## Question 2. Markov Decision Process.
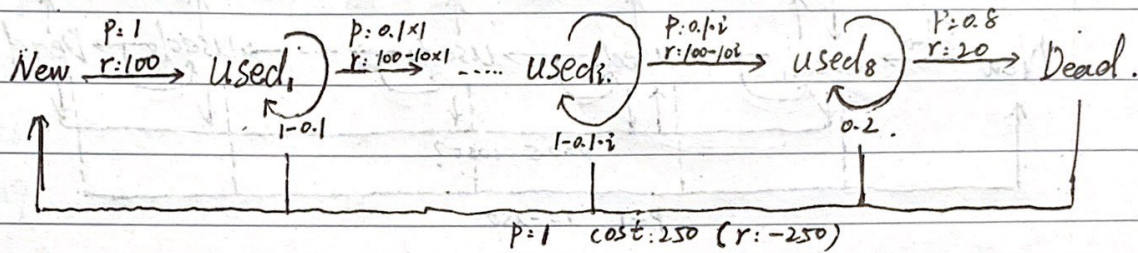
~~The~~ According to the description of the problem, we can transform it to the flow graphy below:



New $\xrightarrow[r:100]{p:1}$ Used$_1$ $\xrightarrow[r:100-10 \times 1]{p:0.1 \times 1}$ .... Used$_i$ $\xrightarrow[r:100-10i]{p:0.1 \cdot i}$ Used$_8$ $\xrightarrow[r:20]{p:0.8}$ Dead.

with loops: $1-0.1$, $1-0.1 \cdot i$, $0.2$

$p:1$ cost:250 $(r:-250)$

[ p represents possibility and r represents reward ]

(a)
(b) $\longrightarrow$ According to the flow graphy above, now we have ~~been~~ update prosses as below given random initial utility: $U(New) = U(used_i) = U(Dead) = 0$

$U_{n+1}(New) = 100 + \beta U_n(U_1)$

for $i = 1$ to $7$

$U_{n+1}(used_i) = \max\left[\left(100 - 10 \cdot i + \beta\left((1 - 0.1 \times i) \cdot U_n(used_i) + 0.1 \times i \cdot U_n(used_{i+1})\right)\right), -250 + \beta U_n(new)\right]$

$U_{n+1}(used_8) = \max\left[\cancel{100} \ 20 + \beta\left(0.2 \cdot U_n(used_8) + 0.8 \cdot U_n(Dead)\right), -250 + \beta U_n(new)\right]$

$U_{n+1}(Dead) = -250 + \beta U_n(New)$
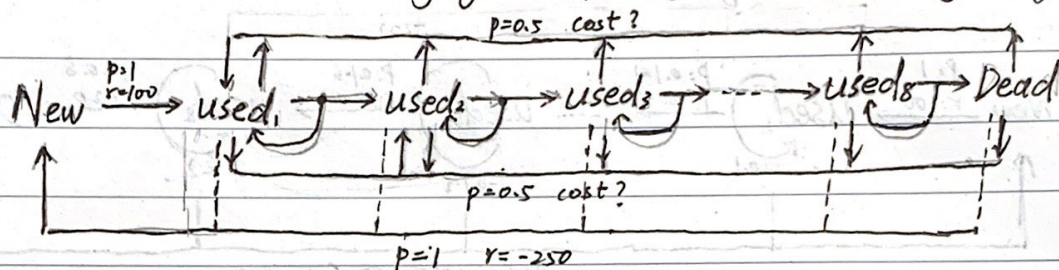
follow this process until results converge.

We define 'converge' like this: When $\max[U_{n+1} \text{ of all states} - U_n \text{ of all state}] < 0.001$

for question (a)(b)

I write a small code to solve this, it ~~stop~~ converges at around $120 \sim 140$ times of iterations. The final utility and optimal policy are shown in the ~~log.txt~~ log_ab.txt. As we can see, when it converges, the policy for 'New' and 'used$_1$'-'used$_5$' is 'use', and the policy for 'used$_6$'-'used$_8$' and 'Dead' is 'replace'.

(c) The option to buy a ~~fresh~~ used_machine ~~$~~ instead of buying a new-machine is changing the flow graph in following way:



According to this graphy, we now have a new update process.

$$U_{n+1}(New) = 100 + \beta U_n(U_1)$$

for $s=1$ to $7$: $U_{n+1}(U_s) = \max\left\{ (100-10s + \beta((100-0.1s)U_n(U_s) + 0.1s \cdot U_n(U_{s+1}))), -250 + \beta U_n(new),\right.$

$$\left. -cost + \beta(0.5 U_n(U_1) + 0.5 U_n(U_2)) \right\}$$

$$U_{n+1}(U_8) = \max\left\{ \cdot(100+0s) + \beta(\overset{0.2}{\cancel{\phantom{xx}}} U_n(U_8) + 0.8 U_n(Dead)), -250 + \beta U_n(New),\right.$$

$$\left. -cost + \beta(0.5 U_n(U_1) + 0.5 U_n(U_2)) \right\}$$

$$U_{n+1}(Dead) = \max\left\{ -250 + \beta U_n(New), -cost + \beta(0.5 U_n(U_1) + 0.5 U_n(U_2)) \right\}$$

We constantly change the cost of buying a used_machine, and we find that the threshold is 169-170. When the cost changes from 169 to 170, the policy for $U_6/U_7/U_8/Dead$ changes from 'Buying used_machine' to 'Buying new-machine. So the highest price shoud be set to 169 for which used_machine is the rational choice. (More data detail in log_c.txt)

(d) we change the β and record the policy when flow converges. And Below is the table of policy.

| β | New | Used1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | Dead |
|---|-----|-------|----|----|----|----|----|----|----|------|
| 0.1~0.8 | u | u | u | u | u | u | u | u | u | re |
| 0.85 | u | u | u | u | u | u | u | re | re | re |
| 0.9 | u | u | u | u | u | u | re | re | re | re |
| 0.93 | u | u | u | u | u | re | re | re | re | re |
| 0.95 | u | u | u | u | u | re | re | re | -re | re |
| [0.96~] | u | u | u | u | re | re | re | re | re | re |

[u: 'use'    re: 'replace']

I ~~use~~ do tests on the base of problem (a) & (b). During the test, we can find that after β gets bigger than 0.96, the policy becomes stable, and the best policy is shown above.

A little thought here: when β is small, it pays more attention to current reward, so each state takes 'use' as optimal policy, for it immediately get rewards. When β is bigger, it gradually takes future ~~~~ utility into accounts. When β is big enough, the future utility compensate for the current reward, and the optimal policy turns from 'use' to 'replace'


Bonus. When the costs of a new machine is 250, ~~as~~ the long term discounted value is 800.53 as we computed in part (a) & (b)

We continue raising the cost of new machine, And we ~~find~~ compute the corresponding utility of new machine. We find

that around ~~749~~, cost of 749, the flow reaches 0 gain (the utility of new machine $\approx$ the cost of a new machine). Below that cost, we are operating at a net gain. Above that cost, we are operating at a net loss.