

# Q LEARNING

Ivan Androš  
Dejan Peretin  
Petra Podolski

25. svibnja 2011.

# Q LEARNING

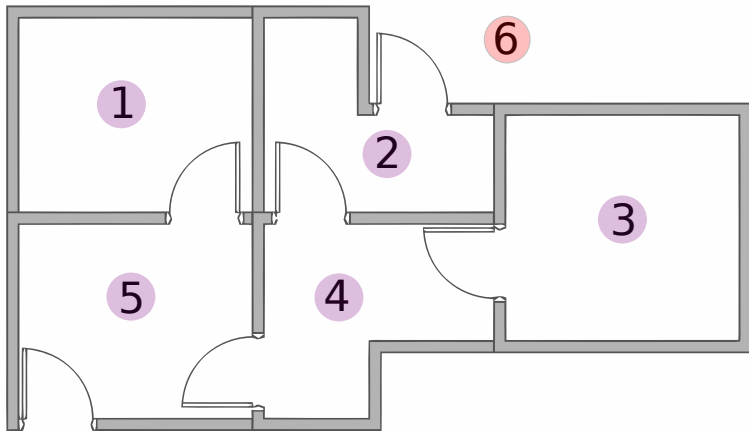
- Tehnika učenja s podrškom.
- Agent uči evaluacijsku funkciju

$$Q : S \times A \rightarrow \mathbb{R}$$

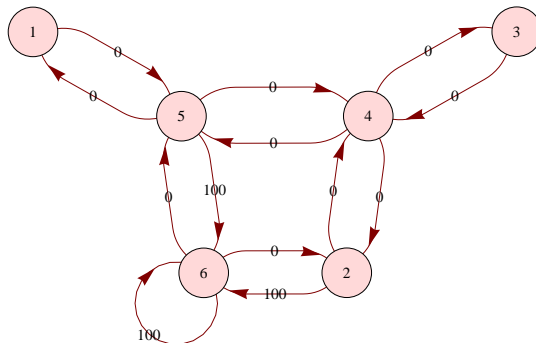
gdje je  $S$  skup stanja, a  $A$  skup akcija.

- Agentu ne mora biti poznat model okoliša.

# KRETANJE ROBOTA



**SLIKA:** Agent se nalazi u jednoj od soba, mora izaći van



SLIKA: Dijagram stanja prethodnog tlocrta

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$\begin{aligned}
 \longrightarrow R = & \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}
 \end{aligned}$$

# UČENJE FUNKCIJE $Q$

$$\begin{array}{c} \longrightarrow \\ R = \end{array} \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$\begin{array}{l}
 \longrightarrow \\
 R = \\
 \longrightarrow
 \end{array}
 \begin{bmatrix}
 - & - & - & - & 0 & - \\
 - & - & - & 0 & - & 100 \\
 - & - & - & 0 & - & - \\
 - & 0 & 0 & - & 0 & - \\
 0 & - & - & 0 & - & 100 \\
 - & 0 & - & - & 0 & 100
 \end{bmatrix}
 \quad
 Q =
 \begin{bmatrix}
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}$$



# UČENJE FUNKCIJE $Q$

$$\begin{array}{l}
 \longrightarrow \\
 R = \\
 \longrightarrow
 \end{array}
 \begin{bmatrix}
 - & - & - & - & 0 & - \\
 - & - & - & 0 & - & 100 \\
 - & - & - & 0 & - & - \\
 - & 0 & 0 & - & 0 & - \\
 0 & - & - & 0 & - & 100 \\
 - & 0 & - & - & 0 & 100
 \end{bmatrix}
 \quad
 Q =
 \begin{bmatrix}
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$\begin{array}{l}
 \longrightarrow \\
 R = \\
 \longrightarrow
 \end{array}
 \begin{bmatrix}
 - & - & - & - & 0 & - \\
 - & - & - & 0 & - & 100 \\
 - & - & - & 0 & - & - \\
 - & 0 & 0 & - & 0 & - \\
 0 & - & - & 0 & - & 100 \\
 - & 0 & - & - & 0 & 100
 \end{bmatrix}
 \quad
 Q =
 \begin{bmatrix}
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}$$

$$Q_{2,6} = R_{2,6} + 0.8 \cdot \max\{Q_{6,2}, Q_{6,5}, Q_{6,6}\} = 100 + 0.8 \cdot 0 = 100$$

