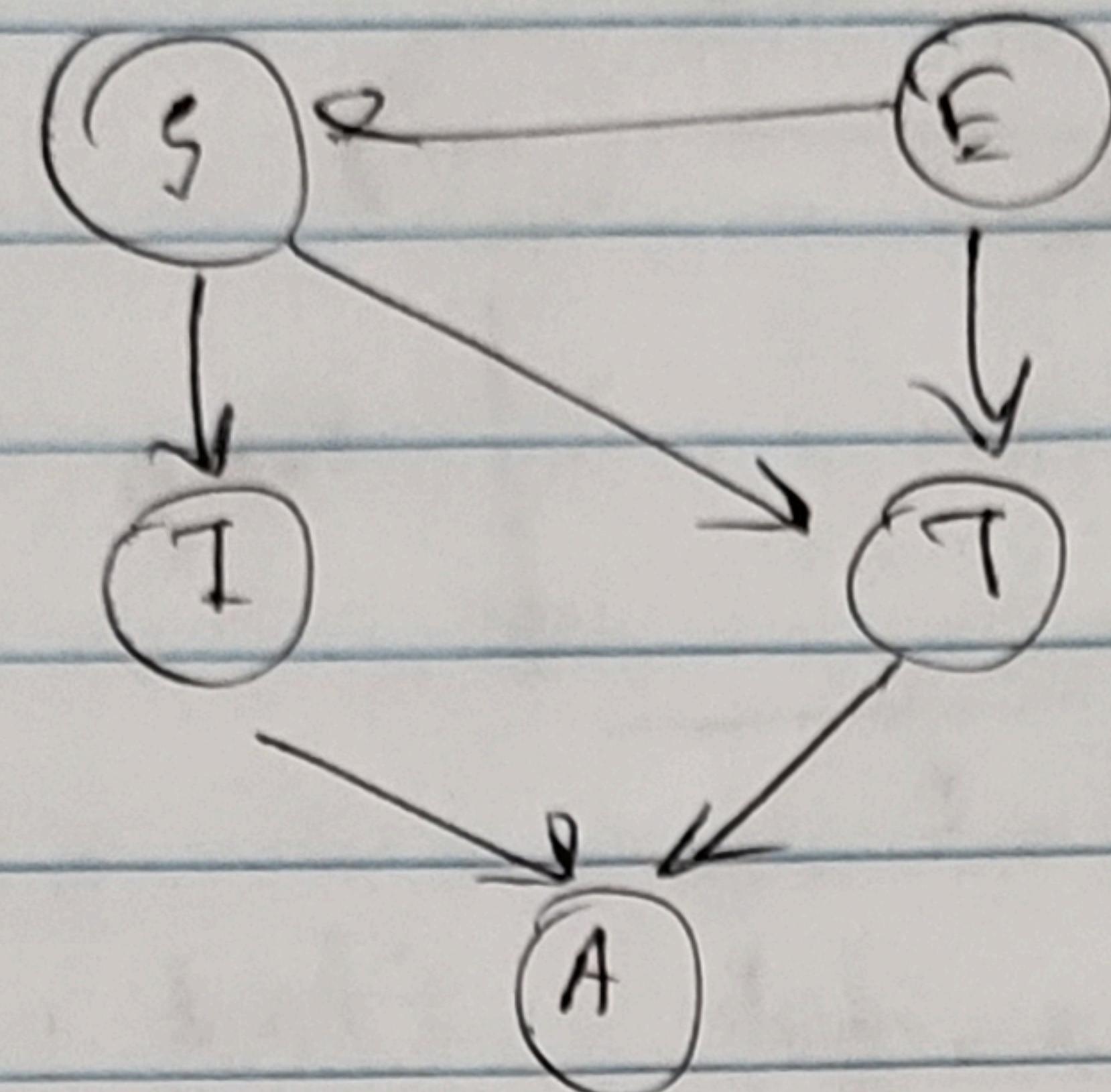


## Bayesian Network



Note: All binary variables

We want to compute  $P(+A, +T)$ .

