

MỤC LỤC

MỤC LỤC	i
DANH MỤC HÌNH ẢNH	ii
DANH MỤC BẢNG	iii
Chương 1 Giới thiệu	1
Chương 2 Kiến thức nền tảng	6
2.1 Các thành phần cơ bản của học tăng cường	6
2.1.1 Hệ thống và môi trường	6
2.1.2 Tổng điểm thưởng	8
2.2 Mô hình Markov Decision Processes	9
2.2.1 Định nghĩa mô hình Markov Decision Processes	9
2.2.2 Chính sách và hàm giá trị	11
2.2.3 Chính sách tối ưu và hàm giá trị tối ưu	14
2.2.4 Phương pháp lặp chính sách	16
2.2.5 Phương pháp lặp giá trị	20
2.3 Tìm chính sách tối ưu bằng phương pháp lặp chính sách trong bài toán thực tế	23
2.3.1 Phương pháp đánh giá chính sách	23
2.3.2 Phương pháp cải thiện chính sách	33
2.4 Tìm chính sách tối ưu bằng phương pháp lặp giá trị trong bài toán thực tế	34
2.4.1 Phương pháp Q-learning	35

Chương 3 Kết hợp học tăng cường với học sâu	38
3.1 Học tăng cường kết hợp với học sâu	39
3.1.1 Học sâu	39
3.1.2 Mạng nơ-ron tích chập	42
3.1.3 Sử dụng mạng nơ-ron để xấp xỉ hàm	46
3.1.4 Học tăng cường kết hợp với xấp xỉ hàm	51
3.2 Kết hợp học tăng cường với học sâu vào bài toán tự động chơi game	53
3.2.1 “Deep Q-Network”	53
3.2.2 Kỹ thuật làm tăng tính ổn định	54
3.2.3 Vấn đề “đánh giá quá cao” của thuật toán “Q-learning” .	58
Chương 4 Kết quả thực nghiệm	63
4.1 Giới thiệu Arcade Learning Environment	63
4.1.1 Hệ máy Atari 2600	63
4.1.2 Interface	64
4.2 Các thiết lập thực nghiệm	65
4.3 Thuật toán “Q-learning”	67
4.4 Thuật toán “Double Q-learning”	70
4.5 Phân tích hàm giá trị hành động	72
Chương 5 Kết luận và hướng phát triển	75
TÀI LIỆU THAM KHẢO	76

DANH MỤC HÌNH ẢNH

1.1	Hình ảnh các game trên hệ máy Atari	3
2.1	Quá trình tương tác giữa hệ thống và môi trường	7
2.2	Đồ thị minh họa chuyển trạng thái cho robot thu gom	11
2.3	Đồ thị minh họa quan hệ giữa những hàm giá trị	13
2.4	Đồ thị minh họa cho hàm giá trị	14
2.5	Đồ thị minh họa quan hệ giữa những hàm giá trị tối ưu	16
2.6	Đồ thị minh họa phương trình Bellman trong hàm giá trị tối ưu	17
2.7	Quy trình tìm chính sách tối ưu bằng phương pháp lặp chính sách	17
2.8	Cập nhật hàm giá trị bằng quy hoạch động	19
2.9	Hội tụ của quy trình tìm kiếm chính sách tối ưu	21
2.10	Cập nhật hàm giá trị bằng phương pháp Monte Carlo	25
2.11	Minh họa phương pháp n -step TD	30
2.12	Minh họa phương pháp Sarsa(0)	33
2.13	Minh họa phương pháp Q-learning	37
3.1	Hình mô phỏng cách hoạt động của mô hình học sâu	41
3.2	Phép tính đặc trưng của CNN	44
3.3	Phép dịch chuyển bộ lọc của CNN	44
3.4	Tầng tích chập với ba bộ lọc	45
3.5	Cấu trúc mạng “Deep Q-Network”	55
3.6	Vấn đề đánh giá quá cao của thuật toán “Q-learning”	61
3.7	Cách giải quyết vấn đề đánh giá quá cao của thuật toán “Double Q-learning”	62
4.1	Hệ máy chơi game Atari 2600	65

4.2	Dồ thị giá trị hành động trung bình và điểm thưởng trung bình	69
4.3	Dồ thị giá trị hành động trung bình và điểm thưởng trung bình của hai thuật toán	71
4.4	Hình ảnh một số “frame” của trò chơi <i>Seaquest</i>	73
4.5	Dồ thị giá trị hành động của một số trạng thái trong một lần chơi game	74

DANH MỤC BẢNG

4.1	Điểm thưởng nhận được của thuật toán “Q-learning”	68
4.2	Điểm thưởng nhận được của thuật toán “Double Q-learning” . .	70

Chương 1

Giới thiệu

Những năm gần đây, **học tăng cường** (Reinforcement learning) liên tục đạt được những thành tựu quan trọng trong lĩnh vực Trí tuệ nhân tạo (Artificial Intelligence). Những đóng góp nổi bật của phương pháp này bao gồm: tự động điều khiển robot di chuyển, điều khiển mô hình máy bay trực thăng, hệ thống chơi cờ vây... Trong số các thành tựu này, hệ thống chơi cờ vây với khả năng chiến thắng những kỳ thủ hàng đầu thế giới là một cột mốc quan trọng của lĩnh vực Trí tuệ nhân tạo. Dù vậy, học tăng cường không phải là một phương pháp mới được phát triển gần đây. Nền tảng lý thuyết của học tăng cường đã được xây dựng từ những năm 1980.

Được xây dựng nhằm mô phỏng quá trình học của con người, ý tưởng chính của học tăng cường là tìm cách lựa chọn hành động *tối ưu* để nhận được **nhiều nhất giá trị điểm thưởng** (Reward). Giá trị điểm thưởng này có ý nghĩa tương tự cảm nhận của con người về môi trường. Khi một đứa trẻ bắt đầu “học” về thế giới xung quanh của mình, những cảm giác như đau đớn (ứng với điểm thưởng thấp) hay vui sướng (điểm thưởng cao) chính là mục tiêu cần tối ưu của việc học. Một điểm quan trọng của học tăng cường là nó được xây dựng với ít giả định nhất có thể về môi trường xung quanh. Hệ thống sử dụng học tăng cường (Agent) không cần biết cách thức hoạt động của môi trường để hoạt động. Ví dụ như để điều khiển robot tìm được đi trong mê cung, hệ thống không cần biết mê cung được xây dựng thế nào hay kích thước là bao nhiêu. Việc hạn chế tối đa những ràng buộc về dữ liệu đầu vào của bài toán học tăng cường giúp cho phương pháp này có thể áp dụng vào nhiều bài toán thực tế.

Học tăng cường được xem là một nhánh trong lĩnh vực máy học ngoài hai nhánh: học có giám sát và học không có giám sát. Trong bài toán học có giám sát, dữ liệu thường được gán nhãn thủ công sẵn và việc chủ yếu của hệ thống là làm sao dự đoán chính xác các nhãn đó với dữ liệu mới. Các nhãn này có thể xem như là sự hướng dẫn trong quá trình học; tính đúng sai của việc học lúc này có thể được xác định dựa vào kết quả dự đoán của hệ thống và nhãn đúng của dữ liệu. Tiếp theo đối với những bài toán học không có giám sát, dữ liệu học thường không được gán nhãn nên công việc của việc học là phải tự tìm ra được cấu trúc “ẩn” bên dưới dữ liệu đó. Khác với hai loại bài toán vừa nêu, trong bài toán học tăng cường, hệ thống *không nhận được nhãn thực sự* (tức hành động tối ưu của tình huống hiện tại) mà chỉ nhận được điểm thưởng từ môi trường. Điểm thưởng lúc này chỉ thể hiện mức độ “tốt/xấu” của hành động vừa chọn chứ không nói lên hành động đó có phải là hành động tối ưu hay không. Điểm thưởng này thông thường rất thưa: ta có thể chỉ nhận được điểm thưởng có ý nghĩa (khác không) sau hàng nghìn hành động. Ngoài ra, giá trị điểm thưởng thường là không định và rất nhiều: cùng một hành động tại cùng một trạng thái, ta có thể nhận được điểm thưởng khác nhau vào hai thời điểm khác nhau. Đây cũng chính là những khó khăn cơ bản của bài toán học tăng cường.

Các trò chơi điện tử thường hay có điểm số mà người chơi cần phải tối ưu hoá. Đặc điểm này trùng với yêu cầu của bài toán học tăng cường, vì vậy các trò chơi này cũng chính là những ứng dụng tự nhiên nhất của phương pháp học tăng cường. Trong luận văn này, chúng em áp dụng phương pháp học tăng cường nhằm xây dựng **hệ thống tự động chơi các game** trên hệ máy Atari. Dữ liệu đầu vào của hệ thống chỉ bao gồm các frame ảnh RGB cùng với điểm số hiện tại. Từ hình ảnh thô này, hệ thống cần tìm cách chơi sao cho điểm số cuối màn chơi (Episode) là lớn nhất có thể. Hệ thống hoàn toàn không biết quy luật của game trước khi bắt đầu quá trình học mà phải tự tìm hiểu quy luật và chiến thuật chơi tối ưu. Lý do luận văn sử dụng game của máy Atari là vì các game này có quy luật chơi tương đối đơn giản nhưng lại rất đa dạng. Mỗi màn chơi thường có độ dài vừa phải (từ 2 - 15 phút) và số hành động có ý nghĩa không quá nhiều (18 hành động). Ngoài ra, các trò chơi này có thể được giả lập trên máy vi tính với tốc độ cao, giúp quá trình học được tăng tốc.



Hình 1.1: Hình ảnh các game trên hệ máy Atari

Một số khó khăn trước mắt có thể thấy ở bài toán tự động chơi game bao gồm:

- Hệ thống không được cung cấp luật chơi của game. Chính vì thế nó cũng không thể biết được hành động nào nên làm hoặc không nên làm ứng với từng tình huống cụ thể.
- Dữ liệu đầu vào là hình ảnh RGB có kích thước 210×160 . Để học được một chiến thuật chơi đơn giản thì hệ thống cũng phải chơi “thử và sai” một số lượng lớn màn chơi (có thể lên đến 10000 frame). Vì vậy, lượng dữ liệu đầu vào cần phải xử lý là rất lớn.
- Các game có hình ảnh, nội dung rất khác nhau. Để có thể học cách chơi của nhiều game khác nhau thì thuật toán học phải mang tính tổng quát cao, không sử dụng các tính chất riêng biệt của từng game.
- Để đạt được điểm số cao (ngang hoặc hơn điểm số của con người) thì phải tìm được chiến thuật chơi mang tính lâu dài. Những phương pháp tham lam, lựa chọn hành động để đạt điểm tối đa trong tương lai gần thường

không tối ưu.

Một trong những tiếp cận đầu tiên cho bài toán tự động chơi game là cuộc thi AAAI General Game Playing được đề xuất từ năm 2005 bởi đại học Stanford. Trong cuộc thi này, các đội thi nhận được mô tả sơ bộ cho những game được sử dụng cuộc thi. Các mô hình chi tiết được thiết kế bên dưới các game này không được mô tả cho các đội chơi. Dựa vào những thông tin nhận được, các đội chơi phải thiết kế một mô hình tổng quát để có thể chơi những game này do trong cuộc thi họ sẽ được nhận ngẫu nhiên một trong những game đã được mô tả [4]. Đội thắng cuộc trong cuộc thi này đã sử dụng mô hình Monte Carlo Tree Search để tìm chiến lược chơi trong những game này.

Trong khoảng ba năm trở lại, tập những game trên hệ máy Atari 2600 trở nên phổ biến trong việc đánh giá khả năng của những phương pháp tiếp cận cho bài toán tự động chơi game. Những game này bao gồm nhiều thể loại khác nhau như: đối kháng, bắn súng... Với sự đa dạng trong cách chơi của những game Atari 2600, cũng như giao diện đơn giản làm cho chúng gây được sự chú ý lớn trong cộng đồng trí tuệ nhân tạo. Nhiều phương pháp đã được đề xuất liên quan tới bài toán tự động chơi game trên hệ máy Atari 2600. [1, 2] thực hiện chia những frame ảnh trong game thành nhiều phần; sau đó xây dựng cấu trúc cây và áp dụng mô hình học Bayesian để mô phỏng lại mô hình của "môi trường" đã được xây dựng trong các game này cho hệ thống. Tuy nhiên, hướng tiếp cận này giả định những phần ảnh lân cận là đủ thông tin để dự đoán những phần ảnh trung tâm, do đó không phù hợp trong một vài game Atari có sự tương tác phức tạp. [7] sử dụng mạng nơ-ron với cấu trúc nông để tạo ra một chính sách có thể tự động chơi một số game Atari. Phương pháp sử dụng thuật toán lập trình tiến hóa để tìm trọng số liên kết tối ưu trong mạng nơ-ron, khác với phương pháp truyền thống sử dụng lan truyền ngược. Mặc dù phương pháp này vượt qua điểm số của con người trong ba game Atari nhưng còn nhiều hạn chế. Một trong những hạn chế đó là đặc tính của lập trình tiến hóa hoạt động giống như một "hộp đen", nghĩa là chúng ta không thể dễ dàng quan sát, cũng như kiểm soát cách tìm kiếm những trọng số của mạng nơ-ron như đã làm với lan truyền ngược [11].

Trong những năm gần đây, học sâu đạt được nhiều bước đột phá trong nhiều

lĩnh vực như Thị giác máy tính (Computer Vision), Nhận diện giọng nói (Speech Recognition), ... Việc kết hợp giữa học sâu và học tăng cường đã dẫn đến một hướng tiếp cận mới cho bài toán tự động chơi game: học tăng cường sâu (Deep reinforcement learning) [10]. Với học sâu, ta có thể học được những đặc trưng cấp cao (high level features) từ hình ảnh thô mà không cần phải tự thiết kế đặc trưng bằng tay (hand-designed features). Khi kết hợp với học tăng cường, ta có một hình “**End-to-end**”: việc học đặc trưng và học chiến thuật chơi được liên kết chặt chẽ với nhau. Trong luận văn này, chúng em thực hiện việc cài đặt lại phương pháp học tăng cường sâu và thử nghiệm mô hình với những tham số khác nhau. Cùng với đó, luận văn thử nghiệm kỹ thuật học chuyển tiếp (Transfer learning) nhằm giảm thời gian huấn luyện cho nhiều game.

Chương 2

Kiến thức nền tảng

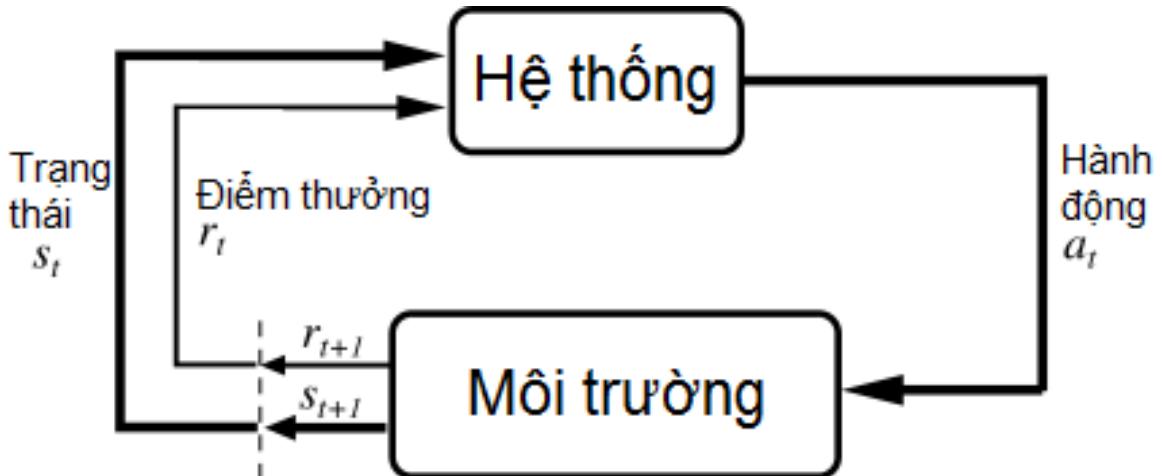
Trong chương này sẽ trình bày những kiến thức nền tảng của học tăng cường. Trong phần đầu tiên chúng em sẽ trình bày định nghĩa của các thành phần cơ bản trong học tăng cường. Tiếp đó sẽ đề cập đến mô hình Markov Decision Processes được áp dụng trong việc đánh giá lý thuyết một số thành phần của bài toán học tăng cường. Cùng với đó sẽ trình bày qui trình tổng quát để đánh giá và cải thiện chính sách trong bài toán. Cuối cùng chúng em sẽ trình bày một số phương pháp phổ biến thường được áp dụng để đánh giá cũng như cải thiện giúp hệ thống có cách giải tốt hơn cho bài toán trên.

2.1 Các thành phần cơ bản của học tăng cường

2.1.1 Hệ thống và môi trường

Trong học tăng cường, đối tượng học và đưa ra quyết định được gọi chung là *hệ thống*. Nó tương tác trực tiếp với một đối tượng được gọi là *môi trường*. Sự tương tác này được diễn ra liên tục. Hệ thống lựa chọn hành động dựa trên những gì nó nhận được từ môi trường. Những thông tin này bao gồm:

- **Trạng thái** (state): Những thông tin hữu ích mà hệ thống có thể cảm nhận được từ môi trường. Ví dụ trong đánh cờ, trạng thái có thể là vị trí những quân cờ đang có trên bàn cờ. Thường được ký hiệu là s .



Hình 2.1: Quá trình tương tác giữa hệ thống và môi trường

- **Điểm thưởng** (reward): Giá trị mà môi trường trả ra tương ứng với trạng thái mà nó vừa đạt được hoặc hành động mà hệ thống vừa thực hiện. Thường được ký hiệu là r . Cũng với ví dụ đánh cờ, điểm thưởng mà hệ thống có thể nhận được từ môi trường trong ví dụ này là: +1 nếu hệ thống thắng, -1 nếu hệ thống thua, và trong quá trình đánh cờ điểm thưởng có thể là 0 cho mỗi trạng thái bàn cờ.

Từ trạng thái và điểm thưởng nhận được, hệ thống dựa vào đó để ra quyết định chọn hành động phù hợp sao cho cố gắng đạt được nhiều điểm thưởng nhất.

Hệ thống và môi trường tương tác theo một chuỗi tuần tự các thời điểm, $t = 0, 1, 2, \dots$. Tại mỗi thời điểm t , hệ thống nhận những mô tả trạng thái của môi trường, $S_t \in \mathcal{S}$, với \mathcal{S} là tập các trạng thái có thể có. Dựa vào những mô tả trạng thái nhận được, hệ thống chọn một hành động, $A_t \in \mathcal{A}(S_t)$, trong đó $\mathcal{A}(S_t)$ là tập các hành động có thể thực hiện tại trạng thái S_t . Tại thời điểm $t+1$ sau đó, hệ thống nhận được giá trị điểm thưởng, $R_{t+1} \in \mathbb{R}$, cùng với trạng thái tiếp theo S_{t+1} . Quá trình tương tác giữa hệ thống và môi trường được mô tả trong hình 2.1

Các thành phần của hệ thống gồm có:

- **Chính sách.** Chính sách, π , xác định khả năng chọn một hành động khi hệ thống nhận được một trạng thái s . Tại thời điểm t , khả năng một hành động a được chọn ở trạng thái s , xác định bằng $\pi_t(a | s) = \mathbb{P}[A_t = a | S_t = s]$. Để đạt được mục tiêu được nhiều điểm thưởng nhất, hệ thống cần có

một chính sách chọn lựa hành động phù hợp mỗi khi gặp một trạng thái. Những phương pháp học tăng cường thường tập trung thay đổi các chính sách của hệ thống sao cho đạt được kết quả tốt trong thực nghiệm.

- **Hàm giá trị.** Hầu hết các thuật toán học tăng cường đều tập trung đánh giá những hàm giá trị, các hàm này đánh giá một trạng thái hoặc hành động là tốt như thế nào cho hệ thống thông qua việc ước lượng tổng điểm thưởng mà hệ thống có thể nhận được ở tương lai. Thông thường, giá trị của một trạng thái s , thực hiện dưới một chính sách π , được ký hiệu $v_\pi(s)$, là tổng điểm thưởng kỳ vọng mà hệ thống có thể nhận được bắt đầu từ trạng thái s về sau; và giá trị của hành động a tại trạng thái s , thực hiện dưới chính sách π , được ký hiệu $q_\pi(s, a)$, là tổng điểm thưởng kỳ vọng mà hệ thống có thể nhận được kể từ sau khi thực hiện hành động a tại trạng thái s .
- **Mô hình.** Trong một số bài toán học tăng cường, hệ thống có thể xây dựng mô hình cho riêng mình để mô phỏng lại môi trường. Qua đó cho phép hệ thống có thể suy luận hoặc dự đoán những thông tin mà nó có thể nhận được từ môi trường trong tương lai.

2.1.2 Tổng điểm thưởng

Tổng điểm thưởng (Return) G_t xác định lượng điểm thưởng mà hệ thống nhận được kể từ thời điểm t đến tương lai. Tổng điểm thưởng thường được xác định bằng nhiều hàm khác nhau, trong đó hàm đơn giản nhất xác định bằng:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T \quad (2.1)$$

ở đây T là thời điểm cuối cùng hệ thống tương tác với môi trường.

Mặt khác, G_t cũng có thể được xác định bằng tổng điểm thưởng đã bị discount qua từng thời điểm. Nó được định nghĩa như sau:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-1} R_T = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (2.2)$$

Trong đó γ là một hệ số với giá trị $0 \leq \gamma \leq 1$. γ cũng được gọi là tỉ lệ discount. Tỉ lệ này xác định độ tin tưởng của hệ thống vào giá trị điểm thưởng ở tương lai. Khi $\gamma \rightarrow 1$, hệ thống có xu hướng quan tâm đến giá trị điểm thưởng tương lai càng nhiều. Đặc biệt với $\gamma = 0$, khi đó hệ thống chỉ quan tâm giá trị điểm thưởng ở hiện tại mà bỏ qua những giá trị điểm thưởng ở tương lai.

2.2 Mô hình Markov Decision Processes

2.2.1 Định nghĩa mô hình Markov Decision Processes

Mô hình Markov Decision Processes (MDP) được sử dụng để mô hình hóa bài toán học tăng cường một cách có hình thức. Cụ thể, MDP là một bộ bao gồm 5 thành phần $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ trong đó:

- \mathcal{S} : tập trạng thái hữu hạn có thể có của môi trường.
- \mathcal{A} : tập những hành động hữu hạn mà hệ thống có thể thực hiện để tương tác với môi trường.
- γ : Hệ số có giá trị thỏa $0 \leq \gamma \leq 1$ thể hiện mức độ tin tưởng về giá trị điểm thưởng nhận được ở tương lai.
- \mathcal{P} : ma trận xác suất chuyển trạng thái. Trong đó $\mathcal{P}_{ss'}^a$ là xác suất chuyển đến trạng thái s' khi hệ thống đang ở trạng thái s và thực hiện hành động a .

$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a] \quad (2.3)$$

- \mathcal{R} : ma trận điểm thưởng theo từng bộ (trạng thái, hành động). \mathcal{R}_s^a là giá trị kỳ vọng điểm thưởng nhận được khi hệ thống thực hiện hành động a ở trạng thái s .

