

H29 ①

- (1) 定義に従って動的計画法により図3を埋めると以下の結果が得られる

-1	0	0
3	0	0
3	2	0

//

- (2) (1)の動的計画法の結果を用いてテーブルを復元することで、最適な経路は以下のようになる

↑	→	→
↑		

//

- (3) $V(s) = \max_{a \in \{UP, DOWN\}} \{ R_{ss'}^a + V(s') \}$ の定義に基づいて

迷路の各マスの $V(s)$ の値を上の行の右のマスから動的計算により求めることを行えばよい。

(図3を見ると更新の向きは右図の通り)

①	→
②	→
...	→

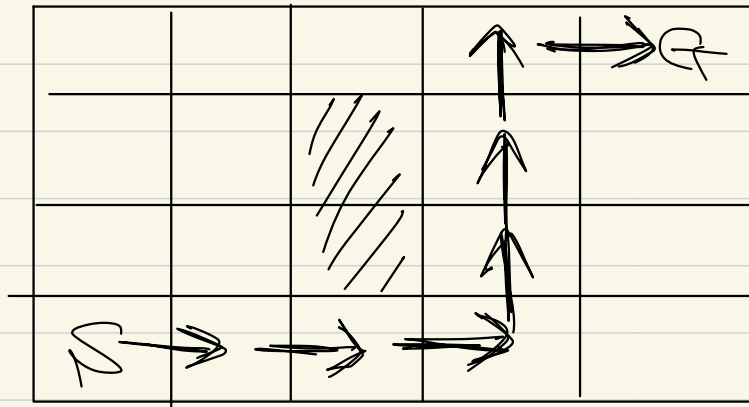
これにより得られる $V(s)$ の値をテーブルにまとめると次ページのように表現できる。

2	0	0	0	0
4	2	////	0	0
6	4	////	10	-5
10	10	10	10	-10

(通行不能な
マスは // 表記)

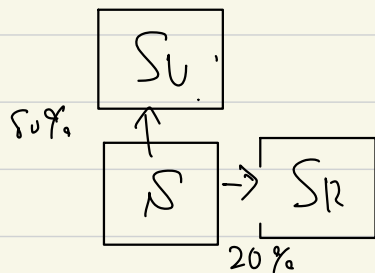
//

- (4) (3)の動的計画法テーブルのデータから得られる
最大利益値は10であり、それを与えるパスを復元するため
始点Sから終点Gに向けてテーブルを戻ると、最適経路
は以下のように復元できる。

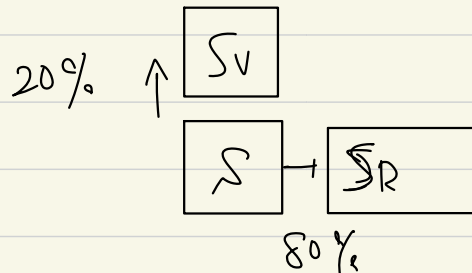


//

以下、今いるマスを P とした時、その上のマスを P_U 、その右
のマスを P_R と定義する。また、便宜上障害物や迷路外にも
アクセスできるものとし、そのような位置 x において、 $V(x) = -\infty$ とする。



① 上を選択



② 右を選択

(但し、どちらかが
不正なルートにある場合
残された方向へ確率1で
移動する)

この時、 $V(S)$ を計算するにあたり、進む方向として上右を選択する方向に
①②のパターンがあらう

① を選択した場合

$$V(s) = \frac{4}{5} \{ R_{s, s_u}^a + V(s_u) \} + \frac{1}{5} \{ R_{s, s_R}^a + V(s_R) \}$$

② を選択した場合

$$V(s) = \frac{1}{5} \{ R_{s, s_u}^a + V(s_u) \} + \frac{4}{5} \{ R_{s, s_R}^a + V(s_R) \}$$

という関係式が成り立つ,

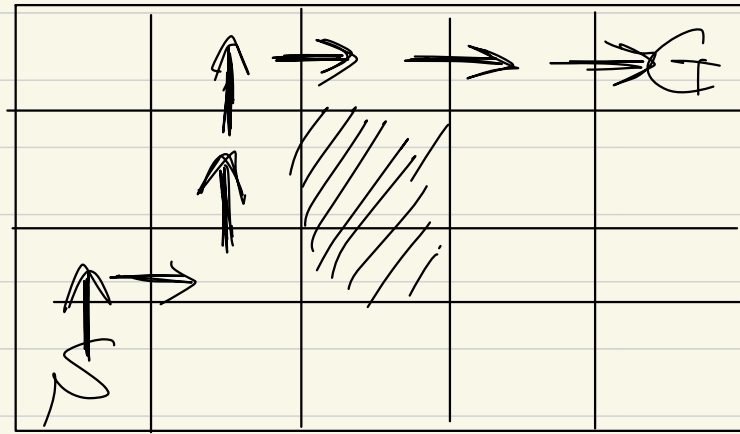
ゆえに どちらにも行くことが可能である場合

$$V(s) = \max \left\{ \frac{4}{5} (R_{s, s_u}^c + V(s_u)) + \frac{1}{5} (R_{s, s_R}^c + V(s_R)), \right. \\ \left. \frac{1}{5} (R_{s, s_u}^c + V(s_u)) + \frac{4}{5} (R_{s, s_R}^c + V(s_R)) \right\}$$

- (6) (5) で ~~定め~~ 仮定と 表において得られる漸化式に基き (そうであるなら (1)~(4)と同じ)
 (3) で更新順序と同じ順序で動的にテーブルを更新
 すると以下の結果を得る、(但し結果は小表第1位までとしている)

2	0	0	0	0
3.2	2	////	-1	0
5.8	4	////	5.2	-5
7.2	5.0	1.1	1.1	-10

- (7) (4) と同様の方法で最適指示路を復元すると、次ページ
 のよう(2)になる。



//

(8) (4) で与えられる最適経路は、それが極めて限定された経路にのみしか達成しえず、そこから外れていけば得られる報酬値は急速に悪化する。

一方、(7) で与えた最適経路では、複数の最大報酬を得ることができる経路が存在するため、仮に経路が指す方向から外れても解の劣化が起りにくい。

その差異が、得られる報酬の期待値として表れたととらえられる。