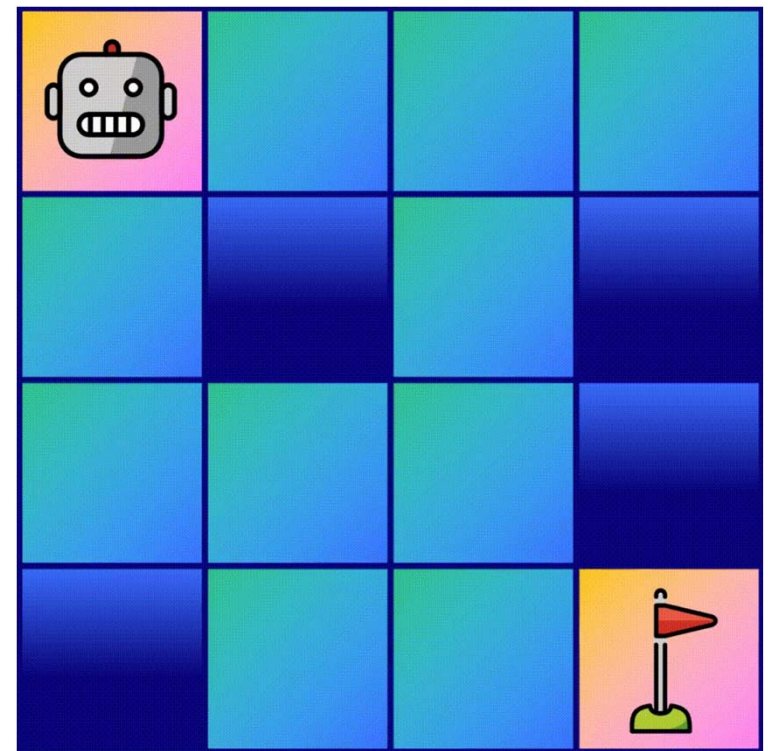


Spodbujevalno učenje – domača naloga

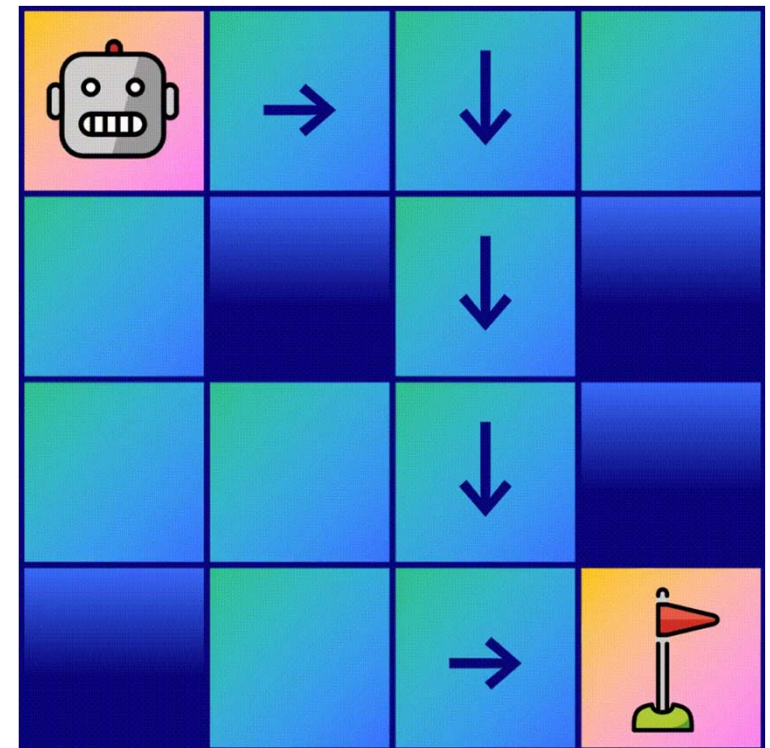
Frozen lake primer

- Začetno stanje - S
- Končno stanje - G
- Dovoljeno mesto - F
- Luknja – H
- Akcije
 - Levo
 - Dol
 - Desno
 - Gor



