

Teză de abilitare:

Planificarea optimistă pentru controlul optimal și pentru sistemele în rețea neliniare

— Rezumat —

Lucian Bușoniu

Activitatea mea de cercetare se concentrează pe obiectivul major de a dezvolta metode de control pentru sisteme complexe, posibil necunoscute. Acest obiectiv este motivat de complexitatea înaltă a sistemelor de control moderne, manifestată prin proprietăți cum ar fi neliniaritatea, stohasticitatea, scara largă, și natura lor distribuită. Mai mult decât atât, aceste sisteme nu pot fi oricând modelate precis, de exemplu fiindcă sunt insuficient înțelese. Pentru acest caz în care modelul este necunoscut, am lucrat în timpul doctoratului la clasa promițătoare a metodelor de învățare prin recompensă (*reinforcement learning*, RL) (Sutton and Barto, 1998; Szepesvári, 2010; Lewis and Liu, 2012), care învață cum să controleze un sistem stohastic neliniar fără a folosi un model. Chiar dacă modelul este cunoscut, problema de control stohastic neliniar rămâne foarte dificilă: metodele de planificare și programare dinamică, analoge învățării prin recompensă dar bazate pe model, pot fi aplicate în acest caz (La Valle, 2006; Bertsekas, 2012; Powell, 2012). Eforturile mele recente, ulterioare doctoratului, s-au concentrat pe acest al doilea caz, și mai exact pe algoritmi de planificare optimistă și pe aplicațiile lor.

Planificarea optimistă (*optimistic planning*, OP) (Munos, 2014) rezolvă probleme de control optimal modelate prin procese de decizie Markov (*Markov decision process*, MDP) (Puterman, 1994). În aceste procese, un controller măsoară la fiecare pas de timp discret starea unui sistem, și aplică o acțiune conform unei legi de control. În urma acestei acțiuni, sistemul trece într-o nouă stare, și un semnal scalar de recompensă este trimis controllerului evaluând calitatea acestei tranziții. Controllerul măsoară noua stare, și întregul ciclu se repetă. Scopul este găsirea unei legi de control optimale, care maximizează recompensa cumulată de-a lungul interacțiunii (returnul). Această paradigmă poate fi aplicată problemelor de control de nivel jos, cum ar fi reglarea stării la o valoare dată, unde recompensele sunt definite de obicei pe baza distanței de la această valoare. Paradigma funcționează însă și pentru probleme de nivel mai înalt, cum ar fi un robot care rezolvă o sarcină într-un mediu necunoscut. În acest caz, recompensele reprezintă succesul sau eșecul în rezolvarea sarcinii dorite.

Tehnicile OP funcționează într-un mod local, prin căutarea acțiunilor la cerere pentru fiecare stare a sistemului întâlnită. Natura locală a acestor tehnici le reduce dependența de dimensiunea vectorului de stare, comparativ cu programarea dinamică sau RL, și le permite să trateze în mod natural variabilele de stare cu valori continue, care sunt esențiale în controlul sistemelor fizice. La fiecare pas, OP rulează o căutare exploratorie în spațiul secvențelor de acțiuni posibile din starea curentă, reprezentând acest spațiu sub forma unui arbore. Apoi, cea mai bună primă acțiune găsită este aplicată, și întregul proces se repetă la pasul următor. Planificarea este așadar un tip foarte general de control predictiv bazat pe model. Cum resursele

de calcul sunt limitate în acest context online, căutarea trebuie să fie eficientă, și o modalitate bună pentru a o eficientiza este aplicarea principiului de *optimism în fața incertitudinii*: date fiind mai multe secvențe de acțiuni cu valori incerte, secvența mai promițătoare este explorată mai întâi. Matematic, secvența cea mai promițătoare este cea cu marginea superioară pe return maximală. Metodele optimiste combină într-un fel inovativ idei din optimizare, planificare, căutarea în grafuri (La Valle, 2006), și RL (Sutton and Barto, 1998), așa-numita teoria bandiților jucând un rol important (Auer et al., 2002) (*bandit theory*, de la numele unui joc de noroc).

