

10 - Reinforcement Learning

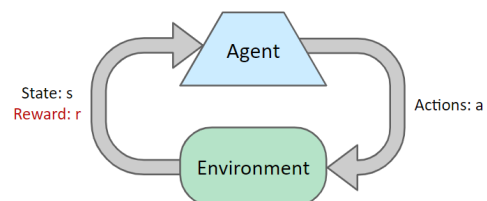
Khái niệm

Important ideas

- Khám phá (exploration): trải nghiệm những hàm động mới để lấy thông tin
- Khai thác (exploitation): sử dụng những trải nghiệm sẵn có → đưa ra hành động tiếp theo
- Hối hận (regret): thử sai
- Lấy mẫu (sampling): vì yếu tố ngẫu nhiên nên phải thực hiện hành động nhiều lần
- Sự khó khăn (difficulty): RL rất khó

Basic ideas

- MDP - offline, RL-online (có sự tương tác với môi trường)
- Tiếp nhận feedback từ môi trường dưới dạng phần thưởng ngay lập tức và trạng thái tiếp theo

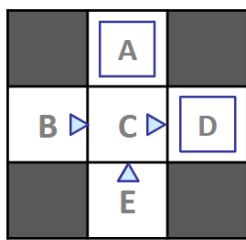


- Agent sẽ quyết định ở trạng thái tiếp theo sẽ làm gì
- Mục tiêu ở mỗi bước sao cho cực đại hóa được tổng điểm thưởng kì vọng khi xuất phát từ trạng thái s trong quá khứ

⇒ có yếu tố Markov nhưng không có P và R

Model-based Learning

- Trong học tăng cường, mục tiêu vẫn là tìm kiếm chiến lược tối ưu, tuy nhiên ta lại không có P (phân phối xác suất) và R (điểm thưởng) → bắt buộc phải tương tác với môi trường, tìm cách xấp xỉ P và R , học ra P và R từ trải nghiệm
- Giả sử có 1 chiến lược π , để xấp xỉ P và R , sẽ dùng chiến lược chơi nhiều ván chơi → dựa trên nhiều ván chơi, xấp xỉ ra P và R → chạy lại policy evaluation dựa trên mô hình vừa được xấp xỉ ⇒ chơi càng nhiều, đánh giá chiến lược càng chính xác
- Ví dụ:



Assume: $\gamma = 1$

giả sử có chiến lược sau

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

thực hiện 4 lần chơi

vì cả 2 lần B → C đều được nên 100%

4 lần C → D nhưng lần 4 thất bại (C → A thay vì D) = 75%

$$\hat{T}(s, a, s')$$

$T(B, \text{east}, C) = 1.00$
 $T(C, \text{east}, D) = 0.75$
 $T(C, \text{east}, A) = 0.25$
...

$$\hat{R}(s, a, s')$$

$R(B, \text{east}, C) = -1$
 $R(C, \text{east}, D) = -1$
 $R(D, \text{exit}, x) = +10$
...

học được P và R sau 4 lần chơi

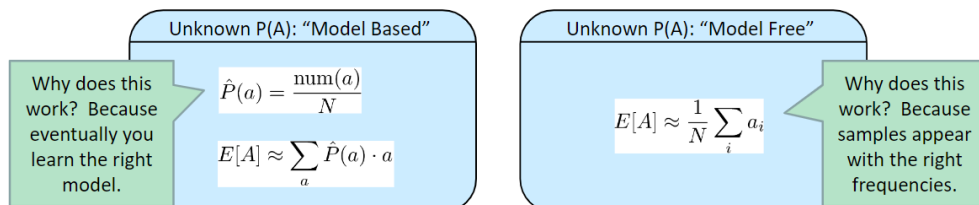
Model-free Learning

-
- Tính giá trị kì vọng trực tiếp không cần thông qua phân bố xs
 - Ví dụ: ước lượng độ tuổi kì vọng của sinh viên

Goal: Compute expected age of CS106 students

Known P(A)
$E[A] = \sum_a P(a) \cdot a = 0.35 \times 20 + \dots$

Without P(A), instead collect samples $[a_1, a_2, \dots, a_N]$



không cần xấp xỉ P như model-based, model-free lấy tổng những mẫu sau đó chia cho N

