# Two-arm bandits

Two actions; A and B.

Action A: always reward $(+6)$

Action B:

if B is of type lucky $(L=l)$ or unlucky $(L=\neg l)$

| B is lucky     CPT |
| --- |
| $P(r=+10 \mid L=l) = 4/5$ |
| $P(r=0 \mid L=l = 1/5$ |

| B is unlucky     CPT |
| --- |
| $P(r=+10 \mid L=\neg l) = 1/5$ |
| $P(r=0 \mid L=\neg l) = 4/5$ |

Discount factor $\gamma = 1$

A-priori $P(L=l) = P(L=\neg l) = 0.5$

## (a) One step time horizon:

MEU for action A: $\boxed{+6}$     [There is no uncertainty in lever A]

MEU for action B:
$$\underset{\underset{\text{reward}}{\uparrow}}{+10} \times \left( \underbrace{\frac{1}{2} \times \frac{4}{5}}_{\underset{\underset{P(r=10 \mid L=l)}{\text{lucky}}}{}} + \underbrace{\frac{1}{2} \times \frac{1}{5}}_{\underset{\underset{P(r=+10 \mid L=\neg l)}{\text{unlucky}}}{}} \right)$$

$$+ \quad 0 \times \left( \frac{1}{2} \times \frac{1}{5} + \frac{1}{2} \times \frac{4}{5} \right)$$
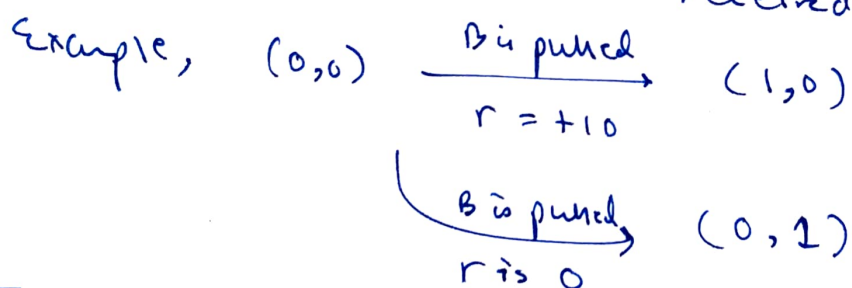
$$= 10 \times \frac{1}{2} = \boxed{+5}$$

MEU(a) > MEU(b) $\Rightarrow$ Action $\underline{\underline{A}}$ should be taken.

) Two-steps in the future

· Note: There is inherent uncertainty in the type of lever B. As we plan ahead, the resulting rewards can change the ~~type~~ belief over the type of lever B. The reward can be treated as an observation that updates the likelihood.

MDP formulation suggested:

State $(m, n)$ where m denotes the # of times lever B was pulled & reward of $(+10)$ was received

n denotes the # of times lever B was pulled & reward of $(0)$ was received.

Example, $(0,0) \xrightarrow[\text{r = +10}]{\text{B is pulled}} (1,0)$

$\xrightarrow[\text{r is 0}]{\text{B is pulled}} (0,1)$

(ii) In state $(0,1)$, take action B, can land up in

· State $(1,1)$ if the reward is +10 with probability ~~$p(L=\ell) p(r=+10|L=\ell) + p(L=7\ell) p(r=+10|L=7\ell)$~~

$$p(L=\ell|r=0)^2 \; p(r=+10|L=\ell) + p(L=7\ell|r=0) \cdot p(r=+10|L=7\ell)$$

Note: previously reward was 0 in the first pull, this obsn will update the bel over the type of lever B is.

$$p(L=\ell|r=0) = \frac{p(r=0|L=\ell)\, p(L=\ell)}{p(r=0|L=\ell)\, p(L=\ell) + p(r=0|L=7\ell)\, p(L=7\ell)}$$

$$= \left[ \frac{1/5 \times \frac{1}{2}}{(1/5 \times 1/2) + (4/5 \times \frac{1}{2})} \right] = \left[ \frac{1}{5} \right]$$