Mai mulți algoritmi de OP erau deja dezvoltati când am abordat acest domeniu, mai ales în cazul acțiunilor discrete (Kocsis and Szepesvári, 2006; Hren and Munos, 2008; Bubeck and Munos, 2010; Walsh et al., 2010), dintre care așa-numitul *Upper Confidence Trees* este probabil cel mai cunoscut (Kocsis and Szepesvári, 2006). Anumiți algoritmi funcționau de asemenea pentru acțiuni continue (Mansley et al., 2011; Weinstein and Littman, 2012). Câteva dintre aceste metode au caracteristici foarte utile: sunt aplicabile la sisteme cu dinamici generale neliniare și la funcții de cost generale, noncuadractice; și furnizează garanții de optimalitate care plasează performanța soluției returnate într-o relație foarte strânsă cu bugetul de calcul investit de algoritm (Hren and Munos, 2008; Bubeck and Munos, 2010). Aceste caracteristici apar însă doar în condiții restrictive: acțiuni discrete; și fie pentru sisteme deterministe fără incertitudini sau perturbații, sau pentru sisteme stohastice dar căutând doar în clasa suboptimală a secvențelor de acțiuni în buclă deschisă. Prima parte majoră a acestei teze va prezenta așadar eforturile mele pentru a elimina aceste limitări.

Mai exact, doi algoritmi sunt discutați în cazul determinist: unul care funcționează pentru acțiuni continue, și altul care adresează perturbații printr-o abordare robustă de tip minimax, tratându-le ca pe acțiunile unui oponent. Apoi, cazul sistemelor stohastice este considerat, care poate de exemplu modela zgomotele și perturbațiile aleatoare cu distribuții de probabilitate cunoscute. Un prim algoritm este furnizat pentru distribuții cu suport discret și finit. Acest algoritm este apoi extins la o clasă de densități continue de probabilitate, prin intermediul discretizării cu puncte sigma (*sigma-point discretization*). Performanța tuturor algoritmilor este caracterizată, cu excepția celui pentru acțiuni continue pentru care analiza este în curs. Noi măsuri de complexitate a problemei sunt propuse, potrivite fiecărei clase de probleme adresate. În plus față de aceste studii analitice, toți algoritmi sunt evaluați în experimente detaliate de simulare.

Pe lângă semnificația lor fundamentală în controlul optimal, generalitatea metodelor de OP le face utile de asemenea pentru a adresa alte provocări în controlul neliniar. În special sistemele în rețea devin extrem de importante în societatea de astăzi: rețelele de comunicare, de energie, de transport, rețelele de calcul distribuite, și rețelele sociale sunt doar câteva exemple de astfel de sisteme care influențează viața de zi cu zi. Așadar, a doua parte majoră a tezei investighează aplicații ale metodelor optimiste la controlul sistemelor în rețea, tratându-le din două puncte de vedere complementare. Prima perspectivă adresează comportamentul coordonat al unor sisteme multiple interconectate numite agenți, sub constrângerile impuse de topologia de interconexiune. În acest context, un prim algoritm este propus folosind optimizarea optimistă a secvențelor de acțiuni de lungime fixă, cu scopul de a atinge un consensus asupra variabilelor de stare ale agenților, presupunând o topologie fixă de comunicație. Apoi, OP în spațiul secvențelor de lungime variabilă este aplicată pentru a obține *flocking* (numit

astfel după comportamentul stolurilor de păsări), unde topologia este dictată de o relație de proximitate între agenți. Pentru acest al doilea algoritm, principalul scop analitic este de a garanta păstrarea topologiei de comunicație dată fiind această constrângere.