⇒ Đa số thuật toán tăng cường là model-free, nhưng các thuật toán sota thì lại là model-based

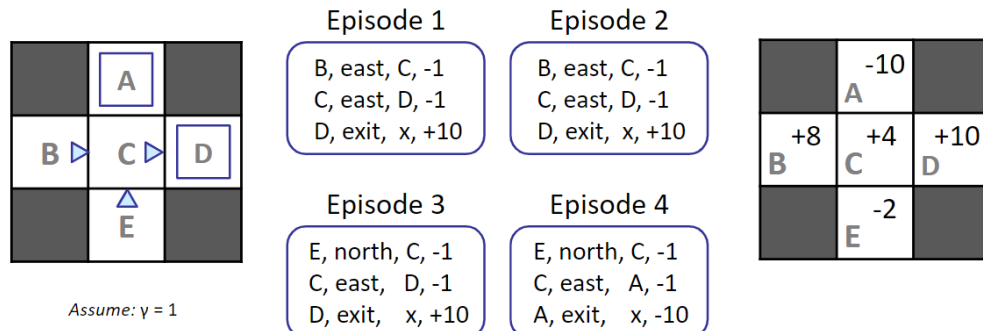
Passive Reinforcement Learning

- Có chiến lược $\pi \rightarrow$ tính $V^\pi(s)$
- Nhưng không biết P và R

Direct Evaluation - Monte Carlo

- Thuật toán: Xấp xỉ trực tiếp ra V^π không cần thông qua P và R dùng model-free

- Áp dụng chiến lược π nhiều lần
- Với mỗi trạng thái s và mỗi lần t mà s được thăm \rightarrow xác định phần thưởng ngay lập tức $r_{t+1}, r_{t+2}, \dots, r_t$
- Lấy mẫu cho trạng thái s ở thời điểm t = tổng giá trị phần thưởng chiết khấu $sample = r_{t+1} + \gamma R_{t+2}(s')$
- Lấy trung bình mẫu
- Ví dụ:



$$\begin{aligned}
 V_{\pi}(A) &= -10 \\
 V_{\pi}(B) &= [(-1-1+10)+(-1-1+10)]/2 = +8 \quad || \quad V_{\pi}(B) = [(-1-1 \times 0.9 + 0.9^2 \times 10) + (-1-1 \times 0.9 + 0.9^2 \times 10)]/2 \\
 V_{\pi}(C) &= [(-1+10)+(-1+10)+(-1+10)+(-1+10)]/4 = +4 \quad || \quad V_{\pi}(C) = [(-1+0.9 \times 10) + (-1+0.9 \times 10) + (-1+0.9 \times 10) + (-1+0.9 \times 10)]/4 \\
 V_{\pi}(D) &= [10 + 10 + 10]/3 = 10 \\
 V_{\pi}(E) &= [(-1-1+10)+(-1-1+10)]/2 = -2
 \end{aligned}$$

- Nhược điểm:
 - Chơi hết 1 ván chơi mới cập nhật \rightarrow lâu
 - Chơi càng nhiều ước lượng càng chính xác \rightarrow tốn chi phí tính toán
 - Dù $B \rightarrow C$ và $E \rightarrow C$ cùng đều tốn 1 bước nhưng $V(B) = 8$ mà $V(E) = -2$ (đáng lẽ phải bằng nhau vì điểm thưởng mỗi bước đi bằng nhau và hệ quả tương lai đều đến được C) \rightarrow do 'thiếu kinh nghiệm'

Mô phỏng policy evaluation khi không có P và R?