$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a] \quad (2.4)$$

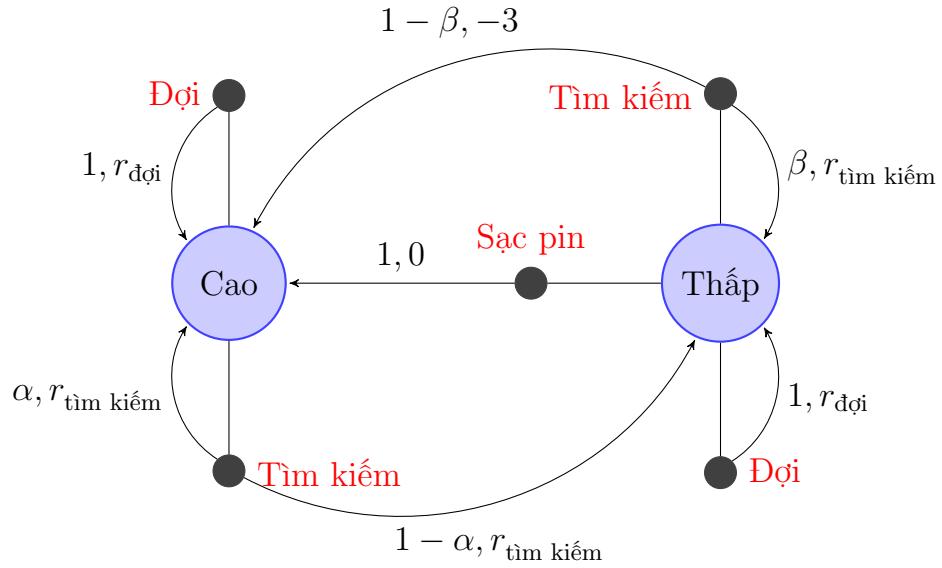
Ví dụ: Mô hình MDP trong robot thu gom Công việc của robot này là thu lượm những lon soda đã được uống hết trong văn phòng. Nó có những cảm biến để xác định những lon soda này, bánh xe và cánh tay để di chuyển và gấp

nhặt những lon này bỏ vào thùng. Robot hoạt động bằng pin sạc. Hệ thống điều khiển của robot có chức năng tiếp nhận những thông tin từ cảm biến từ đó điều khiển bánh xe và cánh tay. Trong ví dụ, chúng em chỉ xét dựa trên mức độ pin hiện tại robot nên quyết định tìm kiếm những lon soda như thế nào? Robot có thể có ba quyết định (1) thực hiện tìm kiếm một lon soda, (2) đứng yên và đợi người khác mang lon soda đến cho nó, (3) quay trở lại nơi sạc pin. Trạng thái của môi trường được xác định là trạng thái của pin hiện tại của robot. Cách tốt nhất để tìm kiếm những lon soda là robot thực hiện hành động tìm kiếm, nhưng việc này sẽ làm giảm dung lượng của pin. Ngược lại nếu robot đứng yên và đợi thì dung lượng pin của nó không giảm. Mỗi khi dung lượng pin của robot ở mức thấp nó sẽ quay lại chỗ sạc pin. Trường hợp xấu nhất có thể xảy ra là robot không đủ dung lượng pin để quay lại nơi sạc khi đó nó sẽ đứng yên và đợi ai đó mang nó đến chỗ sạc. Do đó robot cần có một chiến lược phù hợp để đạt được hiệu năng cao nhất có thể. Hệ thống đưa ra những quyết định của nó dựa trên mức năng lượng pin. Mức năng lượng này có thể được xác định hai mức *cao* và *thấp*. Khi đó tập trạng thái mà hệ thống có thể nhận được $\mathcal{S} = \{\text{cao}, \text{thấp}\}$. Những hành động của hệ thống trong ví dụ này được xét đơn giản gồm ba hành động *đợi*, *tìm kiếm*, và *sạc pin*. Khi dung lượng pin ở trạng thái cao, hệ thống chỉ thực hiện hai hành động: tìm kiếm và đợi. Ngược lại khi ở trạng thái thấp, hệ thống có thể thực hiện ba hành động: tìm kiếm, đợi, và sạc pin.

$$\mathcal{A}(\text{cao}) = \{\text{tìm kiếm}, \text{đợi}\}$$

$$\mathcal{A}(\text{thấp}) = \{\text{tìm kiếm}, \text{đợi}, \text{sạc pin}\}$$

Khi mức năng lượng pin ở mức cao, nếu robot thực hiện tìm kiếm sẽ có xác suất α năng lượng pin vẫn ở mức cao, và $1 - \alpha$ năng lượng của pin sẽ chuyển về mức thấp. Mặt khác, khi mức năng lượng ở mức thấp, nếu robot thực hiện tìm kiếm sẽ có xác suất β năng lượng pin ở mức thấp, và $1 - \beta$ chuyển đến mức cao thường hợp này xảy ra khi dung lượng pin cạn kiệt và cần ai đó mang nó đến chỗ sạc cho đến khi đạt mức năng lượng cao. Ngoài ra, mỗi lần robot thu gom được một lon soda nó sẽ nhận được +1 điểm thưởng và sẽ bị -3 điểm thưởng mỗi khi nó phải cần ai đó mang đến chỗ sạc. $r_{\text{đợi}}$, $r_{\text{tìm kiếm}}$ là số lượng lon soda kỳ



Hình 2.2: Đồ thị minh họa chuyển trạng thái cho robot thu gom. Trong đồ thị có hai loại node: node trạng thái và node hành động. Node trạng thái minh họa những trạng thái có thể có mà hệ thống có thể nhận được, nó được ký hiệu một vòng tròn lớn với tên của trạng thái bên trong. Node hành động tương ứng với cặp (trạng thái, hành động). Việc thực hiện hành động a tại trạng thái s tương ứng trên đồ thị là một cạnh bắt đầu từ node trạng s tới node hành động a . Khi đó môi trường sẽ trả ra trạng thái tiếp theo s' ứng với đích của mũi tên đi từ node hành động a . Xác suất chuyển tới trạng thái s' khi thực hiện hành động a ở trạng thái s $p(s' | s, a)$, và giá trị điểm thưởng kỳ vọng nhận được trong trường hợp này $r(s, a, s')$ tương ứng với ký hiệu trên mũi tên. Ví dụ: khi mức năng lượng pin đăng ở trạng thái *thấp*, hệ thống quyết định thực hiện hành động *sạc pin* khi đó trạng thái tiếp theo mà hệ thống nhận được sẽ là mức năng lượng pin ở trạng thái *cao* và xác suất chuyển tới trạng thái *cao* $p(\text{cao} | \text{thấp}, \text{sạc pin})$ là 1 và giá trị kỳ vọng điểm thưởng tương ứng $r(\text{thấp}, \text{sạc pin}, \text{cao})$ là 0.

vọng mà robot có thể thu gom được trong khi đợi và tìm kiếm. Hình 2.2 minh họa cho mô hình MDP trong robot thu gom.

2.2.2 Chính sách và hàm giá trị

Một chính sách π xác định xác suất mà hệ thống thực hiện hành động a khi nó trong trạng thái s , được ký hiệu $\pi(a | s)$. Có thể nói chính sách như "bộ não" của hệ thống, nó quyết định cách thức mà hệ thống hành động trong những trạng thái cụ thể. Do đó hệ thống có được một chính sách tốt đồng nghĩa với việc khả năng ra quyết định của hệ thống trở nên tốt và thông minh hơn.

Hàm giá trị cho biết những trạng thái hoặc những cặp hành động và trạng thái "tốt" như thế nào khi hệ thống trong một trạng thái hoặc thực hiện những cặp hành động và trạng thái đó. Khái niệm "tốt" ở đây nghĩa là giá trị kỳ vọng tổng điểm thưởng mà hệ thống có thể nhận được kể từ đó cho tới tương lai. Hầu hết các thuật toán trong học tăng cường đều tập trung vào việc đánh giá những hàm giá trị qua đó cải thiện chính sách trở nên tốt hơn. Tổng điểm thưởng mà hệ thống có thể nhận được trong tương lai phụ thuộc vào cách thức chọn hành động mà nó thực hiện. Do đó hàm giá trị chịu ảnh hưởng rất nhiều vào chính sách. Giá trị của trạng thái s dưới một chính sách π , ký hiệu $v_\pi(s)$, là kỳ vọng tổng điểm thưởng mà hệ thống nhận được bắt đầu từ trạng thái s thực hiện theo chính sách π sau đó. Với mô hình MDP, $v_\pi(s)$ được định nghĩa như sau:

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s \right] \quad (2.5)$$

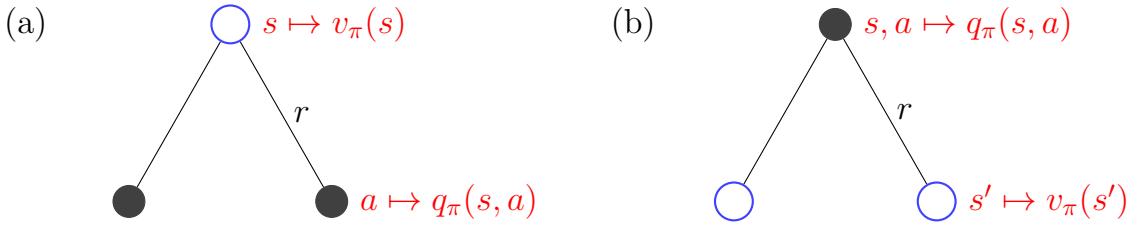
v_π được gọi là hàm giá trị trạng thái dưới chính sách π .

Tương tự, chúng ta định nghĩa giá trị của việc thực hiện hành động a trong trạng thái s dưới chính sách π , được ký hiệu $q_\pi(s, a)$, là giá trị kỳ vọng của tổng điểm thưởng mà hệ thống nhận được bắt đầu từ việc thực hiện hành động a trong trạng thái s theo chính sách π

$$q_\pi(s, a) = \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a \right] \quad (2.6)$$

q_π được gọi là hàm giá trị hành động dưới chính sách π .

Phương trình 2.7 xác định giá trị của một trạng thái bằng giá trị kỳ vọng giá trị của các hành động thực hiện tại trạng thái đó. Hình 2.3a minh họa quan hệ giữa giá trị của một trạng thái s và giá trị của các hành động thực hiện tại trạng thái đó. Hình 2.3b cho thấy từ việc thực hiện hành động a tại trạng thái s , môi trường có thể trả ra nhiều trạng thái tiếp theo s' khác nhau. Do đó giá trị của hành động a ở trạng thái s có thể được xác định bằng giá trị kỳ vọng điểm thưởng nhận được ngay sau đó cộng với tổng giá trị kỳ vọng của các trạng thái tiếp theo đó đã được nhân với hệ số γ . Cách xác định này được biểu diễn



Hình 2.3: Đồ thị minh họa quan hệ giữa hàm giá trị trạng thái và hàm giá trị hành động

trong phương trình 2.8.

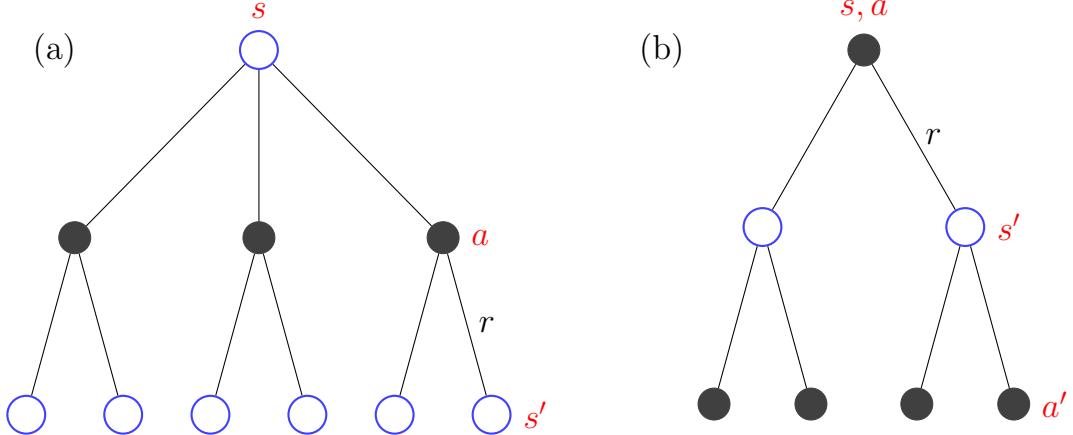
$$v_\pi = \sum_{a \in \mathcal{A}(s)} \pi(a | s) q_\pi(s, a) \quad (2.7)$$

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') \quad (2.8)$$

Hàm giá trị có một tính chất cơ bản thường được áp dụng trong học tăng cường đó là mối quan hệ đệ quy. Cho bất kỳ chính sách π với bất kỳ trạng thái s , hàm giá trị cho một trạng thái được xác định:

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi [G_t | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \end{aligned} \quad (2.9)$$

Phương trình 2.9 được gọi là phương trình Bellman cho v_π . Từ phương trình này ta thấy được mối liên quan giữa giá trị của một trạng s bất kỳ và giá trị của những trạng thái tiếp theo đạt được từ trạng thái đó. Ý tưởng nhìn trước một bước, hay nói cách khác đánh giá trạng thái hiện tại bằng cách nhìn trước tất cả những trạng tái tiếp theo có thể đạt được từ trạng thái đó, được minh họa trong hình 2.4a. Từ một trạng thái, môi trường có thể trả ra nhiều điểm thưởng r và trạng thái tiếp theo s' khác nhau. Phương trình 2.9 sẽ trung bình tất cả các trường hợp có thể đó lại theo xác suất mà chúng xuất hiện. Phương trình



Hình 2.4: Đồ thị minh họa cho (a) v_π và (b) q_π

này cũng cho thấy giá trị của một trạng thái phải bằng tổng giá trị kỳ vọng của những trạng thái tiếp sau đó và giá trị kỳ vọng điểm thưởng nhận được.

$$q_\pi(s, a) = \mathbb{E}_\pi [R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \quad (2.10)$$

Phân tích phương trình 2.6 tương tự như đã làm đối với hàm giá trị hành động, ta có được phương trình 2.10. Hình 2.4b minh họa ý tưởng nhìn trước một bước để đánh giá giá trị của một hành động ở trạng thái hiện tại. Từ một hành động a ở trạng thái s , mỗi trường hợp có thể trả ra nhiều điểm thưởng r và trạng thái s' khác nhau. Trong mỗi trạng thái s' lại có nhiều hành động a' khác nhau có thể thực hiện. Phương trình 2.10 sẽ trung bình tất cả các trường hợp có thể đó lại theo xác suất mà chúng được thực hiện. Hay nói cách khác, phương trình 2.10 cho thấy giá trị của một hành động a tại trạng thái s , $q_\pi(a, s)$ cũng được xác định bằng giá trị kỳ vọng điểm thưởng hệ thống nhận được ngay sau khi thực hiện hành động đó và giá trị kỳ vọng của các hành động trong những trạng thái kế tiếp.

2.2.3 Chính sách tối ưu và hàm giá trị tối ưu

Để giải quyết những vấn đề trong học tăng cường, chúng ta cần tìm một chính sách sao cho hệ thống có thể đạt được nhiều điểm thưởng nhất có thể. Một chính sách π được xác định là tốt hơn hoặc bằng chính sách π' khi với một trạng thái s bất kỳ giá trị của nó theo chính sách π luôn lớn hơn hoặc bằng giá trị của

trạng thái đó theo chính sách π' . Hay có thể định nghĩa theo cách khác:

$$\pi \geq \pi' \iff v_\pi(s) \geq v_{\pi'}(s), \forall s \in \mathcal{S} \quad (2.11)$$

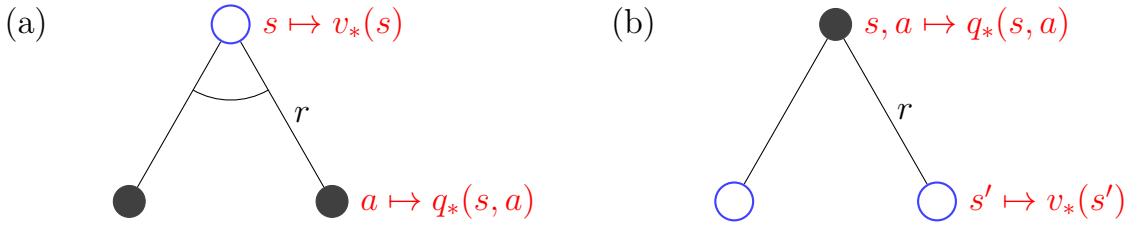
Ta luôn có ít nhất một chính sách tốt hơn hoặc bằng tất cả các chính sách còn lại [13]. Chúng được gọi chung là *chính sách tối ưu* và được ký hiệu π_* . Những chính sách tối ưu đều cùng có chung một hàm giá trị trạng thái và hàm giá trị hành động. Hai loại hàm giá trị này có thể được gọi chung là *hàm giá trị tối ưu*. Chúng ta cũng có thể gọi tách biệt *hàm giá trị trạng thái tối ưu* đối với *hàm giá trị trạng thái* và *hàm giá trị hành động tối ưu* đối với *hàm giá trị hành động*. Do đó, ta có thể dễ dàng tìm được chính sách tối ưu khi có được *hàm giá trị tối ưu* bằng phương pháp tham lam trên *hàm giá trị* này. Và ngược lại khi thực hiện hành động theo chính sách tối ưu, ta cũng dễ dàng xác định được *hàm giá trị tối ưu*. Từ tính chất này, các phương pháp tìm chính sách tối ưu luôn cố gắng tìm *hàm giá trị tối ưu* thay vì tìm trực tiếp chính sách tối ưu. Phương trình 2.12 và 2.13 định nghĩa hình thức cho hai loại hàm này

$$v_*(s) = \max_{\pi} v_{\pi}(s), \forall s \in \mathcal{S} \quad (2.12)$$

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a), \forall s \in \mathcal{S} \text{ và } \forall a \in \mathcal{A}(s) \quad (2.13)$$

Từ hai phương trình 2.12 và 2.13 thấy rằng để xác định *hàm giá trị tối ưu* của mỗi trạng thái s hoặc cặp trạng thái và hành động (s, a) , ta cần thử đánh giá giá trị của chúng theo tất cả các chính sách có thể có và chọn giá trị cao nhất là giá trị tối ưu cho trạng thái s hoặc cặp trạng thái và hành động (s, a) .

Hình 2.5 minh họa quan hệ giữa giá trị trạng thái tối ưu và *hàm giá trị hành động tối ưu*, khi có được *hàm giá trị* này ta dễ dàng có được *hàm còn lại*. Trong hình 2.5a, ta có thể xác định giá trị tối ưu cho trạng thái s dựa trên *hàm giá trị hành động tối ưu* của các hành động có thể thực hiện tại trạng thái đó. Phương trình 2.14 xác định giá trị tối ưu cho trạng thái s bằng cách chọn giá trị hành động tối ưu lớn nhất trong các hành động có thể thực hiện ở trạng thái đó. Tương tự trong hình 2.5b, ta có thể xác định giá trị tối ưu cho hành động a ở trạng thái s , dựa trên *hàm giá trị trạng thái tối ưu* của các trạng thái kế tiếp đạt được từ



Hình 2.5: Đồ thị minh họa quan hệ giữa hàm giá trị trạng thái tối ưu và hàm giá trị hành động tối ưu

hành động đó. Phương trình 2.15 xác định giá trị tối ưu của hành động a tại trạng thái s bằng tổng giá trị kỳ vọng điểm thưởng nhận được từ môi trường và giá trị tối ưu kỳ vọng của những trạng thái kế tiếp đã nhân với hệ số γ .

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_*(s, a) \quad (2.14)$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \quad (2.15)$$

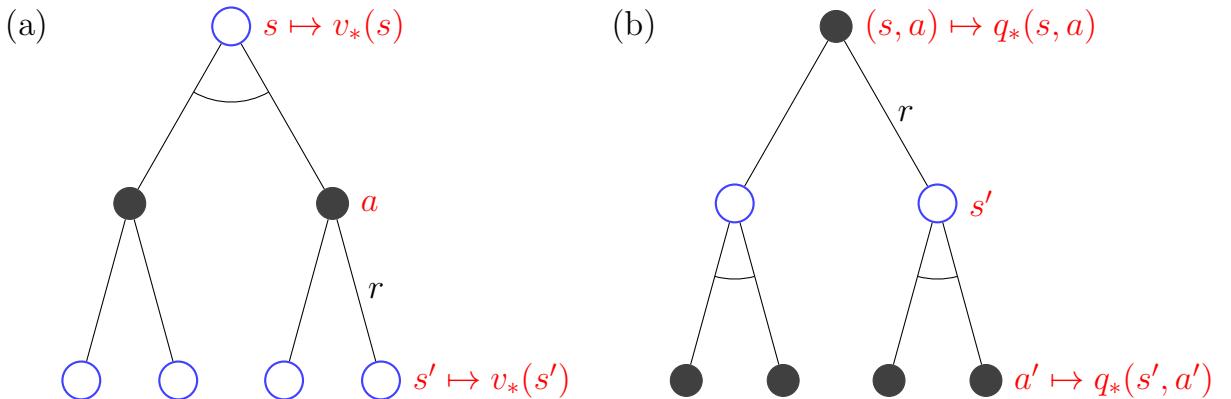
$$v_*(s) = \max_{a \in \mathcal{A}(s)} R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \quad (2.16)$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a' \in \mathcal{A}(s')} q_*(s', a') \quad (2.17)$$

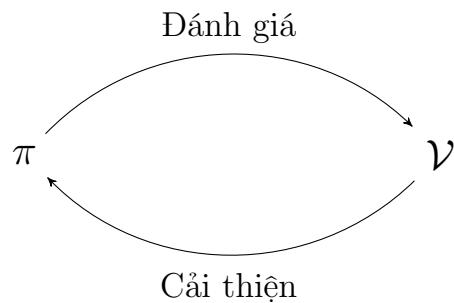
Phương trình 2.16 và 2.17 dễ dàng có được bằng cách thay thế hai phương trình 2.14 và 2.15 qua lại lẫn nhau. Từ hai phương trình này, ta thấy được dạng phương trình Bellman trong hàm giá trị trạng thái tối ưu và hàm giá trị hành động tối ưu. Hình 2.6 minh họa ý tưởng nhìn trước một bước của phương trình Bellman trong hàm giá trị tối ưu. Trong đó hình 2.6a minh họa cách thức xác định giá trị tối ưu cho một trạng thái ứng với phương trình 2.16. Hình 2.6b minh họa cách thức xác định giá trị tối ưu của một hành động ở một trạng thái ứng với phương trình 2.17.

2.2.4 Phương pháp lắp chính sách

Trong các bài toán học tăng cường, mục tiêu chính của ta là tìm được chính sách tối ưu π_* nhằm giúp cho hệ thống giải quyết bài toán tốt nhất có thể. Do



Hình 2.6: Đồ thị minh họa phương trình Bellman trong (a) v_* và (b) q_*



Hình 2.7: Quy trình tìm chính sách tối ưu bằng phương pháp lặp chính sách

đó ta cần có một quy trình để tìm ra chính sách tối ưu. Quy trình lặp chính sách là một trong hai quy trình thường được áp dụng để tìm kiếm chính sách này. Hình 2.7 minh họa quy trình này được chia thành hai giai đoạn:

- **Đánh giá chính sách:** Việc đánh giá một chính sách π được thực hiện bằng cách xác định hàm giá trị trạng thái hoặc hàm giá trị hành động của dưới chính sách đó.
- **Cải thiện chính sách:** Sau khi có được hàm giá trị của một chính sách π , chính sách cải thiện mới π' được tạo ra bằng cách thực hiện tham lam trên hàm giá trị của chính sách π , tức là chỉ chọn thực hiện hành động có giá trị cao nhất dựa trên hàm giá trị hành động có được.

Một chính sách π_1 được cải thiện từ chính sách π_0 dựa trên hàm giá trị trạng thái v_{π_0} . Khi có được chính sách π_1 ta có thể tính được hàm giá trị v_{π_1} qua đó

tiếp tục cải thiện để có được chính sách π_2 . Quá trình này diễn ra cho đến khi đạt được chính sách tối ưu.

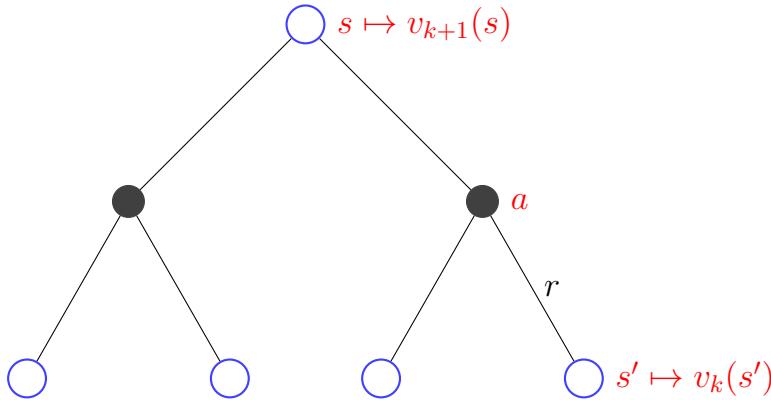
$$\pi_0 \xrightarrow{\text{Đánh giá}} v_{\pi_0} \xrightarrow{\text{Cải thiện}} \pi_1 \xrightarrow{\text{Đánh giá}} v_{\pi_1} \xrightarrow{\text{Cải thiện}} \pi_2 \dots \xrightarrow{\text{Cải thiện}} \pi_* \xrightarrow{\text{Đánh giá}} v_{\pi_*}$$

2.2.4.1 Đánh giá chính sách bằng quy hoạch động (Dynamic Programming)

Quy hoạch động thường được dùng để giải quyết các bài toán tối ưu mà dữ liệu có tính thứ tự, ví dụ như dữ liệu chuỗi hay dữ liệu thời gian. Một bài toán tối ưu có thể được giải quyết bằng quy hoạch động cần có hai đặc điểm:

- Quy tắc tối ưu (Principle of Optimality): các bài toán có thể phân rã thành các bài toán con, và kết quả của bài toán con này đóng góp vào lời giải của bài toán gốc.
- Các bài toán con chồng lấn lên nhau và lặp lại nhiều lần: nhằm tận dụng lại kết quả của những bài toán con đã tính toán trước đó.

Trong nhiều bài toán học tăng cường, kỹ thuật quy hoạch động được dùng để tìm chính sách tối ưu hoặc tối ưu hàm giá trị. Để có thể áp dụng kỹ thuật quy hoạch động, những bài toán này cũng cần phải thỏa yêu cầu là hệ thống có kiến thức đầy đủ về môi trường hay cách khác môi trường có mô hình MDP. Quy hoạch động xác định hàm giá trị của một chính sách bằng cách cập nhật hàm giá trị được khởi tạo bất kỳ ban đầu qua nhiều vòng lặp, dựa vào phương trình Bellman. Ý tưởng của cách xác định này như sau: Ban đầu khởi tạo hàm giá trị v_0 bất kỳ cho tất cả các trạng thái, trừ trạng thái kết thúc được luôn có giá trị là 0. Tiến hành cập nhật hàm giá trị mới v_1 cho chính sách dựa trên hàm giá trị v_0 theo phương trình 2.18. Tương tự cập nhật hàm giá trị mới v_2 dựa trên v_1 . Quá trình lặp cho đến khi độ khác biệt giữa hàm giá trị sau và giá trị trước đó nhỏ hơn một lượng cho trước. Quy trình cập nhật được minh họa trong hình 2.8, trong đó giá trị mới v_{k+1} của trạng thái s được xác định dựa trên giá trị kỳ vọng điểm thưởng nhận được theo chính sách π , và giá trị hiện tại v_k của các trạng thái s' kế tiếp trạng thái s . Tổng thể của việc đánh giá chính sách bằng



Hình 2.8: Đồ thị minh họa cập nhật hàm giá trị bằng quy hoạch động

quy hoạch động được trình bày ở thuật toán 2.1

$$v_{k+1}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a | s) (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s')) \quad (2.18)$$

Thuật toán 2.1 Xác định hàm giá trị bằng quy hoạch động

Đầu vào: Chính sách π cần đánh giá

Đầu ra: Hàm giá trị V xấp xỉ hàm giá trị v_π của chính sách π

Thao tác:

- 1: Khởi tạo ngẫu nhiên $V(s)$ cho tất cả trạng thái s không phải trạng thái kết thúc. Nếu s là trạng thái kết thúc, $V(s) = 0$
- 2: **repeat**
- 3: $\Delta \leftarrow 0 \%$ Tính độ khác biệt giữa hàm giá trị cũ và giá trị mới. Độ lớn của Δ được xác định là độ khác biệt lớn nhất giữa giá trị cũ và giá trị mới của một trạng thái trong tất cả các trạng thái.
- 4: **for** $s \in \mathcal{S}$ **do** %% Với mỗi trạng thái
- 5: $v \leftarrow V(s)$ %% Lưu giá trị hiện tại của trạng thái s
- 6: $V(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a | s) (\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a V(s'))$ %% Tính giá trị mới cho trạng thái s dựa trên giá trị hiện tại của các trạng thái s' kế tiếp của trạng thái s , và giá trị kỳ vọng của các hành động tại trạng thái đó theo chính sách π .
- 7: $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ %% Cập nhật giá trị mới cho Δ
- 8: **end for**
- 9: **until** $\Delta < \theta$ (Một lượng đủ nhỏ)

Mặc dù đã quy hoạch động đã được chứng minh là xấp xỉ tốt hay thậm chí

là tìm được hàm giá trị trạng thái của chính sách π [5], nhưng trong các bài toán học thực tế của học tăng cường đặc biệt là những bài toán lớn thì quy hoạch động trở nên không khả thi do chi phí tính toán cao, trong trường hợp xấu nhất chi phí tính toán thuộc $O(k^n)$ với k là số hành động và n là số trạng thái.

