

2 infcr - p'p'15'n NCR > NF

royleibovici@mail.tau.ac.il ; 206994840 ; 'Tzrif' 171
goltsman@mail.tau.ac.il ; 314645920 ; jn3fcr wde

Question 1

- infcr n73r D'xan uk v12j .1

$$S = \{0, \dots, 2k-1\}$$

$$A = \{CN, CCW\}$$

$$p^{(CCW)}_{ss'} = \begin{cases} 1, & s' = s + 1 \bmod(2k) \\ 0, & s' = s - 1 \bmod(2k) \end{cases}$$

$$p^{(CW)}_{ss'} = \begin{cases} 1, & s' = s - 1 \bmod(2k) \\ 0, & s' = s + 1 \bmod(2k) \end{cases}$$

$$R = \begin{cases} 1, & s = 0 \\ 0, & \text{else} \end{cases}$$

$$s_0 = k$$

after R 3rd (k') π^* random walk π^{3NF} .2
 next $s=0$ 3NF π^{3NF} 2x CCW 1x CW
 , CW SCL CCW 1x CCW SCL CW 3NF π^{3NF} |
 1NF π^{1NF} $s=0$ 3NF π^{3NF} | 2NF π^{2NF} k'3
 . | P 1NF 1x

$$\pi^*(s=s') = \begin{cases} CW, \text{ w.p. } 0.5 \\ CCW, \text{ w.p. } 0.5 \end{cases}, \quad s' = \{s_0, 0\}$$

$$\pi^*(s=s') = CCW, \quad s' = \{s_0 + 1 \bmod(2K), \dots, 2K-1\}$$

$$\pi^*(s=s') = CW, \quad s' = \{1, \dots, s_0 - 1 \bmod(2K)\}$$

$$-\text{gap}, \quad s = \{1, 0\} \quad \text{gap} .3$$

$$V_1(s) = \max_{a \in A} \left\{ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_n(s') \right\}$$

$$V_1(s) = 0 \quad -\text{gap}, \quad s = 0 \quad \text{gap}$$

$$V_1(s=0) = \max_{a \in A} \left\{ r(0, a) + \gamma \sum_{s' \in S} p(s'|0, a) V_n(s') \right\}$$

$$V_1(s=0) = 1$$

$$-\text{dp} \quad S = S \setminus \{0, 1, 2k-1\} \text{ n/a } .4$$

$$V_2(s) = \max_{a \in A} \left\{ \underbrace{r(s, a)}_{=0} + \gamma \sum_{s' \in S} p(s'|s, a) V_n(s') \right\}$$

$$V_2(s) = 0$$

$$-\text{dp} \quad S = \{0, 1, 2k-1\} \text{ n/a }$$

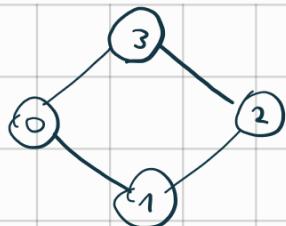
$$V_2(s=0) = \max_{a \in A} \left\{ \underbrace{r(0, a)}_{=1} + \gamma \sum_{s' \in S} p(s'|0, a) V_n(s') \right\} = 1$$

$$V_2(s=1) = \max_{a \in A} \left\{ \underbrace{r(1, a)}_{=0} + \gamma \sum_{s' \in S} p(s'|1, a) V_n(s') \right\} = r$$

$$= \begin{cases} 1, & s'=0 \\ 0, & \text{o.w.} \end{cases}$$

$$V_2(s=2k-1) = \max_{a \in A} \left\{ \underbrace{r(2k-1, a)}_{=0} + \gamma \sum_{s' \in S} p(s'|2k-1, a) V_n(s') \right\} = r$$

$$= \begin{cases} 1, & s'=0 \\ 0, & \text{o.w.} \end{cases}$$



למונטג'ו סט נציגו ערך הולך ועוזר בפתרון

$$\begin{array}{l|l|l} V_2(0) = 1 & V_3(0) = 1 + \gamma^2 & V_4(0) = 1 + \gamma^2 \\ V_2(1) = \gamma & V_3(1) = \gamma & V_4(1) = \gamma(1 + \gamma^2) \\ V_2(2) = 0 & V_3(2) = \gamma^2 & V_4(2) = \gamma^2 \\ V_2(3) = \gamma & V_3(3) = \gamma & V_4(3) = \gamma(1 + \gamma^2) \end{array}$$

$$V_h(0) = 1 + \gamma \cdot V_{h-1}(1)$$

$$V_h(1) = V_h(3) = \gamma \cdot V_{h-1}(0)$$

$$V_h(2) = \gamma \cdot V_{h-1}(1)$$

$$V_h(0) = 1 + \gamma V_h(1)$$

$$V_h(3) = V_h(1) = \gamma V_h(0) = \gamma(1 + \gamma V_h(1))$$

$$V_h(2) = \gamma V_h(1)$$

$$\gamma + \gamma^2 V_h(1) = V_h(1)$$

$$V_h(1) [\gamma^2 - 1] = -\gamma$$

$$V_h(1) = \frac{\gamma}{1 - \gamma^2} = V_h(3)$$

$$V_h(2) = \frac{\gamma^2}{1 - \gamma^2}$$

$$V_h(0) = 1 - \frac{\gamma^2}{1 - \gamma^2} = \frac{1}{1 - \gamma^2}$$

Question 2

1

$$S = \{ \vec{s} = (d_{N+1}, d_N, \dots, d_0)^T \mid 0 \leq i \leq N, d_i \in \{0, 1\}, d_{N+1} \in \{0, \dots, g\} \}$$