# UČENJE FUNKCIJE $Q$

$$\begin{array}{l}
 \longrightarrow \\
 R = \\
 \longrightarrow
 \end{array}
 \begin{bmatrix}
 - & - & - & - & 0 & - \\
 - & - & - & 0 & - & 100 \\
 - & - & - & 0 & - & - \\
 - & 0 & 0 & - & 0 & - \\
 0 & - & - & 0 & - & 100 \\
 - & 0 & - & - & 0 & 100
 \end{bmatrix}
 \quad
 Q =
 \begin{bmatrix}
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 100 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}$$

$$Q_{2,6} = R_{2,6} + 0.8 \cdot \max\{Q_{6,2}, Q_{6,5}, Q_{6,6}\} = 100 + 0.8 \cdot 0 = 100$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & \textcolor{red}{0} & - & \textcolor{red}{100} \\ - & - & - & 0 & - & - \\ - & \textcolor{red}{0} & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \textcolor{red}{0} & 0 & \textcolor{red}{100} \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \textcolor{red}{80} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$Q_{4,2} = R_{4,2} + 0.8 \cdot \max\{Q_{2,4}, Q_{2,6}\} = 0 + 0.8 \cdot 100 = 80$$



# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# UČENJE FUNKCIJE $Q$

$$R = \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \quad Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$Q_{2,6} = R_{2,6} + 0.8 \cdot \max\{Q_{6,2}, Q_{6,5}, Q_{6,6}\} = 100 + 0.8 \cdot 0 = 100$$

# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

PUT : 3

# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & \textcolor{red}{320} & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

PUT : 3  $\rightarrow$  4

# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

PUT : 3  $\rightarrow$  4



# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

PUT : 3  $\rightarrow$  4  $\rightarrow$  2

# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

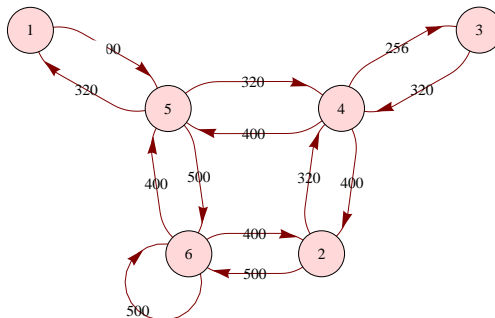
PUT : 3  $\rightarrow$  4  $\rightarrow$  2

