

## 1. Нова верзија на Phone или не?

Замислете дека користите Маркови процеси на одлучување (MDPs) при решавање на проблемот на компанијата Apple, која секои 6 месеци треба да донесе одлука дали да се појави на пазарот со нова верзија на iPhone или да причека. проблемот може да се претстави со 3 состојби **P**, **SS** и **N**, кои одговараат на позитивно, така-така (неутрално) или негативно расположение/сентимент кон компанијата. Дозволените акции се **Нова верзија** и **Причекај**. Наградата, која се добива за било која акција која го води системот во состојбата на позитивен сентимент кон компанијата **P** е +2 (наградата се добива и за премин од **P** во **P**). Наградата за премин во неутралната состојбата **SS** е 0, додека состојбата **N** се смета за непожелна и за премин во неа се добива негативна награда од -1. Факторот на намалување е  $\gamma = 1$ .

Веројатностите за премин од една во друга состојба се прикажани во следната табела.

s	a	T(s, a, s')		
		s' = P	s' = SS	s' = N
P	Нова верзија	0.1	0.9	0
	Причекај	0.2	0.8	0
SS	Нова верзија	0.1	0.9	0
	Причекај	0	0.3	0.7
N	Нова верзија	0.9	0	0.1
	Причекај	0	0.5	0.5

За дадениот проблем чие решавање треба да го дефинирате како MDP, потребно е да одговорите на следните прашања/задачи:

**(а)** Да се пресметаат вредностите на состојбите,  $V(s, a, s')$  &  $Q(s, a, s')$  за првите 3 чекори од постепено рекурзивно пресметување на состојбите (првите 3 итерации) и пополни табелата. Образложете ги пресметките преку формулите кои сте ги користеле при пополнување на табелата.

	P	SS	N
$V^*_0(s)$			
$Q^*_1(s, \text{Нова верзија})$			
$Q^*_1(s, \text{Причекај})$			
$V^*_1(s)$			
$Q^*_2(s, \text{Нова верзија})$			
$Q^*_2(s, \text{Причекај})$			
$V^*_2(s)$			

(6) Која е оптималната политика која Агентот ќе ја преземе доколку се наоѓа во состојбата **SS** и има уште два чекори до крајот на играта? Образложете го вашето решение.

## 2. Во потрага по богатство

Агентот се движи во лавиринт каде е скриено богатство. Акциите кои може да ги преземе агентот се Лево (**L**) или Десно (**R**), но поради стохастичноста може да заврши во поле различно од очекуваното. Ако агентот избере акција Десно со веројатност 0.5 ќе се придвижи во посакуваната насока, но со иста со веројатност 0.5 може да се лизне и падне во бездна каде играта завршува со казна од -4. Ако агентот избере акција Лево се поместува во полето лево со веројатност 1. Воедно, постои и акцијата **Exit** за излез од терминалните полиња, кои се означена со награда. Вредноста на  $\gamma = 1$

	-4	-4	-4	-4	
+100 $s_0$	$s_1$	$s_2$	$s_3$	$s_4$	+100 $s_5$
	-4	-4	-4	-4	

За дадениот проблем чие решавање треба да го дефинирате како MDP, потребно е да одговорите на следните прашања/задачи:

(а) Да се пресметаат вредностите на состојбите,  $V(s, a, s')$  за првите 3 чекори од постепеното рекурзивно пресметување на состојбите ( $i=0, i=1, i=2$ ) и пополни табелата. Образложете ги пресметките преку формулите кои сте ги користеле при пополнување на табелата.

Наоѓате мапа на скриено богатство од претходниот трагач со запис за неговото искуство, т.е политиката која ја преземал и вредноста на состојбата:

$$\begin{aligned} \pi(s_1) &= R & \pi(s_2) &= L & \pi(s_3) &= L & \pi(s_4) &= R \\ V(s_1) &= -4 & V(s_2) &= -4 & V(s_3) &= 8 & V(s_4) &= 10 \end{aligned}$$

(б) Следејќи ја политиката на претходниот трагач, направете евалуација на политиката за првите 3 чекор ( $i=0, i=1, i=2$ ) и пресметајте ги вредностите на сите 4 состојби, кои не се терминални.

(в) Направете подобрување на политиката, врз основа на вредностите на состојбите добиени во претходното барање за чекор  $i=2$ !