Frozen lake primer

- Začetno stanje - S
- Končno stanje - G
- Dovoljeno mesto - F
- Luknja – H
- Akcije
 - Levo
 - Dol
 - Desno
 - Gor

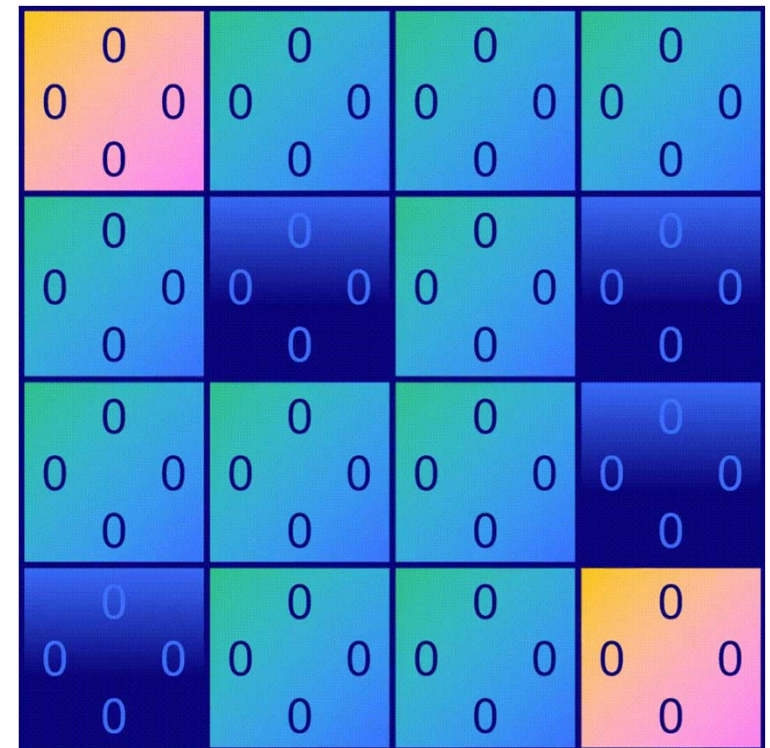
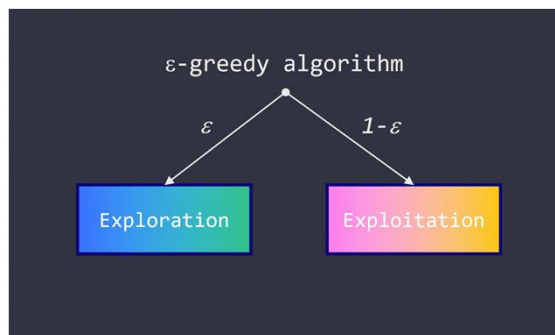


Frozen lake primer

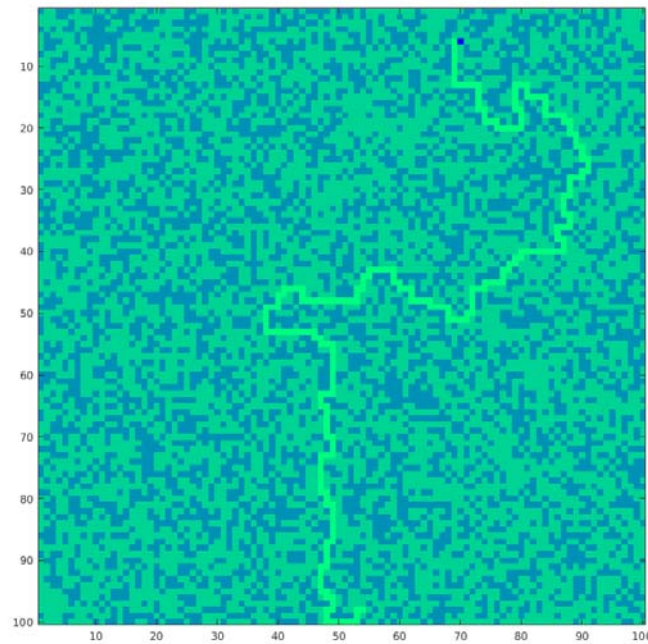
- Reševanje s Q tabelo
 - vrstic: $n \times n$
 - stolpcev: število akcij

$$Q_{new}(s_t, a_t) = Q(s_t, a_t) + \alpha \cdot (r_t + \gamma \cdot \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