2.2.4.2 Cải thiện chính sách bằng phương pháp tham lam (greedy)

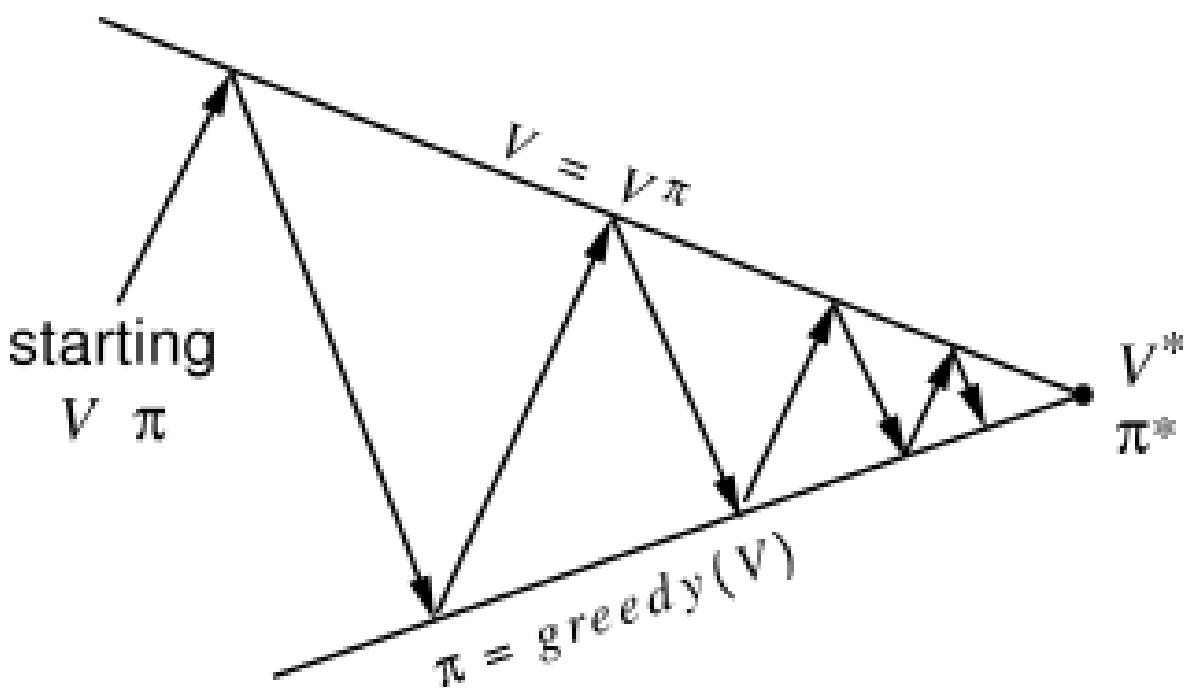
Mục tiêu chúng ta xác định hàm giá trị cho một chính sách là để tìm một chính sách tốt hơn chính sách hiện tại. Giả sử chúng ta có một chính sách π cố định tức là với mỗi trạng thái s , chính sách này luôn chọn thực hiện một hành động cố định, $\pi(s) = a$. Và cũng đã xác định được một hàm giá trị v_π cho một chính sách đó. Với một trạng thái s , câu hỏi đặt ra là chúng ta nên thay đổi một chính sách cố định khác chọn hành động $a' \neq \pi(s)$ không? Chúng ta biết giá trị của trạng thái s theo chính sách hiện tại π , $v_\pi(s)$, tốt như thế nào nhưng liệu chính sách mới π' có tốt hơn hay không nên tê đi? Theo [13], ta có thể cải thiện một chính sách bằng cách chọn hành động có giá trị cao nhất tại mỗi trạng thái s . Chính xác mới π' được xác định trong 2.19.

$$\pi'(a|s) = \begin{cases} 1 & \text{nếu } a = \underset{a \in A}{\operatorname{argmax}} q(s, a) \\ 0 & \text{ngược lại} \end{cases} \quad (2.19)$$

Ngoài ra, việc cải thiện chính sách bằng phương pháp tham lam này, ta đồng thời cũng cải thiện được hàm giá trị [13]. Việc đánh giá và cải thiện chính sách qua nhiều vòng lặp sẽ hội tụ về chính sách tối ưu cũng như hàm giá trị tối ưu. Hình 2.9 minh họa quá trình hội tụ của lặp chính sách.

2.2.5 Phương pháp lặp giá trị

Một nhược điểm của phương pháp lặp chính sách nằm ở giai đoạn đánh giá chính sách của phương pháp này. Bản thân giai đoạn này cũng là một quá trình lặp để tìm ra hàm giá trị của chính sách π đang theo. Ngoài ra, giai đoạn này cũng đã được chứng minh sẽ tìm được chính xác hàm giá trị v_π khi lặp vô hạn lần. Câu hỏi đặt ra là: Để tìm chính sách tối ưu, chúng ta có cần phải xác định chính xác hàm giá trị v_π của chính sách π hiện tại hay không? [12] đã đề xuất



Hình 2.9: Ta có một chính π , đầu tiên ta thực hiện đánh giá chính sách π để có được hàm giá trị theo chính sách này. Khi có được hàm giá trị, ta thực hiện cải thiện chính sách bằng cách lựu chọn tham lam hành động có giá trị lớn nhất dựa trên hàm giá trị đang có. Sau khi có được chính sách mới, ta tiếp tục đánh giá chính sách để có được hàm giá trị theo chính sách đó. Và tiếp tục cải thiện khi đã có được hàm giá trị. Quá trình này được lặp nhiều lần cho đến khi đạt được chính sách tối ưu cũng như hàm giá trị tối ưu.

một phương pháp có thể tìm được chính sách tối ưu mà không cần xác định hàm giá trị v_π , được gọi là phương pháp lặp giá trị. Trong phương pháp lặp chính sách, tại mỗi vòng lặp một ước lượng V cho hàm giá trị v_π được cập nhật nhiều lần để đạt chính xác hàm v_π trước khi thực hiện cải thiện chính sách. Ngược lại, phương pháp lặp giá trị thực hiện duy nhất một lần cập nhật ước lượng kết hợp với tối ưu bằng toán tử max ngay trong mỗi vòng lặp của nó. Chính vì vậy mà phương pháp lặp giá trị không tách biệt hai giai đoạn rõ ràng như phương pháp lặp chính sách và mục tiêu của phương pháp này là tìm được hàm giá trị tối ưu v_* . Phương trình 2.20 mô tả cách thức cập nhật hàm giá trị cũng như tối ưu trong một lần lặp. Ý nghĩa của phương trình này là dùng giá trị hành động lớn nhất cập nhật giá trị cho một trạng thái s .

$$\begin{aligned} v_{k+1}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a [R_{t+1} + \gamma v_k(s')] \end{aligned} \quad (2.20)$$

với $s \in \mathcal{S}$ và v_0 được khởi tạo ngẫu nhiên. Qua nhiều cập nhật chuỗi $\{v_k\}$ sẽ hội tụ về hàm giá trị tối ưu v_* .

Cách thức cập nhật của phương pháp lặp giá trị dựa trên phương trình Bellman của hàm giá trị trạng thái tối ưu 2.16. Mô hình cập nhật của phương pháp lặp giá trị tương tự mô hình cập nhật của phương pháp lặp chính sách 2.8 ngoại trừ nó yêu cầu giá trị cực đại (maximum) được chọn qua các hành động. Thêm vào đó quá trình cập nhật của lặp chính sách để tính hàm giá trị trạng thái v_π của chính sách π hiện tại, trong khi đó quá trình cập nhật của lặp giá trị để tính hàm giá trị trạng thái tối ưu v_* . Các bước thực hiện trong phương pháp lặp giá trị được mô tả trong thuật toán 2.2.

Thuật toán 2.2 Phương pháp lặp giá trị

Đầu vào:

Đầu ra: Hàm giá trị V xấp xỉ hàm giá trị tối ưu v_*

Thao tác:

- 1: Khởi tạo ngẫu nhiên $V(s)$ cho tất cả trạng thái s không phải trạng thái kết thúc. Nếu s là trạng thái kết thúc, $V(s) = 0$
 - 2: **repeat**
 - 3: $\Delta \leftarrow 0$ %% Tính độ khác biệt giữa hàm giá trị cũ và giá trị mới. Độ lớn của Δ được xác định là độ khác biệt lớn nhất giữa giá trị cũ và giá trị mới của một trạng thái trong tất cả các trạng thái.
 - 4: **for** $s \in \mathcal{S}$ **do** %% Với mỗi trạng thái
 - 5: $v \leftarrow V(s)$ %% Lưu giá trị hiện tại của trạng thái s
 - 6: $V(s) \leftarrow \max_a \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a [\mathcal{R}_s^a + \gamma V(s')]$ %% Duyệt và chọn ra giá trị hành động lớn nhất làm giá trị mới của trạng thái s
 - 7: $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
 - 8: **end for**
 - 9: **until** $\Delta < \theta$ (Một lượng đủ nhỏ)
-

2.3 Tìm chính sách tối ưu bằng phương pháp lặp chính sách trong bài toán thực tế

Trong nhiều bài toán thực tế thông thường chúng ta không có kiến thức đầy đủ về môi trường như ma trận chuyển trạng thái \mathcal{P} , ma trận điểm thưởng \mathcal{R} , tập các trạng thái \mathcal{A} . Do đó một yêu cầu được đặt ra là hệ thống phải có khả năng học từ những thông tin mà nó tiếp nhận được qua việc tương tác với môi trường. Các thông tin này thường ở dạng chuỗi (trạng thái, hành động, điểm thưởng) $S_1, A_1, R_2, S_2, A_2, R_3, \dots, S_T$. Với T là thời điểm kết thúc việc tương tác của hệ thống với môi trường.

2.3.1 Phương pháp đánh giá chính sách

Trong phần này, chúng em sẽ trình bày một số phương pháp được áp dụng để đánh giá chính sách.

2.3.1.1 Phương pháp Monte Carlo (MC)

Tương tự với quy hoạch động, Monte Carlo (MC) xác định hàm giá trị của một chính sách bằng cách cập nhật hàm giá trị khởi tạo qua nhiều vòng lặp. Điểm biệt khác với quy hoạch động của phương pháp MC là nó có thể áp dụng để đánh giá chính sách khi hệ thống không có kiến thức đầy đủ về môi trường. MC dựa trên những thông tin mà hệ thống có được qua việc tương tác với môi trường để xấp xỉ hàm giá trị. Thông thường những thông tin này được chia thành các *mẫu thực nghiệm*. Mỗi mẫu thực nghiệm là một chuỗi bắt đầu từ một trạng thái bất kỳ cho đến khi đạt được một trong những trạng thái kết thúc. Khi đó MC chỉ thực hiện cập nhật hàm giá trị khi kết thúc một mẫu thực nghiệm.

Ý tưởng của MC là xác định giá trị của một trạng thái s qua các mẫu thực nghiệm có sự xuất hiện trạng thái đó. Phương pháp MC xác định giá trị của trạng thái s bằng cách trung bình các tổng điểm thưởng mà hệ thống nhận được sau khi quan sát được trạng thái s . Khi quan sát càng nhiều mẫu thực nghiệm có trạng thái s xuất hiện, giá trị trung bình sẽ càng xấp xỉ tốt giá trị thực của trạng thái này theo chính sách π .

Một mẫu thực nghiệm là những thông tin có được trong quá trình hệ thống tương tác với môi trường bằng chính sách π . Giá trị của trạng thái s , $v(s)$, được tính dựa trên những mẫu thực nghiệm có trạng thái s xuất hiện. Một trạng thái s có thể xuất hiện nhiều lần trong một mẫu thực nghiệm. Lần xuất hiện đầu tiên của trạng thái s trong một mẫu thực nghiệm được gọi là first-visit trạng thái đó. Phương pháp first-visit MC xác định giá trị trạng thái s , $v_\pi(s)$, bằng trung bình tất cả các tổng điểm thưởng mà hệ thống nhận sau lần first-visit của trạng thái s trong các mẫu thực nghiệm. Tổng thể của việc đánh giá chính sách bằng first-visit MC được trình bày ở thuật toán 2.3. Hình 2.10 minh họa cách thức cập nhật hàm giá trị trên một mẫu thực nghiệm.

Trong nhiều trường hợp, hệ thống không có được mô hình của môi trường, việc sử dụng hàm giá trị hành động trở nên khả thi hơn hàm giá trị trạng thái. Với việc có được mô hình của môi trường, hàm giá trị trạng thái là đủ để cải thiện một chính sách trở nên tốt hơn; nó đơn giản là nhìn trước trạng thái tiếp theo có thể đến và chọn bất kỳ hành động nào dẫn đến trạng thái đó mà đạt được nhiều điểm thưởng nhất. Ngược lại, nếu không có được mô hình của môi

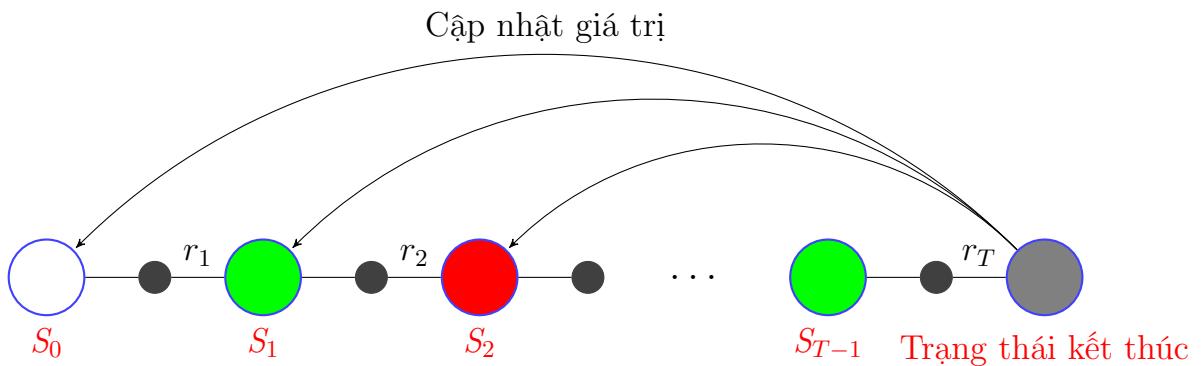
Thuật toán 2.3 Xác định hàm giá trị trạng thái bằng phương pháp first-visit MC

Đầu vào: Chính sách π cần đánh giá

Đầu ra: Hàm giá trị V xấp xỉ hàm giá trị v_π của chính sách π

Thao tác:

- 1: Khởi tạo ngẫu nhiên $V(s)$ cho tất cả trạng thái s không phải trạng thái kết thúc. Nếu s là trạng thái kết thúc, $V(s) = 0$
 - 2: Khởi tạo danh sách rỗng **Returns**(s) cho tất cả trạng thái $s \in \mathcal{S}$ %% Danh sách Returns(s) chứa tất cả các tổng điểm thưởng mà hệ thống nhận được sau lần first-visit của trạng thái s trong các mẫu thực nghiệm.
 - 3: **repeat**
 - 4: Tạo một mẫu thực nghiệm E bằng chính sách π
 - 5: **for** mỗi trạng thái s xuất hiện lần đầu trong E **do**
 - 6: $G \leftarrow$ tổng điểm thưởng nhận được sau lần xuất hiện đầu tiên của s
 - 7: Thêm G vào danh sách Returns(s)
 - 8: $V(s) \leftarrow \text{average}(\text{Returns}(s))$
 - 9: **end for**
 - 10: **until** Thỏa điều kiện dừng
-



Hình 2.10: Đồ thị minh họa cập nhật hàm giá trị trên một mẫu thực nghiệm bằng phương pháp first-visit MC. Hình tròn lớn ký hiệu cho trạng thái xuất hiện. Hình tròn nhỏ ký hiệu cho hành động thực hiện. Màu xác khác nhau giữa các hình tròn biểu thị cho sự khác nhau giữa các trạng thái. Phương pháp first-visit MC chỉ cập nhật giá trị cho các trạng thái khi kết thúc một mẫu thực nghiệm, và mỗi trạng thái chỉ được cập nhật một lần mặc dù trạng thái đó có thể xuất hiện nhiều lần trong cùng một mẫu.

trường, hàm giá trị trạng thái là không đủ do hệ thống không thể xác định được trạng thái tiếp theo là trạng thái gì. Vì vậy, nó cần đánh giá giá trị của mỗi hành động trong mỗi trạng thái để xác định hành động nào nên thực hiện ở mỗi trạng thái qua đó cải thiện chính sách đang thực hiện. Việc xác định hàm giá trị hành động q_π được thực hiện tương tự như đã làm với hàm giá trị trạng thái v_π . Để xác định giá trị của hành động a tại trạng thái s , nó thực hiện tính trung bình các tổng điểm thưởng mà hệ thống nhận được dựa vào các mẫu thực nghiệm có sự xuất hiện của cặp trạng thái, hành động (s, a) . Lần xuất hiện đầu tiên của cặp trạng thái và hành động (s, a) trong một mẫu thực nghiệm được gọi là first-visit của cặp trạng thái và hành động đó. Phương pháp first-visit MC xác định giá trị của hành động a ở trạng thái s , $q_\pi(s, a)$, bằng trung bình tất cả các tổng điểm thưởng nhận được sau lần first-visit của cặp (s, a) trong các mẫu thực nghiệm. Thuật toán 2.4 trình bày cách thức xác định hàm giá trị hành động bằng first-visit MC.

Thuật toán 2.4 Xác định hàm giá trị hành động bằng phương pháp first-visit MC

Đầu vào: Chính sách π cần đánh giá

Đầu ra: Hàm giá trị V xấp xỉ hàm giá trị v_π của chính sách π

Thao tác:

- 1: Khởi tạo ngẫu nhiên $Q(s, a)$ cho tất cả các cặp trạng thái, hành động s, a .
 - 2: Khởi tạo danh sách rỗng **Returns** (s, a) cho tất cả các cặp trạng thái, hành động (s, a) . %% Danh sách Returns (s, a) chứa tất cả các tổng điểm thưởng mà hệ thống nhận được sau lần first-visit của cặp trạng thái, hành động (s, a) trong các mẫu thực nghiệm.
 - 3: **repeat**
 - 4: Tạo một thực nghiệm E bằng chính sách π
 - 5: **for** mỗi cặp trạng thái, hành động (s, a) xuất hiện lần đầu trong E **do**
 - 6: $G \leftarrow$ tổng điểm thưởng nhận được sau lần xuất hiện đầu tiên của cặp trạng thái, hành động (s, a)
 - 7: Thêm G vào danh sách Returns (s, a)
 - 8: $Q(s, a) \leftarrow \text{average}(\text>Returns(s))$
 - 9: **end for**
 - 10: **until** Thỏa điều kiện dừng
-

2.3.1.2 Phương pháp Temporal Difference (TD)

Fương pháp Temporal Difference (TD) kết hợp ý tưởng giữa Monte Carlo và quy hoạch động. Giống như Monte Carlo, phương pháp TD có thể học trực tiếp từ các mẫu thực nghiệm có được qua việc tương tác của hệ thống với môi trường mà không cần có mô hình của môi trường. Một khác tương tự với quy hoạch động, phương pháp TD thực hiện cập nhật giá trị dựa trên những phần đã được xác định trước đó mà không phải đợi đến khi kết thúc một mẫu thực nghiệm như MC.

Giả sử ta có các trung bình μ_1, μ_2, \dots của chuỗi x_1, x_2, \dots có thể được tính như sau:

$$\begin{aligned}\mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\ &= \frac{1}{k} \left(x_k + \sum_{j=1}^{k-1} x_j \right) \\ &= \frac{1}{k} (x_k + (k-1) \mu_{k-1}) \\ &= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})\end{aligned}\tag{2.21}$$

Fương pháp Monte Carlo phải đợi cho đến khi xác định được tổng điểm tưởng sau lần xuất hiện của một trạng thái để thực hiện cập nhật giá trị cho trạng thái đó, và giá trị của một trạng thái được cập nhật qua nhiều vòng lặp. Dựa vào phương trình 2.21 giá trị của mỗi trạng thái có thể được cập nhật như sau:

$$N(S_t) \leftarrow N(S_t) + 1\tag{2.22}$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))\tag{2.23}$$

Khi đó G_t được gọi là mục tiêu cập nhật cho $V(S_t)$. Mặt khác, khi môi trường không ổn định việc cập nhật giá trị trạng thái theo 2.23 thường được cố định bằng hệ số α :

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))\tag{2.24}$$

Khác với phương pháp MC, phương pháp TD chỉ cần đợi tới bước tiếp theo ngay sau đó $t + 1$ để hình thành một cái đích cho việc cập nhật qua việc quan sát điểm thưởng R_{t+1} và giá trị của trạng thái tiếp theo $V(S_{t+1})$. Phương pháp TD đơn giản nhất được gọi là TD(0). Cách thức cập nhật giá trị của một trạng thái trong phương pháp này như sau:

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \quad (2.25)$$

Trong 2.25 ta thấy mục tiêu cập nhật cho $V(S_t)$ trong TD(0) là $R_{t+1} + \gamma V(S_{t+1})$. Vì phương pháp TD thực hiện cập nhật giá trị của một trạng thái dựa một phần vào các giá trị của những trạng thái tiếp theo nên phương pháp này là một phương pháp "bootstrapping", tương tự với quy hoạch động. Như đã định nghĩa trong 2.2.2, giá trị của trạng thái s dưới chính sách π được xác định:

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s] \quad (2.26)$$

$$\begin{aligned} &= \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \\ &= \mathbb{E}_\pi \left[R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \mid S_t = s \right] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s] \end{aligned} \quad (2.27)$$

Qua đó, ta thấy rằng phương pháp MC sử dụng ước lượng của 2.26 là mục tiêu cập nhật; trong khi đó phương pháp quy hoạch động sử dụng ước lượng của 2.27. Mục tiêu cập nhật của MC là một ước lượng vì không biết giá trị kỳ vọng trong 2.26 do đó tổng điểm thưởng trên một mẫu được sử dụng để thay thế cho giá trị kỳ vọng thực sự của nó. Mục tiêu cập nhật của quy hoạch động cũng là một ước lượng không phải vì giá trị kỳ vọng do trong quy hoạch động chúng ta đã giả định hệ thống có mô hình của môi trường; nhưng là vì $v_\pi(S_{t+1})$ là không biết do hiện tại nó đang được đánh giá; do đó $V(S_{t+1})$ được sử dụng để thay thế. Mục tiêu cập nhật trong TD là một ước lượng do cả hai nguyên nhân trên nên TD ước lượng giá trị kỳ vọng trong 2.26 qua mẫu và sử dụng ước lượng của hàm giá trị hiện tại V để thay thế cho hàm giá trị đúng v_π . Vì vậy, phương pháp

TD được cho là phương pháp kết hợp cách lấy mẫu của MC và bootstrapping của quy hoạch động. Từng bước thực hiện cập nhật hàm giá trị bằng phương pháp TD(0) được trình bày trong thuật toán [2.5](#).