- Policy evaluation:

$$\begin{aligned}
 V_0^{\pi}(s) &\leftarrow 0 \\
 V_{k+1}^{\pi}(s) &\leftarrow \sum_{s'} P(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]
 \end{aligned}$$

- Ý tưởng: lấy mẫu kết quả từ s' (bằng cách thực hiện hành động) và lấy trung bình
 - Khởi tạo hệ quả tương lai = 0
 - Lấy mẫu: làm hành động $\pi(s)$ nhiều lần (sau khi có s'_1 thì quay lại s thực hiện lại để ra được s'_n)
 - Hệ quả tương lai chưa biết thì lấy kết quả của vòng lặp trước

$$\begin{aligned} sample_1 &= R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1) \\ sample_2 &= R(s, \pi(s), s'_2) + \gamma V_k^\pi(s'_2) \\ &\vdots \\ sample_n &= R(s, \pi(s), s'_n) + \gamma V_k^\pi(s'_n) \end{aligned}$$

- Lấy trung bình các mẫu:

$$V_{k+1}^\pi(s) \leftarrow \frac{1}{n} \sum_i sample_i$$

⇒ Xấp xỉ được mối quan hệ các trạng thái → cần ít chi phí hơn Monte Carlo

- Nhược điểm:
 - Không thực tế vì sau khi ra được s'_1 , ta cần phải quay lại s để lấy mẫu lần nữa → cần can thiệp vào source code thì mới điều khiển được agent xuất hiện ở bất kì trạng thái nào của game

Temporal Difference Learning

- Cập nhật $V^\pi(s)$ theo diễn biến trò chơi - đi 1 bước, học 1 bước
- Mỗi bước tạo 1 mẫu dữ liệu và dùng mẫu dữ liệu đó để cập nhật
- Mẫu = thông tin chính xác (R) + thông tin không chính xác (hệ quả tương lai) → mẫu là thông tin không chính xác → **Label**

Sample of V(s): $sample = R(s, \pi(s), s') + \gamma V^\pi(s')$

Update to V(s): $V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$

Same update: $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

alpha = learning rate → khi có trải nghiệm mới, muốn học trải nghiệm mới này bao nhiêu % để cập nhật
- Cập nhật $V^\pi(s)$: giá trị kì vọng → **Output**
- Nếu alpha = 1 → $V^\pi(s) = sample$ → tất cả những gì trong quá khứ quên đi hết, chỉ biết trải nghiệm hiện tại
- Nếu alpha = 0 → không học được gì cả

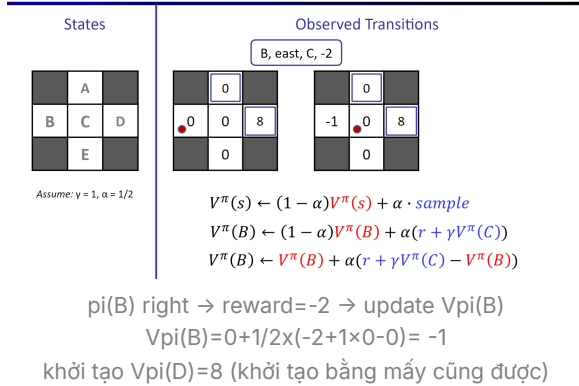
⇒ Phát sinh dữ liệu, gán nhãn trong lúc học nhưng gán nhãn sai

- Thuật toán:

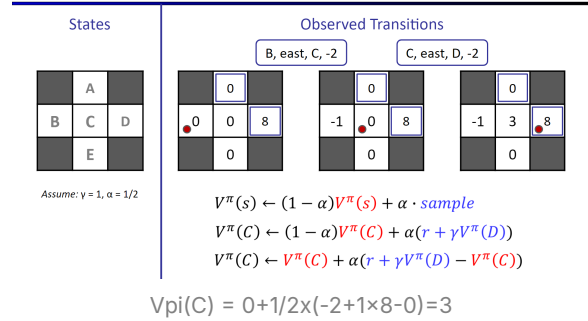
1. **Initialize** each $V^\pi(s)$ with some value.
2. **Observe** experience $(s, \pi(s), r, s')$.
3. **Use** observation in rough estimate of long-term reward $V^\pi(s)$
 $sample = r + \gamma V^\pi(s')$
4. **Update** $V^\pi(s)$ by moving values slightly toward estimate:
 $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

Có thể khởi tạo $V^\pi(s)$ bất kì giá trị nào vì công thức cập nhật giống với momentum (vector gradient trung bình có gán trọng số, thông tin xa thì trọng số khiến nó nhỏ dần đi → quên dần đi)

Example: Temporal Difference Learning



Example: Temporal Difference Learning



• Nhận xét:

- Có khả năng đánh giá được chiến lược bằng cách tính $V^\pi(s)$ nhưng không có khả năng tính được π^*
- Để lấy max được V của các hành động a khác nhau trong 1 trạng thái s, gần như cần phải có source của game, vì sau khi thực hiện a1 trong s thì qua s', cần phải quay lại s để thực hiện a2
- Dù ra được $V^*(s)$ nhưng không thể ra được π^* vì chạy policy extraction (cần thông tin P và R)

\Rightarrow Active Learning: đi tìm Q^* thay vì V^*

Active Reinforcement Learning

Q - Value Iteration

- Chiến lược: đi tìm $Q^*(s, a)$
- Phương trình Bellman tối ưu:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s, a, s') [R(s, a, s') + \gamma V_k^\pi(s')]$$

\Rightarrow Chuyển $V^*(s)$ thành $Q^*(s, a)$

$$Q_{k+1}(s) \leftarrow \sum_{s'} P(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k^*(s', a')]$$

- $V^*(s') \rightarrow Q_k^*(s', \pi^*(s')) \rightarrow \max_{a'} Q_k^*(s', a')$
 - $Q_k^*(s', a')$ đã xác định ở vòng lặp thứ k, nếu vòng lặp đầu thì khởi tạo bất kì giá trị nào
- TD learning: đi 1 bước, cập nhập $V^\pi(s)$
 Q learning: đi 1 bước, cập nhập $Q^*(s, a)$

Q - Learning

- Học ra $Q^*(s,a)$
- Q-table: mỗi hàng là 1 trạng thái, mỗi cột là 1 hành động → mỗi ô là giá trị $Q^*(s,a)$

1. **Initialize** $Q(s,a) = 0$ for each s,a pair
2. **Select an action** and **observe** an experience (s,a,r,s')
3. **Consider** your old estimate: $Q(s,a)$
4. **Use** observation in rough estimate of $Q(s,a)$

$$sample = r + \gamma \max_{a'} Q(s',a')$$
5. **Update** $Q(s,a)$ by moving values slightly toward estimate

$$Q(s,a) \leftarrow Q(s,a) + \alpha(sample - Q(s,a))$$

$$= (1 - \alpha)Q(s,a) + \alpha \cdot sample$$

thường sẽ không có dấu * vì chưa là tối ưu

- Bắt đầu, trạng thái $s_0 \rightarrow$ **chọn hành động** → quan sát (môi trường trả về s' và R) → gán nhãn = điểm thưởng lập tức + hệ quả (hệ quả không có nhưng đã được khởi tạo ở Q table, một lúc sau thì giá trị này khác 0 vì Q table đã được điền) → cập nhật $Q^*(s,a) \rightarrow$ đang ở s' → quay lại bước 2
- Nếu chạy đủ lâu thì sẽ điền được Q table → $Q^*(s,a) \rightarrow$ policy extraction lấy chiến lược tối ưu

ϵ -Greedy Policies

- Chiến lược lựa chọn hành động
- Dựa vào tỉ lệ thắng của những hành động ở những ván chơi trước mà sẽ ưu tiên thực hiện hành động
- Có bảng Q table → lấy $\arg\max_a Q(s,a)$ sẽ ra được hành động có tỉ lệ thắng cao nhất
- Nhưng nếu chỉ $\arg\max$ thì sẽ không thể khám phá được map → phải chọn hành động ngẫu nhiên để điền được hết Q table
- Sau khi học được $Q^*(s,a)$ thì không cần ngẫu nhiên nữa mà lấy $\arg\max$, ý nghĩa của hành động ngẫu nhiên là để học

⇒ Cân bằng giữa exploration (khám phá) và khai phá (exploitation)

$$\pi(a|s) = \begin{cases} \arg\max_a Q(s,a), & \text{with prob. } 1 - \epsilon + \frac{\epsilon}{|A|} \\ a' \neq \arg\max_a Q(s,a), & \text{with prob. } \frac{\epsilon}{|A|} \end{cases}$$

trong đó $|A|$ là số lượng hành động

- Phải chia $|A|$ vì có khả năng rơi trúng trường hợp $\arg\max$

Q - Learning với ϵ -greedy exploration

Với mỗi step trong episode:

1. Initialize $Q(s, a) = 0$ for each s, a pair, $t = 0$
2. Set $\pi_b(s) = \arg \max_{a'} Q(s, a)$ with prob. $1 - \epsilon$, else random.
3. Loop
 1. Sample action a from policy π_b an observe an experience (s, a, r, s')
 2. Consider your old estimate: $Q(s, a)$
 3. Use observation in rough estimate of $Q(s, a)$

$$sample = r + \gamma \max_{a'} Q(s', a')$$
 4. Update $Q(s, a)$ by moving values slightly toward estimate

$$Q(s, a) \leftarrow Q(s, a) + \alpha (sample - Q(s, a)) = Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

$$\pi_b(s) = \arg \max_{a'} Q(s, a) \text{ with prob. } 1 - \epsilon, \text{ else random.}$$
 5. $s \leftarrow s'$

```
q_table = np.zeros((env.observation_space.n, env.action_space.n))
rewards_all = []
for episode in range(num_episodes):
    state, _ = env.reset()
    reward_episode = 0.0
    done = False
    epsilon = min_epsilon + (max_epsilon - min_epsilon) * np.exp(-epsilon_decay_rate*episode)
    for step in range(num_steps_per_episode):
        exploration = random.uniform(0,1)
        if exploration < epsilon:
            action = env.action_space.sample()
        else:
            action = np.argmax(q_table[state, :])
        next_state, reward, terminated, truncated, _ = env.step(action)
        q_table[state, action] = q_table[state, action] * (1 - learning_rate)
            + learning_rate * (reward + gamma * np.max(q_table[next_state, :]))
        reward_episode += reward
        state = next_state
        if terminated or truncated:
            done = True
        if done:
            break
    rewards_all.append(reward_episode)
```

- Sử dụng π_b để chọn hành động $a \rightarrow$ quan sát môi trường trả về s' và $r \rightarrow$ tiến hành gán nhãn dựa vào Q table và $s' \rightarrow$ cập nhật $Q(s, a) \rightarrow \Rightarrow$ xong bước 4 là đang ở s' (gán $s \leftarrow s'$)
- a' chỉ xuất hiện ở argmax gán nhãn, π_b không đung gì vào a' hết

SARSA với ϵ -greedy exploration

Với mỗi step trong episode:

1. Initialize $Q(s, a) = 0$ for each s, a pair, $t = 0$
2. Set $\pi_b(s) = \arg \max_{a'} Q(s, a)$ with prob. $1 - \epsilon$, else random.
3. Sample action a from policy π_b an observe an experience (s, a, r, s')
4. Loop
 1. Sample action a' from policy π_b an observe an experience (s', a', r'', s'')
 2. Consider your old estimate: $Q(s, a)$
 3. Use observation in rough estimate of $Q(s, a)$

$$sample = r + \gamma Q(s', a')$$
 4. Update $Q(s, a)$ by moving values slightly toward estimate

$$Q(s, a) \leftarrow Q(s, a) + \alpha (sample - Q(s, a)) = Q(s, a) + \alpha (r + \gamma Q(s', a') - Q(s, a))$$

$$\pi_b(s) = \arg \max_{a'} Q(s, a) \text{ with prob. } 1 - \epsilon, \text{ else random.}$$
 5. $s \leftarrow s', a \leftarrow a'$

```
q_table = np.zeros((env.observation_space.n, env.action_space.n))
rewards_all = []
for episode in range(num_episodes):
    state, _ = env.reset()
    reward_episode = 0.0
    done = False
    epsilon = min_epsilon + (max_epsilon - min_epsilon) * np.exp(-epsilon_decay_rate*episode)
    exploration = random.uniform(0,1)
    if exploration < epsilon: action = env.action_space.sample()
    else: action = np.argmax(q_table[state, :])
    for step in range(num_steps_per_episode):
        next_state, reward, terminated, truncated, _ = env.step(action)
        exploration = random.uniform(0,1)
        if exploration < epsilon: next_action = env.action_space.sample()
        else: next_action = np.argmax(q_table[next_state, :]) # a'
        q_table[state, action] = q_table[state, action] +
            learning_rate * (reward + gamma * q_table[next_state, next_action] - q_table[state, action])
        reward_episode += reward
        state, action = next_state, next_action
        if terminated or truncated: break
    rewards_all.append(reward_episode)
```

- a' ở Q learning không phụ thuộc π_b vì nó chỉ argmax Q table nhưng a' ở trong vòng lặp 1 Sarsa phụ thuộc π_b vì chọn a' từ π_b
- gán nhãn bằng r của hành động a và hệ quả tương lai (hệ quả này không có max vì a' phụ thuộc π_b chứ không phải lấy argmax) \rightarrow cập nhật $Q(s, a)$
- Ở vòng lặp tiếp theo, Sarsa sẽ cập nhật $Q(s', a')$
- công thức gán nhãn này không còn là Q^* vì a' tạo ra từ $\pi_b(s')$, tức trong trạng thái s làm hành động a , từ s' trở đi tuân theo $\pi_b \rightarrow Q^{\pi_b}(s, a)$

\Rightarrow bản chất: đi tìm $Q^{\pi_b}(s, a)$

- Nhưng Sarsa khởi tạo π_b bằng argmax $Q(s, a)$ với xác suất lớn nên Sarsa không học ra π^* mà học π_b nhưng học đủ lâu thì π_b hội tụ về π^* (để được điều này thì ϵ phải giảm dần để càng về sau nó sẽ ít khi chọn hành động ngẫu nhiên nữa)

Nhận xét

- **Off Policy - Q Learning:** dùng chiến lược π_b phát sinh ra hành động để từ đó tính ra π^* (tính $Q^*(s, a)$) là gì

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

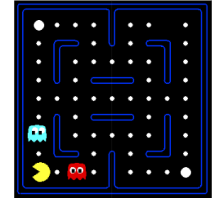
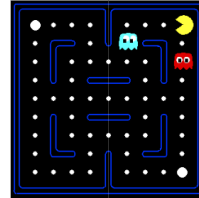
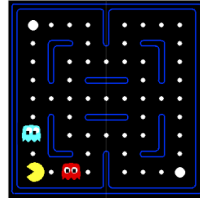
- **On Policy - SARSA:** dùng chiến lược π_b tạo ra hành động để đánh giá π_b (tính $Q^{\pi_b}(s, a)$)

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$$

- Vì giải quyết bài toán dựa trên Q table nên chỉ giải quyết được các bài toán nhỏ, rời rạc

Approximate Q - Learning

- Có những tình huống giống nhau nhưng trạng thái lại khác nhau \rightarrow mất chi phí để học đủ kết quả như nhau



Feature-based Representation

- Biểu diễn trạng thái thông qua các đặc trưng theo mình là quan trọng nhất, như: khoảng cách giữa closest ghost tới pacman, ...
- Có hàm truyền vào trạng thái \rightarrow rút trích đặc trưng

\Rightarrow dùng feature vector

- Biểu diễn giá trị trạng thái và giá trị hành động thành các **hàm tuyến tính**

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Không phải học tăng cường chính xác mà chỉ là xấp xỉ học tăng cường vì ta chỉ dùng những đặc trưng quan trọng nhất chứ không phải toàn bộ đặc trưng \rightarrow approximate \rightarrow đi tìm trọng số $w \rightarrow$ tương tự hồi quy tuyến tính
- Nếu xấp xỉ bằng hàm tuyến tính, rút trích đặc trưng của trạng thái sau khi thực hiện hành động a trong trạng thái s
- Không có nhãn thật \rightarrow không thể làm hồi quy tuyến tính
- Thuật toán:

$$\text{transition} = (s, a, r, s')$$

$$\text{difference} = \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}] \quad \text{Exact Q's}$$

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a) \quad \text{Approximate Q's}$$