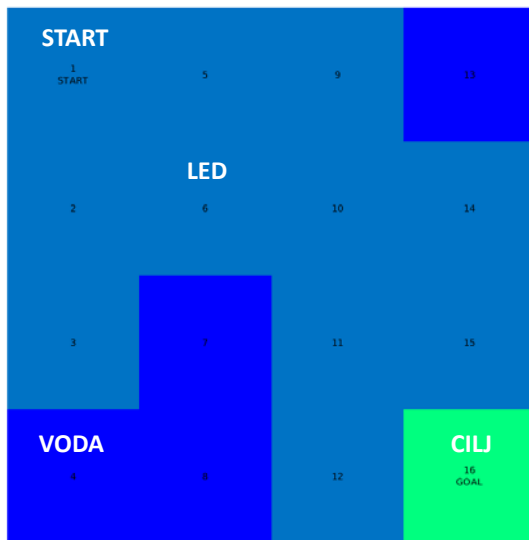
- Izbira akcij



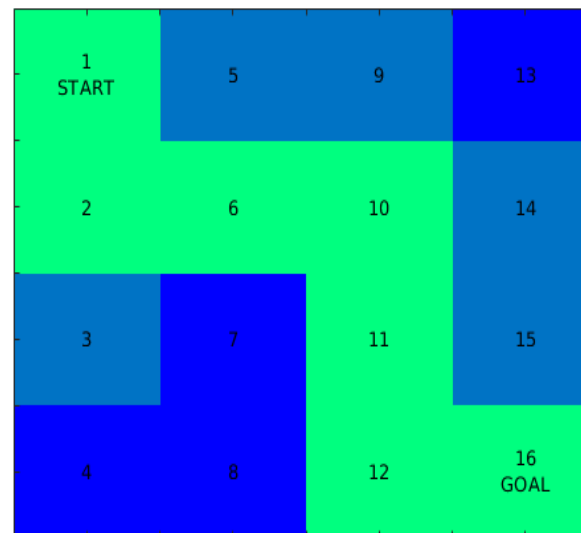
Naloga – „Frozen lake“



Naloga – „Frozen lake“ dimenzije 4 x 4



„Frozen lake“

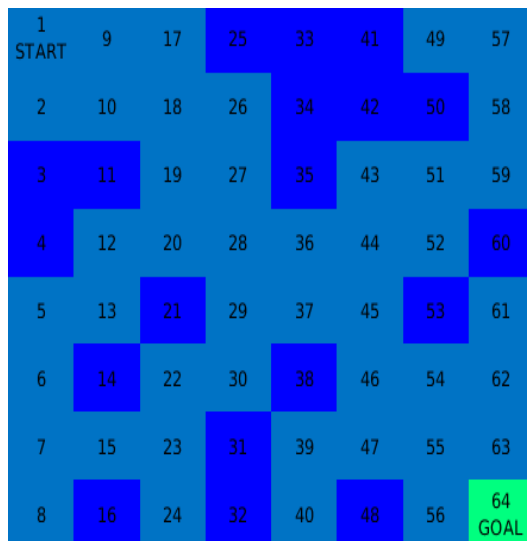


„Frozen lake“ – rešitev problema

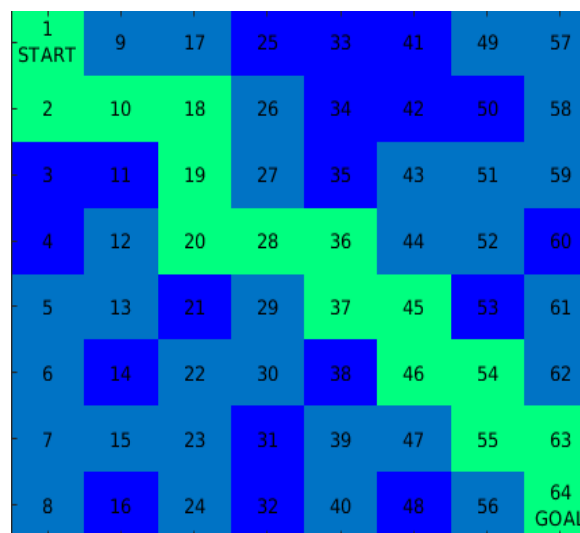
	AKCIJE			
	LEFT	DOWN	RIGHT	UP
1.0000	-2.3578	-1.4293	-1.4293	-2.3578
2.0000	-1.4293	-2.3578	-0.4519	-2.3578
3.0000	-2.3578	-4.0000	-4.0000	-1.4293
4.0000	0	0	0	0
5.0000	-2.3578	-0.4519	-0.4519	-1.4293
6.0000	-1.4293	-4.0000	0.5770	-1.4293
7.0000	0	0	0	0
8.0000	0	0	0	0
9.0000	-1.4293	0.5770	-4.0000	-0.4519
10.0000	-0.4519	1.6600	1.6600	-0.4519
11.0000	-4.0000	2.8000	2.8000	0.5770
12.0000	-4.0000	2.8000	4.0000	1.6600
13.0000	0	0	0	0
14.0000	0.5770	2.8000	1.6600	-4.0000
15.0000	1.6600	4.0000	2.8000	1.6600
16.0000	0	0	0	0

„Frozen lake“ – Q tabela

Naloga – „Frozen lake“ dimenzije 8 x 8



„Frozen lake“



„Frozen lake“ – rešitev problema

AKCIJE				
	LEFT	DOWN	RIGHT	UP
1.0000	-6.3451	-5.6264	-5.6264	-6.3451
2.0000	-5.6264	-8.0000	-4.8699	-6.3451
3.0000	0	0	0	0
4.0000	0	0	0	0
5.0000	-4.9430	-4.9032	-5.0118	-7.9491
6.0000	-4.5211	-4.5135	-7.8540	-4.6217
7.0000	-4.2330	-4.1831	-4.1205	-4.2053
8.0000	-4.0274	-4.0565	-6.7992	-4.0545