לע'ז סדרה s מוגדרת מינימלית כՅ'ן d_{N+1} סטטוס s בז'ן s
 מוגדרת כסדרה סטטוס s בז'ן (s_{N+1}, \dots, s_0)

$$A = \{0, \dots, N\}$$

$$P_{ss'}^{(a)} = \begin{cases} \frac{1}{10}, & d_j' = \begin{cases} d_j, & \forall 0 \leq j \leq N, j \neq a \\ 1, & j = a \\ k, & j = N+1 \end{cases} \\ 0, & \text{else} \end{cases}, \quad a \in A, \quad k \in [0, \dots, g]$$

$$r(s, a) = \frac{d_{N+1} \cdot 10^a}{g \cdot 10^N} = \frac{d_{N+1}}{g} \cdot 10^{a-N}, \quad a \in A$$

הערך המוצע של $r(s, a)$ הוא $\frac{1}{10} \sum_{k=0}^g \frac{d_{N+1}}{g} \cdot 10^{a-k}$

$$s_0 = (0, \dots, 0, k)^T, \quad k \sim \text{Uniform}[0, g]$$

בכדי לסייע לנו בפתרון נזכיר את $V(s)$ ו- $\pi(s)$

$$V(s) = \max_{a \in A} \left\{ r(s, a) + \sum_{s' \in S_{t+1}} p(s'|s, a) \cdot V(s') \right\}$$

לכל פעולה a , (s, a) יופיע בכל חישוב nostro מועד.

אנו אנו יוזם, גנום רצוי - מילוי צלול.

אנו יזק צלול, כדי שוכן השם, ישב בזאת.

אנו אנו רצויים - C.O.O. (Constant)

$$S = (1, 1, 1, \dots, 0, \dots, 1, 1, 0, \dots, 1, 1, d_{n+1})$$

i אמצעי j אמצעי

$a = \{i, j\}$ - יתאפשרו העזים שפער $V(s)$ יהיה זרוף

$$V(s) = \max \left\{ d_{n+1} \cdot 10^i + \sum_{k=0}^9 \frac{1}{10} \cdot V(1, 1, 1, \dots, 1, 1, k), d_{n+1} \cdot 10^j + \sum_{k=0}^9 \frac{1}{10} \cdot V(1, 1, 1, \dots, 1, 1, k) \right\}$$

$$V(1, 1, 1, \dots, 1, 1, d_{n+1}) = \frac{k \cdot 10^j}{9 \cdot 10^n} = \frac{k \cdot b^j}{t}$$

j אמצעי k אמצעי t אמצעי

$$V(s) = \max \left\{ d_{n+1} \cdot \frac{10^i}{t} + \sum_{k=0}^9 \frac{1}{10} \cdot k \cdot \frac{10^j}{t}, d_{n+1} \cdot \frac{10^j}{t} + \sum_{k=0}^9 \frac{1}{10} \cdot k \cdot \frac{10^i}{t} \right\}$$

$$V(s) = \max \left\{ d_{n+1} \cdot \frac{b^i}{t} + \frac{10^j}{t} \cdot 4.5, d_{n+1} \cdot \frac{b^j}{t} + \frac{10^i}{t} \cdot 4.5 \right\}$$

Decision regions

$$d_{n+1} \cdot \frac{10^i}{t} + 4.5 \cdot \frac{10^j}{t} = d_{n+1} \cdot \frac{10^j}{t} + 4.5 \cdot \frac{10^i}{t}$$

$$\cancel{d_{n+1} [10^i - 10^j]} = \cancel{4.5 [10^i - 10^j]}$$

$d_{n+1} = 4.5$

For $i, j > 1$ if d_{n+1} is threshold \rightarrow it is red

Red cell

Non-red cell

$$q = \begin{cases} \max\{i, j\} & d_{n+1} \geq 5 \\ \min\{i, j\} & d_{n+1} < 5 \end{cases}$$

Non-red cell \rightarrow green

Green cell

Red cell \rightarrow red

$Q' = \{q_0, q_1, q_2, \dots, q_n\} \cup \{q_{n+1}\}$ \rightarrow non-red cell

$$V(s) = \max_{q \in Q'} \left\{ r(s_q) + \sum_{s' \in S_{t+1}} p(s'|s_q) \cdot V(s') \right\} =$$

$$\max_{\alpha \in \{\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_n\}} \left\{ \max_{s' \in S_{t+1}} \left\{ r(s, \alpha) + \sum_{s' \in S_{t+1}} P(s'|s, \alpha) \cdot V(s') \right\} \right\},$$

מינימיזציה של הערך הנוכחי
 על מנת לארוך את זמן החיים

$$\max \{ \} = \max \{ \max \{ A \}, \max \{ B \} \} \text{ סמ"כ } A \cup B = C$$

וקטור החלטה $\pi(s)$ מושג על ידי $V(s)$, כלומר $\pi(s) = \frac{1}{C} \sum_{a=0}^{C-1} \pi_a(s)$
 כלומר $\pi_a(s) = \frac{1}{C} \sum_{s' \in S_{t+1}} P(s'|s, a) \cdot V(s')$.

• Initialize: $V(s) = \max_{a \in A} r(s, a)$ for all $s \in S_T$.
 • For $t = T-1, \dots, 1$
 • For all $s \in S_t$, set
 • $V(s) = \max_{a \in A} r(s, a) + \sum_{s' \in S_{t+1}} P(s' | s, a)V(s')$.
 • $\pi(s) = \arg \max_{a \in A} r(s, a) + \sum_{s' \in S_{t+1}} P(s' | s, a)V(s')$.
 • Return π .

