_{best} and get reward r $Sum_{best} \leftarrow Sum_{best} + r; \quad \operatorname{N}_{best} \leftarrow \operatorname{N}_{best} + 1; \quad \operatorname{N} \leftarrow \operatorname{N} + 1$

Why g(N) =
$$\sqrt{2 \log (1 + N \log^2 N)}$$
?
Why g(N) = c $\sqrt{\ln(N)}$?

Why g(N) =
$$\sqrt{2} \log (1 + N \log^2 N)$$
?
Why g(N) = c $\sqrt{\ln(N)}$?
Average reward: $\frac{\sum_{i=1}^{n} Sum_i}{N}$

Why g(N) =
$$\sqrt{2 \log (1 + N \log^2 N)}$$
?
Why g(N) = c $\sqrt{\ln(N)}$?

Average reward:
$$\frac{\sum_{i=1}^{n} \text{Sum}_{i}}{N}$$
Average reward for a policy π : $\mu_{N}^{\pi} = E_{\pi} \left[\frac{\sum_{i=1}^{n} \text{Sum}_{i}}{N} \right]$

Why g(N) =
$$\sqrt{2 \log (1 + N \log^2 N)}$$
?
Why g(N) = c $\sqrt{\ln(N)}$?

Average reward:
$$\frac{\sum_{i=1}^{n} \text{Sum}_{i}}{N}$$
Average reward for a policy π : $\mu_{N}^{\pi} = E_{\pi} \left[\frac{\sum_{i=1}^{n} \text{Sum}_{i}}{N} \right]$

Average expected reward for always picking optimal arm: μ^{best}

Why g(N) =
$$\sqrt{2 \log (1 + N \log^2 N)}$$
?
Why g(N) = c $\sqrt{\ln(N)}$?

Average reward: $\frac{\sum_{i=1}^{n} \text{Sum}_{i}}{N}$ Average reward for a policy π : $\mu_{N}^{\pi} = E_{\pi} \left[\frac{\sum_{i=1}^{n} \text{Sum}_{i}}{N} \right]$

Average expected reward for always picking optimal arm: μ^{best} Regret for a policy: regret_N^{π} = μ^{best} - μ_{N}^{π} How much exploration costs you

Why g(N) =
$$\sqrt{2 \log (1 + N \log^2 N)}$$
?
Why g(N) = c $\sqrt{\ln(N)}$?

Your expected regretKnown Result:will grow at leastRegret_N $^{\pi} = \Omega(\log(N))$ logarithmically with N

Why g(N) = $\sqrt{2} \log (1 + N \log^2 N)$? Why g(N) = c $\sqrt{\ln(N)}$? Your expected regret will grow at least Know<mark>n</mark> Result: Regret_N^{π} = $\Omega(\log(N))$ logarithmically with N $\text{Regret}_{N}^{UCB(g(N))} = O(\log(N))$ UCB with these g(N) functions have regret that grows at worst logarithmically with N

Monte Carlo Tree Search (MCTS) (Section 5.4)

Monte Carlo Tree Search (MCTS) (Section 5.4)

Application of multi-armed bandits

Timeline of Key Ideas in Game Tree Search

- 1948 Alan Turing Look ahead and use an evaluation function
- 1950 Claude Shannon Game tree search
- 1956 John McCarthy Alpha-beta pruning
- 1959 Arthur Samuel Learn evaluation function (Reinforcement learning)

. . .

1997: Deep Blue defeats Gary Kasparov (3½–2½)



Game tree search with alpha-beta pruning plus lots of enhancements

1997: Deep Blue defeats Gary Kasparov (3½–2½) (Not "deep" as in "deep learning") (Deep as in its ancestor, Deep Thought)



Game tree search with alpha-beta pruning plus lots of enhancements



Branching factor is in the 100s

Evaluation function is difficult because payoff may be very far away

Need new ideas

1983 Bruce Ballard Lookahead for probabilistic moves

- 1987 Bruce Abramson Evaluation by expected outcome (repeated simulation)
- 1992 Gerald Tesauro TD-Gammon (reinforcement learning, self-play)
- 1992 Bernd Brugmann Monte Carlo Go (simulated annealing)
- 1999 U of Alberta Simulation in Poker
- 1999 Matt Ginsberg Simulation in Bridge
- 2002 Brian Sheppard Simulation in Scrabble
- 2006 Levente Kocsis and Multi-armed bandits for Monte-Carlo tree search Csaba Szepesvár

2016: AlphaGo defeats Lee Sedol (4-1)



Key ideas of Monte Carlo Tree Search:

1. View move selection as a multi-armed bandit problem

Multi-Armed Bandit



Multi-Armed Bandit for Game Tree Search



Multi-Armed Bandit for Game Tree Search



What move should I try?

Key ideas of Monte Carlo Tree Search:

- 1. View move selection as a multi-armed bandit problem
 - 2. Evaluate moves by simulating games

Multi-Armed Bandit for Game Tree Search



What move should I try on each simulated game?

Monte-Carlo Tree Search (MCTS) Terms

- Leaf node: A state in the game tree that has successors for which no games have been simulated
- Terminal node: End of game state
- Playout/rollout: Simulating a game from a leaf node to a terminal node

• Selection: Make move choices until a leaf node S is reached









- Selection: Make move choices until a leaf node S is reached
- Expansion: Create a new successor state S' for an untried action Simulation: Play a game until you reach a terminal node



- Selection: Make move choices until a leaf node S is reached
- Expansion: Create a new successor state S' for an untried action Simulation: Play a game until you reach a terminal node
- Backpropagation: Update game statistics for the path from S' up to the root



MCTS(state): while TIME-REMAINING() do leaf \leftarrow SELECT(tree) child \leftarrow EXPAND(leaf) result \leftarrow SIMULATE(child) BACKPROPAGATE(result, child) return argmax #playouts(apply(a,state)) aeA Which move gives the game state with most playouts

How do we pick moves?



Remember This? (UCB)

Picking a Move During Selection and Expansion (UCT – Upper Confidence bound applied to Trees)

Sum_i = # of wins N_i = # of times i was tried N = # of simulations thus far (N(parent(i))

$$best \leftarrow \operatorname{argmax}_{1 \le i \le n} \left[\frac{\operatorname{Sum}_i}{\operatorname{N}_i} + c \sqrt{\frac{\ln N}{\operatorname{N}_i}} \right]$$

$$Lets you control how$$

$$much exploration$$



How do we pick moves?

Picking a Move During Simulation

- Light playout: Pick uniformly at random
- Heavy playout: Make a biased selection
 - Simulation statistics
 - Game knowledge

Trade off: Slower run time vs missing a move

Benefits

- Doesn't use an evaluation function!
- Time is linear in depth
- Handles large number of actions
- Let's you make a move when a timer goes off (to manage time) ("anytime algorithm")

AlphaGo / AlphaZero

- Truncated playouts and used (learned) evaluation function
- UCB with additional term for (learned) probability of win