Topics in Reinforcement Learning: Lessons from AlphaZero for (Sub)Optimal Control and Discrete Optimization

> Arizona State University Course CSE 691, Spring 2023

Links to Class Notes, Videolectures, and Slides at http://web.mit.edu/dimitrib/www/RLbook.html

Dimitri P. Bertsekas dbertsek@asu.edu

Lecture 7 Rollout and Approximation in Value Space for Stochastic and Other Problems

Outline

Rollout for Stochastic Problems - One-step Lookahead

- 2 Rollout for Stochastic Problems Multistep Lookahead
- 3 Monte Carlo Tree Search
- 4 Rollout for Deterministic Infinite Spaces Problems
- 5 Stochastic Programming

Stochastic Rollout: A Special Case of Approximation in Value Space

At State
$$x_k$$

$$\underset{u_k \in U_k(x_k)}{\underset{w_k \in U_k($$

 $\widetilde{J}_{k+1,\pi}(x_{k+1})$ is the cost function of some policy π

- The policy π used for rollout is called base policy
- The policy $\tilde{\pi}$ obtained by lookahead minimization is called rollout policy
- Cost improvement property: $J_{k,\tilde{\pi}}(x_k) \leq J_{k,\pi}(x_k)$ for all x_k and k

Approximate variants: Try to approximate $\tilde{J}_{k+1,\pi}(x_{k+1})$

- Possibility of truncated rollout
- Use limited simulation

Bertsekas

Stochastic Rollout for Backgammon (Tesauro, 1996)



- States are the (board position, dice roll) pairs, controls are the different ways to play the dice roll, stochastic disturbance is the dice roll.
- The 1996 version of TD-Gammon uses truncated rollout with cost function approximation provided by a neural network.
- The neural network is trained off-line by a form of approximate policy iteration that used a temporal differences algorithm for policy evaluation.
- The truncated rollout program (1996) plays better than the one without rollout, and better than any human.
- It is too slow for on-line play due to the excessive on-line Monte Carlo simulation.

Cost Improvement Property of the Nontruncated Version of Rollout: $J_{k,\tilde{\pi}}(x_k) \leq J_{k,\pi}(x_k)$ for all x_k and k

We prove this inequality by induction. Clearly it holds for k = N, since $J_{N,\pi} = J_{N,\pi} = g_N$. Assuming that it holds for index k + 1, we have for all x_k ,

$$\begin{aligned} J_{k,\tilde{\pi}}(x_k) &= E\left\{g_k(x_k,\tilde{\mu}_k(x_k),w_k) + J_{k+1,\tilde{\pi}}\left(f_k(x_k,\tilde{\mu}_k(x_k),w_k)\right)\right\} \\ &\leq E\left\{g_k(x_k,\tilde{\mu}_k(x_k),w_k) + J_{k+1,\pi}\left(f_k(x_k,\tilde{\mu}_k(x_k),w_k)\right)\right\} \\ &= \min_{u_k \in U_k(x_k)} E\left\{g_k(x_k,u_k,w_k) + J_{k+1,\pi}\left(f_k(x_k,u_k,w_k)\right)\right\} \\ &\leq E\left\{g_k(x_k,\mu_k(x_k),w_k) + J_{k+1,\pi}\left(f_k(x_k,\mu_k(x_k),w_k)\right)\right\} \\ &= J_{k,\pi}(x_k), \end{aligned}$$

where:

- The first equality is the DP equation for the rollout policy $\tilde{\pi}$.
- The first inequality holds by the induction hypothesis.
- The second equality holds by the definition of the rollout algorithm.
- The final equality is the DP equation for the base policy π .

• Given x_k , we compute for each $u_k \in U_k(x_k)$ the Q-factor

$$Q_{k,\pi}(x_k, u_k) = E \Big\{ g_k(x_k, u_k, w_k) + J_{k+1,\pi} \big(f_k(x_k, u_k, w_k) \big) \Big\}$$

and minimize over u_k (equivalently compare Q-factor differences).

- This requires that for each u_k , we generate many sample disturbance trajectories $(w_k, w_{k+1}, \ldots, w_{N-1})$ and we obtain the Q-factor as their average cost.
- In practice the number of samples is finite, so the calculated values $\hat{Q}_{k,\pi}(x_k, u_k)$ are approximate and involve stochastic variance.
- We should aim to reduce the variance of the calculated Q-factor differences

$$\hat{Q}_{k,\pi}(x_k,u_k)-\hat{Q}_{k,\pi}(x_k,u_k')$$

for control pairs (u_k, u'_k) .