Deterministic Markov policy

$$V_3((d_3, 1, 1, 1), a) = 0 \quad \text{if } a \in \{0, \dots, 9\} \leftarrow \begin{array}{l} \text{for } p(s) \text{ in } S \\ \text{if } p(s) \neq 0 \text{ then } \\ \quad \cdot 10^2 \text{ for } k \\ \quad \cdot 10^1 \text{ for } k \end{array}$$

$$V_2((d_3, 0, 1, 1), 2) = \frac{d_3}{g} \cdot 10^{-2} = \frac{d_3}{g} \Rightarrow \pi_2^*((d_3, 0, 1, 1)) = 2$$

$$V_2((d_3, 1, 0, 1), 1) = \frac{d_3}{g} \cdot 10^{-1} \Rightarrow \pi_2^*((d_3, 1, 0, 1)) = 1$$

$$V_2((d_3, 1, 1, 0), 0) = \frac{d_3}{g} \cdot 10^{-2} \Rightarrow \pi_2^*((d_3, 1, 1, 0)) = 0$$

• π_1^* for $(d_3, 0, 0, 1)$, $t=1$

$$\pi_1^*((d_3, 0, 0, 1)) = \arg \max_{a \in [0, 9]} \left[r((d_3, 0, 0, 1), a) + \sum_{s' \in S_{t+1}} p((s', 1, 0, 1) | (d_3, 0, 0, 1), a) \cdot V_2(s', 1) \right]$$

$$r((d_3, 0, 0, 1), a) + \sum_{s' \in S_{t+1}} p((s', 1, 0, 1) | (d_3, 0, 0, 1), a) \cdot V_2(s', 1) =$$

$$= \arg \max_{a \in [0, 9]} \left[\frac{d_3}{g} + \sum_{k=0}^9 \frac{1}{10} \cdot \frac{K}{g} \cdot 10^{-1}, \frac{d_3}{g} \cdot 10^{-1} + \sum_{k=0}^9 \frac{1}{10} \cdot \frac{k}{g} \right] =$$

$$= \arg \max_{a \in \{0,1\}} \left[\frac{d_3}{g} + 0.05, \quad \frac{d_3}{90} + 0.5 \right] \Rightarrow \begin{cases} \frac{d_3}{g} + 0.05 & \geq \\ & \leq \\ & \leq \end{cases} \frac{d_3}{90} + 0.5 \Rightarrow$$

$$\Rightarrow \frac{a_3}{10} > 0.45 \Rightarrow a_3 > 4.5$$

$$\mathcal{P}_1^*(d_3, 1, 0, 0) = \arg \max_{a \in \{1, 0\}} \left[r(d_3, 1, 0, 0, 1) + \sum p(k, 1, 0) | (d_3, 1, 0, 0, 1) V_2(k, 1, 1, 0) - r(d_3, 1, 0, 0, 0) + \right. \\ \left. + \sum p(k, 1, 0, 1) | (d_3, 1, 0, 0, 0) V_2(k, 1, 0, 1) \right] =$$

$$= \underset{a \in [1, \infty)}{\operatorname{arg\,max}} \left[\frac{d_3}{g} \cdot 10^{-1} + \sum_{k=0}^y \frac{1}{10} \cdot \frac{k}{g} \cdot 10^{-2}, \quad \frac{d_3}{g} \cdot 10^{-2} + \sum_{k=0}^y \frac{1}{10} \cdot \frac{k}{g} \cdot 10^{-1} \right] =$$

$$= \arg \max_{a \in [1, 0]} \left[\frac{d_3}{90} + \frac{1}{200}, \frac{d_3}{900} + \frac{1}{20} \right] \Rightarrow \frac{d_3}{90} + \frac{1}{200} > \frac{1}{900} + \frac{1}{20} \Rightarrow$$

$$\Rightarrow \frac{d_3}{100} > \frac{9}{200} \Rightarrow d_3 < 4.5 \Rightarrow$$

$$\mathcal{T}_1^*(d_3, 0, 1, 0) = \underset{a \in \{0, 1\}}{\operatorname{argmax}} \left[r(d_3, 0, 1, 0, 1) + \sum p(k, 1, 1, 0) | (d_3, 0, 1, 0, 1) V_k(k, 1, 1, 0) - r(d_3, 0, 1, 0, 0) \right]$$

$$+ \sum p((k, o_{1,1}) | (d_{2,1}, o_{1,1}, o)) V_2(k, o_{1,1}) \Big] =$$

$$= \underset{a \in [2, 6]}{\operatorname{argmax}} \left[\frac{d_3}{g} + \sum_{k=0}^g \frac{1}{10} \cdot \frac{k}{g} \cdot 10^{-2}, \frac{d_3}{g} \cdot 10^{-2} + \sum_{k=0}^g \frac{1}{10} \cdot \frac{k}{g} \right] = \underset{a \in [2, 6]}{\operatorname{argmax}} \left[\frac{d_3}{g} + \frac{1}{200}, \frac{d_3}{g} + \frac{1}{2} \right] \Rightarrow$$

$$\Rightarrow \frac{d_3}{g} + \frac{1}{200} \geq \frac{d_3}{900} + \frac{1}{2} \Rightarrow \frac{11}{100} \cdot d_3 \geq \frac{99}{200} \Rightarrow d_3 \geq 4.5$$

- 0'70

$$\pi_1^*(d_3, 0, 0, 1) = \begin{cases} 1, & 0 \leq d_3 \leq 4 \\ 2, & 5 \leq d_3 \leq 9 \end{cases} \quad V_1(d_3, 0, 0, 1) = \begin{cases} \frac{d_3}{90} + \frac{1}{2}, & 0 \leq d_3 \leq 4 \\ \frac{d_3}{9} + \frac{1}{200}, & 5 \leq d_3 \leq 9 \end{cases}$$

$$\pi_1^*(d_3, 1, 0, 0) = \begin{cases} 0, & 0 \leq d_3 \leq 4 \\ 1, & 5 \leq d_3 \leq 9 \end{cases} \quad V_1(d_3, 1, 0, 0) = \begin{cases} \frac{d_3}{900} + \frac{1}{20}, & 0 \leq d_3 \leq 4 \\ \frac{d_3}{9} + \frac{1}{200}, & 5 \leq d_3 \leq 9 \end{cases}$$

$$\pi_1^*(d_3, 0, 1, 0) = \begin{cases} 0, & 0 \leq d_3 \leq 4 \\ 2, & 5 \leq d_3 \leq 9 \end{cases} \quad V_1(d_3, 0, 1, 0) = \begin{cases} \frac{d_3}{900} + \frac{1}{2}, & 0 \leq d_3 \leq 4 \\ \frac{d_3}{9} + \frac{1}{200}, & 5 \leq d_3 \leq 9 \end{cases}$$