9.0000	-6.4069	-4.8699	-4.8923	-5.6488
10.0000	-5.6264	-8.0000	-4.0736	-5.6264
11.0000	0	0	0	0
12.0000	-7.9959	-5.6764	-2.3530	-7.9997
13.0000	-5.2995	-7.8378	-7.9629	-4.7347
14.0000	0	0	0	0

STANJA

...

„Frozen lake“ – Q tabela

Algoritem učenja

Vhod: *strategija π , uint no_epizod, faktor učenja α , potek ϵ*
Izhod: optimalna *funkcija vrednosti Q* (če je število epizod dovolj veliko)

```
Inicializacija:  $q(s, a) = 0 \forall s \in \mathcal{S} \wedge a \in \mathcal{A}$ , in  $q(\text{končno stanje}, \cdot) = 0$ 
for  $i = 1$  to  $no\_epizod$ 
    Izberemo vrednost  $\epsilon$ 
    Opazujemo stanje  $S_0$ 
     $t \leftarrow 0$ 
    repeat
        Izberemo akcijo  $A_t$  na osnovi strategije iz  $Q$  (na primer  $\epsilon$  – požrešna strategija)
        Izvedemo akcijo  $A_t$  in opazujemo nagrado  $R_{t+1}$  in novo stanje  $S_{t+1}$ 
        Posodobimo  $q(S_t, A_t) \leftarrow q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a q(S_{t+1}, a) - q(S_t, A_t)]$ 
         $t \leftarrow t + 1$ 
    until  $S_t$  je končno stanje
end
```

Parametri za učenje

no_epizod – število epizod učenja
 α – parameter hitrosti učenja
 γ – parameter zniževanja vrednosti nagrade
 ϵ – izbira ϵ -požrešne strategije

