

a) optimal search strategy: visit each grid exactly once.

	a	b	c	d	e
1	24	9	10	11	12
2	23	8	1	2	13
3	22	7	*	3	14
4	21	6	5	4	15
5	20	19	18	17	16

Empty

Suppose the prize appears in each cell with the same possibility, so $p = \frac{1}{24}$

If the prize is at the i th cell, then the discounted reward will be $\gamma^i \cdot 1 = \gamma^i$

Thus, the average discounted award will be:

$$\begin{aligned}
 \sum_{i=1}^{24} p \cdot \gamma^i &= \frac{1}{24} \sum_{i=1}^{24} \gamma^i \\
 &= \frac{1}{24} \left(\frac{1 - \gamma^{25}}{1 - \gamma} - \gamma^0 \right) \\
 &= \frac{1}{24} \left(\frac{\gamma - \gamma^{25}}{1 - \gamma} \right) \\
 &= \frac{1}{24} \left(\frac{0.95 - 0.95^{25}}{1 - 0.95} \right) \\
 &\approx 0.5605
 \end{aligned}$$

b)

	a	b	c	d	e
1				3	4
2			1	2	5
3			*		
4					
5					

Four room

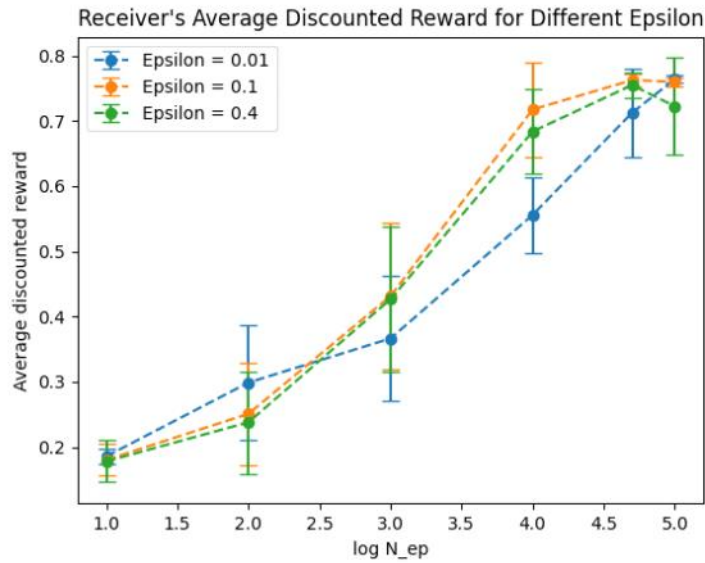
WLOG, suppose the agent choose the top right room. The prize still appears in the empty cells uniformly randomly: $p = \frac{1}{18}$

Similarly: the average reward will be:

$$\begin{aligned}
 &p (\gamma + \gamma^2 + \gamma^3 + \gamma^4 + \gamma^5) \\
 &= \frac{1}{18} \left(\sum_{i=1}^5 \gamma^i \right) \\
 &= \frac{1}{18} \left(\frac{\gamma - \gamma^6}{1 - \gamma} \right) = \frac{1}{18} \left(\frac{0.95 - 0.95^6}{1 - 0.95} \right) \approx 0.2388
 \end{aligned}$$

The agent can randomly choose any room, but the average discounted reward will be the same: 0.2388 (We can also take the average for four rooms, but they have identical results so give the same answer)

C) All log in the following questions are base 10



Messages:

```
[[0 0 0 3 3]
 [0 0 3 3 3]
 [0 0 0 0 0]
 [1 1 1 2 2]
 [1 1 0 2 2]]
```

So:

Msg 0 means the prize is at the top left room

Msg 1 means the prize is at the bottom left room

Msg 2 means the prize is at the bottom right room

Msg 3 means the prize is at the top right room

Actions for message 0					
	a	b	c	d	e
0	↓	←		↓	↓
1	→	↑	←	←	←
2			↑		
3	→	→	↑	←	↓
4	↑	↑		↑	←
Actions for message 1					
	a	b	c	d	e
0	↓	↓		←	↓
1	→	→	↓	←	←
2			↓		
3	→	↓	←	←	←
4	↑	←		↑	↑
Actions for message 2					
	a	b	c	d	e
0	←	↓		↓	←
1	→	→	↓	←	←
2			↓		
3	→	→	→	↓	←
4	→	↑		→	↑
Actions for message 3					
	a	b	c	d	e
0	↓	↓		↓	←
1	→	→	→	→	↑
2			↑		
3	→	→	↑	←	←
4	→	↑		↑	↑