. 10'10 10'61 0'6 , t=0 7'70

$$\pi_0^*(d_3, 0, 0, 0) = \arg \max_{a \in \{0, 1, 2\}} R((d_3, 0, 0, 0), a) + \sum p((k, 1, 0, 0) | (d_3, 0, 0, 0), a) V_1(k, 1, 0, 0),$$

$$R((d_3, 0, 0, 0), 1) + \sum p((k, 0, 1, 0) | (d_3, 0, 0, 0), 1) V_1(k, 0, 1, 0),$$

$$R((d_3, 0, 0, 0), 0) + \sum p((k, 0, 0, 1) | (d_3, 0, 0, 0), 0) V_1(k, 0, 0, 1) =$$

$$= \arg \max_{a \in \{0, 1, 2\}} \left[\frac{d_3}{900} + \frac{1}{10} \left(\sum_{k=0}^4 \left(\frac{k}{900} + \frac{1}{20} \right) + \sum_{k=5}^9 \left(\frac{k}{90} + \frac{1}{200} \right) \right), \right.$$

$$\frac{d_3}{90} + \frac{1}{10} \left(\sum_{k=0}^4 \left(\frac{k+1}{90} + \frac{1}{2} \right) + \sum_{k=5}^9 \left(\frac{k+1}{9} + \frac{1}{200} \right) \right),$$

$$\left. \frac{d_3}{900} + \frac{1}{10} \cdot \left(\sum_{k=0}^4 \left(\frac{d_3 + k}{90} + \frac{1}{2} \right) + \sum_{k=5}^9 \left(\frac{d_3 + k}{9} + \frac{1}{20} \right) \right) \right] =$$

$$= \underset{a \in [2, 1, 0]}{\operatorname{argmax}} \left[\frac{d_3}{9} + \frac{27}{400}, \frac{d_3}{90} + \frac{257}{400}, \frac{d_3}{900} + \frac{27}{40} \right]$$

$$\frac{d_3}{9} + \frac{27}{400} \stackrel{?}{\leq} \frac{d_3}{90} + \frac{257}{400}$$

$$d_3 \cdot \frac{1}{10} \stackrel{?}{\geq} \frac{23}{400}$$

$$d_3 \stackrel{?}{\geq} 5.75$$

$$\frac{d_3}{90} + \frac{257}{400} \stackrel{?}{\leq} \frac{d_3}{900} + \frac{27}{40}$$

$$d_3 \cdot \frac{1}{100} \stackrel{?}{\geq} \frac{13}{400}$$

$$d_3 \stackrel{?}{\geq} 3.25$$

$$\frac{d_3}{900} + \frac{27}{40} \stackrel{?}{\leq} \frac{d_3}{900} + \frac{27}{40}$$

$$\frac{d_3}{9} + \frac{27}{400} \stackrel{?}{\leq} \frac{d_3}{900} + \frac{27}{40}$$

$$d_3 \cdot \frac{11}{100} \stackrel{?}{\geq} \frac{243}{400}$$

$$d_3 \stackrel{?}{\geq} 5.52$$

- CPO

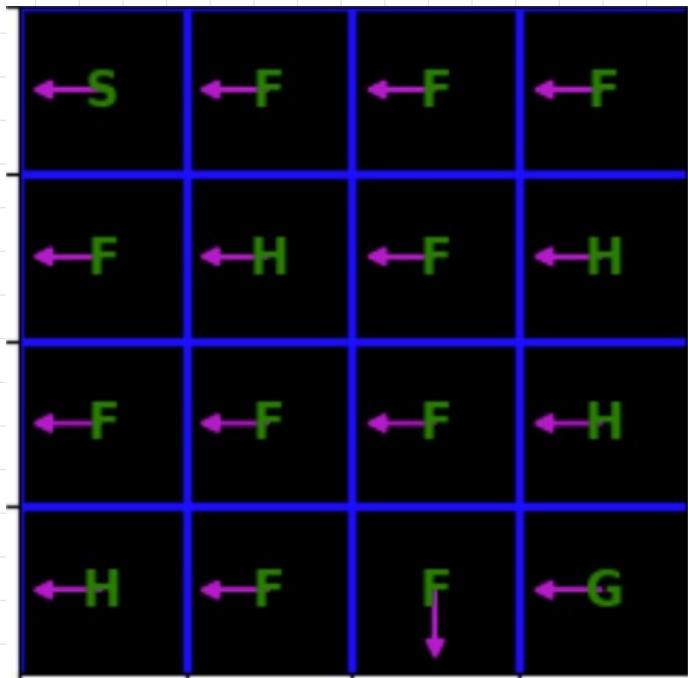
$$\pi_o^*(d_3, 0, 0, 0) = \begin{cases} 0, & 0 \leq d_3 \leq 3 \\ 1, & 4 \leq d_3 \leq 5 \\ 2, & 6 \leq d_3 \leq 9 \end{cases}$$

