

TD072

Phương pháp Học Tăng Cường Sâu (Deep Reinforcement Learning) trong Điều Khiển Robot Công Nghiệp

Lần ban hành: 1

1. Giới thiệu tổng quan

Học Tăng Cường Sâu (Deep Reinforcement Learning – DRL) là sự kết hợp giữa Học Tăng Cường (Reinforcement Learning – RL) và Mạng Nơ-ron Sâu (Deep Neural Networks). Đây là một hướng đi đột phá cho phép hệ thống học từ **tương tác trực tiếp với môi trường**, đưa ra các quyết định tối ưu mà không cần mô hình hóa chính xác toàn bộ động lực học.

Trong bối cảnh robot công nghiệp, DRL cho phép:

- Robot tự khám phá chiến lược điều khiển để tối ưu hóa hiệu suất.
- Thích ứng linh hoạt khi điều kiện môi trường hoặc cấu hình hệ thống thay đổi.
- Giảm phụ thuộc vào các chuyên gia thiết kế thuật toán điều khiển thủ công.

Cốt lõi của DRL là **quy trình học qua thử-sai**: hệ thống tự động điều chỉnh chính sách để tối đa hóa phần thưởng (reward) tích lũy trong thời gian dài.

2. Kiến trúc hệ thống DRL

Một hệ thống DRL tiêu chuẩn trong robot công nghiệp thường bao gồm:

- Agent (Tác nhân): Bộ não của robot, thực thi thuật toán DRL để chọn hành động.
- Environment (Môi trường): Không gian làm việc, có thể là thực tế (nhà máy, kho hàng) hoặc mô phỏng (Gazebo, MuJoCo).
- State (Trạng thái): Dữ liệu cảm biến gồm vị trí, vận tốc, dữ liệu hình ảnh, tín hiệu lực, v.v.
- Action (Hành động): Tín hiệu điều khiển như xoay khớp tay, dịch chuyển, kẹp hoặc nhả vật.

22.37.13_AI Race

TD072

Phương pháp Học Tăng Cường Sâu (Deep Reinforcement Learning) trong Điều Khiến Robot Công Nghiệp

Lần ban hành: 1

Reward (Phần thưởng): Con số định lương để đánh giá hiệu quả của hành đông, ví du "hoàn thành nhiệm vụ trong thời gian ngắn" hoặc "giữ đô chính xác cao".

Quy trình lặp: Agent quan sát trạng thái \rightarrow thực hiện hành động \rightarrow nhận phần thưởng \rightarrow cập nhật chính sách \rightarrow lặp lại.

Thuật toán DRL tiêu biểu và phân tích 3.

Các thuật toán DRL phổ biến:

3.1 Deep Q-Network (DQN)

- Kết hợp Q-learning với mạng nơ-ron sâu.
- Xấp xỉ hàm giá tri Q(s,a) để tìm hành đông tối ưu.
- Dùng replay buffer để ổn định huấn luyện.

Policy Gradient Methods (ví du: REINFORCE, A2C, A3C)

- Tối ưu trưc tiếp chính sách $\pi(a|s)$.
- Thích hợp cho không gian hành động liên tục.

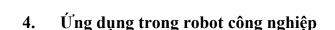
3.3 Actor-Critic

- Actor chọn hành động, Critic đánh giá chất lượng.
- Cải thiện tốc độ và tính ổn định so với chỉ dùng Policy Gradient.

3.4 Proximal Policy Optimization (PPO)

- Cân bằng giữa hiệu suất và độ ổn định.
- Phổ biến trong các ứng dụng robot phức tạp.

125-09-23 22.37.13_AI Race Các thuật toán này có thể được kết hợp với kiến trúc mang nơ-ron hiện đại như CNN (xử lý hình ảnh), LSTM/GRU (xử lý chuỗi tín hiệu) để nâng cao khả năng nhận diện trạng thái.