Thuật toán 2.5 Xác định hàm giá trị trạng thái bằng TD(0)

Đầu vào: Chính sách π cần đánh giá

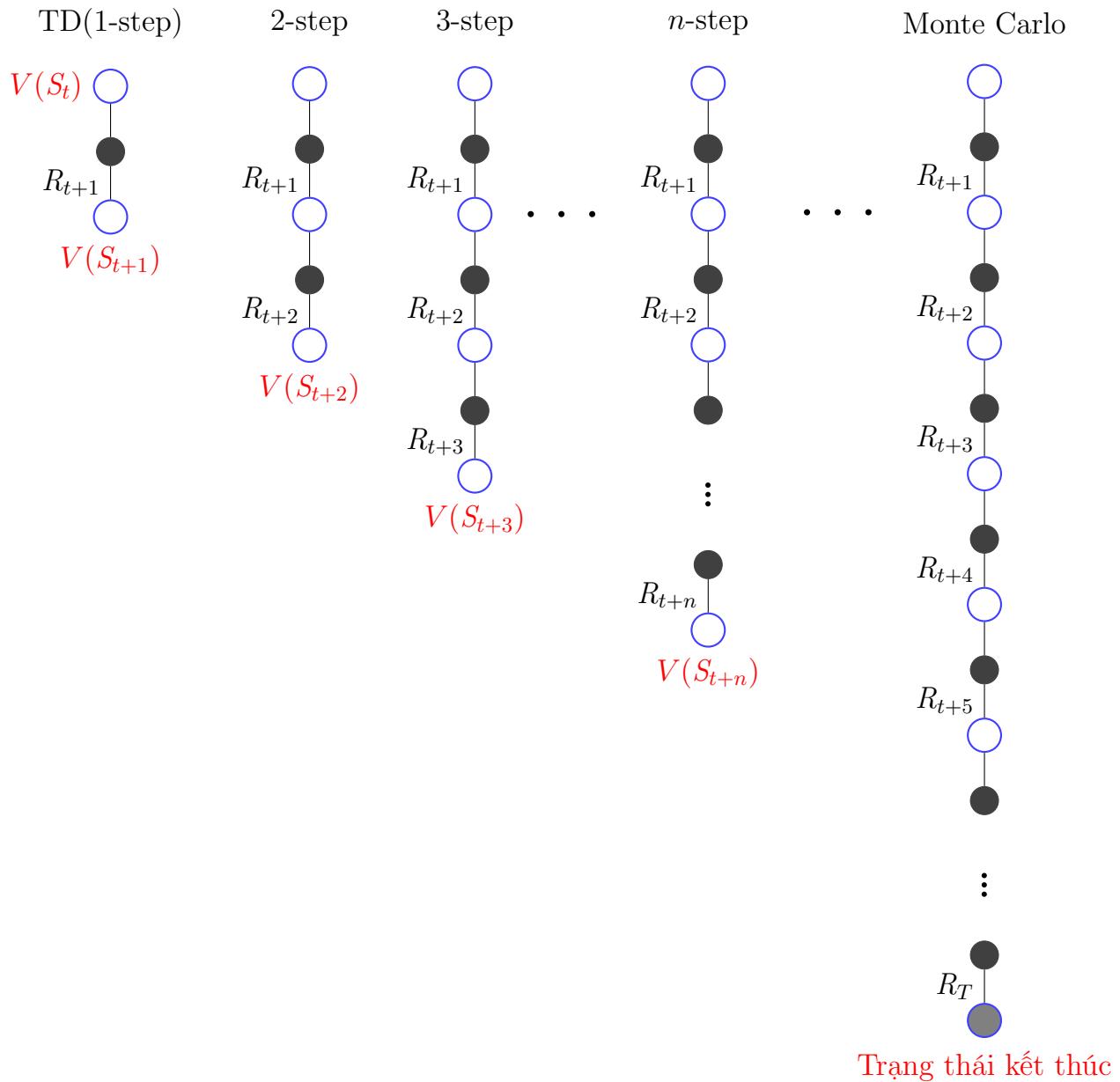
Đầu ra: Hàm giá trị V xấp xỉ hàm giá trị v_π của chính sách π

Thao tác:

- 1: **repeat**
 - 2: Tạo một mẫu thực nghiệm E bằng chính sách π
 - 3: Khởi tạo trạng thái s
 - 4: **for** mỗi bước trong E **do**
 - 5: $A \leftarrow$ hành động được chọn theo chính π tại s
 - 6: Thực hiện hành động A ; quan sát điểm thưởng r nhận được, và trạng thái tiếp theo s'
 - 7: $V(s) \leftarrow V(s) + \alpha [r + \gamma V(s') - V(s)]$ %% Thực hiện cập nhật giá trị cho trạng thái s
 - 8: $s \leftarrow s'$
 - 9: **end for**
 - 10: **until** thỏa điều kiện dừng
-

Với bất kỳ chính sách π cố định. Việc xác định hàm giá trị bằng phương pháp TD đã được chứng minh hội tụ về hàm giá trị v_π theo luật số lớn. Trong thực nghiệm, phương pháp TD thường hội tụ về hàm giá trị v_π nhanh hơn phương pháp MC.

Phương pháp TD(0) xem giá trị của các trạng thái kế tiếp từ một trạng thái s là giá trị đại diện cho lượng điểm thưởng mà trạng thái s có thể nhận được ở tương lai; và dựa vào giá trị đại diện này và điểm thưởng nhận được ngay trạng thái s để xác định giá trị của trạng thái đó. Tổng quát cho phương pháp TD là n-step TD. Để xác định giá trị của một trạng thái s , phương pháp n-step TD xem giá trị của trạng thái thứ n sau đó là giá trị đại diện cho lượng điểm thưởng mà hệ thống có thể nhận được từ bước thứ n trở về sau; và xác định giá trị của trạng thái s dựa trên giá trị đại diện này cùng với điểm thưởng đã nhận được ở n bước sau đó. Hình [2.11](#) minh họa cách xác định hàm giá trị trạng thái bằng phương pháp n-step TD. Phương pháp này thực hiện cập nhật cho một trạng thái ở bước thứ n sau khi trạng thái s xuất hiện trong mẫu thực nghiệm.



Hình 2.11: Đồ thị bên trái ngoài cùng minh họa xác định hàm giá trị bằng phương pháp TD(0), trong khi đó đồ thị bên phải ngoài cùng minh họa cho phương pháp Monte Carlo. Các đồ thị ở giữa minh họa phương pháp n -step TD ứng với từng giá trị của n

Xét một chuỗi trạng thái, điểm thưởng $S_t, R_{t+1}, S_{t+1}, R_{t+2}, \dots, S_T$. Như chúng ta đã biết, Monte Carlo thực hiện cập nhật ước lượng giá trị của trạng thái s chỉ khi tính được tổng điểm thưởng nhận được kể từ trạng thái đó:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

trong đó T là thời điểm cuối cùng trong một mẫu thực nghiệm. Ngược với MC, TD(0) thực hiện cập nhật cho một trạng thái dựa trên điểm thưởng vừa nhận được ngay trạng thái đó và giá trị hiện tại của các trạng thái kế tiếp sau đó mà không cần đợi đến khi tính được tổng điểm thưởng:

$$G_t^{(1)} = R_{t+1} + \gamma V_k(S_{t+1})$$

với V_k là giá trị hiện tại của trạng thái sau khi cập nhật k lần.

Tổng quát, phương pháp n -step TD thực hiện cập nhật giá trị ước lượng của một trạng thái s sau n bước kể từ lúc trạng thái s xuất hiện trong mẫu thực nghiệm; và mục tiêu cập nhật được xác định:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_k(S_{t+n}), \forall n \geq 1 \quad (2.28)$$

Sau khi xác định mục tiêu cập nhật, việc thực hiện cập nhật giá trị ước lượng của một trạng thái tại lần lặp thứ $k+1$ tương tự như đã thực hiện ở hai phương pháp trên:

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t^{(n)} - V(S_t)]$$

2.3.1.3 Phương pháp Sarsa

Fương pháp Sarsa được dùng để xác định hàm giá trị hành động q_π thay vì giá trị trạng thái v_π cho chính sách π . Trong Sarsa, chúng ta quan tâm chuyển từ cặp trạng thái, hành động này sang cặp trạng thái, hành động khác và học giá trị của những cặp trạng thái, hành động; thay vì chỉ quan tâm đến sự chuyển tiếp trạng thái, cũng như giá trị của chúng như trong TD. Tương tự với TD, Sarsa dựa trên phương pháp "bootstrapping" để xác định giá trị điểm thưởng mà hệ thống có thể nhận được ở tương lai từ một thời điểm xác định, đồng thời xác

định giá trị của một cặp trạng thái và hành động qua nhiều lần lặp cập nhật.

Phương pháp đơn giản nhất trong Sarsa được gọi là Sarsa(0). Cách thức cập nhật giá trị của một cặp trạng thái và hành động bằng Sarsa(0) được thực hiện:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)] \quad (2.29)$$

Giá trị của hành động A_t tại trạng thái S_t được cập nhật dựa trên điểm thưởng từ môi trường ứng với hành động đó và giá trị của hành động ở trạng thái kế tiếp sau đó. Nếu trạng thái tiếp theo S_{t+1} là trạng thái kết thúc khi đó giá trị của các hành động tại trạng thái đó đều có giá trị là không; tức là $Q(S_{t+1}, A_{t+1}) = 0$. Cách thức cập nhật của Sarsa(0) được minh họa trong hình 2.12. Cập nhật này sử dụng một bộ gồm 5 phần tử cho cập nhật $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$.

Thuật toán 2.6 Xác định hàm giá trị hành động bằng Sarsa(0)

Đầu vào: Chính sách π cần đánh giá

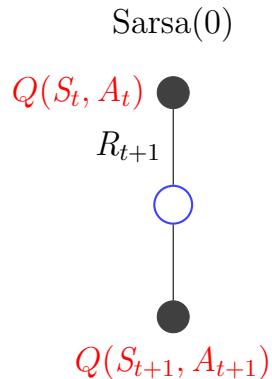
Đầu ra: Hàm giá trị Q xấp xỉ hàm giá trị q_π của chính sách π

Thao tác:

- 1: **repeat**
- 2: Tạo một mẫu thực nghiệm E bằng chính sách π
- 3: Khởi tạo trạng thái S
- 4: Chọn một hành động A tại trạng thái S theo chính sách π
- 5: **for** mỗi bước trong E **do**
- 6: Thực hiện hành động A ; quan sát điểm thưởng R nhận được, và trạng thái tiếp theo S'
- 7: Chọn hành động A' ở trạng thái S' theo chính sách π
- 8: $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ %% Thực hiện cập nhật giá trị cho cặp trạng thái, hành động S, A
- 9: $S \leftarrow S'$
- 10: $A \leftarrow A'$
- 11: **end for**
- 12: **until** thỏa điều kiện dừng

Phương pháp n -step Sarsa là phương pháp tổng quát cho việc đánh giá chính sách bằng cách xác định hàm giá trị hành động. Mục tiêu cập nhật của n -step Sarsa cho một cặp trạng thái và hành động (S_t, A_t) được xác định:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}, A_{t+1}) \quad (2.30)$$



Hình 2.12: Đồ thị minh họa xác định hàm giá trị hành động bằng phương pháp Sarsa(0)

Sau khi xác định được mục tiêu cập nhật, giá trị của một cặp trạng thái và hành động được cập nhật tương tự như cách cập nhật trong 2.29 của Sarsa(0):

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[G_t^{(n)} - Q(S_t, A_t) \right] \quad (2.31)$$

2.3.2 Phương pháp cải thiện chính sách

Phương pháp ϵ -greedy

Trong những bài toán học tăng cường thực tế, chúng ta thường không có mô hình của môi trường như: tập các trạng thái có thể có, xác suất chuyển từ trạng thái này sang trạng thái khác, và điểm thưởng nhận được cho tất cả các trạng thái. Do đó, câu hỏi đặt ra là chúng ta nên thực hiện *khai thác* những kiến thức đã học để tìm lời giải tối ưu hay *khám phá* những kiến thức mới để có cơ hội tìm được lời giải tối ưu hơn cho bài toán. Ví dụ: trong bài toán chọn nhà hàng.

- Khai thác: Chọn nhà hàng đã từng vào mà mình thích nhất.
- Khám phá: Thử chọn một nhà hàng mới.

Nếu chỉ khai thác những kiến thức đã học mà không thực hiện khám phá kiến thức mới thông thường chúng ta chỉ tìm được lời giải tối ưu cục bộ trong những kiến thức đã học. Ngược lại nếu chỉ thực hiện khám phá kiến thức mới mà không

thực hiện khai thác kiến thức đã học, chúng ta không thể tìm được lời giải tối ưu.

Ý tưởng của phương pháp ϵ -greedy đảm bảo luôn khám phá kiến thức mới trong quá trình khai thác. Trong bài toán cải thiện chính sách, giả sử chúng có hàm giá trị hành động Q , phương pháp ϵ -greedy đảm bảo các hành động luôn có khả năng được chọn thực hiện tại mỗi trạng thái. Phương pháp ϵ -greedy sẽ thực hiện hoàn toàn ngẫu nhiên mà không quan tâm đến hàm giá trị hành động với xác suất ϵ tại mỗi trạng thái; ngược lại ϵ -greedy sẽ chọn thực hiện hành động dựa trên tham lam hàm giá trị hành động với xác suất $1 - \epsilon$. Qua đó, ta thấy được nếu tại một trạng thái s có thể thực hiện được m hành động khác nhau thì xác xuất để chọn một hành động bất kỳ được chọn mà không quan tâm đến giá trị của chúng là ϵ/m . Từ 2.32, tại một trạng thái s , hành động có giá trị lớn nhất sẽ có xác suất được chọn thực hiện là $\epsilon/m + 1 - \epsilon$, ngược lại một hành động không có giá lớn nhất cũng có xác suất được chọn là ϵ/m .

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{nếu } a = \underset{a \in A}{\operatorname{argmax}} Q(s, a) \\ \epsilon/m & \text{ngược lại} \end{cases} \quad (2.32)$$

Câu hỏi đặt một chính sách π' thực hiện ϵ -greedy trên hàm giá trị hành động của chính sách π có tốt hơn chính sách π không? Theo định lý tối ưu: Với chính sách thực hiện theo phương pháp ϵ -greedy π bất kỳ, một chính sách khác π' thực hiện ϵ -greedy trên hàm giá trị hành động của chính sách π , q_{π} , luôn là một chính sách tốt hơn hoặc bằng chính sách π , $v_{\pi'}(s) \geq v_{\pi}(s), \forall s \in \mathcal{S}$.

2.4 Tìm chính sách tối ưu bằng phương pháp lắp giá trị trong bài toán thực tế

Như đã nói ở phần trước, những bài toán thực tế không có tính chất MDP của môi trường hay nói cách khác những thông tin như: những trạng thái có thể có, xác suất chuyển trạng thái, và những giá trị kỳ vọng điểm thưởng ứng với từng hành động tại mỗi trạng thái, chúng ta không biết hoặc chỉ biết một phần. Chính vì lý do đó, ngoài việc các phương pháp có khả năng tìm chính sách tối

ưu trên, phương pháp lặp giá trị áp dụng trong bài toán thực tế sẽ có gắng xấp xỉ hàm giá trị hành động tối ưu q_* khác với mô hình MDP.

2.4.1 Phương pháp Q-learning

On-policy và off-policy

Ý tưởng của on-policy là dựa trên những kinh nghiệm thực tế của chính hệ thống, nó có thể tự cải thiện trở nên tốt hơn. Những phương pháp thực hiện theo *on-policy* là những phương pháp đánh giá và cải thiện chính sách π dựa trên những mẫu dữ liệu có được qua việc tương tác với môi trường theo chính chính sách đó. Các phương pháp quy hoạch động, MC, TD, Sarsa là những phương pháp thực hiện theo on-policy. Ngược lại với on-policy, ý tưởng của off-policy là hệ thống có thể cải thiện chính sách của nó dựa trên những kinh nghiệm thực tế của một hệ thống có một chính sách thực hiện khác. Những phương pháp thực hiện theo *off-policy* đánh giá và cải thiện chính sách π dựa trên những mẫu dữ liệu có được qua việc tương tác với môi trường theo một chính sách π' khác. Ngoài ra, những phương pháp thực hiện theo off-policy có thể tận dụng lại những mẫu dữ liệu cũ để tiếp tục cải thiện chính sách.

Phương pháp off-policy Q-learning

Một trong những đột phá quan trọng nhất trong học tăng cường được phát triển dựa trên phương pháp TD được biết đến chính là Q-learning [13]. Tương tự như những phương pháp ở trên, phương pháp Q-learning xác định giá trị của một cặp trạng thái và hành động bằng cách cập nhật giá trị ước lượng của nó qua nhiều lần lặp:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right] \quad (2.33)$$

Mục tiêu cập nhật giá trị của cặp trạng thái, hành động (S_t, A_t) bằng phương pháp Q-learning dựa trên giá trị điểm thưởng nhận được R_{t+1} và giá trị của hành động lớn nhất ở trạng thái kế tiếp S_{t+1} . Hình 2.13 minh họa cho cách thức xác định mục tiêu cập nhật của Q-learning. Một điểm khác biệt của phương

pháp Q-learning so với các phương pháp trên là nó xấp xỉ trực tiếp hàm giá trị hành động tối ưu q_* , mà không phải phụ thuộc quá nhiều vào chính sách mà nó đang theo. Ảnh hưởng của chính sách mà hệ thông đang theo dõi với phương pháp này là: nó xác định cặp trạng thái và hành động nào được xuất hiện trong các mẫu thực nghiệm. Tuy nhiên để đảm bảo xác định được giá trị hành động tối ưu, một yêu cầu tối thiểu là các cặp trạng thái và hành động đều được xuất hiện trong quá trình đánh giá chính sách và giá trị các cặp trạng thái, hành động đều được cập nhật liên tục. Từng bước thực hiện của phương pháp Q-learning được trình bày trong thuật toán 2.7.

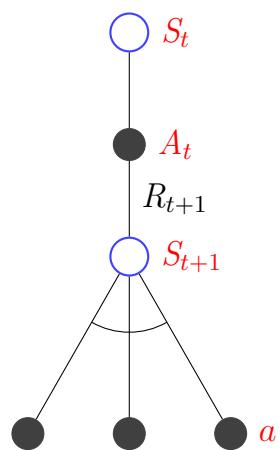
Thuật toán 2.7 Xác định hàm giá trị hành động tối ưu bằng Q-learning

Đầu vào:

Đầu ra: Hàm giá trị Q xấp xỉ hàm giá trị q_*

Thao tác:

- 1: **repeat**
 - 2: Tạo một mẫu thực nghiệm E bằng chính sách π' (thực hiện ϵ -greedy theo hàm giá trị Q hiện tại)
 - 3: Khởi tạo trạng thái S
 - 4: **for** mỗi bước trong E **do**
 - 5: Chọn một hành động A tại trạng thái S theo chính sách π'
 - 6: Thực hiện hành động A ; quan sát điểm thưởng R nhận được, và trạng thái tiếp theo S'
 - 7: $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ %% Thực hiện cập nhật giá trị cho cặp trạng thái, hành động S, A
 - 8: $S \leftarrow S'$
 - 9: **end for**
 - 10: **until** thỏa điều kiện dừng
-



Hình 2.13: Đồ thị minh họa cập nhật hàm giá trị hành động bằng phương pháp Q-learning. Q-learning xác định mục tiêu cập nhật cho giá trị của cặp trạng thái, hành động (S_t, A_t) bằng tổng giữa điểm thưởng R_{t+1} nhận được khi thực hiện hành động A_t tại trạng thái S_t và giá trị hành động a lớn nhất tại trạng thái kế tiếp S_{t+1} đã nhân với hệ số γ .

Chương 3

Kết hợp học tăng cường với học sâu

Những thành công gần đây của học sâu (Deep learning) trong các bài toán như xử lý ngôn ngữ tự nhiên, nhận diện đối tượng trong ảnh... đặt ra vấn đề: liệu các kỹ thuật trong học sâu có thể áp dụng vào học tăng cường? Để trả lời câu hỏi đó, chương này trình bày về hướng tiếp cận kết hợp học sâu với học tăng cường để áp dụng vào bài toán “Tự động chơi game”. Hướng tiếp cận mới mẻ này của lĩnh vực học tăng cường này mang tên “Học tăng cường sâu” (Deep reinforcement learning)

Chương này trình bày hai phần:

- Nguyên nhân cần sử dụng học sâu và kiến thức cơ bản về học sâu
- Áp dụng học tăng cường sâu vào bài toán “Tự động chơi game”

3.1 Học tăng cường kết hợp với học sâu

3.1.1 Học sâu

Lý do cần áp dụng học sâu

Những thuật toán học tăng cường được trình bày trong chương trước đều tìm chính sách tối ưu dựa vào hàm giá trị. Việc tính **đúng** và **nhanh** hàm giá trị ảnh hưởng rất nhiều đến kết quả của bài toán. Các thuật toán học tăng cường cổ điển như “Monte Carlo” (MC) hay “Temporal-Difference” (TD) đều đã được chứng minh là luôn hội tụ trong những điều kiện nhất định [13]. Ngoài ra, khi áp dụng vào các bài toán kinh điển của học tăng cường thì các thuật toán này đều hội tụ khá nhanh.

Tuy nhiên, với những bài toán thực tế với số trạng thái rất lớn thì việc lưu véc-tơ hàm giá trị trạng thái v_π (hoặc ma trận hàm giá trị hành động q_π) là việc không thể. Ví dụ như “frame hình” của bài toán tự động chơi game có kích thước $210 \times 160 \times 3 = 100800$ điểm ảnh; mỗi điểm ảnh có giá trị trong khoảng $[0, 127]$ nên số trạng thái có thể có lên đến 128^{100800} . Vì vậy, việc lưu trữ hàm giá trị dưới dạng bảng là không khả thi về mặt bộ nhớ. Còn về mặt tốc độ tính toán thì các thuật toán học tăng cường trên đều tính hàm giá trị *rồi rạc* cho từng trạng thái. Với số trạng thái quá lớn như trên thì ta không thể duyệt lần lượt từng trạng thái để tính được.

Những lý do trên dẫn đến việc sử dụng một phương pháp xấp xỉ hàm là bắt buộc cho các bài toán học tăng cường với số trạng thái lớn. Một trong những tiếp cận rất tự nhiên đó là sử dụng các mô hình học có giám sát như là một phương pháp xấp xỉ hàm giá trị. Đặc biệt, với những đột phá gần đây của học sâu trong lĩnh vực xử lý ảnh, video... thì việc áp dụng các mô hình phổ biến của học sâu vào bài toán tự động chơi game là đầy hứa hẹn.

Giới thiệu học sâu

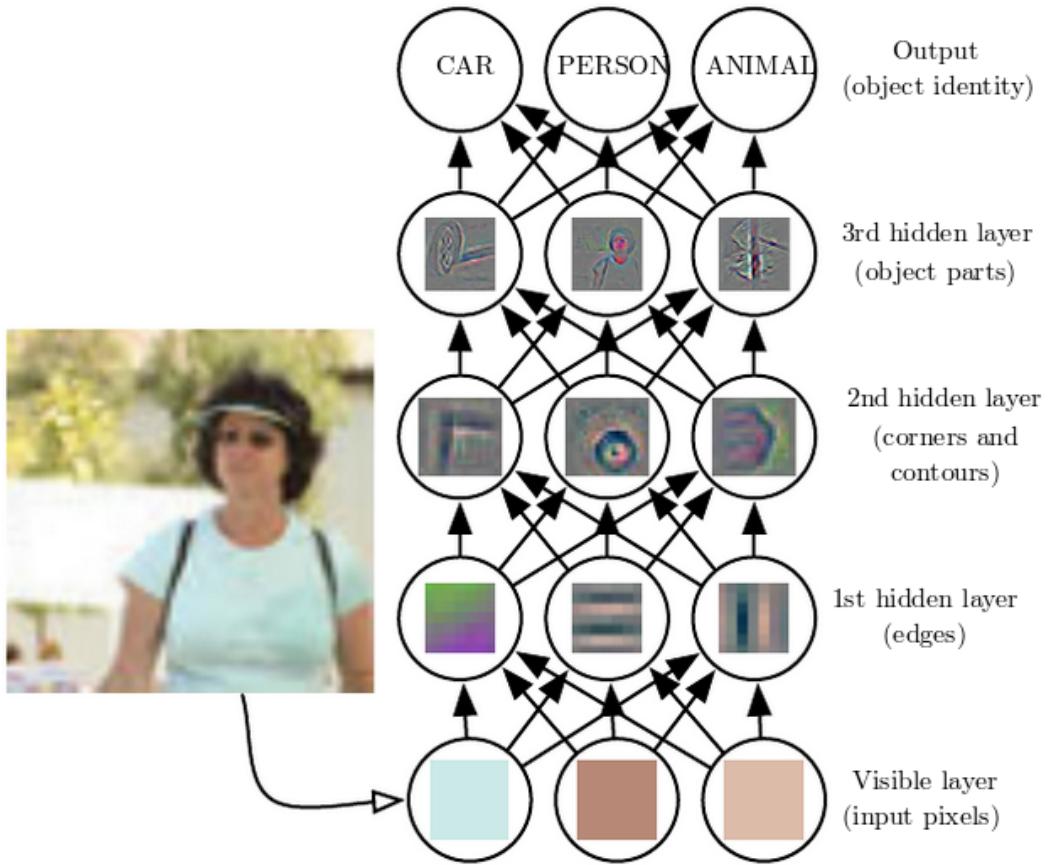
Các mô hình truyền thống trong lĩnh vực máy học như “Linear regression”, “Bayesian learning”... thông thường đều hoạt động trên các đặc trưng được rút

trích một cách thủ công (hand-designed features). Với dữ liệu thô thu thập được từ thực tế, các nhà khoa học xây dựng các phương pháp rút trích ra những “thông tin hữu ích” (thường được gọi là đặc trưng) để cung cấp cho các mô hình máy học. Kết quả nhận được từ các mô hình này phụ thuộc rất lớn vào cách biểu diễn dữ liệu. Ví dụ như trong bài toán nhận diện người nói từ một đoạn âm thanh, các đặc trưng có thể bao gồm: độ lớn âm thanh, tần số trung bình của đoạn âm,... Nếu các đặc trưng này không đủ “mạnh” (như có hai người nói đoạn âm thanh nhưng lại có chung độ lớn, tần số...) thì các mô hình học sẽ không thể phân biệt.

Để giải quyết vấn đề thiết kế đặc trưng, các mô hình máy học thuộc loại “Học biểu diễn” (*Representation learning*) ra đời. Các mô hình này có khả năng tự động học luôn các đặc trưng cần thiết cho quá trình phân tích dữ liệu. Nhờ vậy, các mô hình này có thể áp dụng được dễ dàng hơn vào các bài toán thực tế mà không cần con người phải can thiệp. Tuy nhiên, các đặc trưng cần thiết lại có thể rất phức tạp. Việc học ra các đặc trưng này có thể khó ngang với việc giải bài toán gốc. Ví dụ như trong bài toán nhận diện người nói trên, ta có thể sử dụng thông tin về giọng địa phương của người nói (accent). Đặc trưng này rất trừu tượng, không dễ phân tích bằng máy tính mặc dù con người có thể nhận biết một cách khá dễ dàng. Vì vậy, nếu các đặc trưng cần thiết quá khó để học thì các mô hình máy học thuộc loại “Học biểu diễn” cũng không thể cho kết quả tốt.