— QVN PFS

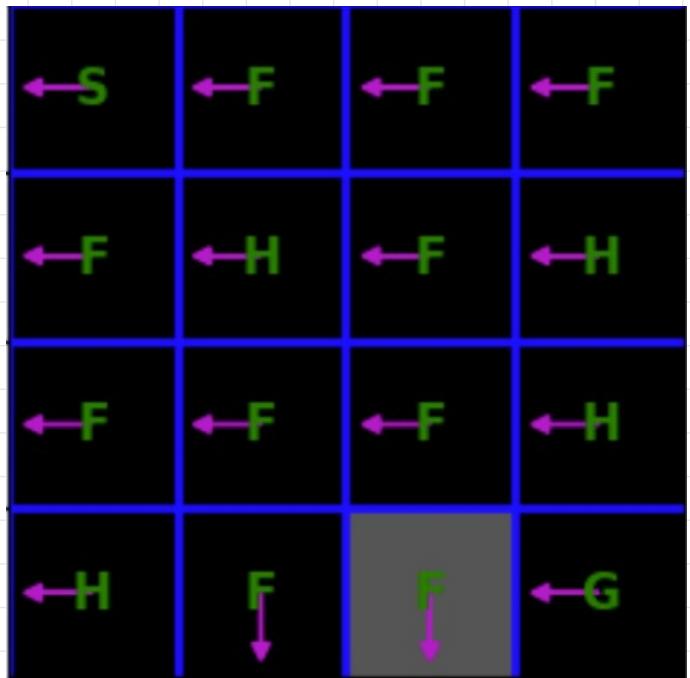
Question 1 - Value Iteration

of "Frozen Lake" and start with zero value for all states
• frozen trap & Bellman Value- \rightarrow to black