Ri

VIETTEL AI RACE

TD072

Phương pháp Học Tăng Cường Sâu (Deep Reinforcement Learning) trong Điều Khiến Robot Công Nghiệp

Lần ban hành: 1

DRL đã được ứng dụng thành công trong nhiều mảng:

- Lắp ráp tự động:
 Cánh tay robot học cách điều chỉnh lực và vị trí để lắp ráp các chi tiết cực kỳ nhỏ với sai số dưới 0.1 mm.
- Quản lý kho thông minh:
 Robot tự tìm đường tối ưu, tránh va chạm, tiết kiệm thời gian di chuyển trong kho diện tích lớn.
- Hàn và sơn tự động:
 Học đường đi và áp suất phun để đảm bảo lớp sơn/hàn đồng đều trên bề mặt phức tạp.
- Bảo trì và kiểm tra: Robot gắn camera học cách di chuyển trong không gian hẹp, kiểm tra lỗi sản phẩm mà không cần lập trình tỉ mỉ.

5. Quy trình triển khai thực tế

- 5.1 Mô phỏng trước khi thử nghiệm thực tế Sử dụng các nền tảng như Gazebo, PyBullet, hoặc Isaac Gym để giảm chi phí và rủi ro. Quá trình huấn luyện có thể kéo dài hàng trăm nghìn bước, nên mô phỏng là lựa chọn tiết kiệm.
- 5.2 Transfer Learning (Chuyển giao học tập)
 Sau khi huấn luyện trên mô phỏng, mô hình được tinh chỉnh để thích
 ứng với robot thực, khắc phục chênh lệch giữa môi trường ảo và thực
 tế.

5.3 Thiết kế phần thưởng (Reward Shaping)

- Chia nhỏ nhiệm vụ thành nhiều giai đoạn.
- Phạt nặng các hành động nguy hiểm (quá tốc độ, va chạm).
- Thưởng dần theo mức độ tiến bộ để học nhanh hơn.

22.37.13_AI Race



TD072

Phương pháp Học Tăng Cường Sâu (Deep Reinforcement Learning) trong Điều Khiển Robot Công Nghiệp

Lần ban hành: 1

5.4 Đảm bảo an toàn và ràng buộc vật lý Áp dụng giới hạn tốc độ, lực, vùng hoạt động. Cảm biến lực và công tắc khẩn cấp cần được tích hợp để ngăn sự cố.

6. Thách thức kỹ thuật

- Chi phí tính toán:
 Huấn luyện DRL đòi hỏi GPU/TPU mạnh và thời gian tính toán dài.
- Khả năng tổng quát hóa:
 Robot cần hoạt động tốt trong nhiều tình huống khác nhau, đòi hỏi dữ liệu đa dạng.
- Tích hợp cảm biến đa dạng:

 Camera, LiDAR, cảm biến lực, IMU... tạo ra dòng dữ liệu khổng lồ cần xử lý theo thời gian thực.
- Quản lý dữ liệu và nhật ký huấn luyện:
 Lưu trữ, phân tích log để tối ưu thuật toán.

7. Hướng phát triển tương lai

- World Model + DRL:

 Robot tự xây dựng mô hình thế giới nội tại để dự đoán trạng thái tương lai, giảm nhu cầu tương tác thật.
- Học tự giám sát (Self-Supervised RL):
 Sử dụng dữ liệu không gán nhãn để khởi tạo chính sách ban đầu.
- Edge Computing & Triển khai gọn nhẹ: Lượng tử hóa mô hình (model quantization) giúp chạy trên phần cứng hạn chế như bộ điều khiển nhúng.





TD072

Phương pháp Học Tăng Cường Sâu (Deep Reinforcement Learning) trong Điều Khiển Robot Công Nghiệp

Lần ban hành: 1

Kết hợp với Học Liên Kết (Federated RL):
 Nhiều robot học chung mà không cần chia sẻ dữ liệu thô, tăng tính bảo mật.

22.37.13_AI Race

2025-09-23 22:37:13 AI Race

2025-09-23 22.37.13_AI Race

2012.05