Học sâu (Deep learning) được ra đời nhằm giải quyết các vấn đề trên. Học sâu là một nhánh của “Học biểu diễn” nên vẫn có thể tự động học ra các đặc trưng hữu ích. Học sâu được thiết kế để học ra các đặc trưng có quan hệ với nhau theo nhiều tầng (layer). Các đặc trưng ở tầng phía sau được xây dựng dựa vào các đặc trưng ở tầng phía trước. Các tầng đầu tiên bao gồm những đặc trưng đơn giản và các tầng tiếp theo ngày càng trừu tượng, ngày càng phức tạp hơn. Mỗi tầng chỉ cần học cách xây dựng đặc trưng từ những đặc trưng ở tầng phía trước (đã có sẵn độ trừu tượng nhất định) thay vì học từ dữ liệu thô ban đầu; điều này giúp cho học sâu có khả năng học được những đặc trưng rất phức tạp.

Một trong những mô hình học sâu phổ biến nhất và cũng cho kết quả rất



Hình 3.1: Hình mô phỏng cách hoạt động của mô hình học sâu cho bài toán nhận diện đối tượng trong ảnh. Dữ liệu đầu vào là hình ảnh RGB chứa đối tượng cần xác định. Tầng đầu tiên của mô hình là tầng “input” tiếp nhận thông tin này dưới dạng ma trấn số. Các tầng tiếp theo ngoại trừ tầng cuối cùng được gọi là tầng ẩn “hidden layer” vì đặc trưng học được tại đây con người không quan sát được. Các tầng ẩn học các đặc trưng ngày càng trừu tượng dựa vào đặc trưng ở tầng phía trước. Tầng ẩn đầu tiên học được các đặc trưng về cạnh bằng cách so sánh độ sáng giữa các điểm ảnh gần nhau. Tầng ẩn thứ hai học được các đặc trưng về đường cong bằng cách tổng hợp đặc trưng về cạnh ở tầng trước đó. Tầng ẩn thứ ba học được các đặc trưng về bộ phận của đồ vật như khuôn mặt, bánh xe... dựa vào các đặc trưng về đường cong ở tầng trước. Tầng cuối cùng được gọi là tầng “output” có nhiệm vụ tìm kiếm các bộ phận và trả về lớp đối tượng tương ứng. Bằng cách học đặc trưng ngày càng trừu tượng hơn, các mô hình học sâu có khả năng tự động học được các đặc trưng từ đơn giản đến phức tạp; tất cả đều nhằm hỗ trợ cho quá trình phân lớp dữ liệu (hình được chỉnh sửa từ [3])

tốt với dữ liệu ảnh đó là mạng nơ-ron tích chập (Convolutional Networks). Với bài toán tự động chơi game, dữ liệu hệ thống nhận được từ môi trường là ảnh RGB. Chính vì vậy, mạng nơ-ron tích chập là một mô hình rất phù hợp để kết hợp với các thuật toán học tăng cường.

3.1.2 Mạng nơ-ron tích chập

Mạng nơ-ron tích chập (Convolutional Networks - CNN) là một mô hình học sâu được áp dụng rộng rãi trong các bài toán liên quan đến ảnh, video hoặc âm thanh... CNN được thiết kế để tận dụng thông tin về không gian (spatial structures) của các loại dữ liệu nêu trên. Ví dụ như trong bài toán nhận diện mặt người trong ảnh, thông tin về vị trí của mắt, mũi, miệng... là rất quan trọng. Nếu ta coi các điểm ảnh đều có ý nghĩa tương tự nhau thì ta đã bỏ quên thông tin về vị trí của những đặc trưng này. Ví dụ như khi ta tìm được hai mắt và miệng trong bức ảnh, ta có thể xác định vị trí tương đối của mũi là nằm ở giữa hai đặc trưng này. Trong ví dụ trên, nếu ta áp dụng các mô hình không quan tâm đến loại dữ liệu (coi từng thuộc tính dữ liệu đầu vào là độc lập) thì ta sẽ không tận dụng được thông tin về không gian trong đó.

Để tận dụng thông tin về không gian trong dữ liệu ảnh, CNN được thiết kế để học các đặc trưng trên một vùng nhỏ của ảnh. Các đặc trưng này được lưu trữ dưới dạng những bộ lọc (filter) thường có kích thước nhỏ hơn nhiều so với kích thước ảnh gốc. Các đặc trưng sau khi được học được áp dụng trên toàn bộ ảnh gốc để kiểm tra xem vị trí nào của ảnh xuất hiện đặc trưng này. CNN thực hiện phép kiểm tra này bằng cách “trượt” các bộ lọc này trên toàn bộ ảnh gốc. Phép “trượt” được thực hiện lần lượt từ trái qua phải và từ trên xuống dưới. Một cách tổng quát, phép “trượt” này có thể di chuyển không đồng đều theo hai chiều: ta có thể chỉ di chuyển qua phải một điểm ảnh nhưng lại di chuyển xuống dưới hai điểm ảnh. Tại mỗi vị trí, bộ lọc có nhiệm vụ kiểm tra thử đặc trưng được học có xuất hiện (hoặc mức độ rõ ràng của đặc trưng - đặc trưng xuất hiện nhiều hay ít) tại vị trí này không. Kết quả của phép kiểm tra này được lưu trữ lại dưới dạng một “bức ảnh” nhỏ hơn; giá trị mỗi “điểm ảnh” lúc này chính là kết quả của phép kiểm tra đặc trưng của bộ lọc. Do phép “trượt” được thực hiện theo thứ tự được nêu ở trên, các “điểm ảnh” kết quả vẫn mang thông tin

tương đối về vị trí: “điểm ảnh” bên trái ứng với kết quả của vùng nằm bên trái trong ảnh gốc và tương tự với điểm ảnh bên phải.

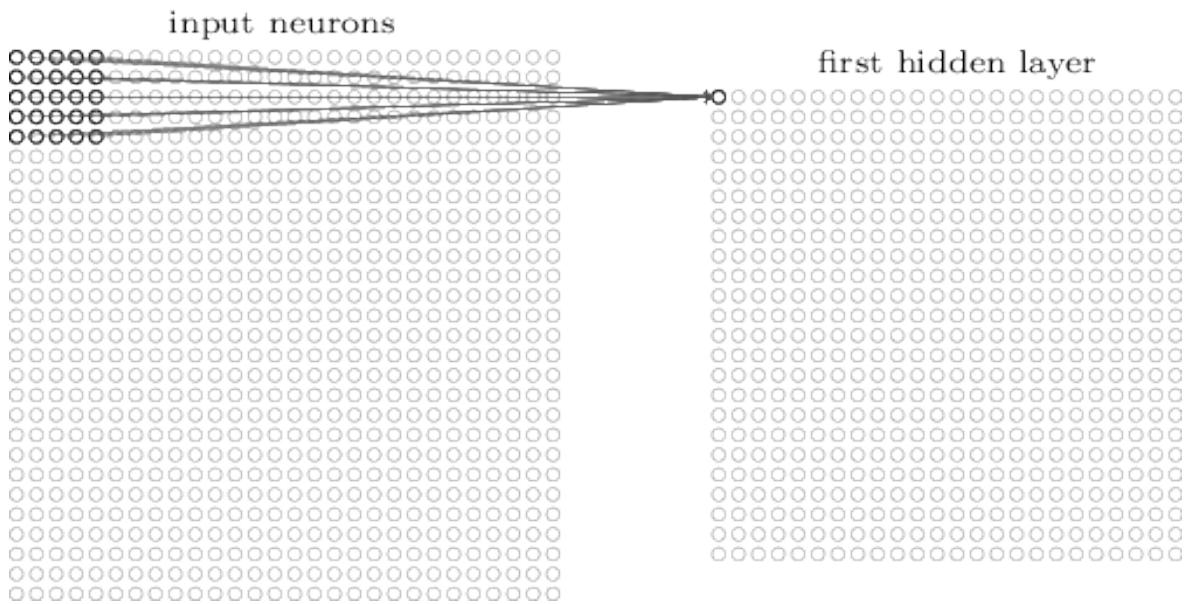
Phép “kiểm tra” đặc trưng của bộ lọc được thực hiện thông qua phép toán **tích chập** (convolution):

$$a_{i,j} = \sigma \left(b + \sum_{u=0}^{height} \sum_{v=0}^{width} w_{u,v} x_{i+u, j+v} \right) \quad (3.1)$$

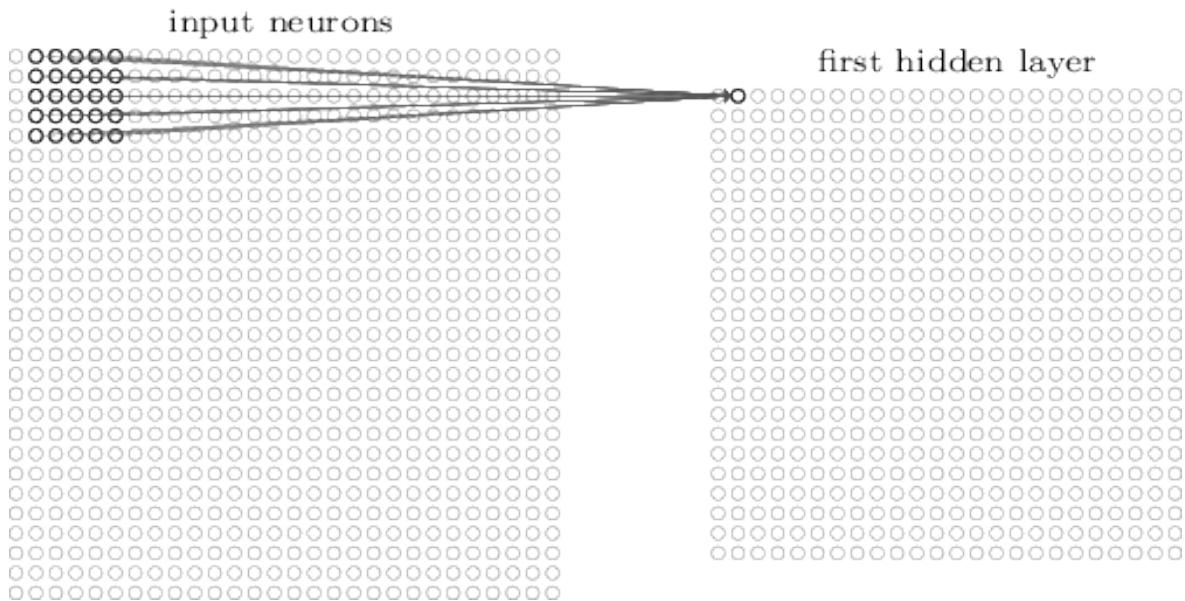
Trong đó:

- $a_{i,j}$ là kết quả của phép kiểm tra tại vùng có góc trái trên là i, j trong ảnh gốc.
- b là số thực gọi là “bias” của bộ lọc.
- $width, height$ lần lượt là chiều rộng và chiều cao của bộ lọc (5×5 trong hình (3.2)).
- w là ma trận trọng số có kích thước $height \times width$ của bộ lọc. $w_{u,v}$ là thành phần dòng u cột v của ma trận.
- $x_{i+u, j+v}$ là giá trị điểm ảnh đầu vào tại vị trí dòng $i+u$ và cột $j+v$.
- σ là một hàm phi tuyến được gọi là hàm kích hoạt (activation function).

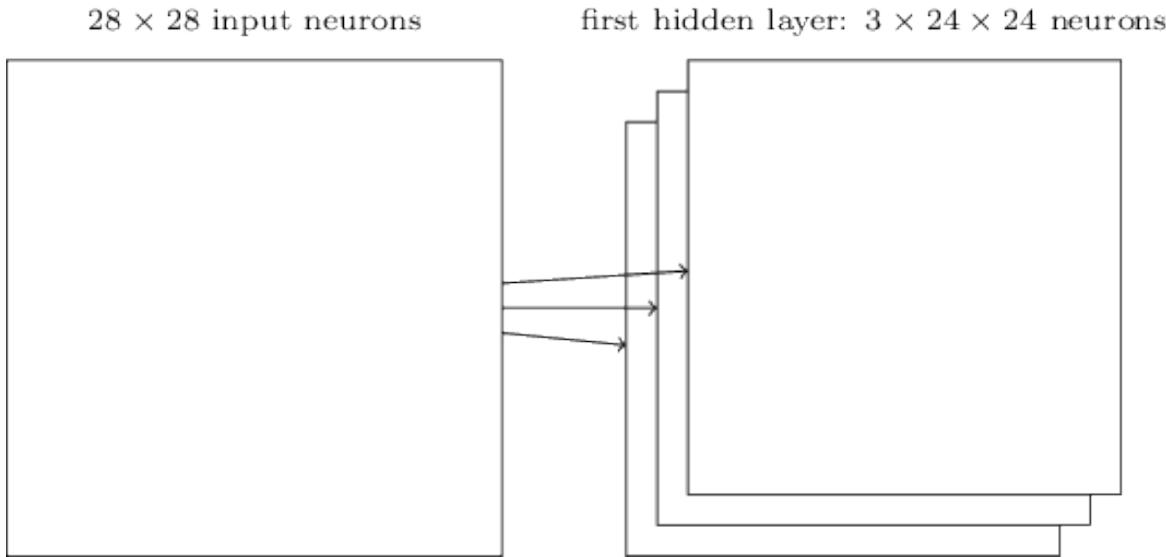
Phép toán tích chập này chỉ đơn giản là một tổng gồm các tích của trọng số và giá trị điểm ảnh. Hàm kích hoạt sau đó nhận giá trị tổng của phép tích chập để thực hiện phép biến đổi phi tuyến tính; nhờ có hàm kích hoạt, bộ lọc có thể học được những đặc trưng phi tuyến tính. Một trong những hàm phi tuyến hay được áp dụng đó là hàm “Rectified linear”: $\sigma(x) = \max(0, x)$. Công thức tích chập trên chứa hai thành phần cần được “học” đó là “bias” b và ma trận trọng số w . Hai thành phần này lưu trữ thông tin về đặc trưng mà bộ lọc học được. Một trong những đặc điểm quan trọng nhất là việc ma trận trọng số w và “bias” b được sử dụng chung cho việc “kiểm tra” đặc trưng tại mọi vị trí của ảnh gốc. Nhờ việc sử dụng chung ma trận trọng số và “bias” này, số lượng trọng số phải học của CNN được giảm đi rất nhiều. Cụ thể hơn, với ảnh gốc kích thước



Hình 3.2: Hình mô tả phép “kiểm tra” đặc trưng của bộ lọc có kích thước 5×5 tại vị trí đầu tiên (góc trái trên) của ảnh đầu vào. Kết quả của phép “kiểm tra” này được lưu lại thành một “điểm ảnh” trong “bức ảnh kết quả”. Lưu ý ảnh đầu vào trong ví dụ ở đây có kích thước 28×28 . Do bộ lọc chỉ kiểm tra những vùng nằm hoàn toàn trong ảnh nên kích thước của “bức ảnh” đầu ra là 24×24



Hình 3.3: Bộ lọc được trượt sang một điểm ảnh về phía bên phải để kiểm tra vùng cục bộ 5×5 bên cạnh. Kết quả của phép “kiểm tra” này được lưu lại thành một “điểm ảnh” bên phải của kết quả của vùng cục bộ trước đó. Thực hiện lần lượt quá trình “trượt” này đến hết ảnh ta sẽ có “bức ảnh kết quả” cuối cùng. “Bức ảnh kết quả” này mang thông tin về một đặc trưng cụ thể tại những vùng cục bộ liên tiếp nhau của ảnh gốc.



Hình 3.4: Hình mô tả tầng tích chập với ba bộ lọc học đặc trưng từ ảnh đầu vào. Mỗi bộ lọc lúc này học một bộ trọng số w và giá trị “bias” b khác nhau. Do các bộ lọc có cùng kích thước (5×5) nên “bức hình” kết quả cũng có cùng kích thước. Ta coi mỗi “bức hình” kết quả như một kênh (channel) khác nhau của bức hình và ghép chúng lại thành một bức hình lớn hơn gồm ba kênh. Các kênh này cũng giống như ba kênh màu khác nhau của ảnh RGB.

$28 \times 28 = 764$ như trong hình (3.2), nếu áp dụng mạng nơ-ron truyền thẳng thông thường để học đặc trưng từ toàn bộ ảnh, số lượng trọng số sẽ lên đến 764 cho một đặc trưng. Trong khi đó, CNN chỉ gồm khoảng $5 \times 5 = 25$ trọng số cần học.

Một bộ lọc của CNN chỉ học được một đặc trưng cụ thể. Tuy nhiên trong bài toán thực tế, ta cần nhiều đặc trưng khác nhau trên ảnh. Ví dụ như để nhận diện khuôn mặt trên ảnh, ta cần đặc trưng về mắt, miệng... Để học được nhiều đặc trưng khác nhau với CNN, ta chỉ việc thiết kế thêm nhiều bộ lọc học đặc trưng song song với nhau từ ảnh gốc. Với mỗi bộ lọc, ta cần học một ma trận trọng số và giá trị “bias” khác nhau. Ta gọi tập hợp các bộ lọc này là một tầng tích chập (convolution layer). Hình (3.4) mô tả tầng tích chập với ba bộ lọc.

CNN là một mô hình học sâu. Vì vậy để học được những đặc trưng có tính trừu tượng cao, CNN áp dụng phương pháp tổng hợp đặc trưng phức tạp từ những đặc trưng đơn giản hơn. Để tổng hợp đặc trưng, ta chỉ việc coi “bức hình kết quả” như là bức ảnh đầu vào và thêm tầng tích chập mới học đặc trưng từ bức hình này. Do kết quả của tầng tích chập trước mang thông tin về đặc trưng

được học bởi các bộ lọc, tầng tích chập tiếp theo thực hiện việc học từ các đặc trưng đã học ở tầng trước. Nhờ vậy, việc thêm vào các tầng tích chập tiếp theo sẽ giúp mô hình học được những đặc trưng ngày càng phức tạp hơn.

Phần tiếp theo sẽ trình bày cách sử dụng CNN như là một công cụ xấp xỉ hàm hỗ trợ cho các thuật toán học tăng cường. Do có khả năng học đặc trưng tự động, ta có thể thiết kế một công cụ xấp xỉ hàm “đa năng”: vừa học đặc trưng từ hình ảnh, vừa xấp xỉ hàm mong muốn từ những hình ảnh đó.

3.1.3 Sử dụng mạng nơ-ron để xấp xỉ hàm

Việc sử dụng mạng nơ-ron tích chập (hoặc mạng nơ-ron nói chung) để xấp xỉ hàm giá trị mang lại ba lợi ích quan trọng:

- Mạng nơ-ron có khả năng học được những đặc trưng phức tạp từ dữ liệu thô.
- Mạng nơ-ron có thể xấp xỉ hàm giá trị phức tạp.
- Mạng nơ-ron có thể tổng quát hoá từ trạng thái đã biết sang trạng thái chưa biết.

Mạng nơ-ron là một mô hình học sâu nên có thể học được những đặc trưng phức tạp từ dữ liệu thô như đã nói ở phần trên. Cụ thể hơn trong bài toán tự động chơi game, ta có thể đưa dữ liệu đầu vào là các “frame hình” RGB thẳng vào “input” của mạng nơ-ron để tính ra giá trị tương ứng. Nhờ vậy, ta không cần phải thiết kế các đặc trưng bằng tay để biểu diễn cho từng trạng thái của game. Ngoài ra, do qua trình học đặc trưng là hoàn toàn tự động, thuật toán học tăng cường lúc này có thể áp dụng cho bất kỳ game nào mà không cần phải thay đổi cách rút trích đặc trưng. Với một mô hình cố định lúc này, ta có thể học chơi được nhiều game. Đây cũng chính là mục đích của bài toán tự động chơi game: xây dựng mô hình có khả năng tự động học chơi *tốt* nhiều game chứ không chỉ chơi “*hoàn hảo*” một game.

Với bài toán tự động chơi game, giá trị của một trạng thái bất kỳ rất khó xác định do một màn game thường rất dài. Vì thế, hàm giá trị của bài toán này là một hàm phi tuyến phức tạp và không liên tục. Công cụ xấp xỉ hàm giá trị

phải có khả năng xấp xỉ những hàm phức tạp như vậy thì thuật toán học tăng cường mới đạt được hiệu quả. Với cách nhìn nhận mạng nơ-ron như là một công cụ xấp xỉ hàm, ta có thể thấy mạng nơ-ron rất linh hoạt với khả năng xấp xỉ hàm đích bất kỳ.

Một tính chất quan trọng khác của mạng nơ-ron chính là khả năng tổng quát hoá (generalization) từ trạng thái đã biết sang trạng thái chưa biết. Nhờ đặc điểm này mà quá trình học tăng cường được tăng tốc đáng kể. Thay vì phải duyệt qua từng trạng thái (thậm chí phải duyệt nhiều lần) để tính hàm giá trị tại đó, mạng nơ-ron có khả năng “dự đoán” giá trị của một trạng thái chưa từng thấy dựa vào những trạng thái đã thấy. Như trong bài toán tự động chơi game thì các “frame hình” liên tiếp thường rất giống nhau và các trạng thái này thường cũng có giá trị tương đương nhau. Vì vậy, khi học xong cách chơi một game nào đó, thuật toán học tăng cường vẫn hoạt động tốt trong quá trình kiểm thử khi gặp những tình huống game chưa từng thấy trong lúc huấn luyện.

Các thuật toán học tăng cường ở chương 2 đều lưu hàm giá trị dưới dạng bảng (lookup table). Để sử dụng mạng nơ-ron như một công cụ xấp xỉ hàm, lúc này ta coi hàm giá trị là một hàm có tham số (parameterized function) và tìm các tham số này:

$$\hat{v}(s; \theta) \approx v_\pi(s) \quad (3.2)$$

$$\hat{q}(s, a; \theta) \approx q_\pi(s, a) \quad (3.3)$$

θ là bộ trọng số của mạng nơ-ron mà ta cần học. Để học được bộ trọng số xấp xỉ tốt hàm đích ($v_\pi(s)$ hoặc $q_\pi(s, a)$), ta cung cấp các mẫu dữ liệu (data sample). Mỗi mẫu bao gồm dữ liệu đầu vào của mạng nơ-ron (tức trạng thái s hoặc bộ trạng thái, hành động s, a) cùng với giá trị đích mong muốn (tức $v_\pi(s)$ hoặc $q_\pi(s, a)$). Khi ta cung cấp đủ nhiều mẫu dữ liệu cho mạng nơ-ron, các bộ tham số sẽ được thay đổi để mạng xấp xỉ được hàm đích mong muốn. Số lượng mẫu dữ liệu càng lớn thì mạng nơ-ron càng “thấy” được nhiều giá trị tại nhiều vị trí khác nhau của hàm đích hơn, khi đó mạng nơ-ron càng có khả năng xấp xỉ hàm đích tốt hơn. Để xét xem mạng nơ-ron có xấp xỉ tốt hàm đích hay chưa, ta có thể tính độ “khác biệt” của giá trị đích với giá trị xấp xỉ trên cả không gian đầu

vào:

$$J(\theta) = \mathbb{E}_{s \sim \pi}[(v_\pi(s) - \hat{v}(s; \theta))^2] \quad (3.4)$$

$$J(\theta) = \mathbb{E}_{s, a \sim \pi}[(q_\pi(s, a) - \hat{q}(s, a; \theta))^2] \quad (3.5)$$

Kỳ vọng $\mathbb{E}_{s \sim \pi}$ ý chỉ kỳ vọng với biến ngẫu nhiên là trạng thái s được lấy từ phân bố do chính sách π tạo nên. Ví dụ như ta cần xấp xỉ hàm giá trị của một chính sách chỉ luôn “đi qua trái” thì s sẽ là những trạng thái mà khi đi theo chính sách này, ta có thể gặp được s . Trong khi đó nếu chính sách đang xét là “đứng yên” (không thay đổi trạng thái) thì s chỉ có thể có một trạng thái duy nhất đó là trạng thái bắt đầu. Giá trị nằm trong kỳ vọng là độ lỗi bình phương giữa giá trị xấp xỉ và giá trị đích. Với hàm lỗi bình phương, có thể thấy khi $J(\theta)$ càng nhỏ thì bộ trọng số θ giúp cho mạng nơ-ron xấp xỉ hàm đích càng tốt. Lý do chính ta chọn độ lỗi bình phương (ví dụ thay vì độ lỗi theo trị tuyệt đối) là vì việc tính toán (như đạo hàm) trên hàm bình phương khá dễ dàng.

Lưu ý là ở đây, hàm lỗi $J(\theta)$ là hàm theo θ ; tức là lúc này, ta cố định bộ trọng số θ để tìm ra “sai số” mà bộ trọng số này gây ra. Nếu ta có hai bộ trọng số θ_1 và θ_2 , ta có thể so sánh khả năng xấp xỉ hàm đích của chúng bằng cách so sánh giá trị $J(\theta_1)$ và $J(\theta_2)$ tương ứng.