- For variance reduction purposes, it is often best to use the same sample disturbance trajectories (w_k, w_{k+1},..., w_{N-1}) for all u_k (see the class notes).
- Example: Calculate the difference $q_1 q_2$ by subtracting two simulation samples $s_1 = q_1 + w_1$ and $s_2 = q_2 + w_2$. Var $(s_1 s_2)$ decreases as correlation of w_1 and w_2 increases (it is zero when $w_1 = w_2$).



Consider the pure case (no truncation, no terminal cost approximation)

- Additional cost improvement is obtained with longer lookahead
- But the necessary simulation increases rapidly with the length of the lookahead
- The big issue: How do we save in simulation effort?
- One possibility is to use the certainty equivalence approximation (fix w_{k+1},..., w_{N-1} to nominal values)
- Another possibility is Monte Carlo Tree Search (MCTS)

Certainty Equivalence Approximation (Requires Much Less Simulation)



Fix w_{k+1}, \ldots, w_{N-1} at some nominal values $\overline{w}_{k+1}, \ldots, \overline{w}_{N-1}$, and Monte Carlo average many trajectories of the form ($w_k, \overline{w}_{k+1}, \ldots, \overline{w}_{N-1}$). A two-fold benefit: deterministic rather than stochastic simulation, and fewer applications of the base policy.

Monte Carlo Tree Search - A Stochastic Form of Pruning Motivation: Save Simulation Effort

We assumed equal effort for evaluation of Q-factors of all controls at a state x_k

Drawbacks:

- Some controls may be clearly inferior to others and may not be worth as much sampling effort.
- Some controls that appear to be promising may be worth exploring better through multistep lookahead.

Monte Carlo Tree Search (MCTS) is a form of approximate multistep lookahead minimization that tries to economize in simulation time

- MCTS involves adaptive simulation (simulation effort adapted to the perceived quality of different controls).
- Aims to balance exploitation (extra simulation effort on controls that look promising) and exploration (adequate exploration of the potential of all controls).
- MCTS does not directly improve performance; it just tries to save in simulation effort. But the saving allows longer lookahead for a given computational budget.

MCTS - One-Step Approximation in Value Space



MCTS provides an economical sampling policy to estimate the Q-factors

$$\tilde{Q}_k(x_k, u_k) = E\Big\{g_k(x_k, u_k, w_k) + \tilde{J}_{k+1}\big(f_k(x_k, u_k, w_k)\big)\Big\}, \qquad u_k \in U_k(x_k)$$

Simulation scheme: Pick a control u (in some way) and generate a single sample of its Q-factor

- After the *n*th sample we have Q_{u,n}, the empirical mean of the Q-factor of each control *u* (total sample value divided by total number of samples corresponding to *u*). We can view Q_{u,n} as an exploitation index (a measure of quality of *u*).
- We could use the estimates $Q_{u,n}$ to select the control to sample next ... but how do we make sure that we do not overlook some less explored controls.

Bertsekas

Reinforcement Learning

MCTS Based on Statistical Tests



Main idea: To balance exploitation (sample controls that seem most promising, i.e., a small $Q_{u,n}$) and exploration (sample controls with small sample count).

- A popular strategy: Sample next the control *u* that minimizes the sum $Q_{u,n} + R_{u,n}$ where $R_{u,n}$ is an exploration index.
- $R_{u,n}$ is based on a confidence interval formula and depends on the sample count S_u of control u (which comes from analysis of multiarmed bandit problems).
- The UCB rule (upper confidence bound) sets $R_{u,n} = -c\sqrt{\log n/S_u}$, where *c* is a positive constant, selected empirically (values $c \approx \sqrt{2}$ are suggested, assuming that $Q_{u,n}$ is normalized to take values in the range [-1, 0]).
- MCTS with UCB rule has been extended to multistep lookahead ... but AlphaZero has used a different (semi-heuristic) rule.

