

# INTELIGENȚĂ , ARTIFICIALĂ



**Sisteme inteligente**

Învățare prin întărire

Laura Dioșan

# Sumar

---

## A. Scurtă introducere în Inteligența Artificială (IA)

## B. Rezolvarea problemelor prin căutare

- Definirea problemelor de căutare
- Strategii de căutare
  - Strategii de căutare neinformate
  - Strategii de căutare informate
  - Strategii de căutare locale (Hill Climbing, Simulated Annealing, Tabu Search, Algoritmi evolutivi, PSO, ACO)
  - Strategii de căutare adversială

## C. Sisteme inteligente

- Sisteme care învață singure
  - Arbori de decizie
  - Rețele neuronale artificiale
  - Mașini cu suport vectorial
  - Algoritmi evolutivi
  - *Q-learning*
- Sisteme bazate pe reguli
- Sisteme hibride

# Materiale de citit și legături utile

---

- ❑ capitolul IV.20 din *S. Russell, P. Norvig, Artificial Intelligence: A Modern Approach, Prentice Hall, 1995*
- ❑ capitolul 6 din *H.F. Pop, G. Șerban, Inteligență artificială, Cluj Napoca, 2004*
- ❑ capitolul 13 din *T. M. Mitchell, Machine Learning, McGraw-Hill Science, 1997*
- ❑ capitolul 17 din *C. Groșan, A. Abraham, Intelligent Systems: A Modern Approach, Springer, 2011*
- ❑ capitolul 9 din *Adrian A. Hopgood, Intelligent Systems for Engineers and Scientists, CRC Press, 2001*

# Învățare automată

---

- Învățarea automată, în funcție de experiența pe care se bazează, poate fi:
  - Învățare supervizată
  - Învățare nesupervizată
  - Învățare cu întărire

# Învățare prin întărire (reinforcement learning)

---

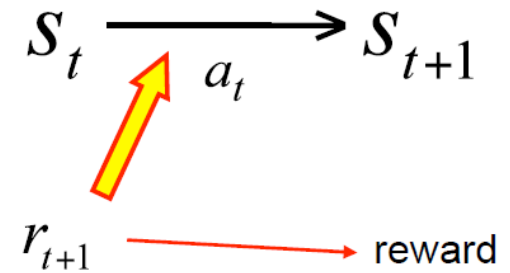
## □ Ideea de bază

- Învățarea din interacțiuni
- Învățarea unui mod de acțiune pentru a maximiza o recompensă (numerică)
  - Ex. *Ce recompensa primesc daca fac acest lucru?*
- Unul sau mai mulți agenți care
  - Învăță permanent conform principiului "încercare și eroare"
  - Planifică permanent
  - Afectează mediul înconjurător
  - Execută o mulțime de sarcini (acțiuni)

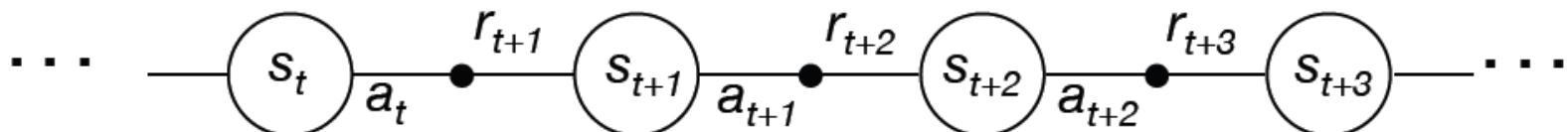
# Învățare prin întărire (reinforcement learning)

## □ Definirea problemei

- Un mediu definit printr-un set de stări posibile
  - $S = \{s_1, s_2, \dots\}$
  - Stări temporale
- Un set de acțiuni posibil de efectuat
  - Într-o anumită stare  $A(s_t) = \{a_1, a_2, \dots\}$
  - În toate stările  $A = \cup_{s_t \in S} A(s_t)$
- O recompensă  $r_{t+1}$  obținută prin deplasarea dintr-o stare  $s_t$  într-o stare  $s_{t+1}$  executând o acțiune  $a_t$



- Agentul și mediul interacționează la diferite momente de timp  $t$ ,  $t = 0, 1, 2, \dots$ 
  - Agentul observă starea la momentul  $t$ :  $s_t$
  - Efectuează o acțiune la momentul  $t$ :  $a_t$
  - Primește o recompensă:  $r_{t+1}$



# Învățare prin întărire (reinforcement learning)

---

## □ Aplicații

- Robotică
- Planificare
- Jocuri
- Sisteme de control

# Învățare prin întărire (reinforcement learning)

---

## □ Învățare supervizată vs. Învățare cu întărire

### ■ Învățarea supervizată

- Informațiile de antrenare → date brute + scopul urmărit
  - Caracteristici ale datelor și clasele din care fac parte
- Modul de învățare
  - Minimizarea erorii
    - Diferența între ieșirea corectă și ieșirea calculată de algoritm
  - Fără interacțiuni cu mediul

### ■ Învățarea cu întărire

- Informațiile de antrenare → date brute + evaluări (recompense/penalizări)
  - Acțiuni posibile și recompensele corespunzătoare lor
- Modul de învățare
  - Optimizarea evaluărilor primite
  - Interacțiuni cu mediul



# Învățare prin întărire (reinforcement learning)

---

## □ Învățare nesupervizată vs. Învățare cu întărire

### ■ Învățarea nesupervizată

- Informațiile de antrenare → date brute, fără rezultate (scop sau recompensă/pedeapsă)
  - Caracteristici ale datelor (fără clasele din care fac parte)
- Modul de învățare
  - Minimizarea diferențelor dintre elementele aceleiași clase
  - Maximizarea diferențelor dintre elementele claselor diferite
  - Fără interacțiune cu mediul

### ■ Învățarea cu întărire

- Informațiile de antrenare → date brute + evaluări (recompense/penalizări)
  - Acțiuni posibile și recompensele corespunzătoare lor
- Modul de învățare
  - Optimizarea evaluărilor primite
  - