

Reinforcement Learning for Wumpus

1. Cara kerja

Q-Learning (Off-Policy)

Model ini bekerja dengan menyimpan tabel Q yang berisi “nilai” dari setiap kombinasi state dan aksi. Secara umum, Q-Learning berjalan dengan dua tahap, yaitu interaksi dengan lingkungan dan update tabel Q.

Interaksi dengan lingkungan:

- a. Pertama, tabel Q dibuat dengan nilai nol untuk semua state dan aksi.
- b. Agen memulai dari state awal.
- c. Agen memilih aksi menggunakan metode epsilon-greedy:
 - Dengan probabilitas ϵ , agen mencoba aksi acak (eksplorasi).
 - Dengan probabilitas $1-\epsilon$, agen memilih aksi dengan Q terbesar (eksploitasi).
- d. Agen melakukan aksi, lalu lingkungan mengembalikan reward, state baru, dan tanda selesai atau tidak.

Update tabel Q:

- a. Dari state baru, agen mencari nilai Q tertinggi untuk semua aksi yang mungkin ($\max_a Q(s',a')$).
- b. Nilai Q pada state lama dan aksi yang dipilih diperbarui dengan rumus:

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

c. Parameter:

- α (learning rate): seberapa besar update.
- γ (discount factor): seberapa jauh agen mempertimbangkan reward masa depan.
- d. Ulangi langkah ini sampai episode selesai, lalu ulangi lagi untuk banyak episode.

SARSA (On-Policy)

Model ini hampir sama dengan Q-Learning, tapi cara update Q sedikit berbeda. SARSA memperbarui nilai Q berdasarkan aksi yang benar-benar diambil oleh agen, bukan aksi terbaik secara teori.

Interaksi dengan lingkungan:

- a. Sama seperti Q-Learning, tabel Q diinisialisasi nol.
- b. Agen memulai dari state awal.
- c. Agen memilih aksi dengan epsilon-greedy.
- d. Agen menjalankan aksi, lalu lingkungan mengembalikan reward dan state baru.
- e. Agen langsung memilih aksi berikutnya a' untuk state baru, tetap dengan epsilon-greedy.

Update tabel Q:

a. Nilai Q diperbarui menggunakan reward yang diterima dan aksi yang benar-benar dipilih di state berikutnya:

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

b. Ulangi langkah ini untuk semua transisi sampai episode selesai.

Q-Learning (Off-Policy) meng-update Q menggunakan aksi terbaik (max), sehingga cenderung lebih optimistik, cenderung agresif dalam mencari jalur dengan nilai tinggi, meskipun berisiko. Di lain pihak, SARSA (On-Policy) meng-update Q menggunakan aksi yang benar-benar diambil sehingga lebih realistis, mempertimbangkan eksplorasi, hasilnya lebih aman/konservatif.

2. Analisis Perbandingan Hasil

Berdasarkan eksperimen, didapat hasil sbb:

=== Evaluasi Konvergensi Q-Learning ===

Rata-rata reward setiap 500 episode:

Episode 1 - 500: -293.45

Episode 501 - 1000: -89.48

Episode 1001 - 1500: 273.97

Episode 1501 - 2000: 934.88

Episode 2001 - 2500: 1309.49

Episode 2501 - 3000: 1316.66

Episode 3001 - 3500: 1288.61

Episode 3501 - 4000: 1257.44

Episode 4001 - 4500: 1268.10

Episode 4501 - 5000: 1227.60

Q-Learning dianggap mulai konvergen sekitar episode 473 (rata-rata reward ≥ -100)

Stabilitas akhir: std reward 500 episode terakhir = 955.05

=== Evaluasi Konvergensi SARSA ===

Rata-rata reward setiap 500 episode:

Episode 1 - 500: -266.98

Episode 501 - 1000: 38.19

Episode 1001 - 1500: 723.50

Episode 1501 - 2000: 1145.56

Episode 2001 - 2500: 1475.39

Episode 2501 - 3000: 1449.31

Episode 3001 - 3500: 1331.75

Episode 3501 - 4000: 1267.95

Episode 4001 - 4500: 1256.11

Episode 4501 - 5000: 960.47

SARSA dianggap mulai konvergen sekitar episode 314 (rata-rata reward ≥ -100)

Stabilitas akhir: std reward 500 episode terakhir = 904.95

Policy final (Q-Learning):

```
[2 2 2 1 1 2 2 0 0 1 0 0 3 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 2 2 1 1
2 0 0 0 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 0 0 1 0 1 1 2 0 0
1 0 1 0 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 1 2 0 0 0 0 0 0 0 0
0 3 0 3 3 0 0 0 0 0 0 0 0 0 0 0 0 0 1 2 0 0 0 4 1 4 0 0 0 0 0 0 0 4 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0]
```

Policy final (SARSA):

```
[2 0 1 1 0 1 2 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 2 0 0 1
2 0 1 1 0 0 0 4 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 0 1 1 0 0 1 2 3 0]
```

Perbedaan ini juga tercermin pada jalur (path) optimal yang ditempuh. Q-Learning memang berhasil menemukan jalur menuju emas dan kembali ke permukaan, tetapi perjalanannya sedikit lebih panjang dengan beberapa aksi yang tampak berulang atau tidak efisien. Sementara itu, SARSA menghasilkan jalur yang lebih singkat dan lebih rapi, tanpa banyak langkah tambahan yang tidak perlu. Ini menunjukkan bahwa meskipun Q-Learning lebih cepat belajar, SARSA menghasilkan solusi yang cenderung lebih aman dan efisien.

Secara keseluruhan, analisis ini menunjukkan bahwa Q-Learning unggul dari segi kecepatan, tetapi SARSA lebih realistis dan hati-hati dalam menyusun jalur akhir. Perbedaan ini sejalan dengan teori dasar keduanya: Q-Learning bersifat off-policy sehingga update berdasarkan asumsi aksi terbaik, sedangkan SARSA bersifat on-policy sehingga update berdasarkan aksi nyata yang diambil.