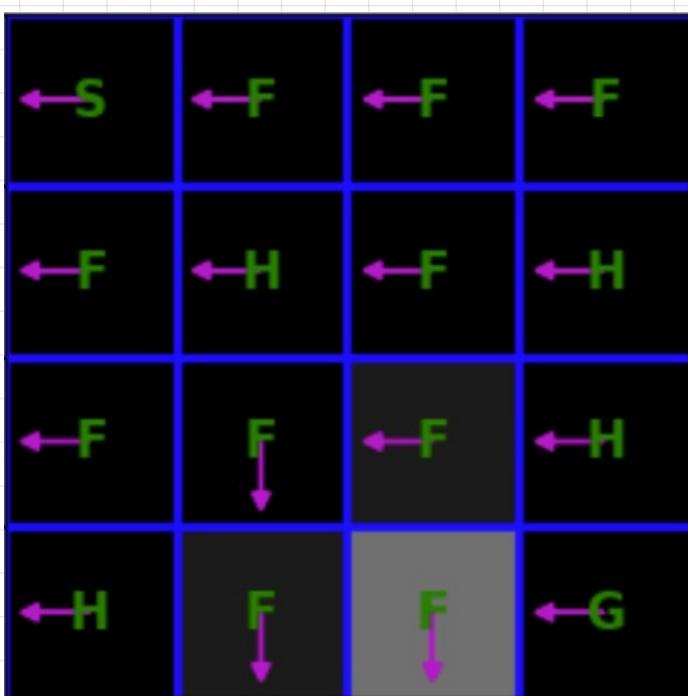
• initial state of frozen lake



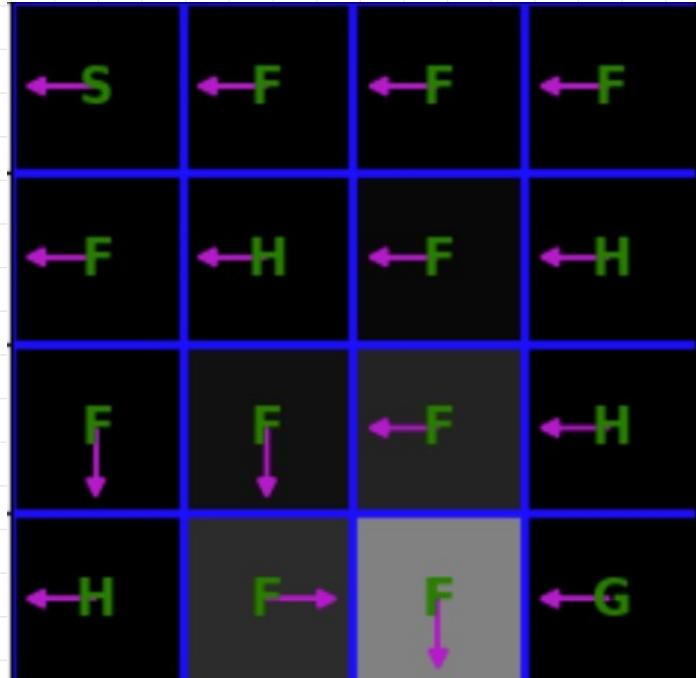
Iteration 1



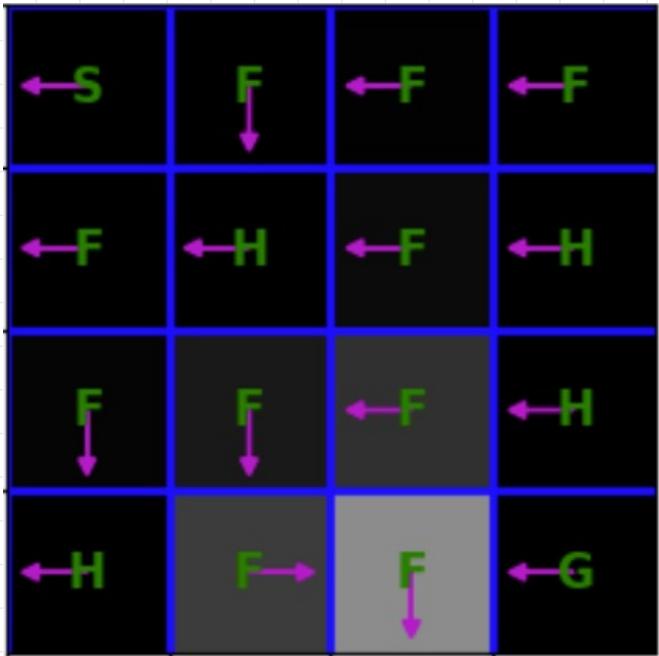
Iteration 2



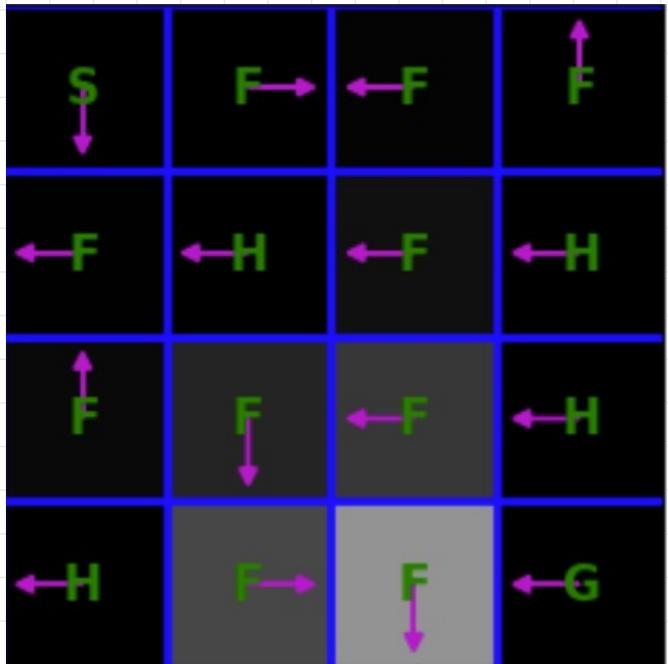