So, the free variables are E, S, I.  
For variable elimination, which  
of the following orders are  
easy to compute?

① E  $\triangleright$  S  $\triangleright$  I

$$P(+A, +T)$$

$$= \sum_E P(E) \sum_S P(S|E) P(+T|S, E) \\ \sum_I P(I|S) P(+A|I, +T)$$

② I  $\triangleright$  E  $\triangleright$  S

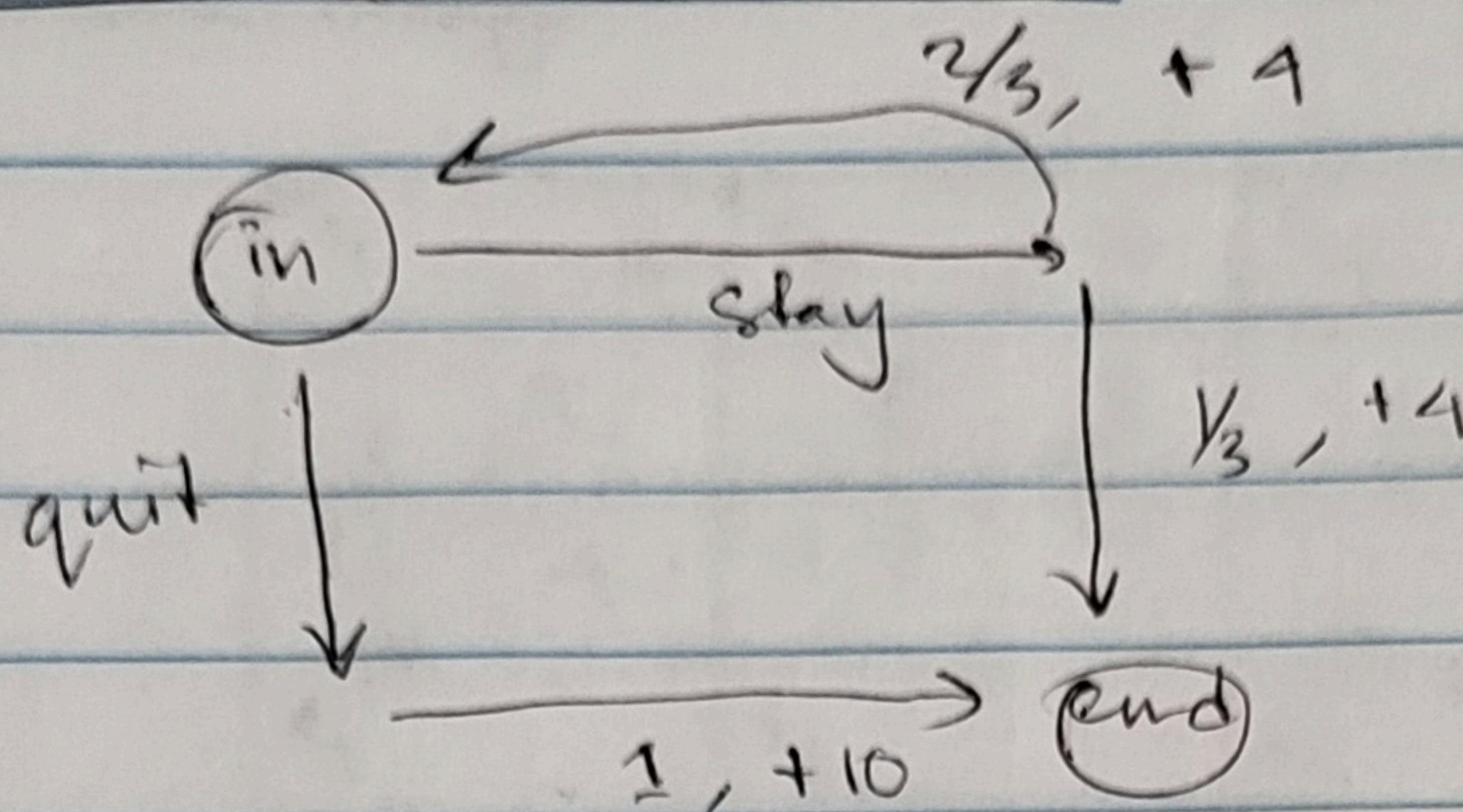
$$P(+A, +T)$$

$$= \sum_I P(+A|I, +T) \sum_E P(E) \\ \sum_S P(I|S) P(S|E) P(+T|S, E)$$

Out of these two, ① has  
three terms and more variables  
in the last part.

So, ① seems easier.

MDP (Dice Game)



Let's define this MDP in tabular form.

$s$	$a$	$s'$	$T(s,a,s')$	$R(s,a,s')$
in	stay	in	2/3	+4
in	stay	end	1/3	+4
in	quit	end	1	+10

\* No rows for  $s=\text{end}$  (why?)

$v_{in} = v_{end} = 0, \gamma = 0.5$

let's do VALUE ITERATION

Iteration 1

$$Q^1(\text{in}, \text{stay}) = \frac{1}{3} [+4 + \gamma \cdot 0] + \frac{2}{3} [+4 + \gamma \cdot 5] = +9$$

$$Q^1(\text{in}, \text{end}) = \cancel{1} \cdot 1 \cdot [-10 + \gamma \cdot 0] = +10$$

$$v^2_{in} = \max_a Q^1(\text{in}, a) = +\underline{\underline{10}}$$

$$\pi^2(\text{in}) = \text{avg} \max_a Q^1(\text{in}, a) = \underline{\underline{\text{end}}}$$