Matlab predloga

- predloga *frozen_lake_tmplt.m*
 - v spremenljivko *vpisna_stevilka* vpišete vašo vpisno številko in dodate še eno cifro
 - ustvari tabelo *lake* s frozen lake okoljem in nagradami
 - prikaz okolja
- vpišete vaš algoritem za spodbujevalno učenje
 - na koncu skripte še klic funkcije *visualization_Q4.p* za prikaz končne rešitve
 - zapis *num_steps = visualization_Q_arrows4(Q, lake)* izriše akcije v obliki puščic

```

1 %% Create frozen lake env
2 vpisna_stevilka = 649901670;
3 rng(vpisna_stevilka)
4 n = 4;
5
6 lake = -1*ones(n,n);
7
8 for i=1:n
9     for j=1:n
10         if (rand() < 0.25)
11             lake(i,j) = -n;
12         end
13     end
14 end
15 lake(1,1) = -1;
16 lake(1,2) = -1;
17 lake(2,1) = -1;
18 lake(2,2) = -1;
19 lake(n-1,n-1) = -1;
20 lake(n,n-1) = -1;
21 lake(n-1,n) = -1;
22 lake(n,n) = n;
23
24 % Render environment
25 disp(lake)
26
27 fh = figure;
28 imagesc(lake);
29 colormap(winter);
30
31 for i=1:n
32     for j=1:n
33
34         if (i==1) && (j==1)
35             text(1,1,'1','START','HorizontalAlignment','center');
36         elseif (i==n) && (j==n)
37             text(n,n,{num2str(n*n)},'GOAL','HorizontalAlignment','center')
38         else
39             text(j,i,num2str(i+n*(j-1)),'HorizontalAlignment','center')
40         end
41     end
42 end
43
44 axis off
45
46 %%
47 %Vaša koda
48
49 %%
50 % Vizualizacija rešitve
51 indexQ = int32([(1:(n*n))]);
52 visQ = table(indexQ,Q)
53
54 num_steps = vizualizacija_Q(Q, lake);
55
56
57

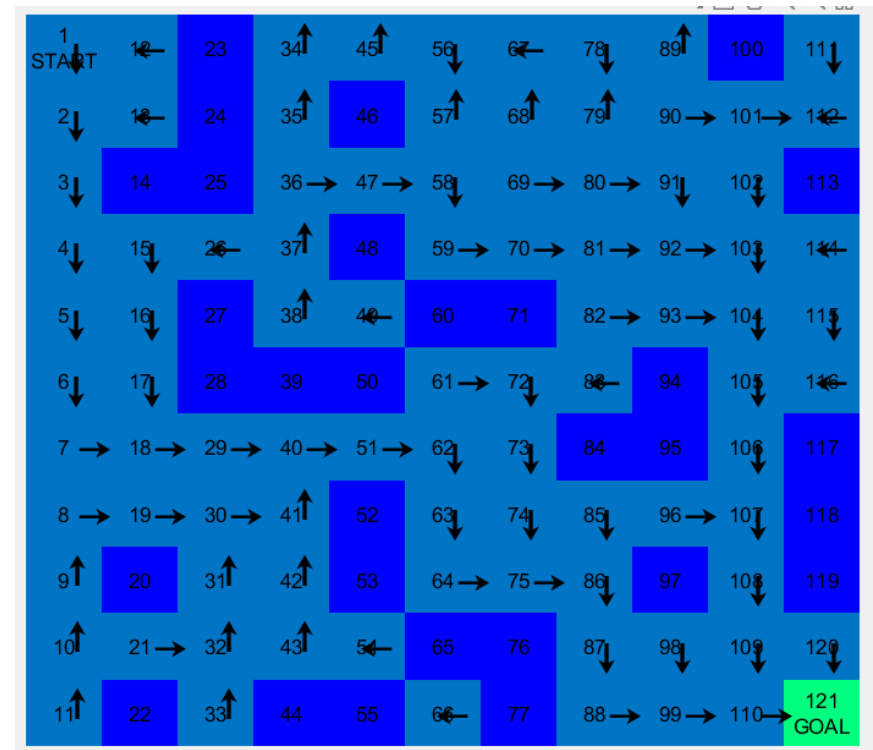
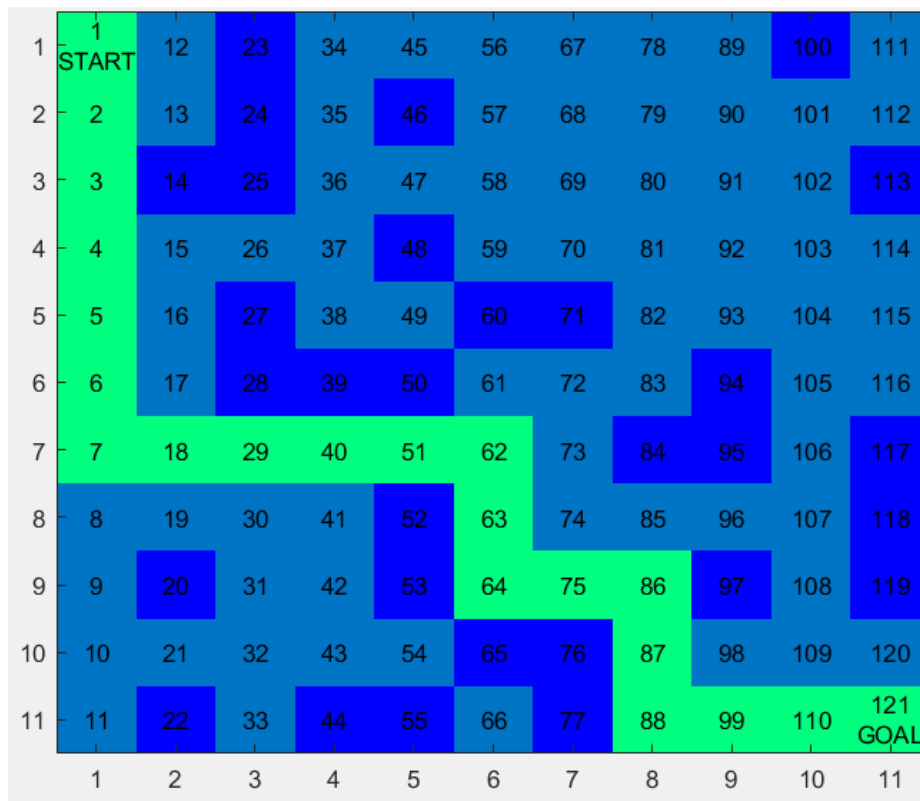
```

