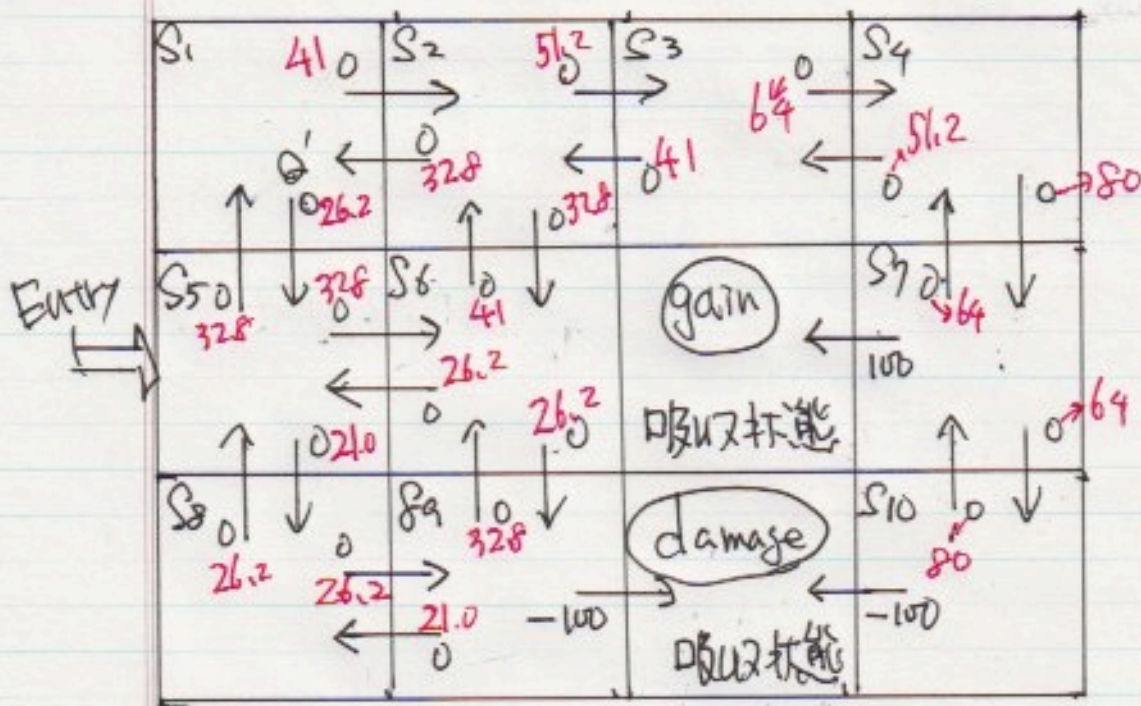


AI2. 12th week

$$Q(s,a) \equiv r(s,a) + \gamma \max_{a'} Q(s,a')$$



$$\gamma = 0.8$$

$$Q = -7.14$$

状態	行動	Q'
S1	→	0
S1	↓	0
S2	→	0
S2	←	0
S2	↓	⋮
⋮	⋮	⋮
S10	↑	0

- Q'
- 41
- 26.2
- 51.2
- 32.8
- ⋮
- ⋮
- 80.

- $\pi(S_1) = \text{right}$
- 学習の結果
- $\pi(S_5) = \uparrow$
- $\pi(S_7) = \rightarrow$
- ⋮
- ⋮
- $\pi(S_9) = \uparrow$
- ⋮

最適行動ポリシー  $\pi$

- (↑, →, →, →, ↓, ←)
- (→, ↑, →, →, ↓, ←)