Tuy nhiên ta không thể tính độ lỗi trên mọi điểm dữ liệu đầu vào theo công thức (3.4) và (3.5) được vì số lượng trạng thái là rất lớn. Vậy ta có thể xấp xỉ giá trị $J(\theta)$ bằng cách chỉ xét độ lỗi trên một “tập huấn luyện” (hay còn gọi là “Batch”) các trạng thái mà ta biết được giá trị đích. Khi đó, kỳ vọng $\mathbb{E}_{s \sim \pi}$ được thay thế bằng giá trị trung bình độ lỗi trên từng mẫu dữ liệu của “batch”:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (v_\pi(S_i) - \hat{v}(S_i; \theta))^2 \quad (3.6)$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (q_\pi(S_i, A_i) - \hat{q}(S_i, A_i; \theta))^2 \quad (3.7)$$

Trong đó:

- N là số mẫu dữ liệu trong “batch”.

- S_i, A_i tương ứng là trạng thái và hành động của mẫu dữ liệu thứ i của “batch”.

Với một tập huấn luyện, ta mong muốn tìm được bộ trọng số giúp cho mạng nơ-ron xấp xỉ tốt hàm đích trên các mẫu dữ liệu thuộc tập huấn luyện này. Một thuật toán đơn giản và hay được sử dụng để tìm bộ trọng số này đó là “Batch Gradient Descent” (BGD). Để cực tiểu hóa hàm $J(\theta)$, thuật toán BGD thực hiện lặp lại nhiều “bước đi” nhỏ để thay đổi bộ trọng số θ dần dần; mỗi bước đi sẽ giúp cho hàm $J(\theta)$ giảm đi một ít. Để chọn “hướng đi” (tức cách cập nhật θ) thì BGD sẽ “nhìn” xung quanh vị trí hiện tại và đi theo hướng nào giúp giảm $J(\theta)$ nhiều nhất có thể. Hướng đi này chính là ngược hướng véc-tơ đạo hàm riêng (tức “gradient”) của hàm $J(\theta)$ tại θ . Như vậy, thuật toán BGD thực hiện lặp lại nhiều lần việc cập nhật bộ trọng số θ theo công thức:

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} J(\theta) \quad (3.8)$$

Trong đó:

- θ_t là bộ trọng số tại bước thứ t .
- α là hệ số học (learning rate); giá trị này dùng để điều khiển độ lớn của “bước đi”.
- $\nabla_{\theta} J(\theta)$ là véc-tơ đạo hàm riêng của hàm $J(\theta)$ tại vị trí θ

Do ta chỉ “nhìn” cục bộ tại ví trí hiện tại nên nếu ta đi một bước quá dài (hệ số học α lớn) thì giá trị hàm lỗi J tại điểm đến sẽ không chắc là sẽ giảm; còn nếu bước đi quá ngắn thì mỗi bước chỉ giảm J một ít, khi đó ta sẽ tốn rất nhiều thời gian để J đạt giá trị cực tiểu.

Một thuật toán cải tiến của BGD hay được sử dụng đó là “Stochastic Gradient Descent” (SGD). Điểm yếu của thuật toán BGD là ta cần phải tính véc-tơ đạo hàm riêng cho **tất cả** các mẫu trong “batch” để cập nhật được một lần cho bộ trọng số. Thuật toán SGD khắc phục điểm yếu này bằng cách chọn ngẫu nhiên một số mẫu dữ liệu trong “batch” (gọi là “mini-batch”), tính véc-tơ đạo hàm riêng trung bình trên “mini-batch” này và thực hiện cập nhật bộ trọng số. Lúc

này, véc-tơ đạo hàm riêng trung bình trên “mini-batch” có thể coi là một xấp xỉ của véc-tơ đạo hàm riêng trung bình trên toàn bộ tập huấn luyện. Do véc-tơ này chỉ tính trên “mini-batch” nên khi cập nhật, giá trị $J(\theta)$ có thể tăng. Tuy nhiên, khi cập nhật nhiều lần thì xu hướng chung là hàm lỗi $J(\theta)$ sẽ giảm. Khi tập huấn luyện càng lớn thì lợi thế của thuật toán SGD càng rõ. Với tập huấn luyện gồm 1000 mẫu thì thuật toán BGD chỉ cập nhật trọng số được **một** lần sau khi duyệt qua hết dữ liệu; trong khi đó, thuật toán SGD với kích thước “mini-batch” là 10 có thể cập nhật được 100 lần. Kích thước của “mini-batch” lúc này ảnh hưởng đến độ chính xác của véc-tơ đạo hàm riêng của $J(\theta)$. Nếu kích thước quá nhỏ thì véc-tơ đạo hàm riêng trung bình trên “mini-batch” sẽ không xấp xỉ tốt véc-tơ đạo hàm riêng trung bình trên toàn bộ tập huấn luyện. Nếu kích thước quá lớn thì lợi thế về tốc độ của SGD so với BGD sẽ không còn cao.

Như đã đề cập ở chương 2, để cải tiến chính sách khi không có đầy đủ thông tin về môi trường (tức không có các ma trận của MDP) ta cần xấp xỉ q_π thay vì v_π . Áp dụng thuật toán SGD để cực tiểu hóa hàm lỗi (3.7), công thức cập nhật bộ trọng số tại thời điểm t có dạng:

$$\theta_{t+1} = \theta_t - \Delta\theta_t \quad (3.9)$$

$$= \theta_t - \alpha \nabla_{\theta_t} J(\theta_t) \quad (3.10)$$

$$= \theta_t - \alpha \nabla_{\theta_t} \left(\frac{1}{B} \sum_{i=1}^B (q_\pi(S_i, A_i) - \hat{q}(S_i, A_i; \theta_t))^2 \right) \quad (3.11)$$

$$= \theta_t - \alpha \frac{1}{B} \sum_{i=1}^B (q_\pi(S_i, A_i) - \hat{q}(S_i, A_i; \theta_t)) \nabla_{\theta_t} \hat{q}(S_i, A_i; \theta_t) \quad (3.12)$$

Ở đây:

- $\Delta\theta$ là giá trị cập nhật tại cho một “mini-batch”.
- B là kích thước của “mini-batch”
- $q_\pi(S_i, A_i)$ là giá trị **thật sự** của hành động A_i tại trạng thái S_i khi thực hiện theo chính sách π
- $\hat{q}(S_i, A_i; \theta_t)$ là giá trị **xấp xỉ** của hành động A_i tại trạng thái S_i khi thực

hiện theo chính sách π

- $\nabla_{\theta_t} \hat{q}(S_i, A_i; \theta_t)$ là véc-tơ đạo hàm riêng của hàm \hat{q} tại trạng thái S_i và hành động A_i

Công thức trên bao gồm một tổng các số hạng có thể được tính riêng lẻ cho từng mẫu dữ liệu. Vậy để đơn giản hóa công thức, ta viết lại công thức trên cho duy nhất một mẫu dữ liệu; khi cần thiết lập công thức tổng quát cho một “mini-batch”, ta chỉ việc tính giá trị trung bình của nhiều mẫu. Lúc này, công thức mới sẽ ứng với trường hợp kích thước “mini-batch” đúng bằng một nên ta có thể bỏ chỉ số i của từng mẫu dữ liệu:

$$\theta_{t+1} = \theta_t - \alpha(q_\pi(S, A) - \hat{q}(S, A; \theta_t)) \nabla_{\theta_t} \hat{q}(S, A; \theta_t) \quad (3.13)$$

3.1.4 Học tăng cường kết hợp với xấp xỉ hàm

Áp dụng thuật toán SGD để tối ưu hàm lỗi J theo công thức (3.12) thì ta sẽ cực tiểu hoá được độ khác biệt giữa hai hàm q_π và \hat{q} . Tuy nhiên, giá trị đích thực sự mà ta mong muốn là $q_\pi(s, a)$ ta không biết được mà chỉ có một số **mẫu** S_t, A_t khi tương tác với môi trường. Các thuật toán học tăng cường chính là kỹ thuật giúp ta có được một ước lượng đơn giản của giá trị đích này. Như ở chương 2, ta có hai thuật toán để ước lượng hàm giá trị: thuật toán “Monte-Carlo” (MC) và thuật toán “Temporal Difference” (TD).

Với thuật toán MC, ta chỉ việc thay thế $q_\pi(s, a)$ bằng một ước lượng không chêch của giá trị này. Theo định nghĩa:

$$q_\pi(s, a) = \mathbb{E}_\pi [R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T \mid S_t = s, A_t = a] \quad (3.14)$$

Vậy ta có thể lấy $R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$ làm một ước lượng cho $q_\pi(s, a)$. Giá trị này có thể dễ dàng có được bằng cách cho hệ thống tương tác với môi trường đến khi kết thúc một màn (tức một “episode”). Sau đó với mỗi trạng thái S_t của “episode”, ta chỉ cần tính tổng điểm thưởng đến cuối “episode” để có giá trị ước lượng mong muốn. Đây là một ước lượng không chêch do kỳ vọng của biểu thức này bằng đúng $q_\pi(s, a)$. Vậy công thức cập nhật bộ trọng số trong

(3.13) được thay bằng:

$$\theta_{t+1} = \theta_t - \alpha(G_t - \hat{q}(S, A; \theta_t)) \nabla_{\theta_t} \hat{q}(S, A; \theta_t) \quad (3.15)$$

G_t ở đây là tổng điểm thưởng (tức “Returns”) nhận được khi thực hiện hành động A tại trạng thái S . Kết hợp thuật toán MC với thuật toán SGD để cực tiểu hóa hàm lỗi của mạng nơ-ron, ta có được một thuật toán học tăng cường kết hợp học sâu hoàn chỉnh để giải các bài toán có số trạng thái lớn.

Với ý tưởng tương tự, để sử dụng mạng nơ-ron để xấp xỉ hàm giá trị cho thuật toán “Q-learning”, ta sử dụng tổng điểm thưởng được “bootstrap” để ước lượng $q_\pi(s, a)$. Cụ thể hơn, ta sẽ thay thế giá trị $q_\pi(s, a)$ trong công thức (3.14) bằng $R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a; \theta)$. Có thể thấy rằng, giá trị “bootstrap” này là một ước lượng chêch của $q_\pi(s, a)$. Tuy vậy, ước lượng này chỉ gồm tổng của hai số hạng R_{t+1} và $\gamma \max_a \hat{q}(S_{t+1}, a; \theta)$ nên có phương sai thấp hơn nhiều so với ước lượng của MC (gồm tổng của nhiều số hạng R_{t+1}, R_{t+2}, \dots). Đây cũng chính là một sự “thoả hiệp” (trade-off) giữa hai giá trị “bias” và “variance” của hai thuật toán MC và “Q-learning”. Thuật toán MC sử dụng một ước lượng không chêch nên có “bias” bằng không nhưng “variance” (tức phương sai) càng cao; khi đó tổng điểm thưởng của cùng một bộ (S, A) sẽ thay đổi rất nhiều. Khi các giá trị này thay đổi quá nhanh thì mạng nơ-ron sẽ học chậm hơn. Trong khi đó “Q-learning” sử dụng một ước lượng chêch nên có “bias” lớn nhưng phương sai lại nhỏ; khi đó giá trị “bootstrap” của cùng một bộ (S, A) sẽ ít thay đổi trong những lần duyệt đến khác nhau. Tương tự như thuật toán MC, công thức cập nhật bộ trọng số trong (3.13) được thay bằng:

$$\theta_{t+1} = \theta_t - \alpha(R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a; \theta) - \hat{q}(S, A; \theta_t)) \nabla_{\theta_t} \hat{q}(S, A; \theta_t) \quad (3.16)$$

Mã giả của thuật toán “Q-learning” với xấp xỉ hàm được trình bày ở (3.1).

Thuật toán 3.1 “Q-learning” kết hợp với xấp xỉ hàm

Đầu vào: Số “episode” cần thực hiện để cập nhật bộ trọng số mạng nơ-ron

Đầu ra: Bộ trọng số θ của mạng nơ-ron

Thao tác:

- 1: Khởi tạo ngẫu nhiên bộ trọng số θ của mạng nơ-ron
 - 2: **repeat**
 - 3: Tương tác với môi trường dựa vào chính sách có hàm giá trị được xấp xỉ bởi $\hat{q}(\cdot; \theta)$ đến khi kết thúc “episode” để có được tập các mẫu dữ liệu $S_1, A_1, R_2, S_2, A_2, R_3, \dots, S_{T-1}, A_{T-1}, R_T$
 - 4: Chia tập dữ liệu trên thành các “mini-batch” gồm B mẫu dữ liệu có dạng S_i, A_i, R_{i+1}
 - 5: **for** mỗi “mini-batch” của tập dữ liệu **do**
 - 6: Cập nhật θ theo công thức (3.16) cho toàn bộ các phần tử trong “mini-batch”
 - 7: **end for**
 - 8: **until** Thực hiện đủ số “episode”
-

3.2 Kết hợp học tăng cường với học sâu vào bài toán tự động chơi game

Phần này trình bày cách kết hợp thuật toán học tăng cường với học sâu vào bài toán tự động chơi game. Đầu tiên, chúng em sẽ trình bày về cấu trúc mạng “Deep Q-Network” [10] - cấu trúc mạng nơ-ron tích chập kết hợp với mạng nơ-ron truyền thẳng được thiết kế riêng biệt cho bài toán tự động chơi game. Phần tiếp theo trình bày hai kỹ thuật quan trọng giúp tăng tính ổn định của quá trình học. Phần cuối cùng đề cập đến vấn đề “đánh giá quá cao” (overestimation) ảnh hưởng thế nào lên kết quả của hệ thống cũng như cách thức giải quyết vấn đề này.

3.2.1 “Deep Q-Network”

Để có được cấu trúc mạng “Deep Q-Network” [10] hoàn chỉnh và hoạt động tốt cho bài toán tự động chơi game, ta kết hợp thuật toán “Q-learning” với mạng nơ-ron tích chập có cấu trúc phù hợp với bài toán. Trong bài toán tự động chơi game, dữ liệu đầu vào mà hệ thống nhận được từ môi trường tại mỗi thời điểm

là một ảnh RGB có kích thước 210×160 . Ta có thể đưa cả hình ảnh này làm dữ liệu đầu vào cho mạng nơ-ron tích chập; tuy nhiên với kích thước khá lớn như vậy, việc huấn luyện hệ thống sẽ tốn nhiều thời gian. Để tăng tốc độ huấn luyện lên, ta có thể thực hiện việc thu nhỏ (scale) ảnh về kích thước nhỏ hơn. Việc tiền xử lý ảnh bằng cách thu nhỏ có thể làm mất mát thông tin, tuy nhiên thực nghiệm cho thấy hệ thống vẫn có độ chính xác cao.

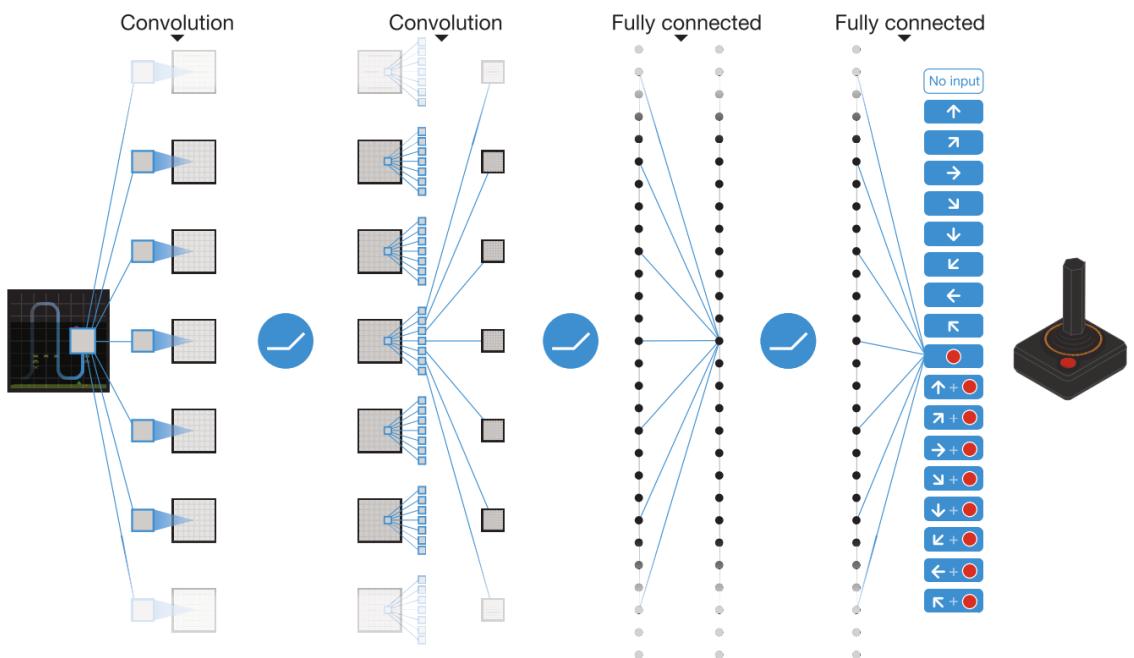
Một điểm quan trọng trong bài toán tự động chơi game là hình ảnh tại mỗi thời điểm không mô tả hết thông tin cần thiết. Ví dụ như khi có hình của một quả bóng, ta không biết hướng và vận tốc hiện tại của quả bóng. Để giải quyết điều này, ta có thể ghép hình ảnh của nhiều thời điểm liên tiếp lại theo thứ tự thời gian. Khi đó một trạng thái S_t sẽ gồm nhiều hình ảnh của các thời điểm liên nhau và chứa luôn cả thông tin về thời gian.

Với cách thiết kế mạng nơ-ron thông thường, ta cần cung cấp trạng thái S và hành động A để tính được giá trị xấp xỉ $\hat{q}(S, A; \theta)$. Cách thiết kế này không phù hợp cho thuật toán “Q-learning”: ta cần phải lan truyền tiến nhiều lần để tính giá trị $\max_a \hat{q}(S, a; \theta)$ trong công thức (3.16). Để khắc phục nhược điểm này, ta chỉnh lại cấu trúc mạng để “input” là trạng thái S và “output” là giá trị của mọi hành động tại trạng thái này: $\hat{q}(S, a; \theta), \forall a$. Với cấu trúc này, ta chỉ việc lan truyền tiến một lần và lấy max giá trị “output” của mạng nơ-ron. Như vậy, số nơ-ron của tầng “output” là số lượng hành động có thể có.

Để tổng hợp các đặc trưng được học từ các tầng tích chập, các tầng ẩn cuối cùng của mạng nơ-ron sẽ là các tầng “fully-connected”. Các tầng tích chập có chức năng tìm kiếm các đặc trưng cục bộ cần thiết còn tầng “fully-connected” có nhiệm vụ tổng hợp các đặc trưng đó trên **toàn ảnh**. Tầng “output” cũng là một tầng “fully-connected” với số nơ-ron là số hành động có thể có. Lưu ý là do tầng này trả về kết quả là giá trị của từng hành động nên ta không áp dụng hàm kích hoạt tại đây. Hình (3.5) mô tả cấu trúc mạng “Deep Q-Network”.

3.2.2 Kỹ thuật làm tăng tính ổn định

Đặc điểm của thuật toán “Q-learning” khi áp dụng xấp xỉ hàm đó là thuật toán không chắc sẽ hội tụ [13]. Hội tụ ở đây ý chỉ sau một số bước cập nhật hữu hạn, giá trị hành động của trạng thái sẽ hội tụ về giá trị nhất định. Do đặc điểm này,



Hình 3.5: Hình mô tả cấu trúc mạng “Deep Q-Network” [10]. Hai tầng ẩn đầu tiên là tầng tích chập với hàm kích hoạt “rectified linear”. Tầng ẩn tiếp theo là tầng “fully-connected” có nhiệm vụ tổng hợp đặc trưng trên toàn ảnh của tầng trước. Tầng “output” cũng là tầng “fully-connected” trả về kết quả là giá trị của từng hành động ứng với trạng thái đầu vào. Tầng “output” không có hàm kích hoạt.

nếu ta huấn luyện mạng nơ-ron bằng thuật toán SGD thông thường thì nhiều khả năng “Q-learning” sẽ không hội tụ dẫn đến chính sách tương ứng sẽ không được tốt. Để giải quyết vấn đề này, chúng em áp dụng hai kỹ thuật giúp tăng tính ổn định và khả năng hội tụ: kỹ thuật “Experience replay” [8] và kỹ thuật cố định “Q-target” [9].

“Experience replay”

Thuật toán “Q-learning” kết hợp với xấp xỉ hàm cho phép ta học “online”: sau mỗi bước tương tác với môi trường và nhận được dữ liệu mới, ta có thể cập nhật ngay bộ trọng số mạng nơ-ron. Tuy nhiên, cách học này có thể gây ra vấn đề không hội tụ. Lý do là các mẫu dữ liệu liên tiếp nhau thường có tương quan (correlation) với nhau rất lớn. Ngoài ra, các mẫu dữ liệu phía sau bị ảnh hưởng bởi chính sách (tương ứng với bộ trọng số mạng nơ-ron) ngay trước đó. Ví dụ như chính sách hiện tại là “luôn qua trái” thì các mẫu dữ liệu liên tiếp sẽ lấy từ phân bố “qua trái”. Khi chính sách thay đổi sang “luôn qua phải” thì tất cả các mẫu dữ liệu mới sẽ lấy từ phân bố “qua phải”. Khi các mẫu dữ liệu có tương quan lớn và phân bố dữ liệu thay đổi nhanh thì “Q-learning” sẽ khó hội tụ [10].

Kỹ thuật “experience replay” [8] giải quyết vấn đề này bằng cách lưu trữ lại một tập các mẫu dữ liệu nhận được gần đây nhất. Sau đó, mỗi lần muốn cập nhật bộ trọng số, ta chỉ việc lấy một ngẫu nhiên một “mini-batch” từ tập dữ liệu này để học. Ý tưởng của kỹ thuật này giống như việc ta “gọi nhớ” lại các kinh nghiệm đã có từ trước đó và học lại từ chúng. Cách làm này cũng mang tính “cứng cố” lại những kinh nghiệm mà hệ thống học được trong quá khứ. Do việc lấy ngẫu nhiên nên các mẫu dữ liệu trong “mini-batch” sẽ không còn tương quan lớn với nhau. Ngoài ra, việc lưu trữ lại giúp các mẫu dữ liệu có thể được học lại nhiều lần. Điều này giúp hệ thống học được nhiều hơn mà không phải tương tác với môi trường quá nhiều.

Cố định “Q-target”

Mỗi khi mạng nơ-ron được cập nhật bộ trọng số thì chính sách hiện tại sẽ thay đổi. Khi chính sách thay đổi thì hàm q_π mục tiêu cũng thay đổi theo. Thực tế cho thấy, việc thay đổi nhỏ của bộ trọng số sẽ gây ra thay đổi lớn của chính

sách [9]. Để thấy được điều này, ta có thể thay đổi trọng số ở những tầng đầu tiên của mạng nơ-ron. Khi đó, trọng số này sẽ ảnh hưởng dây chuyền đến những tầng sau làm cho “output” của mạng nơ-ron bị thay đổi nhiều. Do vậy, mỗi khi cập nhật bộ trọng số, hàm mục tiêu q_π tại một bộ trạng thái s và hành động a nhiều khả năng sẽ thay đổi lớn. Do giá trị mục tiêu của việc xấp xỉ hàm thay đổi quá nhiều và liên tục, mạng nơ-ron có thể không xấp xỉ được hàm giá trị như ta mong muốn.

Một kỹ thuật được đề xuất bởi [10] có tên **cố định “Q-target”** (fix Q-target) nhằm giải quyết vấn đề này. Thay vì thay đổi giá trị mục tiêu một cách liên tục, ta có thể cố định lại hàm mục tiêu q_π trong một số bước cập nhật trọng số để mạng nơ-ron xấp xỉ tốt hàm này. Sau một khoảng thời gian, ta mới thay đổi hàm mục tiêu thành hàm ứng với chính sách hiện tại. Các lần cập nhật trọng số tiếp theo lại tiếp tục cố gắng xấp xỉ tốt hàm mục tiêu mới. Quá trình này được lặp lại nhiều lần để hàm mục tiêu cuối cùng mà ta có cũng chính là hàm giá trị tối ưu.

Để có cái nhìn trực quan hơn về kỹ thuật này, ta có thể lấy ví dụ về việc xây một căn nhà. Do quá trình thiết kế ban đầu chưa tính trước đến những trường hợp phát sinh (như thiếu vật liệu, thợ không đủ trình độ...) trong lúc xây dựng nên ta phải điều chỉnh thiết kế lại liên tục. Ta có thể vừa xây dựng căn nhà (ứng với việc thay đổi trọng số mạng nơ-ron) vừa điều chỉnh lại thiết kế căn nhà cho phù hợp (ứng với việc thay đổi chính sách mục tiêu). Tuy nhiên, khi thiết kế căn nhà thay đổi quá nhanh và liên tục thì việc xây dựng theo sẽ gặp nhiều khó khăn và có khả năng khi hoàn thành kết quả sẽ không tốt. Kỹ thuật cố định “Q-target” tương tự với việc ta cố gắng xây dựng hoàn tất một phần của căn nhà (tức xấp xỉ tốt hàm giá trị một chính sách cũ) trước khi thay đổi thiết kế. Cứ như vậy, mỗi khi hoàn tất một phần, ta thay đổi thiết kế (ứng với việc thay đổi hàm mục tiêu) rồi lại xây dựng tiếp. Cách làm này giống như tạo dựng một số các “điểm kiểm soát” (checkpoint) với những mục tiêu xác định. Do các “điểm kiểm soát” này cố định nên việc thay đổi kiến trúc căn nhà sẽ không quá nhiều, dẫn đến quá trình xây dựng sẽ dễ dàng đạt được mục tiêu hơn.

Tuy nhiên, kỹ thuật này dẫn đến một vấn đề là ta cần phải thay đổi hàm mục tiêu sau bao nhiêu lần cập nhật trọng số mạng nơ-ron (tức tần suất thay

đổi hàm mục tiêu)? Nếu tần suất quá cao (tức thay đổi hàm quá nhiều) thì kỹ thuật này không mang lại được nhiều hiệu quả. Nếu tần suất quá thấp (tức ít khi thay đổi hàm mục tiêu) thì việc học sẽ rất chậm do ta phải cập nhật trọng số rất nhiều để xấp xỉ một hàm giá trị cũ; khi cập nhật hàm mục tiêu thì chính sách mới chỉ tốt lên một chút so với chính sách cũ. Vì vậy, ta cần phải chọn một giá trị vừa phải, phù hợp với bài toán tự động chơi game.