Bertsekas

Classical Control Problems - Infinite Control Spaces



On-Line Rollout for Deterministic Infinite-Spaces Problems



Suppose the control space is infinite (so the number of Q-factors is infinite)

- One possibility is discretization of $U_k(x_k)$; but the number of Q-factors is excessive.
- Another possibility is to use optimization heuristics that look $(\ell 1)$ steps ahead.
- Seemlessly combine the *k*th stage minimization and the optimization heuristic into a single *l*-stage deterministic optimization (under favorable circumstances).
- Can solve it by nonlinear programming/optimal control methods (e.g., quadratic programming, gradient-based). Constraints can be readily accommodated.
- Possibility of a terminal cost approximation.
- This is the idea underlying model predictive control (MPC).

A Supply Chain Example



- System: $x_{k+1}^1 = x_k^1 + u_k^1 u_k^2$, $x_{k+1}^2 = x_k^2 + u_{k-\tau}^2 d_k$, (*d_k* is given)
- Objective: Minimize sum of costs for production (u¹_k), transportation (u²_k), and excess/shortage inventory (|x²_k d_k|) over N stages.
- The delay requires state augmentation, so the problem is intractable by DP.
- One possibility: Rollout with an $(\ell 1)$ -stages optimization heuristic that may be solvable by nonlinear programming.
- The optimization variables are the 2ℓ states (x_{k+1}^1, x_{k+1}^2) and the 2ℓ controls (u_k^1, u_2^k) , with $k = 0, ..., \ell 1$.
- This approach readily handles constraints and on-line replanning. Generalizes to integer constraints and multiple products.
- Bears close similarity to MPC.
- The computations per stage are simpler than solving the original problem. Using a different/simpler type of base heuristic and discretized DP is an alternative.

Generic resource allocation over time:

- System: $x_{k+1} = A_k x_k + B_k u_k$, (x_k and u_k are vectors, A_k , B_k are given matrices)
- Objective: Minimize a linear cost

$$c_N' x_N + \sum_{k=0}^{N-1} (c_k' x_k + d_k' u_k)$$

over N stages (c_k , d_k are given vectors, prime denotes transpose).

- Constraints: Linear on x_k and u_k (possibly some additional integer constraints).
- For large *N* and/or integer constraints this is a hard problem.
- One possibility: Rollout with an $(\ell 1)$ -stages linear programming-based heuristic.
- Readily handles on-line replanning.
- Generalizes to integer constraints, making use of integer programming software.
- Using a different/simpler type of base heuristic and discretized DP is an alternative, but does not exploit the linear programming structure of the problem.

Stochastic Programming - Two-Stage Case



Classical two-stage stochastic programming problem:

- In the first stage we choose a vector $u_0 \in U_0$ with cost $g_0(u_0)$.
- Then an uncertain event will occur, represented by a random variable *w*₀, which takes one of the values *w*¹,..., *w*^m with probabilities *p*¹,...,*p*^m.
- Once w_0 occurs, we will know its value w^i , and then we choose a vector $\mu_1(u_0, w^i) \in U_1(u_0, w^i)$ at a cost $g_1(\mu_1(u_0, w^i), w^i)$.
- The objective is to minimize the expected cost $g_0(u_0) + \sum_{i=1}^m p^i g_1(\mu_1(u_0, w^i), w^i)$
- Can be viewed as a nonlinear programming problem, whose optimization variables are u_0 , $\mu_1(u_0, w^i)$, i = 1, ..., m (five vectors in the figure).

Bertsekas

In multistage stochastic programming, the decision u_k at the *k*th stage is a function of the history $(u_0, w_0, u_1, w_1, \dots, u_{k-1}, w_{k-1})$, the state of the *k*th stage.

Similar formulation to the two-stage case ... but exact solution by DP or NLP gets rapidly out of hand as the number of stages increases.

We view this as a special case of finite horizon stochastic optimal control

- Rollout with or without truncation applies.
- Base heuristic could be based on two-stage stochastic programming with terminal cost approximation.
- Alternative base heuristics can be based on certainty equivalence approximations (only w_0 is stochastic and subsequent disturbances are fixed at nominal values).

- Review of multiagent problems.
- Multiagent rollout demos.

Please review our discussion of multiagent problems in Lecture 4 and in Chapter 2 of the class notes.

Recommended videolecture at https://www.youtube.com/watch?v=eqbb6vVIN38.

Last homework to be announced next week