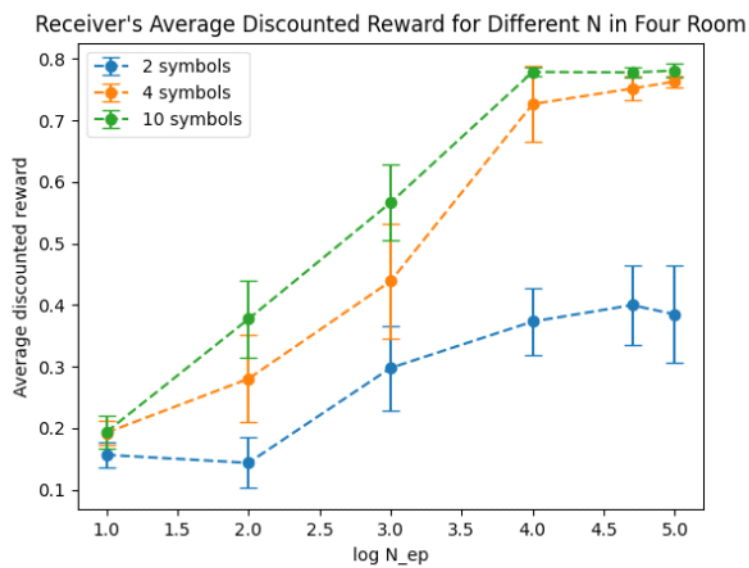
Receiver's policy for msg 0, top left room

Receiver's policy for msg 1, bottom left room

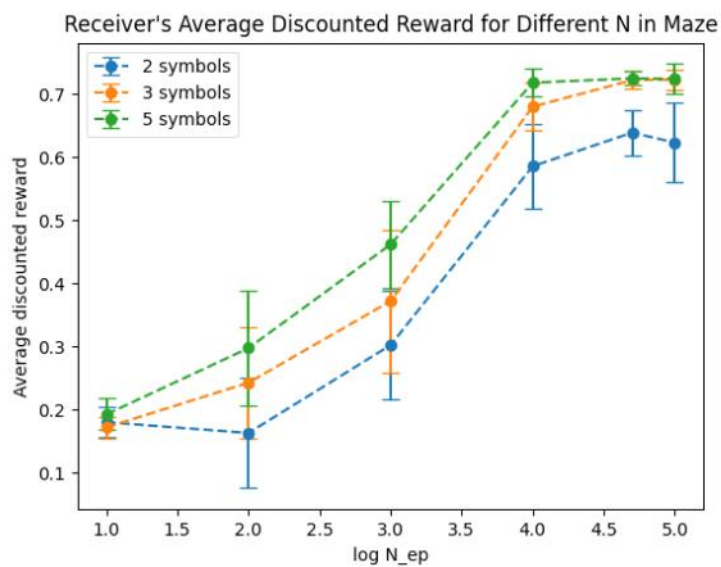
Receiver's policy for msg 2, bottom right room

Receiver's policy for msg 3, top right room

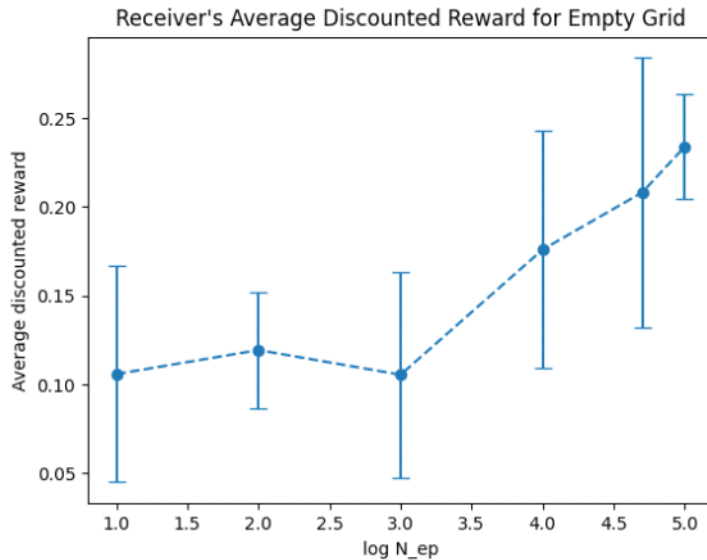
D)



E)



F)



G)

Assuming that the average discounted reward is approximately linear to $\log N_{ep}$ according to the graph we found in f)

By interpolating:

$$\begin{aligned} \text{reward} &\approx 0.106 & \text{at } N=10 & \log N_{ep} = 1 \\ \text{reward} &= 0.234 & \text{at } N=100000 & \log N_{ep} = 5 \end{aligned}$$

$$\text{Slope } k = \frac{0.234 - 0.106}{5 - 1} \approx 0.032$$

from a): the optimal reward we found is 0.5605

$$\begin{aligned} \text{so: } \log N_{opt} &= 5 + \frac{0.5605 - 0.234}{0.032} \approx 15.203 \\ \Rightarrow N_{opt} &\approx 10^{15.203} \approx 1.6 \times 10^{15} \end{aligned}$$

Although from the graph in f), it looks like the reward grows slowly for small N_{ep} , and faster for large N_{ep} , which makes it like linear from $\log N_{ep} = 3$ to $\log N_{ep} = 5$. If we only use that part for interpolation,

$$k = \frac{0.234 - 0.106}{5 - 3} \approx 0.063, \log N_{opt} = \left(3 + \frac{0.5605 - 0.234}{0.063} \right) \approx 10.18, N_{opt} \approx 10^{10}$$

From the result for c, d, e, the curve all goes like a "sigmoid" shape, meaning growing slow at the beginning and the end. So this might be an overestimate.

So I probably need to run some iterations between 10^{10} and 10^{15} to get close enough to the optimal