

① 強化学習の数理モデル

⇒ 離散時間マルコフ意思決定過程
(Discrete MDP)

・ 状態 : $s_t \in S \leftarrow$ 全状態集合
(state) 時間, $t = 0, 1, 2, \dots$

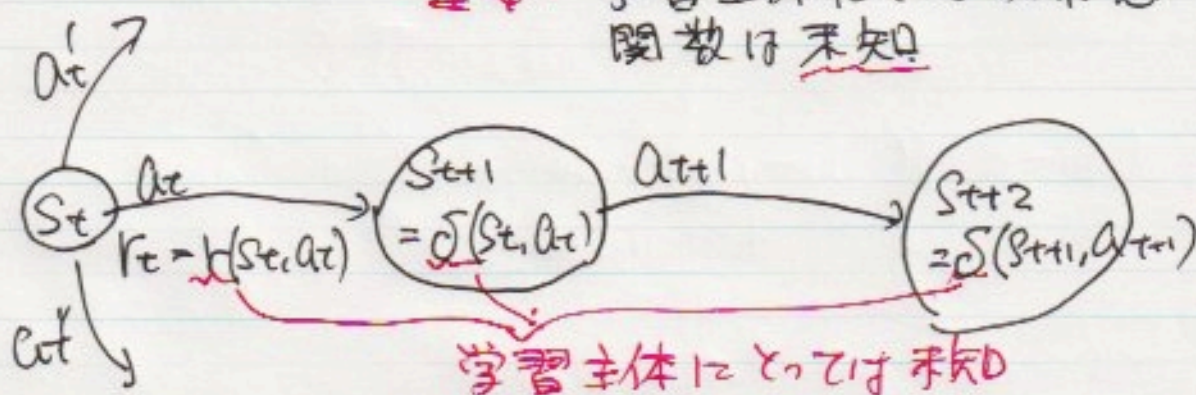
・ 行為 : $a_t \in A \leftarrow$ 全行為集合
(action)

・ 報酬 : $r_t = r(s_t, a_t)$
(reward) ↳ 報酬関数

重要: 学習主体にとっては報酬関数は未知

次状態 $s_{t+1} = \gamma(s_t, a_t)$
↳ 次状態関数

重要: 学習主体にとっては次状態関数は未知



マルコフ性 : $\gamma(s_t, a_t), r(s_t, a_t)$ は s_t と a_t だけ (T) に依存している。

② 強化学習の目標 : 報酬を最大化する

行動ポリシー を学習する
(制御)

$\pi(s_t) = a_t$

制御ポリシーを $\pi(s_t) = a_t$ と t_t と τ

累積報酬 $V^\pi(s_t)$

$$\underline{V^\pi(s_t)} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

$$= \sum_{i=0}^{\infty} \gamma^i r_{t+i}$$

γ : 減衰係数 ($0 \leq \gamma < 1$)

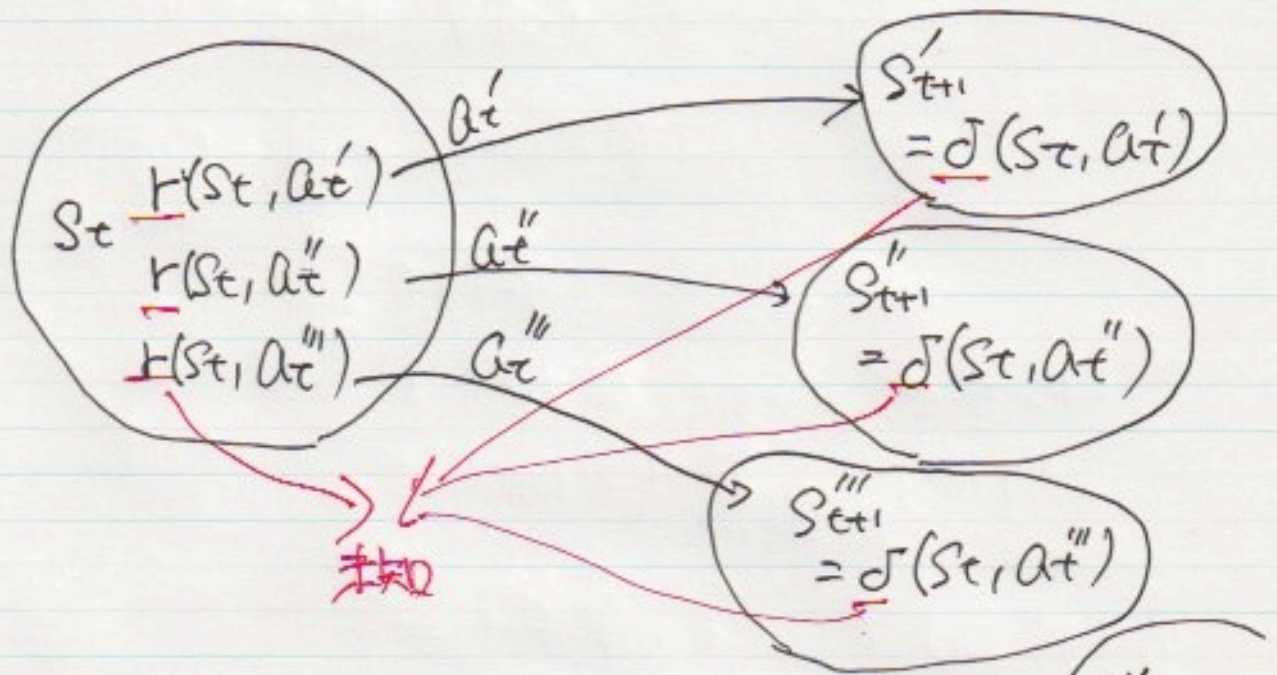
技術的理由: $V^\pi(s_t)$ を有限の値にするために導入する.

(\Rightarrow 社会学における未来係数と)
関係がある.)

\rightarrow 強化学習の目的は $V^\pi(s_t)$ を最大化すること

② Q学習: $V^*(s_t)$ の最大値を推定

3



$$r(s, a) + \gamma V^*(\gamma(s, a))$$

未知なので推定ね

A* 最適な
経路を
探索

推定値 Q を導入

$$Q(s, a) \stackrel{\text{def}}{=} r(s, a) + \gamma \overset{*}{V}(\gamma(s, a))$$

最適値

$$\text{つまり } V^*(\gamma(s, a)) = \max_{a'} Q(\gamma(s, a), a')$$

$$\text{つまり } Q(s, a) = r(s, a) + \gamma \max_{a'} Q(\gamma(s, a), a')$$