3.2.3 Vấn đề “đánh giá quá cao” của thuật toán “Q-learning”

Thuật toán “Q-learning” cải tiến từ thuật toán “Sarsa(0)” bằng cách lấy max giá trị các hành động của trạng thái tiếp theo. Nhắc lại công thức cập nhật của “Q-learning” là:

$$\theta_{t+1} = \theta_t - \alpha \left[R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a; \theta) - \hat{q}(S, A; \theta_t) \right] \nabla_{\theta_t} \hat{q}(S, A; \theta_t) \quad (3.17)$$

Nhờ việc lấy max mà thuật toán “Q-learning” hội tụ nhanh hơn [13]: ta luôn ước lượng giá trị của trạng thái tiếp theo dựa vào hành động tốt nhất tại đó. Đây có thể hiểu là một cách nhìn “lạc quan” về tương lai: nếu có nhiều trường hợp tốt, xấu khác nhau thì ta luôn giả sử điều tốt nhất sẽ xảy ra. Tuy việc giả sử này có thể không đúng trong mọi tình huống, do các trạng thái tiếp theo sau đó cũng được cập nhật lại giá trị nên thuật toán “Q-learning” vẫn hội tụ chính xác về chính sách tối ưu.

Rõ ràng là cách nhìn “lạc quan” này cũng có điểm xấu: ta sẽ có cái nhìn quá “lạc quan” (optimistic) về tương lai và dễ dẫn đến khả năng đánh giá quá cao (overestimation) giá trị của trạng thái kế tiếp. Nhất là khi áp dụng xấp xỉ hàm, giá trị ước lượng hiện tại của mạng nơ-ron chưa được chính xác. Khi đó nếu mạng nơ-ron ước lượng cao hơn giá trị thật sự thì vấn đề đánh giá quá cao sẽ xảy ra nhiều hơn. Mặc dù vậy, ta vẫn cần đặt ra câu hỏi: liệu việc đánh giá quá cao có gây ảnh hưởng đến kết quả của chính sách cuối cùng hay không? Rõ ràng là nếu tất cả hành động đều bị đánh giá quá cao (ví dụ như tăng một lượng bằng nhau C) thì việc lựa chọn hành động không bị ảnh hưởng (do ta vẫn chọn hành động tốt nhất theo thuật toán greedy). Tuy nhiên, nghiên cứu [?] chỉ ra

rằng thuật toán “Q-learning” với xấp xỉ hàm gặp hiện tượng đánh giá quá cao không đồng đều. Điều này làm cho kết quả của một số bài toán cụ thể bị ảnh hưởng lớn. Trong nghiên cứu này, các tác giả đã giả sử sai số do xấp xỉ hàm là do công cụ xấp xỉ hàm không đủ linh hoạt. Sau đó, nghiên cứu của [6] cho thấy sai số do môi trường gây ra cũng gây ảnh hưởng đến việc xấp xỉ hàm và gây nên hiện tượng đánh giá quá cao. Các nghiên cứu trên đều cho thấy vấn đề này có ảnh hưởng lớn đến chính sách tìm được của thuật toán “Q-learning”.

Để khắc phục vấn đề đánh giá quá cao, thuật toán “Double Q-learning” [6] ra đời và được áp dụng vào bài toán tự động chơi game [16]. Đây là một thuật toán cải tiến của “Q-learning” bằng cách thay đổi lại cách ước lượng giá trị của trạng thái tiếp theo. Cụ thể hơn, giá trị này theo “Q-learning” là:

$$R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a; \theta) \quad (3.18)$$

Ta có thể tách việc lấy max giá trị ra thành hai bước: chọn hành động tối ưu và đánh giá hành động đó. Bước chọn hành động được thực hiện bằng cách chọn hành động cho giá trị cao nhất tại trạng thái: $a' = \arg \max_a \hat{q}(S_{t+1}, a; \theta)$. Bước đánh giá hành động được thực hiện bằng cách lấy lại giá trị của hành động tốt nhất: $\hat{q}(S_{t+1}, a'; \theta)$. Ghép hai bước này lại ta có công thức ước lượng giá trị của thuật toán “Q-learning”:

$$R_{t+1} + \gamma \hat{q}(S_{t+1}, \arg \max_a \hat{q}(S_{t+1}, a; \theta); \theta) \quad (3.19)$$

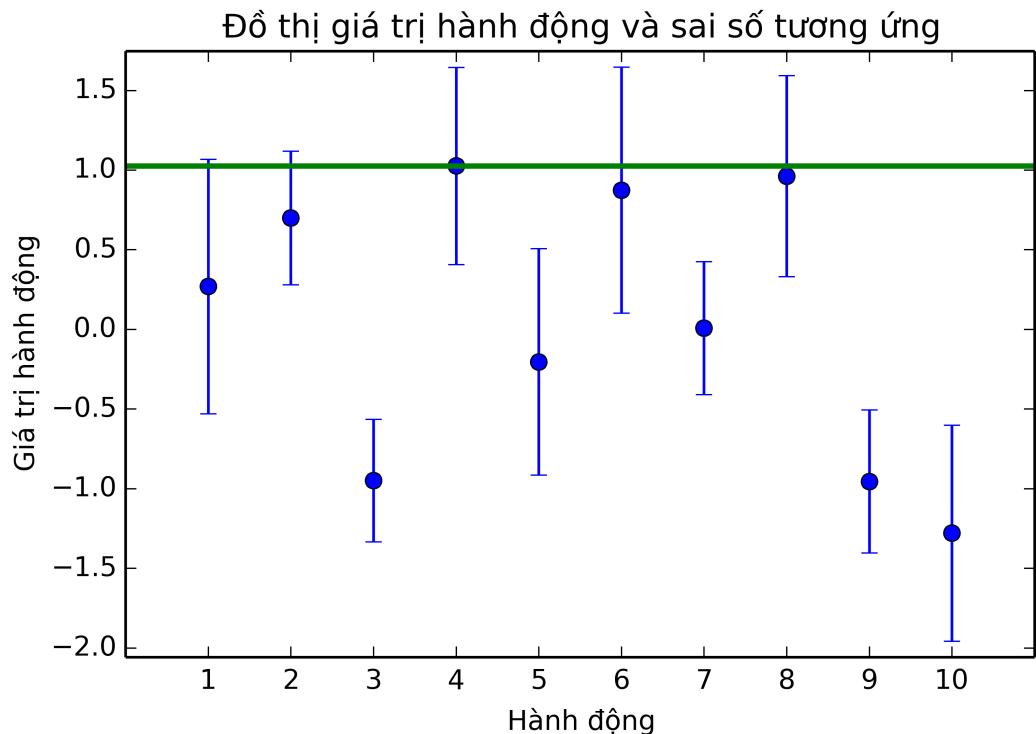
Lưu ý rằng công thức (3.19) chỉ là một cách nhìn khác của công thức (3.18) và hai công thức này hoàn toàn tương đương nhau.

Thuật toán “Double Q-learning” giảm thiểu vấn đề đánh giá quá cao bằng cách đánh giá hành động tốt nhất bằng hàm giá trị của chính sách khác với chính sách hiện tại. Cụ thể hơn, ta cần học hai mạng nơ-ron với hai bộ tham số θ và θ^- . Khi cần cập nhật bộ tham số θ của mạng đầu tiên, ta áp dụng công thức:

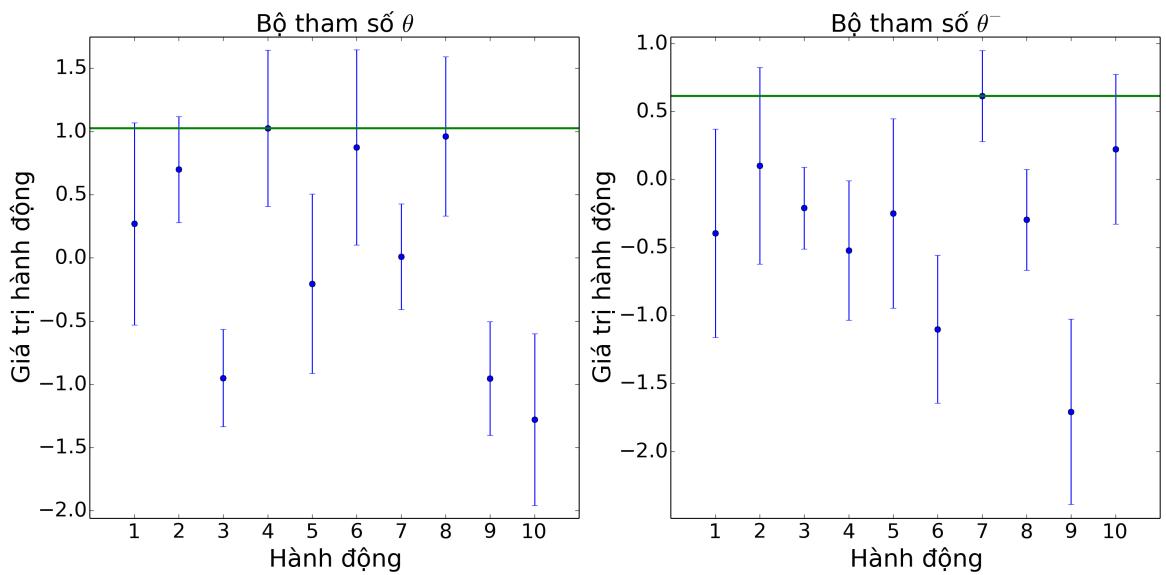
$$R_{t+1} + \gamma \hat{q}(S_{t+1}, \arg \max_a \hat{q}(S_{t+1}, a; \theta); \theta^-) \quad (3.20)$$

Công thức trên chỉ khác với công thức (3.19) ở chỗ ta đánh giá hành động

được chọn bằng một mạng nơ-ron khác. Nhờ việc “phân tách” (decouple) hai bước trên, “Double Q-learning” làm giảm khả năng bị đánh giá quá cao mà “Q-learning” gặp phải ([6]). Vậy tại sao cách cập nhật này lại không gặp vấn đề nêu trên? Trong “Q-learning”, ta cần tính max của giá trị của nhiều hành động. Do mỗi giá trị hành động tính được đều là xấp xỉ nên sẽ có sai số nhất định (tức sai số giữa giá trị xấp xỉ và giá trị thực $\hat{q}(s, a; \theta) - q_\pi(s, a)$). Giả sử sai số này là “uniform” trong đoạn $[-\epsilon, \epsilon]$ với ϵ là một giá trị dương nhỏ. Khi lấy max giá trị nhiều hành động lại thì nhiều khả năng hành động được chọn sẽ có sai số dương (do các hành động có sai số âm thì giá trị xấp xỉ bị giảm so với giá trị đúng nên khả năng được chọn thấp hơn). Khi càng có nhiều hành động thì khả năng hành động đó có sai số dương lại càng cao hơn. Điều này làm cho “Q-learning” gặp vấn đề đánh giá quá cao. Trong khi với thuật toán “Double Q-learning”, gọi hành động được chọn theo bộ tham số θ là a_* . Khi đó, xác suất để a_* cũng có sai số dương theo bộ tham số θ^- là thấp (do hai bộ tham số học từ dữ liệu khác nhau). Nhờ vậy, việc phân tách hai bước chọn và đánh giá hành động giúp làm giảm vấn đề đánh giá quá cao. Xem hình (3.6) và (3.7) mô tả một ví dụ đơn giản để hiểu hơn về vấn đề đánh giá quá cao trong hai thuật toán “Q-learning” và “Double Q-learning”.



Hình 3.6: Hình mô tả vấn đề đánh giá quá cao của thuật toán “Q-learning”. Trục hoành thể hiện mười hành động và trục tung thể hiện giá trị tương ứng. Chấm tròn thể hiện giá trị thực sự của từng hành động và các đoạn song song với trục tung thể hiện sai số do quá trình xấp xỉ. Đường thẳng nằm ngang thể hiện max của giá trị thực sự của tất cả các hành động. Đây là giá trị đích mong muốn của “Q-learning”. Có nhiều hành động có giá trị xấp xỉ với sai số dương lớn hơn giá trị đích này (hành động 4, 6 và 8) nên khi lấy max theo “Q-learning” thì ta gặp tình trạng đánh giá quá cao.



Hình 3.7: Hình mô tả cách giải quyết vấn đề đánh giá quá cao của thuật toán “Double Q-learning”. Đồ thị bên trái thể hiện giá trị hành động được xấp xỉ theo bộ tham số θ ; đồ thị bên phải là giá trị hành động theo bộ tham số θ^- . Khi chọn hành động theo bộ tham số θ , nhiều khả năng ta sẽ chọn hành động 4, 6 hoặc 8. Tuy nhiên, theo bộ tham số θ^- thì các hành động này lại không phải là những hành động có giá trị xấp xỉ cao nhất. Nhờ việc lựa chọn và đánh giá hành động một cách độc lập như vậy mà thuật toán “Double Q-learning” không gặp vấn đề đánh giá quá cao.

Chương 4

Kết quả thực nghiệm

Trong chương này, chúng em trình bày kết quả thực nghiệm của bài toán tự động chơi game. Đầu tiên, môi trường giả lập của bài toán tự động chơi game cùng với các thiết lập thực nghiệm được giới thiệu. Phần tiếp theo, chúng em trình bày về kết quả thực nghiệm của thuật toán “Q-learning” để chứng minh tính khả thi của việc kết hợp học tăng cường với học sâu vào bài toán tự động chơi game. Kế tiếp, chúng em phân tích mức độ ảnh hưởng của vấn đề “đánh giá quá cao” và kết quả của thuật toán “Double Q-learning”. Phần cuối cùng bao gồm một số thực nghiệm chứng minh rằng mô hình mạng nơ-ron xấp xỉ tốt hàm giá trị hành động.

4.1 Giới thiệu Arcade Learning Environment

Arcade Learning Environment (ALE) là một framework được thiết kế giúp dễ dàng trong việc phát triển những hệ thống có thể tự động chơi bất kỳ những game trên hệ máy Atari 2600.

4.1.1 Hệ máy Atari 2600

[4.1](#) là hệ máy chơi game Atari 2600 tại gia đình được phát triển và sử dụng phổ biến vào những năm 1970. Hơn 500 game đã được phát triển cho hệ máy này; một số game vẫn tiếp tục được phát triển cho tới ngày nay. Độ phân giải của những game này là 210×160 . Mặc dù vậy nhiều game được phát triển cho hệ

máy này, nhưng cấu hình phần cứng của chúng khá đơn giản bao gồm một CPU 1.19 Mhz, một bộ nhớ ROM 2-4KB để lưu giữ code của game, và dung lượng RAM của máy cũng chỉ là 128 bytes. Đồ họa của các game được cố định ở độ phân giải 160×210 , với ảnh màu RGB. Người chơi tương tác với game qua một thiết bị được gọi là joystick, có thể thực hiện tối đa 18 hành động.

4.1.2 Interface

ALE được xây dựng trên nền Stella, nó giả lập máy hệ máy Atari 2600. Qua đó cho phép người dùng tương tác với Atari 2600 bằng cách tiếp nhận những di chuyển của joystick, gửi thông tin của game cho người dùng như hình ảnh, điểm số nhận được. ALE cũng cung cấp một *lớp sử lý* để chuyển đổi những thông tin trong game theo chuẩn của bài toán học tăng cường như điểm số đạt được, kết thúc game chưa. Mặc định, lớp này biểu diễn mỗi trạng thái trong game bằng một mảng 7-bit pixels 2 chiều, và không gian hành động gồm 18 hành động tương ứng với các hành động có thể thực hiện trong joystick. Lớp sử lý cung cấp thông tin tập những hành động cụ thể có thể thực hiện trong một game. ALE có thể tạo ra 60 frame hình ảnh game trong một giây theo mặc định, và số frame ảnh nhiều nhất trong một giây nó có thể tạo ra 6000 frame. Điểm số nhận được tại từng thời điểm được định nghĩa theo từng game khác nhau. Một mẫu thực nghiệm có thể được tạo ra từ frame đầu tiên của game cho đến khi game kết thúc, ngoài ra lớp sử lý cũng cho phép người người tạo mẫu thực nghiệm với số lượng frame cố định.

ALE hơn nữa cũng cung cấp chức năng *sao lưu* và *khôi phục* trạng thái của game. Khi lệnh sao lưu được thực thi, ALE lưu lại tất cả các dữ liệu liên quan tại thời điểm đó trong game như nội dung RAM, registers. Ngược lại khi lệnh khôi phục được thực thi, nó sẽ khởi tạo lại trạng thái game đã lưu lại trước đó.



Hình 4.1: Hệ máy chơi game Atari 2600

4.2 Các thiết lập thực nghiệm

Chúng em thực hiện huấn luyện hệ thống trên ba trò chơi: *Assault*, *Breakout* và *Seaquest*. Việc thực nghiệm trên ba trò chơi là để chứng minh hướng tiếp cận kết hợp học sâu với học tăng cường có thể áp dụng vào nhiều trò chơi khác nhau. *Assault* là trò chơi bắn máy bay truyền thống: hệ thống điều khiển máy bay bay qua trái, phải hoặc bắn tên lửa để hạ máy bay địch. *Breakout* là trò chơi “phá gạch”: hệ thống điều khiển một mái chèo (paddle) qua trái, phải để giữ cho trái banh không rơi xuống dưới; phía trên là những viên gạch bị phá huỷ và cho điểm thưởng khi chạm phải trái banh. Cuối cùng là trò chơi *Seaquest* có nội dung tương tự như “*Assault*” nhưng là điều khiển tàu ngầm theo cả tám hướng xung quanh. Tuy các trò chơi này có nội dung tương tự nhau nhưng luật chơi cụ thể lại rất khác biệt và đòi hỏi chiến thuật mang tính “dài hơi”. Ví dụ như màn hình trò chơi *Seaquest* có một thanh không khí thể hiện mức độ không khí còn lại trong tàu ngầm. Khi thanh này cạn thì ta cần điều khiển tàu ngầm nổi lên mặt nước để lấy không khí. Vì vậy, hệ thống cần học được đặc trưng thể hiện mức không khí còn lại và dựa vào đó để điều khiển tàu ngầm nổi lên mặt nước. Ngoài ra, để chơi tốt và đạt nhiều điểm, hệ thống cần học được chiến thuật đợi không khí còn lại ít nhất rồi mới nổi lên mặt nước một lần để tiết kiệm thời gian. Ví dụ này cho thấy các trò chơi này đều đòi hỏi các đặc trưng học được phải rất trừu tượng và chiến thuật chơi cũng phải khá phức tạp. Hình () thể

hiện màn hình một khung ảnh của ba trò chơi trên.

Để giảm bớt tính toán của mô hình, khung ảnh RGB đầu vào được chuyển về ảnh mức xám (grayscale) và được giảm kích thước về 84×84 . Để lưu thông tin về sự di chuyển và tốc độ của các đối tượng trong ảnh, ta ghép bốn khung ảnh liên tiếp nhau lại thành một trạng thái. Vậy một trạng thái lúc này sẽ gồm bốn ảnh mức xám 84×84 và được đưa trực tiếp vào “input” của mạng nơ-ron. Ngoài ra, do các khung hình gần nhau thường không thay đổi nhiều, một kỹ thuật đơn giản để tăng tốc độ huấn luyện đó là ta lặp lại một hành động cho k khung hình liên tiếp. Nhờ kỹ thuật này, ta có thể huấn luyện nhiều hơn k lần so với khi không áp dụng. Trong thực nghiệm này, chúng em chọn k bằng 4 theo đề xuất của ([9]).

Cấu trúc mạng nơ-ron tích chập được dùng để xấp xỉ hàm giá trị là “Deep Q-Network” ([10]). Mạng này bao gồm bốn tầng ẩn với ba tầng đầu là tầng tích chập và tầng cuối là tầng “Fully-connected”. Tầng tích chập đầu tiên gồm 32 bộ lọc với kích thước 8×8 và bước dịch chuyển (stride) là 4×4 . Tầng tích chập thứ hai gồm 64 bộ lọc với kích thước 4×4 và bước dịch chuyển là 2×2 . Tầng tích chập cuối cùng gồm 64 bộ lọc với kích thước 3×3 và bước dịch chuyển là 1×1 . Đầu ra của tầng tích chập được đưa vào một tầng “Fully-connected” gồm 512 nơ-ron ẩn để tổng hợp đặc trưng. Tầng “output” của mạng cũng là một tầng “Fully-connected” với số nơ-ron ẩn là số hành động của trò chơi (6 cho “Assault”, 4 cho “Breakout” và 18 cho “Seaquest”). Tầng này sẽ trả về giá trị hành động tương ứng cho trạng thái đầu vào. Hàm kích hoạt của tất cả các tầng ẩn là hàm “Rectified”: $\sigma(x) = \max(0, x)$.

Để huấn luyện mạng nơ-ron này, chúng em áp dụng thuật toán “RMSprop” ([15]) để tối thiểu hàm lỗi bình phương. Đây là một thuật toán cải tiến của SGD với ý tưởng chọn một hệ số học khác nhau cho từng trọng số của mạng nơ-ron. Hệ số học của mỗi trọng số sẽ phụ thuộc vào giá trị đạo hàm của hàm lỗi đối với trọng số này ở các bước cập nhật trước. Cụ thể hơn, nếu trọng số này được cập nhật ít ở những lần cập nhật trước thì hệ số học sẽ được tăng lên; trọng số được cập nhật nhiều thì sẽ có hệ số học giảm đi. Nhờ thuật toán này mà mức độ thay đổi của các trọng số là tương đương nhau.

Chúng em sử dụng “Python” (kết hợp với framework “Theano” ([14])) làm

ngôn ngữ lập trình chính. Do việc huấn luyện đòi hỏi sức mạnh tính toán rất lớn nên chúng em thực hiện cài đặt tính toán song song trên GPU để tăng tốc. GPU được dùng là NVIDIA GTX 980.

Với các thiết lập thực nghiệm trên, chúng em tiến hành huấn luyện ba trò chơi với hai thuật toán “Q-learning” và “Double Q-learning”. Mỗi trò chơi với mỗi thuật toán được huấn luyện 50 triệu hành động. Thời gian huấn luyện là khoảng 50 giờ cho mỗi trò chơi với mỗi thuật toán.

4.3 Thuật toán “Q-learning”

Để thực nghiệm khả năng học và tự động chơi game, chúng em thực hiện thử nghiệm thuật toán “Q-learning” kết hợp với học sâu. Chúng em sử dụng chính sách ϵ -greedy để làm chính sách khám phá không gian trạng thái và không gian hành động. ϵ -greedy được chọn ở đây là vì chính sách này đơn giản cho việc cài đặt mà vẫn đảm bảo khả năng khám phá không gian cũng như khai thác kiến thức đã học của mô hình. Do ban đầu hệ thống chưa có nhiều hiểu biết về môi trường, chúng ta nên khám phá nhiều hơn; sau khi có một số hiểu biết nhất định thì hệ thống lại cần khai thác những hiểu biết đó. Ý tưởng này có thể dễ dàng được cài đặt bằng cách giảm dần hệ số ϵ trong khi đang huấn luyện. Với thuật toán “Q-learning”, chúng em thực hiện giảm dần ϵ từ 1.0 về 0.1 trong 1 triệu hành động đầu tiên.

Để tiện cho việc quan sát kết quả và thử nghiệm hệ thống, chúng em chia quá trình huấn luyện 50 triệu hành động ra thành 200 “epoch” (tức mỗi “epoch” bao gồm 250000 hành động). Sau mỗi “epoch”, chúng em kiểm thử bộ trọng số học được trên 30 màn chơi và lưu lại điểm số trung bình. Kết quả điểm số cuối cùng được báo cáo của mỗi trò chơi là giá trị điểm số lớn nhất của 200 epoch. Các thiết lập trên tương đồng với ([10]) để việc so sánh là công bằng nhất.