A doua perspectivă tratează constrângerile de comunicație induse de rețeaua interpusă între un singur sistem și controlerul său. Două strategii de control folosind OP sunt propuse pentru a reduce numărul de transmisii în rețea. În prima strategie, secvențe de acțiuni sunt transmise la sistem cu o perioadă fixă. În a doua strategie, algoritmul lucrează cu un buget de calcul fix și decide momentul următoarei transmisii pe baza ultimei stări măsurate, ducând la o lege auto-declanșată (*self-triggered*). Garanții analitice de performanță sunt furnizate pentru ambii algoritmi.

Pe lângă aceste analize teoretice, toți algoritmi pentru sisteme în rețea sunt evaluați în exemple numerice detaliate.

Cele două direcții discutate mai sus, implicând dezvoltări fundamentale în algoritmi OP pe de o parte, și aplicațiile lor la controlul neliniar pe de alta, formează principala mea direcție de cercetare după doctorat. Direcții adiționale în RL, aplicații, și subiecte secundare de planificare și control sunt trecute în revistă în capitole separate. Cercetările rezultate direct din activitatea doctorală nu sunt discutate, chiar dacă au fost efectuate sau publicate după data doctoratului. În mod similar, nu sunt menționate activitățile la care am participat dar fără a juca un rol important.

În rezumat, contribuțiile importante prezentate în această teză sunt, în ordinea în care sunt discutate:

- Algoritmul *simultaneous optimistic optimization for planning*, pentru sisteme cu acțiuni continue (Buşoniu et al., 2013a).
- Algoritmul *optimistic minimax search*, pentru abordarea robustă a problemelor cu perturbații (Buşoniu et al., 2014).
- Algoritmul *optimistic planning for Markov decision processes*, pentru sisteme stohastice cu zgomote discrete (Buşoniu and Munos, 2012).
- Extensia acestuia numită *optimistic planning with sigma-point discretization*, pentru sisteme stohastice cu zgomote continue (Buşoniu and Tamas, 2014).
- O metodă bazată pe optimizarea optimistă pentru consensus în sisteme multiagent decentralizate (Buşoniu and Morarescu, 2014).
- O metodă bazată pe planificarea optimistă pentru flocking în sisteme multiagent (Buşoniu and Morarescu, 2015, 2013).
- Două strategii bazate pe planificarea optimistă pentru controlul în rețea (Buşoniu et al., 2013b).
- Garanții analitice de performanță pentru toate metodele cu excepția celei pentru acțiuni continue.
- Validări numerice detaliate pentru toate metodele.

În viitor, voi porni de la cercetările de planificare și control prezentate în această teză și voi integra noi direcții de control, învățare automată și RL, pentru a mă apropia de obiectivul meu general al unui *cadru algoritmic pentru controlul prin învățare și planificare al sistemelor complexe*. O componentă importantă va fi combinarea OP cu idei de RL din experiența mea anterioară, pentru a obține algoritmi hibridi cu avantajele ambelor clase de metode. Garanțiile de stabilitate pentru soluția obținută vor fi de asemenea esențiale. Direcțiile aplicative vor fi continuate și noi aplicații vor fi investigate, atât în contextul general al problemelor de control neliniar, cât și în domeniul specific al roboticii asistive. Toate aceste rezultate vor crea o platformă solidă pornind de la care se vor explora noi direcții în control și luarea deciziilor pe de o parte, și pe de alta în învățarea automată și inteligența artificială.