Iteration 3



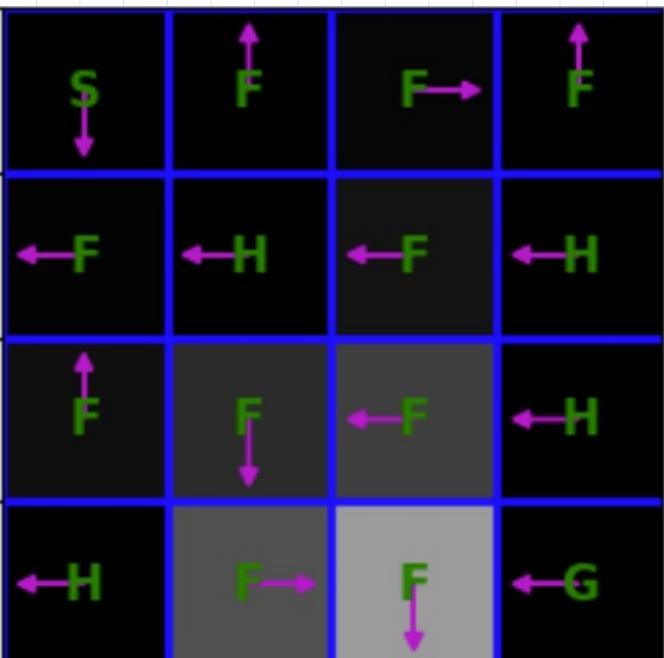
Iteration 4



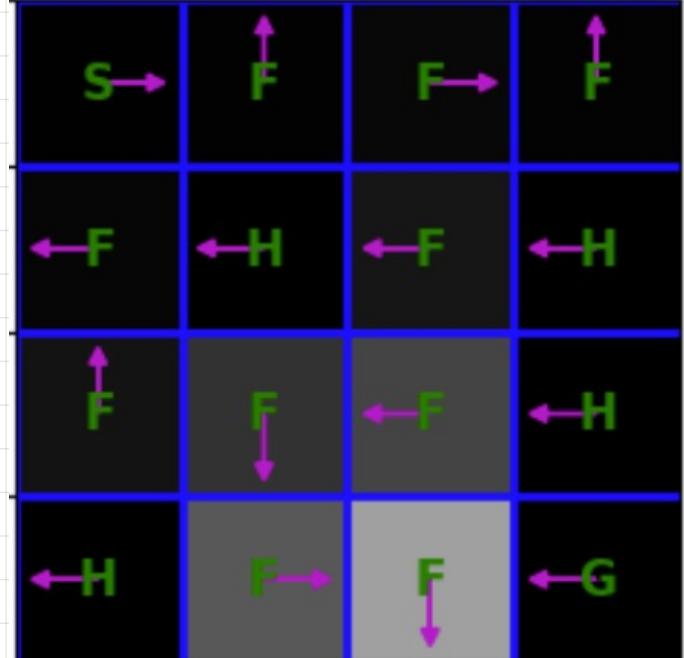
Iteration 5



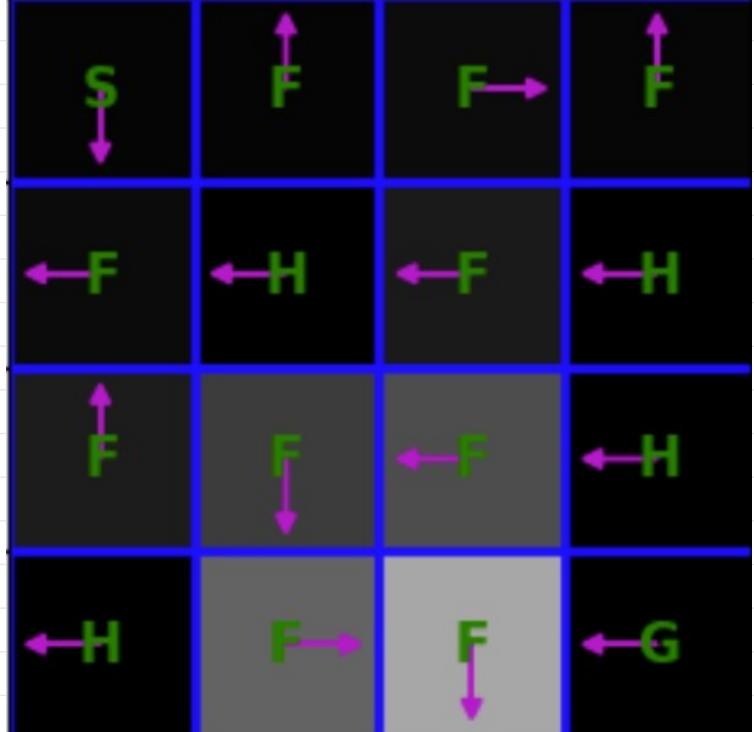
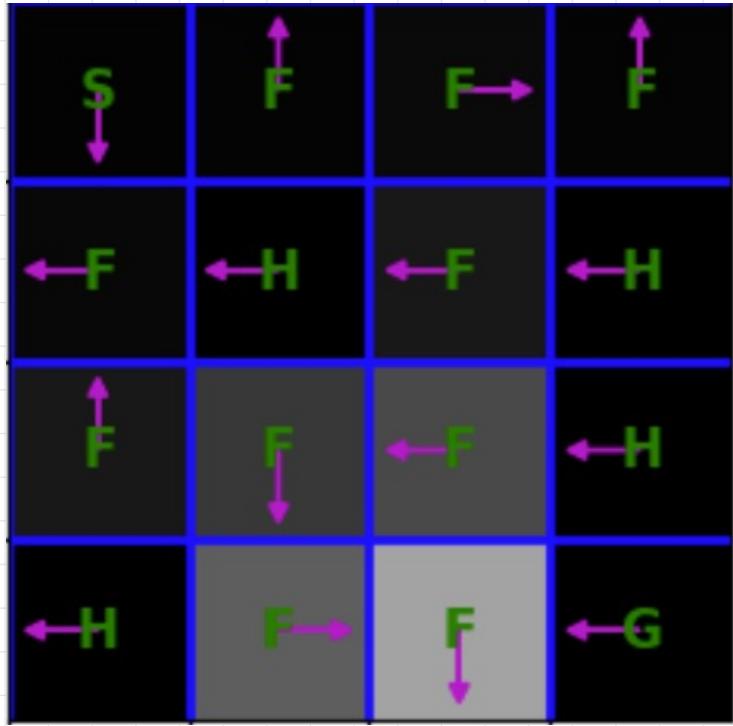
Iteration 6