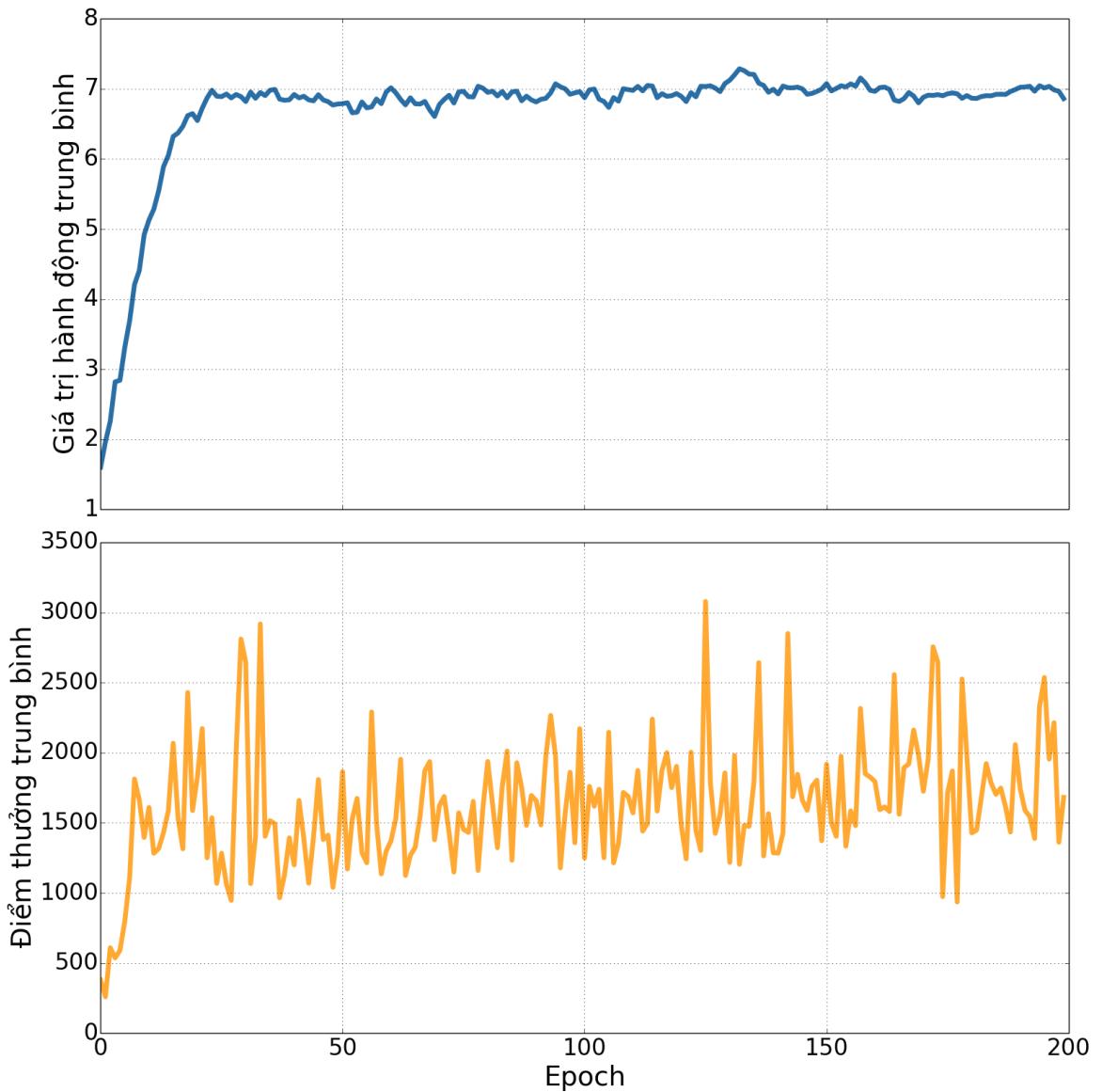
Một đặc điểm của bài toán tự động chơi game đó là kết quả điểm thưởng sau từng “epoch” thường rất nhiều. Giá trị của các trọng số mạng nơ-ron có thể thay đổi ít sau mỗi “epoch” nhưng hàm giá trị tương ứng sẽ thay đổi nhiều. Điều này dẫn đến chính sách tương ứng cũng thay đổi nhiều làm cho điểm thưởng nhận được khi chơi game của các chính sách này sẽ rất khác nhau. Vì vậy nếu

Bảng 4.1: Kết quả *Ngẫu nhiên* nhận được bằng cách chơi với chính sách ngẫu nhiên (chọn hành động theo phân bố đều). Kết quả *Con người* nhận được bằng cách cho một số người học và chơi các game này. Kết quả gốc là kết quả thực nghiệm của bài báo [10]. Kết quả *Thực nghiệm* là kết quả do nhóm thực hiện cài đặt lại thuật toán “Q-learning”. Giá trị được in đậm là giá trị lớn nhất của mỗi dòng.

	<i>Ngẫu nhiên</i> [10]	<i>Con người</i> [10]	<i>Kết quả gốc</i> [10]	<i>Thực nghiệm</i>
<i>Assault</i>	222.4	742.0	4,280.4	3,078.8
<i>Breakout</i>	1.7	30.5	385.5	377.6
<i>Seaquest</i>	68.4	42,054.7	5,860.6	6,340.7

ta chỉ quan sát đồ thị điểm thưởng của quá trình kiểm thử sau mỗi “epoch”, ta sẽ không biết chắc là hệ thống có đang chơi tốt lên hơn không. Một cách khác để quan sát quá trình học được đề xuất bởi ([9]) đó là ta quan sát đồ thị giá trị của một tập trạng thái sau từng “epoch”. Đầu tiên, ta thực hiện chơi game một cách ngẫu nhiên và lưu lại một tập các trạng thái gặp được trong quá trình chơi này (gọi là tập “hold-out”). Trong quá trình huấn luyện, sau mỗi “epoch”, ta tính trung bình giá trị hành động tốt nhất của các trạng thái trong tập “hold-out”. Một chính sách chơi tốt thì thường sẽ đánh giá các hành động của các trạng thái trong tập “hold-out” cao hơn so với các chính sách chơi không tốt. Vì vậy, quan sát đồ thị này là một cách đơn giản để theo dõi quá trình học của hệ thống.

Hình (4.2) thể hiện đồ thị điểm số và giá trị trạng thái trên tập “hold-out” của trò chơi *Assault*. Điểm thưởng kết quả của ba game được báo cáo trong bảng (4.1). Cả ba game được thực nghiệm cho kết quả vượt xa kết quả chơi ngẫu nhiên. Điều này cho thấy hệ thống thật sự học được cách chơi để đạt được nhiều điểm số dù ban đầu không hề biết quy luật chơi. Hai trong số ba game còn cho kết quả hơn cả kết quả do con người chơi. Tuy nhiên, bài báo [10] chỉ cho người chơi học thử khoảng 2 giờ cho mỗi game, ít hơn nhiều so với việc huấn luyện 50 triệu hành động (tức khoảng 38 ngày chơi) của hệ thống. Vì vậy, điểm thưởng nhận được của con người trong bảng (4.1) chỉ mang tính tham khảo. Dù chúng em thực hiện cài đặt lại thuật toán được đề xuất bởi bài báo [10] nhưng kết quả cũng không bằng nhau (cao hơn ở game *Seaquest* và thấp hơn ở hai game còn lại). Lý do chính của điều này là vì cách cài đặt khác nhau và chương trình giả lập game Atari khác nhau.



Hình 4.2: Đồ thị kết quả của trò chơi *Assault*. Đồ thị ở trên thể hiện trung bình giá trị hành động trên tập “hold-out” theo “epoch” huấn luyện; đồ thị ở dưới thể hiện điểm thưởng trung bình theo “epoch” huấn luyện. Điểm thưởng trung bình thay đổi rất nhiều sau từng “epoch”. Nếu sử dụng đồ thị này để theo dõi quá trình huấn luyện thì ta sẽ không biết là hệ thống có đang học được cách chơi game hay không. Trong khi đó, đồ thị giá trị hành động lại ít nhiễu hơn và hội tụ dần đến một mức nhất định. Quan sát đồ thị này, ta thấy những “epoch” đầu hệ thống học rất nhanh và sau khoảng 50 “epoch” thì giá trị hành động không tăng cao hơn nữa. Lúc này, ta có thể dừng việc học sớm. Tuy nhiên, một số game có thể cần học lâu hơn. Vì vậy, chúng em cố định số “epoch” để huấn luyện được nhiều trò chơi mà không phải thay đổi siêu tham số.

4.4 Thuật toán “Double Q-learning”

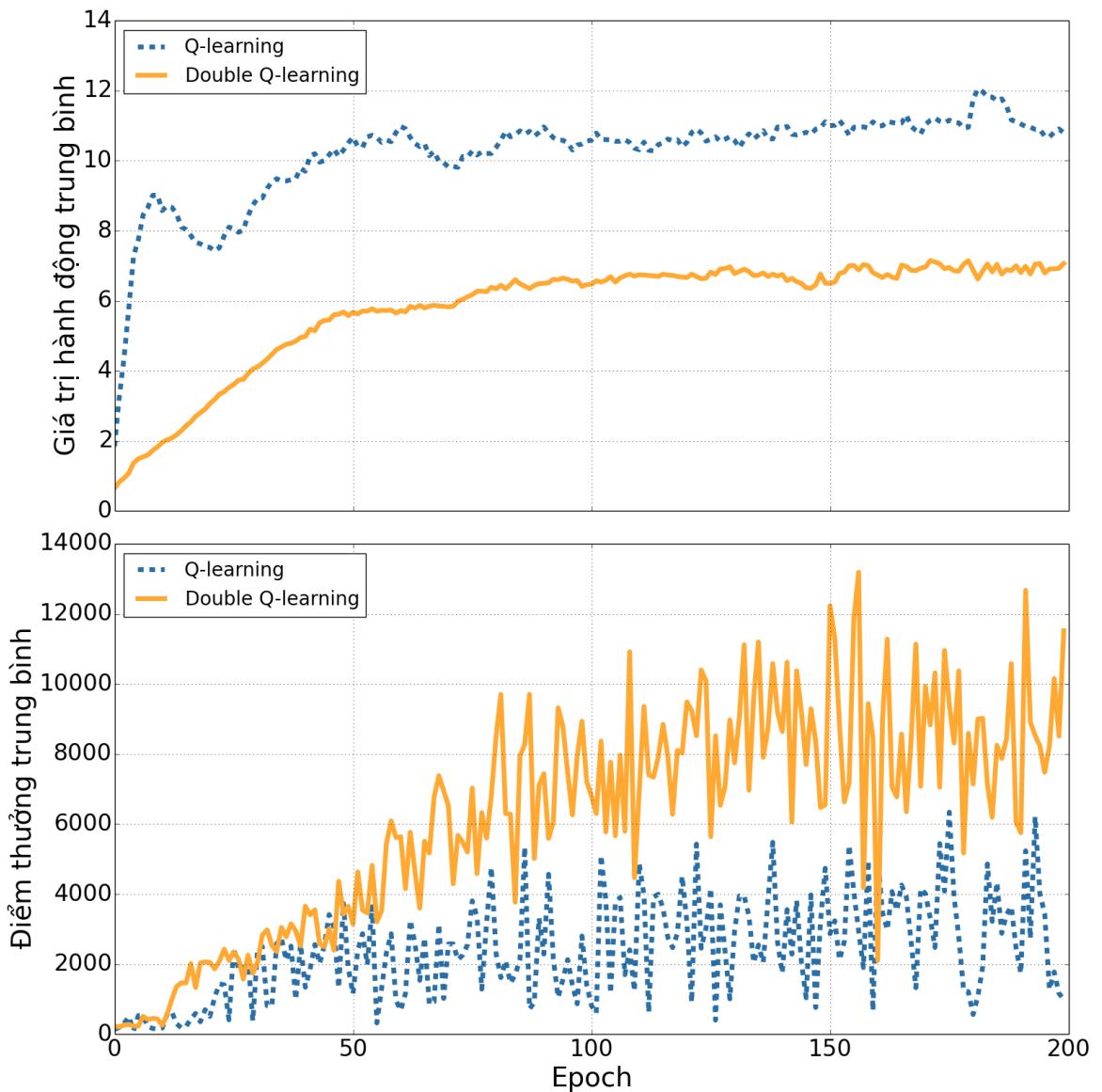
Chúng em thực hiện cài đặt và thử nghiệm lại thuật toán “Double Q-learning” [6], [16] để giải quyết vấn đề đánh giá quá cao. Các thiết lập thực nghiệm đều giống như khi thực nghiệm thuật toán “Q-learning”. Tuy nhiên, để tương đồng với cài đặt của bài báo [16], chúng em giảm hệ số ϵ từ 1.0 xuống 0.01 (thay vì 0.1 như ở phần thực nghiệm “Q-learning”) trong 1 triệu hành động đầu tiên của quá trình huấn luyện.

Bảng 4.2 cho biết kết quả của thuật toán “Double Q-learning” trên ba game. Thuật toán “Double Q-learning” cho kết quả tốt hơn hẳn điểm số của “Q-learning” trên cả ba game và cao hơn gấp đôi ở game *Seaquest*. Lý do “Double Q-learning” hoạt động tốt như vậy trên game *Seaquest* là vì game này có nhiều hành động (18 so với 4 của *Breakout* và 6 của *Assault*). Số lượng hành động càng nhiều thì khả năng bị đánh giá quá cao lại càng tăng. Vì vậy, thuật toán “Q-learning” bị ảnh hưởng rất lớn bởi vấn đề đánh giá quá cao và thuật toán “Double Q-learning” cho kết quả tốt hơn hẳn khi giải quyết được vấn đề này. Các kết quả này cho thấy thuật toán “Double Q-learning” là một cải tiến rất tốt cho bài toán tự động chơi game.

Để thấy rõ hơn việc hạn chế vấn đề đánh giá quá cao của “Double Q-learning”, ta có thể coi đồ thị giá trị hành động trung bình trên tập “hold-out” và đồ thị điểm thưởng trung bình của cả hai thuật toán ở hình (4.3). Thuật toán “Double Q-learning” cho kết quả điểm thưởng tốt hơn hẳn “Q-learning” vì đánh giá hành động chính xác hơn (không bị vấn đề đánh giá quá cao).

Bảng 4.2: Kết quả *Ngẫu nhiên* và *Con người* được trích từ bài báo [16]. Do môi trường giả lập game khác nhau nên chúng em chỉ so sánh “Double Q-learning” với phiên bản cài đặt “Q-learning” của nhóm thay vì so sánh với phiên bản “Double Q-learning” của bài báo [16]. Giá trị được in đậm là giá trị lớn nhất của mỗi dòng.

	<i>Ngẫu nhiên</i> [10]	<i>Con người</i> [10]	<i>Q-learning</i>	<i>Double Q-learning</i>
<i>Assault</i>	222.4	742.0	3,078.8	4864.6
<i>Breakout</i>	1.7	30.5	377.6	431.2
<i>Seaquest</i>	68.4	42,054.7	6,340.7	13,186

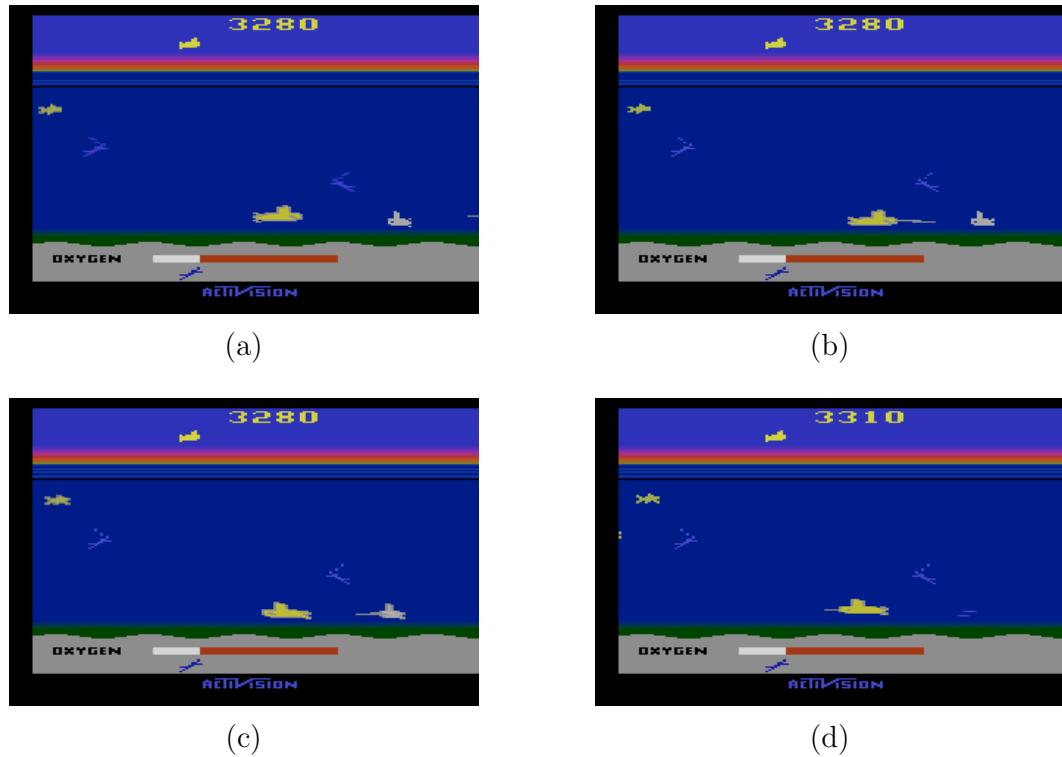


Hình 4.3: Đồ thị kết quả của trò chơi *Seaquest*. Hình phía trên là đồ thị giá trị hành động trung bình theo “epoch”; hình phía dưới là đồ thị điểm thưởng trung bình theo “epoch”. Ở đồ thị phía trên, thuật toán “Q-learning” đánh giá rất cao các hành động của trạng thái thuộc tập “hold-out”. Tuy nhiên, khi chơi thử game thì thuật toán “Q-learning” lại cho điểm thưởng rất thấp. Điều này cho thấy thuật toán “Q-learning” gặp vấn đề đánh giá quá cao và vấn đề này ảnh hưởng đến kết quả chơi game của hệ thống. Trong khi đó, thuật toán “Double Q-learning” tuy đánh giá các hành động của trạng thái thuộc tập “hold-out” thấp hơn nhưng khi chơi game thì lại được nhiều điểm thưởng hơn. Thuật toán “Double Q-learning” đánh giá chính xác hơn giá trị hành động và vì thế, chính sách chơi game cũng tốt hơn.

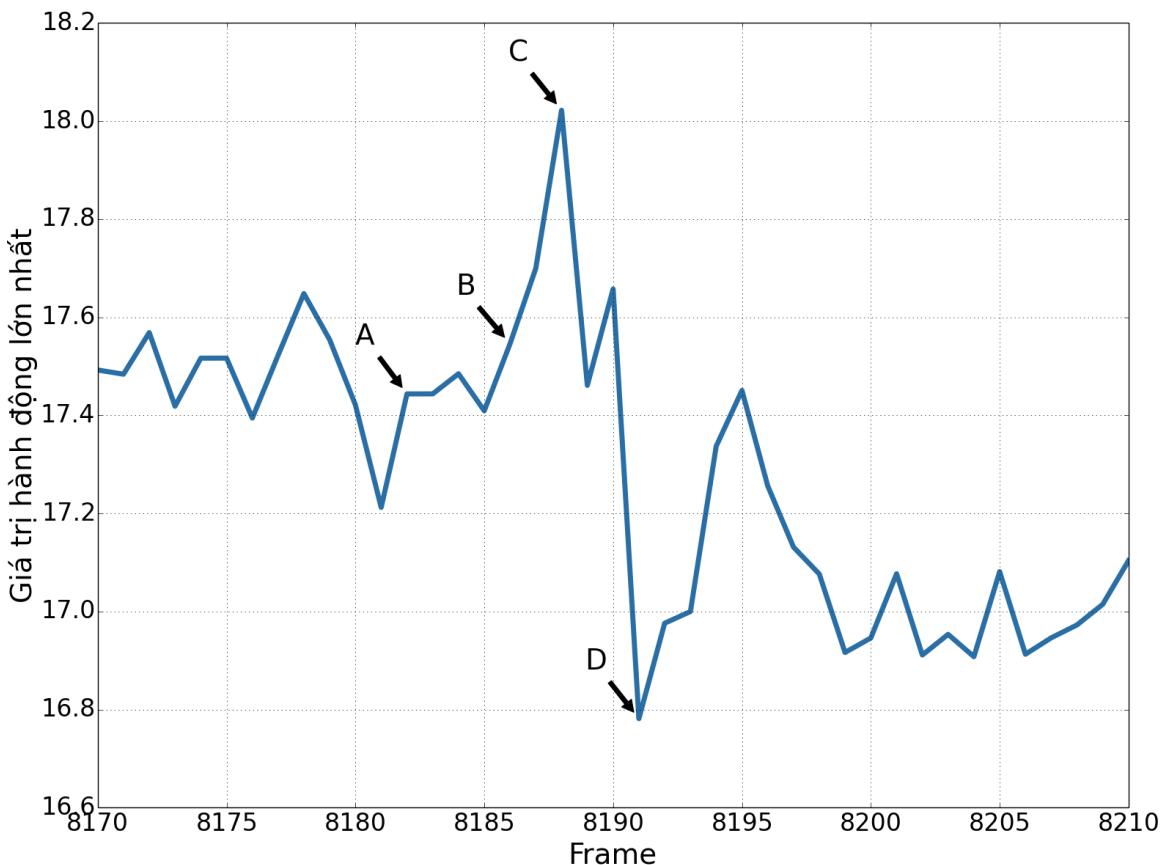
4.5 Phân tích hàm giá trị hành động

Để thấy rõ hơn khả năng xấp xỉ hàm giá trị của mạng nơ-ron, ta có thể xem xét giá trị của một số trạng thái trong quá trình chơi game. Cụ thể hơn, chúng em thực hiện lưu một số “frame” ảnh trong quá trình chơi cùng với giá trị của chúng được tính từ mạng nơ-ron (với thuật toán “Q-learning”).

Hình 4.4 và đồ thị 4.5 thể hiện các “frame” ảnh và giá trị của các trạng thái tương ứng. Quan sát giá trị tương đối của các trạng thái liên tiếp nhau trong một lần chơi, ta có thể thấy mạng nơ-ron học được một hàm giá trị phức tạp và có ý nghĩa cao. Nhờ hàm giá trị được xấp xỉ tốt, chính sách chơi game mới lựa chọn hành động hợp lý để nhận được điểm thưởng cao nhất có thể.



Hình 4.4: Hình ảnh 4 “frame” của trò chơi *Seaquest*. Ở “frame” đầu tiên (hình 4.4a), tàu ngầm do hệ thống điều khiển ở giữa và tàu ngầm địch ở bên phải. Ở “frame” kế tiếp (hình 4.4b), hệ thống thực hiện hành động phóng ngư lôi. Tiếp đó (hình 4.4c), ngư lôi trúng tàu ngầm địch. Cuối cùng (hình 4.4d), tàu ngầm địch bị phá huỷ cho điểm thưởng. Các “frame” này ứng với các giá trị trong đồ thị ở hình 4.5.



Hình 4.5: Đồ thị thể hiện giá trị hành động lớn nhất tại một số trạng thái; các trạng thái này ứng với các “frame” trò chơi *Seaquest*. Tại “frame” được đánh dấu *A*, giá trị hành động lớn nhất đang ở mức trung bình ứng với việc tàu ngầm của ta chưa thực hiện hành động gì ảnh hưởng đến điểm số. Tại “frame” *B*, tàu ngầm của ta bắn ngư lôi về phía tàu ngầm địch. Mạng nơ-ron khi nhận được “frame” ảnh này “hiểu” được rằng nếu ngư lôi chạm vào địch thì ta sẽ nhận được điểm thưởng. Vì vậy, đồ thị giá trị bắt đầu tăng lên đến “frame” *C*. Tại “frame” *C*, ngư lôi trúng tàu ngầm địch nên việc nhận được điểm là chắc chắn; hàm giá trị biết được điều này nên đồ thị giá trị đạt đỉnh điểm tại đây. Sau khi nhận được điểm thưởng và tàu ngầm địch bị phá huỷ, mạng nơ-ron nhận thấy rằng xung quanh không có yếu tố nào có thể giúp tăng điểm số nên đánh giá “frame” *D* có giá trị thấp.

Chương 5

Kết luận và hướng phát triển

TÀI LIỆU THAM KHẢO

- [1] M. Bellemare, J. Veness, and M. Bowling, “Bayesian learning of recursively factored environments,” in *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 2013, pp. 1211–1219. [4](#)
- [2] M. Bellemare, J. Veness, and E. Talvitie, “Skip context tree switching,” in *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, T. Jebara and E. P. Xing, Eds. JMLR Workshop and Conference Proceedings, 2014, pp. 1458–1466. [Online]. Available: <http://jmlr.org/proceedings/papers/v32/bellemare14.pdf> [4](#)
- [3] I. G. Y. Bengio and A. Courville, “Deep learning,” 2016, book in preparation for MIT Press. [Online]. Available: <http://www.deeplearningbook.org> [41](#)
- [4] M. Genesereth, N. Love, and B. Pell, “General game playing: Overview of the aaai competition,” *AI magazine*, vol. 26, no. 2, p. 62, 2005. [4](#)
- [5] G. J. Gordon, “Stable function approximation in dynamic programming,” in *Proceedings of the twelfth international conference on machine learning*, 1995, pp. 261–268. [20](#)
- [6] H. V. Hasselt, “Double q-learning,” in *Advances in Neural Information Processing Systems*, 2010, pp. 2613–2621. [59](#), [60](#), [70](#)
- [7] M. Hausknecht, J. Lehman, R. Miikkulainen, and P. Stone, “A neuroevolution approach to general atari game playing,” *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 6, no. 4, pp. 355–366, 2014. [4](#)

- [8] L.-J. Lin, “Reinforcement learning for robots using neural networks,” DTIC Document, Tech. Rep., 1993. [56](#)
- [9] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013. [56](#), [57](#), [66](#), [68](#)
- [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, pp. 529–533, 2015. [5](#), [53](#), [55](#), [56](#), [57](#), [66](#), [67](#), [68](#), [70](#)
- [11] S. Risi and J. Togelius, “Neuroevolution in games: State of the art and open challenges,” 2014. [4](#)
- [12] R. D. Smallwood and E. J. Sondik, “The optimal control of partially observable markov processes over a finite horizon,” *Operations Research*, vol. 21, no. 5, pp. 1071–1088, 1973. [20](#)
- [13] R. S. Sutton and A. G. Barto, *Introduction to reinforcement learning*. MIT Press Cambridge, 1998, vol. 135. [15](#), [20](#), [35](#), [39](#), [54](#), [58](#)
- [14] Theano Development Team, “Theano: A Python framework for fast computation of mathematical expressions,” *arXiv e-prints*, vol. abs/1605.02688, May 2016. [Online]. Available: <http://arxiv.org/abs/1605.02688> [66](#)
- [15] T. Tielemans and G. Hinton, “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude,” *COURSERA: Neural Networks for Machine Learning*, vol. 4, p. 2, 2012. [66](#)
- [16] H. Van Hasselt, A. Guez, and D. Silver, “Deep reinforcement learning with double q-learning,” *arXiv preprint arXiv:1509.06461*, 2015. [59](#), [70](#)