

(1A). (a). As the prob. are not given, we can assume equal probabilities and hence arm 3 is most preferred and arm 1 is least preferred.

(b). $E[r_1] \approx 5.17$; $E[r_2] = 4.88$; $E[r_3] = 4.32$

(c). Arm 1 is best as it has highest expected reward.

(d). Write either ϵ -greedy / UCB discussed in class.

(2A).

(a).

State	Action	Next state	Prob.	Reward
R	M	R	0.6	$-20 + 100$
R	M	B	0.4	-20
R	NM	R	0.3	100
R	NM	B	0.7	0
B	Re	R	0.6	$-40 + 100$
B	Re	B	0.4	-40
B	RP	R	1	$-150 + 100$

Here, Running - R ; broken - B ; maintenance - M, no maintenance - NM, repair - Re ; replace - RP.

(b).
$$V^*(R) = \max_{\{M, NM\}} \begin{cases} 0.6 [80 + \gamma V^*(R)] + 0.4 [-20 + \gamma V^*(B)], \\ 0.3 [100 + \gamma V^*(R)] + 0.7 [0 + \gamma V^*(B)] \end{cases}$$

$$V^*(B) = \max_{\{Re, RP\}} \begin{cases} 0.6 [60 + \gamma V^*(R)] + 0.4 [-40 + \gamma V^*(B)], \\ 1 [-50 + \gamma V^*(R)] \end{cases}$$

(3A). (a). $G_0 = 1 + 2x + (2x)^2 + (2x)^3 + \dots = \frac{1}{1-2x} = 5.$

$$G_1 = 2 + 4x + 8x^2 + 16x^3 + \dots$$

$$= 2(1 + 2x + (2x)^2 + (2x)^3 + \dots) = 2\left(\frac{1}{1-2x}\right) = 10.$$

(b). $G_0 = 2 + 4x + 2x^2 + 4x^3 + 2x^4 + \dots$

$$= 2(1 + x^2 + x^4 + \dots) + 4(x + x^3 + \dots)$$

$$= 2\left(\frac{1}{1-x^2}\right) + 4\left(\frac{x}{1-x^2}\right) = \frac{5.2}{0.36} \approx 14.45$$

$$G_1 = 4 + 2x + 4x^2 + \dots$$

$$= 4(1 + x^2 + \dots) + 2(x + x^3 + \dots)$$

$$= 4\left(\frac{1}{1-x^2}\right) + 2\left(\frac{x}{1-x^2}\right) = \frac{5.6}{0.36} \approx 15.56.$$