$p(L=7\ell|r=0) = 1 - 1/5 = 4/5$

w, from state $(0,1)$ taking action $B$ can be in

- State $(1,1)$ if reward is $+10$ win probability

$$p(L=\ell \mid r=0) \, p(r=+10 \mid L=\ell) + p(L=7\ell \mid r=0)$$
$$p(r=+10 \mid L=7\ell)$$

$$= \frac{1}{5} \times \frac{4}{5} + \frac{4}{5} \times \frac{1}{5} = \frac{8}{25}$$

- State $(0,2)$ if reward is $(+0)$ with probability

$$p(L=\ell \mid r=0) \, p(r=0 \mid L=\ell) + p(L=7\ell \mid r=0) \, p(r=0 \mid L=7\ell)$$

$$= \frac{1}{5} \times \frac{1}{5} + \frac{4}{5} \times \frac{4}{5} = \frac{17}{25}$$

==

(ii) Now, from state $(1,0)$, taking action $B$ can be in

- State $(2,0)$ if reward is $+10$ with probability

$$p(L=\ell \mid r=0) \, p(r=+10 \mid L=\ell) + p(L=7\ell \mid r=+10)$$
$$p(r=10 \mid L=7\ell)$$

We require $p(L=\ell \mid r=10) = \dfrac{p(r=10 \mid L=\ell)\, p(L=\ell)}{p(r=10 \mid L=\ell)\, p(L=\ell) + p(r=10 \mid L=7\ell)\, p(L=7\ell)}$

$$= \frac{\frac{4}{5} \times \frac{1}{2}}{\frac{4}{5} \times \frac{1}{2} + \frac{1}{5} \times \frac{1}{2}} = \frac{4}{5}$$

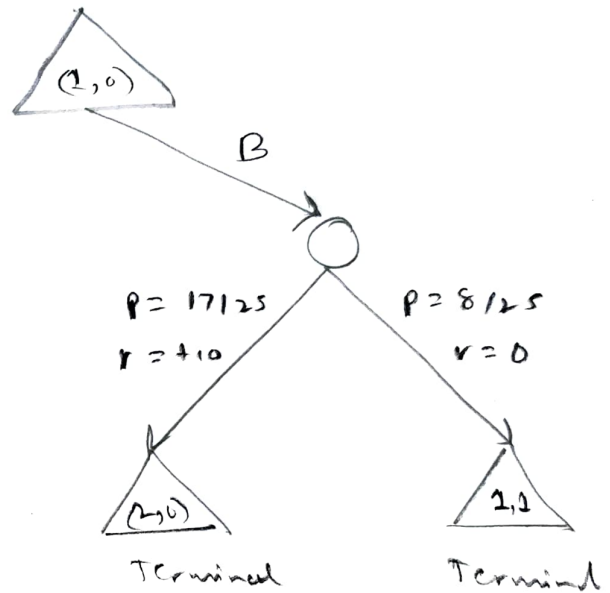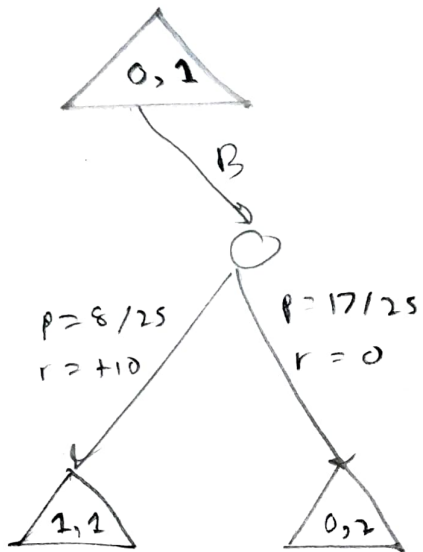$$p(L=7\ell \mid r=+10) = \frac{1}{5}$$

New state probability is:

$$\frac{4}{5} \times \frac{4}{5} + \frac{1}{5} \times \frac{1}{5} = \frac{17}{25}$$