- Ở Q learning, tìm $\max Q(s', a')$ ta dựa vào Q table \rightarrow không có Q table, lấy hàm $Q(s, a)$ thể từng hành động a' xong lấy max \rightarrow lấy sample - $Q(s, a)$ (sự khác biệt giữa label và output của mô

hình, label chỉ có 1 phần thông tin chính xác là điểm thưởng lập tức, còn phần hệ quả không là thông tin chính xác vì lấy ra từ mô hình đang học) → vì không có Q table nên thay vì cập nhật $Q(s,a)$, ta sẽ cập nhật trọng số w

- Khi đã hội tụ (tìm được $Q^*(s,a)$) → difference = 0 nhưng nếu đang học thì không thể bằng 0 và do đang sử dụng những đặc trưng quan trọng nhất chứ không phải toàn bộ đặc trưng

⇒ tìm w sao cho difference càng nhỏ càng tốt

- Tại sao phải nhân với $f_i(s, a)$ vì liên quan đến đạo hàm hàm lỗi bình phương

$$\text{error}(w) = \frac{1}{2} \left(y - \sum_k w_k f_k(x) \right)^2$$

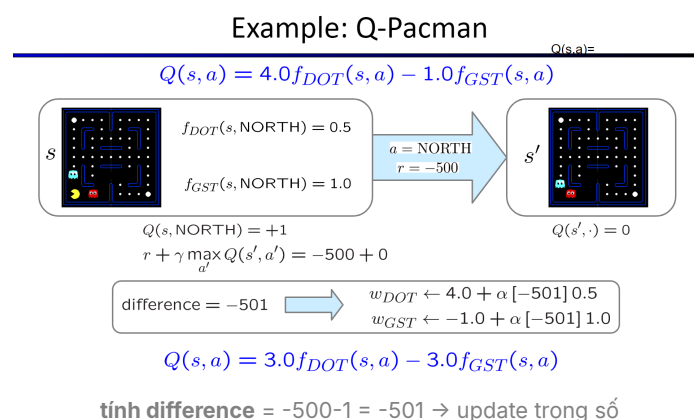
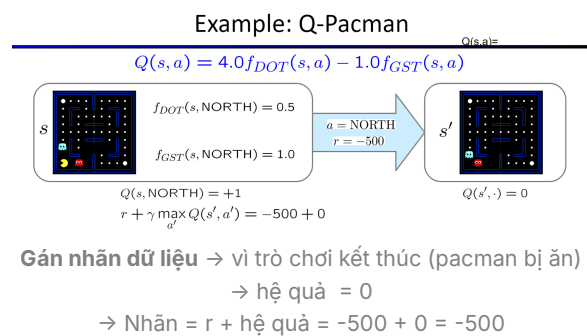
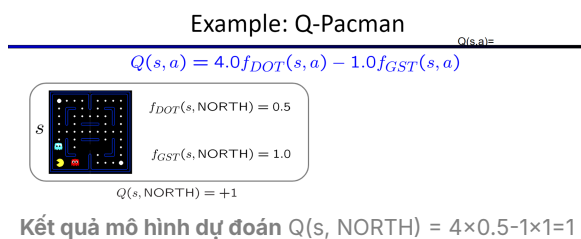
$$\frac{\partial \text{error}(w)}{\partial w_m} = - \left(y - \sum_k w_k f_k(x) \right) f_m(x)$$

$$w_m \leftarrow w_m + \alpha \left(y - \sum_k w_k f_k(x) \right) f_m(x)$$

$$w_m \leftarrow w_m + \alpha \left[\underbrace{r + \gamma \max_{a'} Q(s', a')}_{\text{"target"}} - \underbrace{Q(s, a)}_{\text{"prediction"}} \right] f_m(s, a)$$

- Nếu chạy đủ lâu thì nhãn sẽ được gán chính xác dần

- Ví dụ:



Từng ra thì

- Nếu chọn số đặc trưng đúng bằng số ô trong Q table thì khi đó approximate Q learning sẽ thành exact Q learning → chạy approximate sẽ hội tụ về đúng Q^* . Chuyển Q Learning từ dạng bảng về dạng tuyến tính → chạy đủ lâu thì sẽ hội tụ về Q^*