Matlab predloga

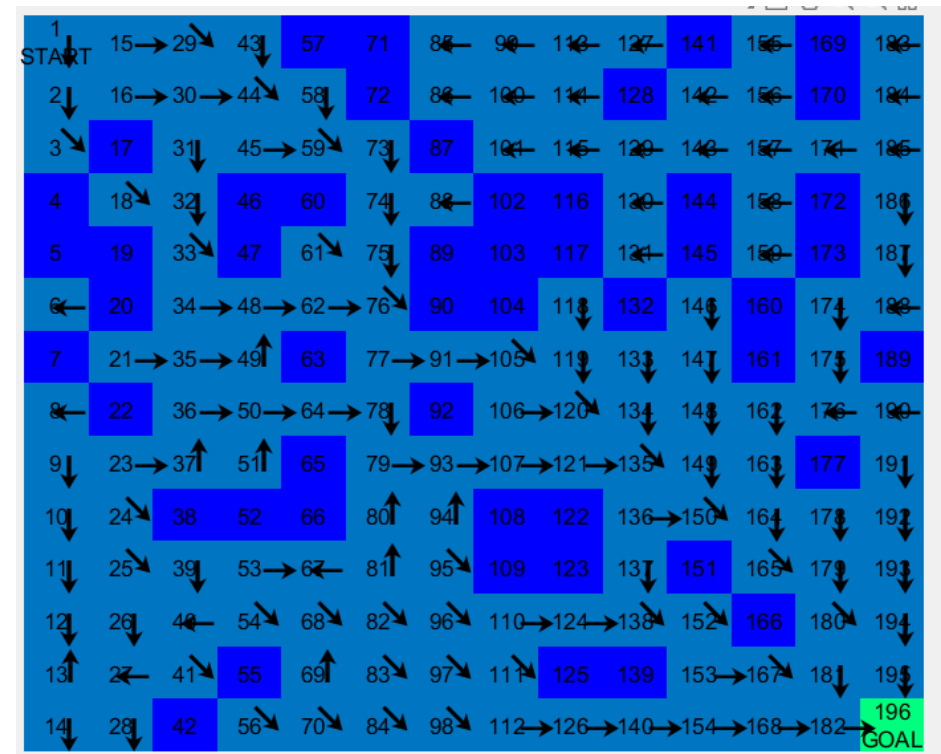
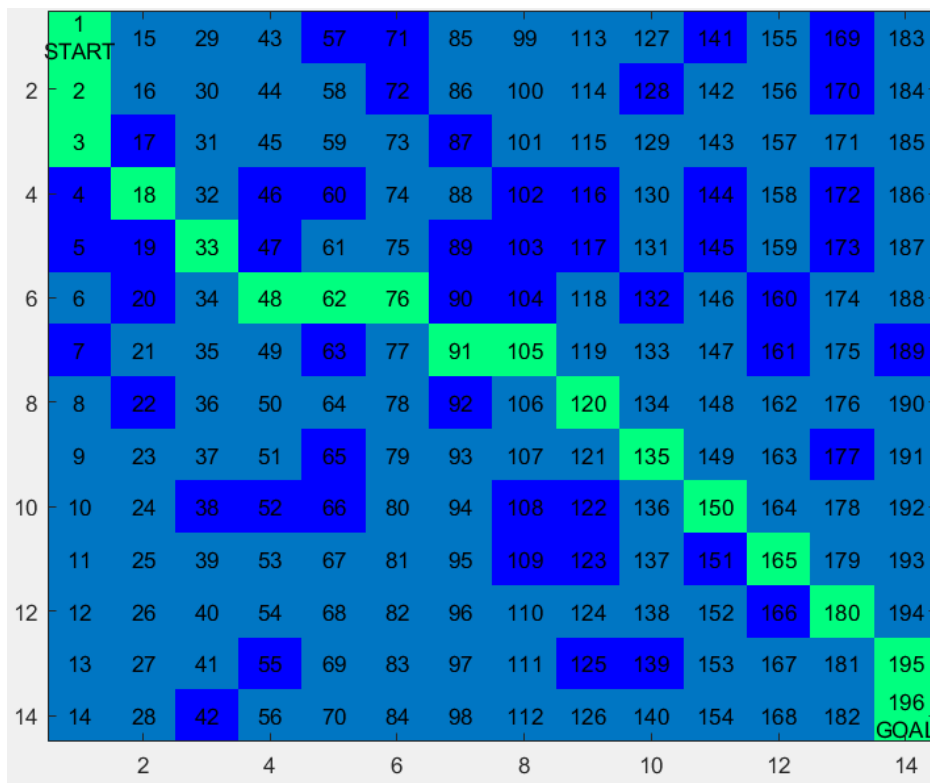
- vizualizacija rešitve s funkcijo *visualization_Q4.p*
- Q tabela mora imeti $n \times n$ vrstic, ter 4 stolpce za 4 akcije:
 - 1. stolpec za akcijo LEFT
 - 2. stolpec za akcijo DOWN
 - 3. stolpec za akcijo RIGHT
 - 4. stolpec za akcijo UP
- vizualizacija rešitve s funkcijo *visualization_Q5.p*
- Q tabela mora imeti $n \times n$ vrstic, ter 5 stolpcev za 5 akcij:
 - 1. stolpec za akcijo LEFT
 - 2. stolpec za akcijo DOWN
 - 3. stolpec za akcijo RIGHT
 - 4. stolpec za akcijo UP
 - 5. stolpec za akcijo RIGHT-DOWN

```
49  
50 %%  
51 % Vizualizacija rešitve  
52 indexQ = int32([(1:(n*n))]');  
53 visQ = table(indexQ,Q)  
54  
55 num_steps = vizualizacija_Q(Q, lake);  
56  
57
```

Primeri rešitev



Primeri rešitev



Vrednosti stanj

- Predavanje
 - Ovrednotenje naključne strategije v „majhni mreži“ (stran 4, zgornja prosojnica)
 - Deterministično iteriranje vrednosti (stran 5 , spodnja prosojnica)
 - *11 - Spodbujevalno učenje - planiranje in predikcija.pdf*