- State $(1,1)$ if round is to win probability

$$P(L=\ell \mid r=10) \times P(r_2 \text{ } 10 \mid L=\ell) + P(L=7\ell \mid r=+10) \times$$
$$P(r=0 \mid L=7\ell)$$

$$= \frac{4}{5} \times \frac{1}{5} + \frac{1}{5} \times \frac{4}{5} = \frac{8}{25}$$
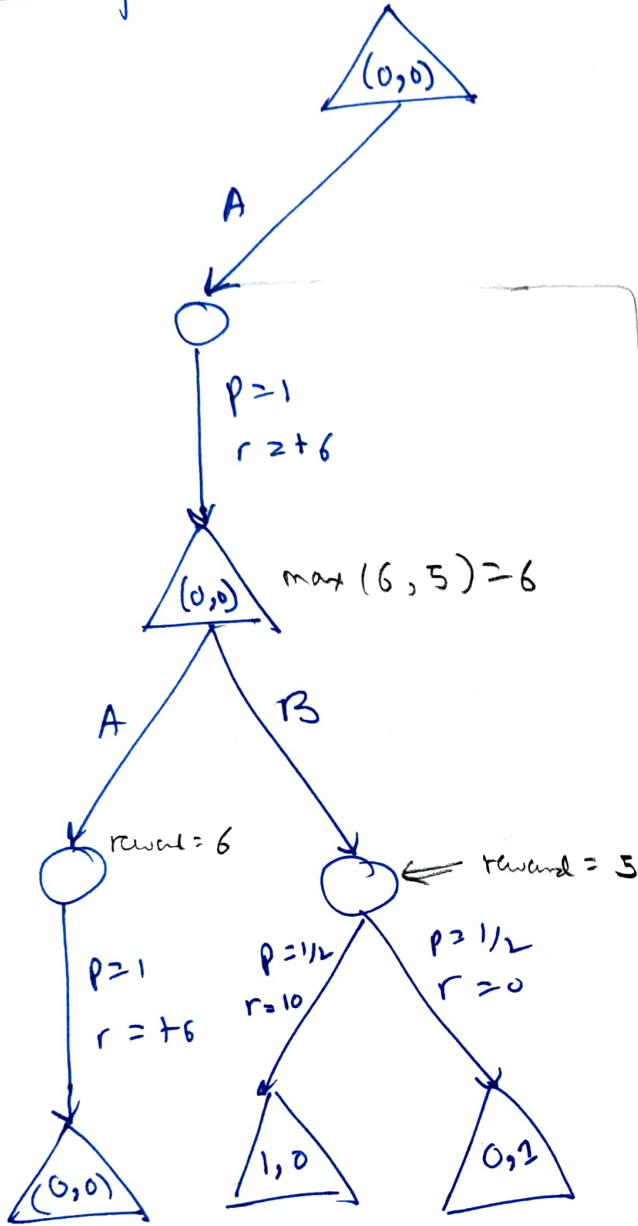
These transitions as trees

iii) Computation Tree

) Which action is optimal A or B or both by MEU.

; For left sub-tree



(0,0)

A

$p = 1$
$r = +6$

(0,0)    max $(6, 5) = 6$

A          B

reward = 6          ← reward = 5

$p = 1$        $p = 1/2$        $p = 1/2$
$r = +6$       $r = 10$         $r = 0$