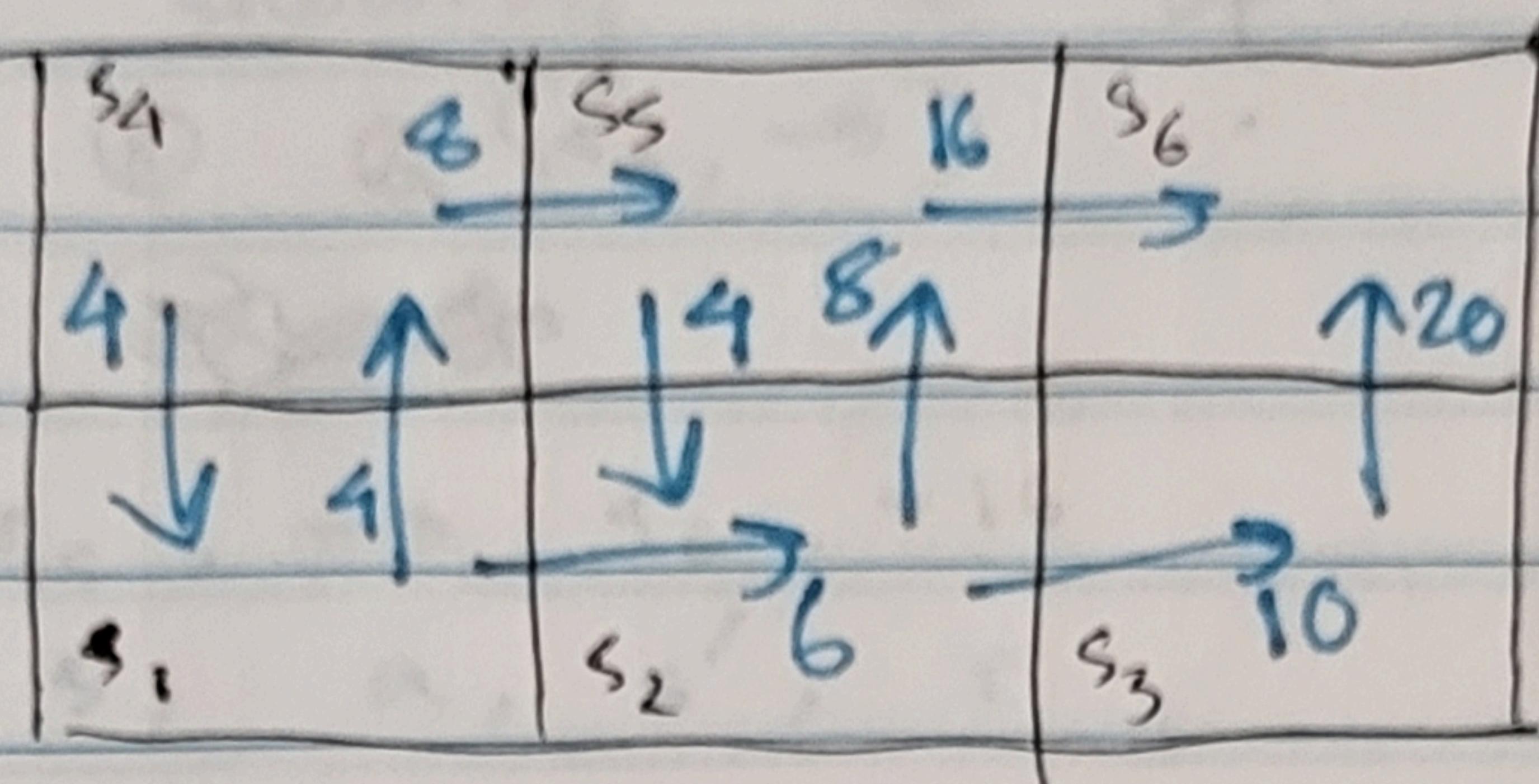
Iteration 2

$$Q^2(\text{in}, \text{stay}) = \frac{1}{3} [+9 + 0.5 \times 0] + \frac{2}{3} [+9 + 0.5 \times +9] = 16/3$$

$$Q^2(\text{in}, \text{end}) = 1 \cdot (+10 + 0) = +10$$

$$v^3_{in} = 10, \pi^3_{in} = \underline{\underline{\text{end}}}$$

## Q-learning



Arrows indicate Q-values (after a few iterations)  
Let's list all of them.

$$Q(s_1, \uparrow) = 4, Q(s_1, \rightarrow) = 6$$

$$Q(s_2, \uparrow) = 8, Q(s_2, \rightarrow) = 10$$

$$Q(s_3, \uparrow) = 20$$

$$Q(s_4, \rightarrow) = 8, Q(s_4, \downarrow) = 4$$

$$Q(s_5, \rightarrow) = 16, Q(s_5, \downarrow) = 4$$

Now, assume the following episode.

$$(r=0) \quad (r=0) \quad (r=+16)$$

$$s_1, \uparrow, s_4, \rightarrow, s_5, \rightarrow, s_6$$

Update Q-values. ( $\gamma = 0.3, \alpha = 0.8$ )

$$\textcircled{1} \quad (s_1, a, s', r) \\ (s_1, \uparrow, s_4, \rightarrow)$$

$$\text{sample} = r + \gamma \max_{a'} Q(s', a')$$

$$= 0 + 0.3 \times 8 = 2.4$$

$$Q(s_1, \uparrow) = (1-\alpha) Q(s_1, \uparrow) + \alpha \cdot \text{sample} \\ = 0.2 \times 4 + 0.8 \times 2.4 \\ = 2.72$$

Not showing the updates for

$$\textcircled{11} \quad Q(s_4, \rightarrow)$$

~~sample~~

$$\textcircled{11} \quad s_5, \rightarrow, s_6, +16$$

$$s, a, s', r$$

$$\begin{aligned}\text{sample} &= r + \gamma \max_{a'} Q(s', a') \\ &= +16 + 0.3 \times \underline{6} \\ &= +16\end{aligned}$$

$$\begin{aligned}Q(s_5, \rightarrow) &= 0.2 \times 16 + 0.8 \times 16 \\ &= 16\end{aligned}$$