

1) State Space: {warm, cool, Off}

Actions: {fast, slow}

Transition Model: {

$P(\text{cool} \mid \text{cool}, \text{slow}) = 1$   
 $P(\text{warm} \mid \text{cool}, \text{slow}) = 0$   
 $P(\text{cool} \mid \text{cool}, \text{fast}) = \frac{1}{4}$   
 $P(\text{warm} \mid \text{cool}, \text{fast}) = \frac{3}{4}$   
 $P(\text{cool} \mid \text{warm}, \text{slow}) = \frac{1}{4}$   
 $P(\text{warm} \mid \text{warm}, \text{slow}) = \frac{3}{4}$   
 $P(\text{cool} \mid \text{warm}, \text{fast}) = 0$   
 $P(\text{warm} \mid \text{warm}, \text{fast}) = \frac{7}{8}$   
 $P(\text{off} \mid \text{warm}, \text{fast}) = \frac{1}{8}$   
 $P(\text{off} \mid \text{off}, \text{slow}) = 1$   
 $P(\text{off} \mid \text{off}, \text{fast}) = 1$

}

Reward Function: {

off  $\rightarrow$  any action  $\rightarrow$  off = 0  
warm or cold  $\rightarrow$  slow  $\rightarrow$  warm or cold = 4  
warm or cold  $\rightarrow$  fast  $\rightarrow$  warm or cold or off = 10

}

2) With the policy of always going fast:

$U(\text{off}) = 0$

$U(\text{warm}) = \frac{7}{8}(10 + U(\text{warm})) + \frac{1}{8}(10 + U(\text{off}))$

$U(\text{cool}) = \frac{1}{4}(10 + U(\text{cool})) + \frac{3}{4}(10 + U(\text{warm}))$

Solving gives  $U(\text{warm}) = 80$ ,  $U(\text{cool}) = 93.33$

With the policy of always going slow:

$U(\text{off}) = 0$

$U(\text{warm}) = \frac{3}{4}(4 + U(\text{warm})) + \frac{1}{4}(4 + U(\text{cool}))$

$U(\text{cool}) = 1 * (4 + U(\text{cool}))$

This solution is undefined because the  $U(\text{cool})$  is actually infinite. If we drive slowly then we will stay cool forever and continue to get 4 reward points

3) First policy eval. I just added a discount factor to the first policy above and resolved the system.  $U(\text{off}) = 0$ ,  $U(\text{warm}) = 47.06$ ,  $U(\text{cool}) = 53.89$

Because the order of states to iterate doesn't matter, I'll do cool then warm then off.

State cool, action slow  $\rightarrow 1 * 53.89 = 53.89$

State cool, action fast  $\rightarrow \frac{1}{4}(53.89) + \frac{3}{4}(47.06) = 48.76$

Argmax(actions) means update policy cool to be slow,

State warm, action slow  $\rightarrow \frac{3}{4}(47.06) + \frac{1}{4}(53.89) = 48.77$

State warm, action fast  $\rightarrow \frac{7}{8}(47.06) + \frac{1}{8} * 0 = 41.18$

Argmax(actions) means update policy warm to be slow

State off, action slow  $\rightarrow 0$

State off, action fast  $\rightarrow 0$

Argmax keeps it the same, when off try to go fast

Next  $U = \text{policyeval}$

$U(\text{off}) = 0$

$U(\text{warm}) = \frac{3}{4} * (4 + .9 * U(\text{warm})) + \frac{1}{4} * (4 + .9 * U(\text{cold}))$

$U(\text{cool}) = 1 * (4 + .9 * U(\text{cool}))$

Solving gives  $U(\text{cool}) = 40$ ,  $U(\text{warm}) = 40$ ,  $U(\text{off}) = 0$

Now iterate over actions

State cool, action slow  $\rightarrow 1 * 40 = 40$

State cool, action fast  $\rightarrow \frac{1}{4} * 40 + \frac{3}{4} * 40 = 40$

If argmax breaks ties by keeping the current policy then stay slow on cool

State warm, action slow  $\rightarrow \frac{3}{4} * (40) + \frac{1}{4} * (40) = 40$

State warm, action fast  $\rightarrow \frac{7}{8} * (40) + \frac{1}{8} * 0 = 35$

Argmax keeps it the same

State off doesn't really matter, they all end up as zero, policy doesn't matter either.

Policy after second iteration depends on how argmax breaks ties. If it sticks with the current policy then the policy is cool  $\rightarrow$  slow, warm  $\rightarrow$  slow, off  $\rightarrow$  fast