(0,0)          1,0              0,1
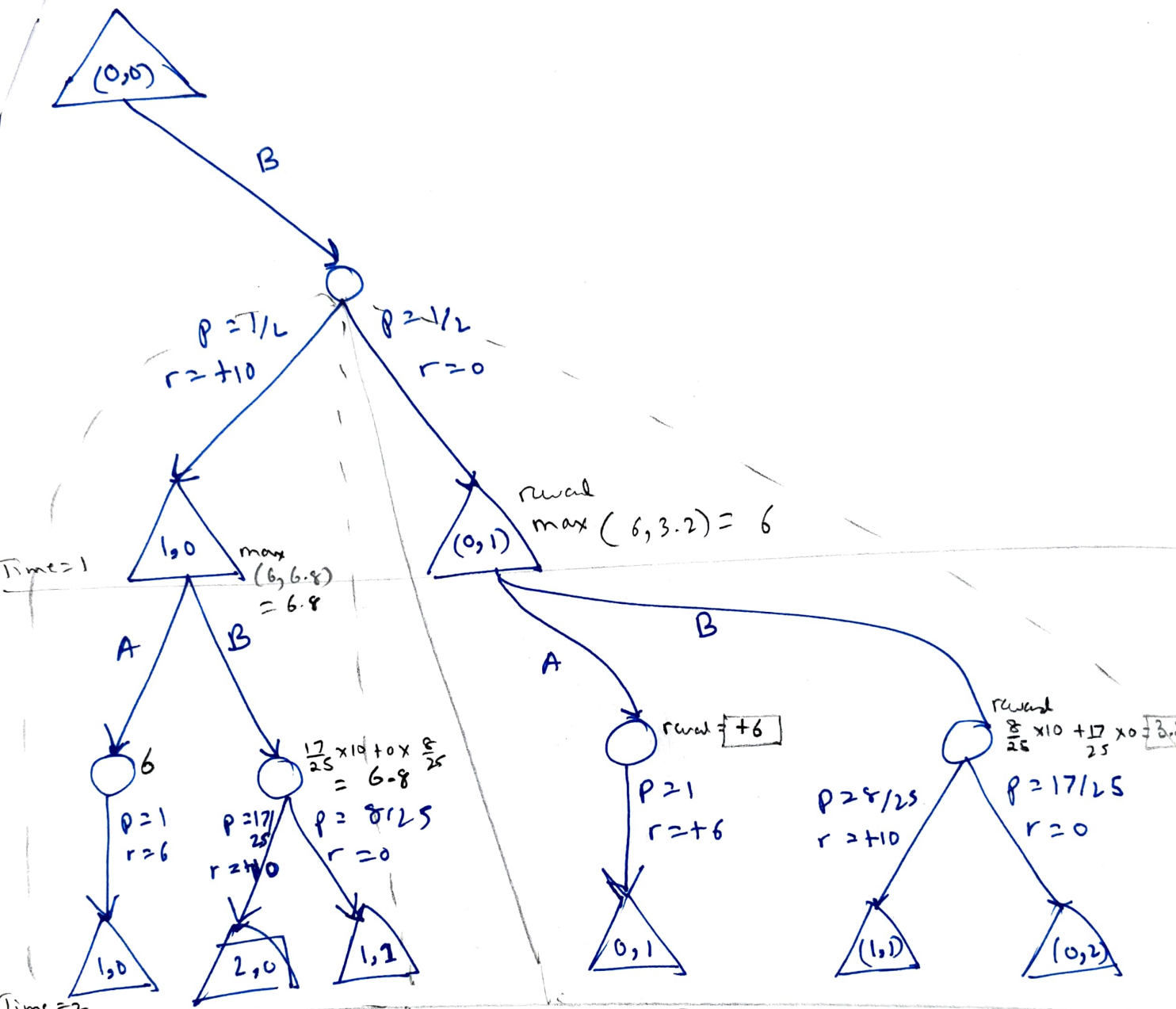
The total reward for the sub-tree after taking action A in state (0,0) is

$6 + \gamma \cdot 6 = \boxed{12}$

as $\gamma = 1$

For the right sub-tree



(0,0)

B

$p = 1/2$
$r = +10$

$p = 1/2$
$r = 0$

Time = 1

1,0   max
(6, 6.8)
= 6.8

A   B

6
$p = 1$
$r = 6$

$p = 17/25$
$r = +10$

$\frac{17}{25} \times 10 + 0 \times \frac{8}{25}$
$= 6.8$

$p = 8/25$
$r = 0$

Time = 2

1,0   2,0   1,2

(0,1)   reward
max (6, 3.2) = 6

B

A

$p = 1$
$r = +6$   reward = +6

$\frac{8}{25} \times 10 + \frac{17}{25} \times 0 = 3.2$   reward

$p = 8/25$
$r = +10$

$p = 17/25$
$r = 0$

0,1   (1,1)   (0,2)

Left
Reward $= +10 + 8 \cdot 6.8$

$= 10 + 6.8$

$= \boxed{16.8}$

Right
Reward $= 0 + 8 \cdot 6 = \boxed{6}$
as $r = 1$

expected $= \frac{1}{2} \times 16.8 + \frac{1}{2} \times 6$

$= 8.4 + 3 = \boxed{11.4}$

MEU (A) > MEU (B)   still action A is optimal.