# NALAŽENJE NAJKRAĆEG PUTA

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

PUT : 3  $\rightarrow$  4  $\rightarrow$  2  $\rightarrow$  6

# NALAŽENJE NAJKRAĆEG PUTA

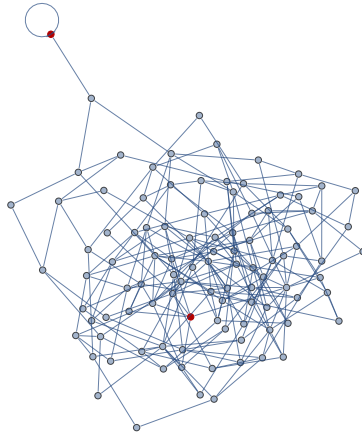


SLIKA: Dijagram stanja iz perspektive funkcije  $Q$

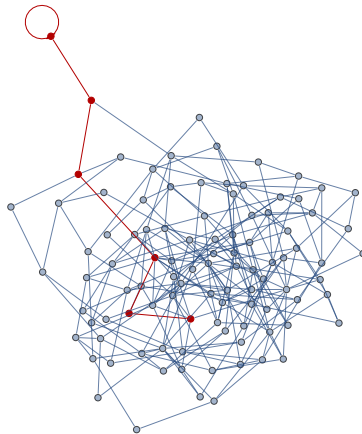
# PSEUDOKOD

- 1 učitaj parametra  $\gamma$  i matricu  $R$
- 2 inicijaliziraj vrijednosti matrice  $Q$  na 0
- 3 ponavljaj za svaku epizodu
  - na slučajan način izaberi inicijalno stanje
  - ponavljaj dok ne dođeš u ciljno stanje
    - izaberi jedno od mogućih akcija za trenutno stanje
    - $Q_{s,a} = R_{s,a} + \gamma \cdot \max_i \{Q_{a,a_i}\}$
    - postavi sljedeće stanje za trenutno stanje

učitaj parametar  $\gamma$  i matricu  $R$   
inicijaliziraj vrijednosti matrice  $Q$  na 0  
**while** nema konvergencije **do**  
    na slučajan način izaberi inicijalno stanje  
    **while** nismo u završnom stanju **do**  
        izaberi jedno od mogućih akcija za trenutno stanje  
         $Q_{s,a} = R_{s,a} + \gamma \cdot \max_i \{Q_{a,a_i}\}$   
        postavi sljedeće stanje za trenutno stanje  
    **end while**  
**end while**

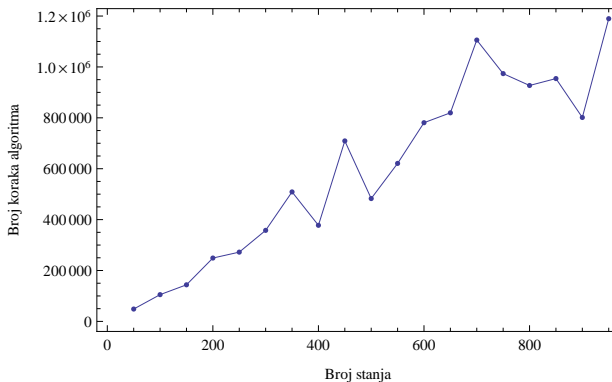


SLIKA: Dijagram sa 100 stanja

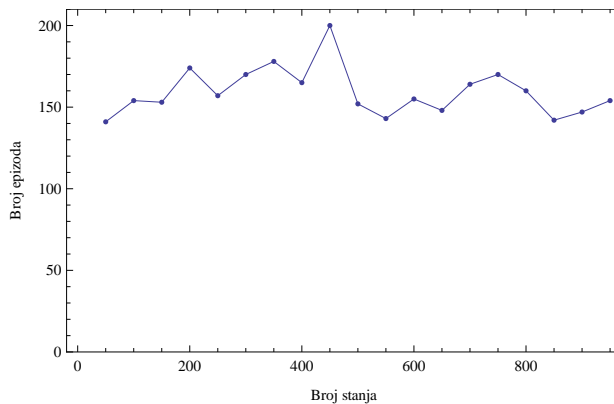


SLIKA: Dijagram sa 100 stanja

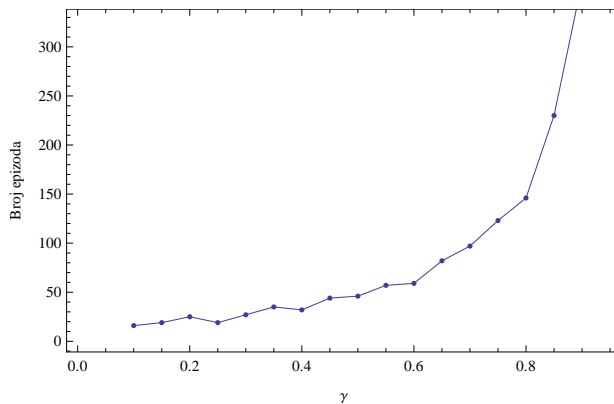




**SLIKA:** Broj koraka algoritma u odnosu na broj stanja



**SLIKA:** Broj epizoda u odnosu na broj stanja



**SLIKA:** Broj epizoda u odnosu na vrijednost  $\gamma$