Iteration 7



Iteration 8



Iteration 12

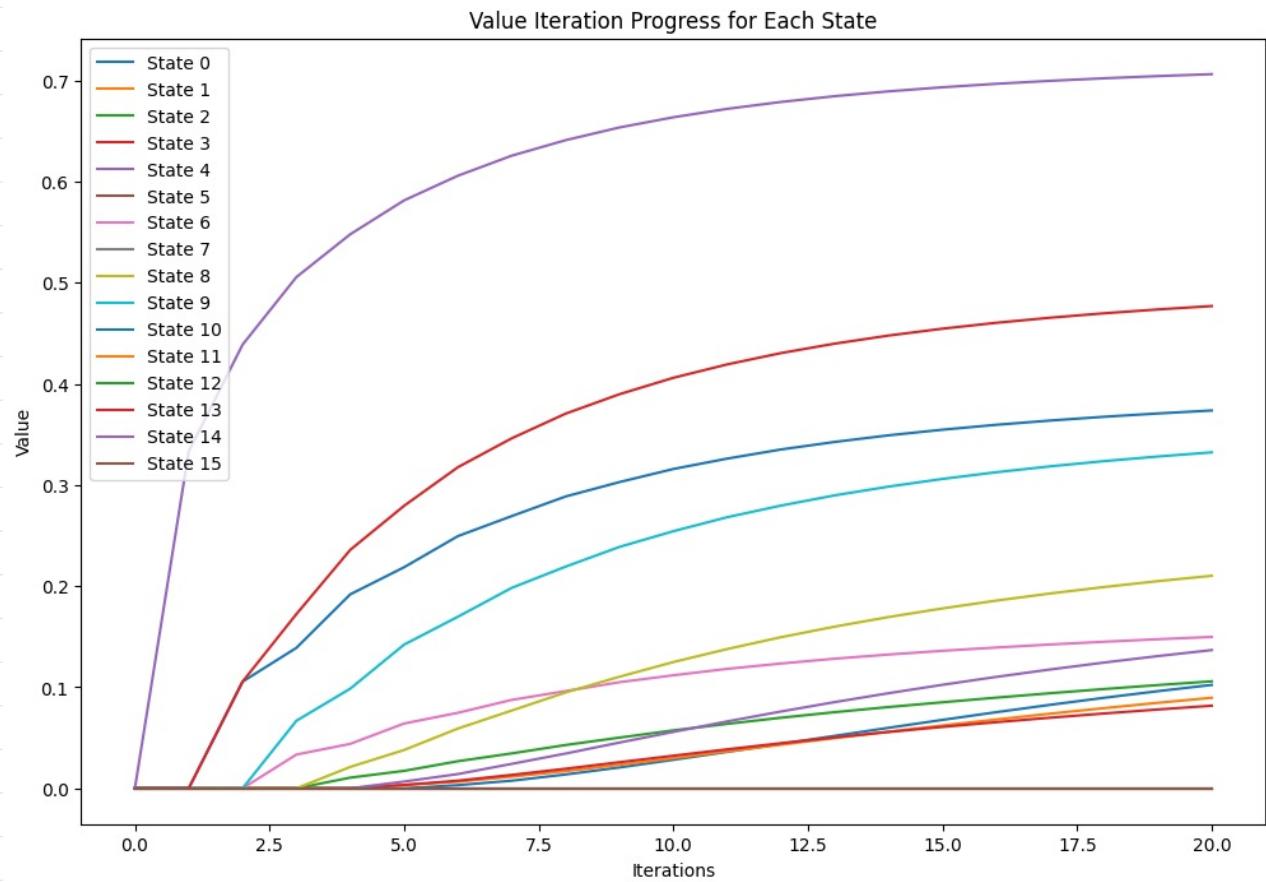
Iteration 14 ↑

mark selection without
no DNN

Iteration Tab

Iteration	max V-Vprev	# chg actions	V[θ]
0	0.33333	N/A	0.000
1	0.10556	1	0.000
2	0.06685	1	0.000
3	0.06351	2	0.000
4	0.04357	1	0.000
5	0.03821	4	0.003
6	0.02857	2	0.008
7	0.02437	1	0.014
8	0.01952	1	0.021
9	0.01624	0	0.028
10	0.01384	0	0.036
11	0.01173	0	0.044
12	0.01047	1	0.052
13	0.00948	0	0.060
14	0.00852	1	0.068
15	0.00782	0	0.075
16	0.00733	0	0.083
17	0.00694	0	0.090
18	0.00656	0	0.096
19	0.00618	0	0.102

→ 23% for 100%
23% for Bellman Value →
→ 70% for 37% CIG 70%

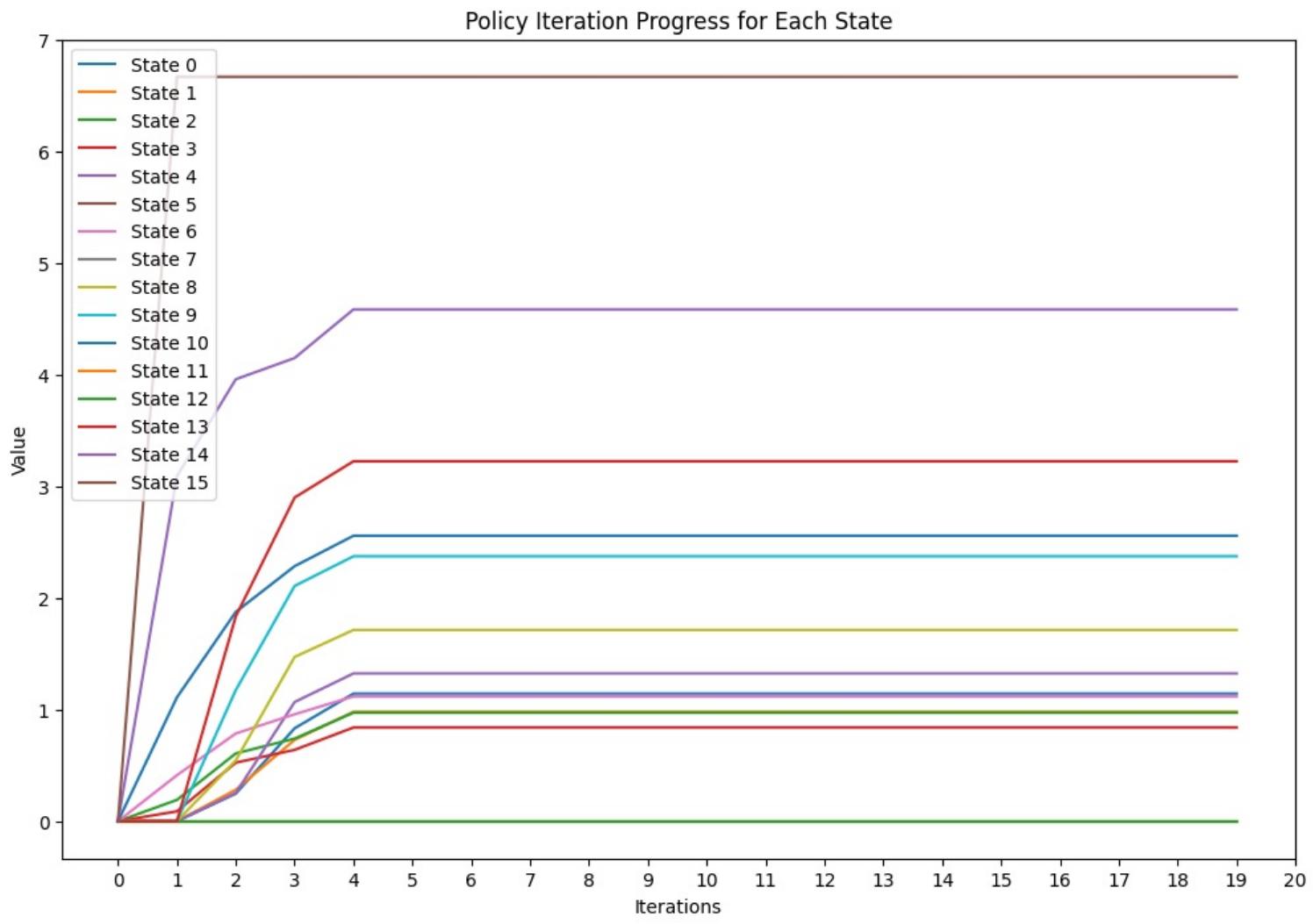


Question 2 — Policy Iteration

Iteration	# chg actions	v[0]
0	1	0.00000
1	9	0.00000
2	5	0.24722
3	3	0.83222
4	0	1.14299
5	0	1.14299
6	0	1.14299
7	0	1.14299
8	0	1.14299
9	0	1.14299
10	0	1.14299
11	0	1.14299
12	0	1.14299
13	0	1.14299
14	0	1.14299
15	0	1.14299
16	0	1.14299
17	0	1.14299
18	0	1.14299
19	0	1.14299

rek / yea 21 12/10 is ok for
 "red" "dark" "QND" "rewn"
 "frozen lake"
 "van" "spa" "RN D"
 "37(k)" "ac" "pop"

pan of & Bellman Value -> at 23'nd freq prob
- 11'2nd 6 27C's 17sra



Policy Iteration

⑥ Value Iteration