Bibliografie

- Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256.
- Bertsekas, D. P. (2012). *Dynamic Programming and Optimal Control*, volume 2. Athena Scientific, 4th edition.
- Bubeck, S. and Munos, R. (2010). Open loop optimistic planning. In *Proceedings 23rd Annual Conference on Learning Theory (COLT-10)*, pages 477–489, Haifa, Israel.
- Buşoniu, L., Daniels, A., Munos, R., and Babuška, R. (2013a). Optimistic planning for continuous-action deterministic systems. In *2013 IEEE International Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL-13)*, Singapore.
- Buşoniu, L. and Morarescu, C. (2013). Optimistic planning for consensus. In *Proceedings American Control Conference 2013 (ACC-13)*, pages 6735–6740, Washington, DC.
- Buşoniu, L. and Morarescu, C. (2014). Consensus for black-box nonlinear agents using optimistic optimization. *Automatica*, 50(4):1201–1208. .
- Buşoniu, L. and Morarescu, C. (2015). Topology-preserving flocking of nonlinear agents using optimistic lanning. *Control Theory and Technology*, 13(1):333–344.
- Buşoniu, L. and Munos, R. (2012). Optimistic planning for Markov decision processes. In *Proceedings 15th International Conference on Artificial Intelligence and Statistics (AISTATS-12)*, volume 22 of *JMLR Workshop and Conference Proceedings*, pages 182–189, La Palma, Canary Islands, Spain.
- Buşoniu, L., Páll, E., and Munos, R. (2014). An analysis of optimstic, best-first search for minimax sequential decision making. In *2014 IEEE International Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL-14)*, Orlando.
- Buşoniu, L., Postoyan, R., and Daafouz, J. (2013b). Near-optimal strategies for nonlinear networked control systems using optimistic planning. In *Proceedings American Control Conference 2013 (ACC-13)*, pages 3020–3025, Washington, DC.
- Buşoniu, L. and Tamas, L. (2014). Optimistic planning for the near-optimal control of general nonlinear systems with continuous transition distributions. In *19th IFAC World Congress (IFAC-14)*, Cape Town, South Africa.

- Hren, J.-F. and Munos, R. (2008). Optimistic planning of deterministic systems. In *Proceedings 8th European Workshop on Reinforcement Learning (EWRL-08)*, pages 151–164, Villeneuve d’Ascq, France.
- Kocsis, L. and Szepesvári, C. (2006). Bandit based Monte-Carlo planning. In *Proceedings 17th European Conference on Machine Learning (ECML-06)*, pages 282–293, Berlin, Germany.
- La Valle, S. M. (2006). *Planning Algorithms*. Cambridge University Press.
- Lewis, F. and Liu, D., editors (2012). *Reinforcement Learning and Adaptive Dynamic Programming for Feedback Control*. Wiley.
- Mansley, C., Weinstein, A., and Littman, M. L. (2011). Sample-based planning for continuous action Markov decision processes. In *Proceedings 21st International Conference on Automated Planning and Scheduling*, pages 335–338, Freiburg, Germany.
- Munos, R. (2014). The optimistic principle applied to games, optimization and planning: Towards foundations of Monte-Carlo tree search. *Foundations and Trends in Machine Learning*, 7(1):1–130.
- Powell, W. B. (2012). *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. Wiley, 2 edition.
- Puterman, M. L. (1994). *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. Wiley.
- Sutton, R. S. and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. MIT Press.
- Szepesvári, Cs. (2010). *Algorithms for Reinforcement Learning*. Morgan & Claypool Publishers.
- Walsh, T. J., Goschin, S., and Littman, M. L. (2010). Integrating sample-based planning and model-based reinforcement learning. In *Proceedings 24th AAAI Conference on Artificial Intelligence (AAAI-10)*, Atlanta, US.
- Weinstein, A. and Littman, M. L. (2012). Bandit-based planning and learning in continuous-action Markov decision processes. In *Proceedings 22nd International Conference on Automated Planning and Scheduling (ICAPS-12)*, São Paulo, Brazil.

Despre autor

Lucian Buşoniu deține din 2014 poziția de conferențiar universitar la Universitatea Tehnică din Cluj-Napoca, România, unde a fost și șef de lucrări în perioada 2011-2014. El a primit doctoratul (*cum laude*) în 2009 de la Universitatea Tehnică din Delft (*Technische Universiteit Delft*), Olanda, și diploma de inginer (ca șef de promoție) în 2003 la Universitatea Tehnică din Cluj-Napoca. A activat pe poziții de cercetător la TUDelft, INRIA Lille, și Centrul de Cercetare în Automatică din Nancy, Franța. Interesele lui de cercetare includ planificarea optimistă, controlul neliniar optimal, controlul prin învățare, și sistemele în rețea.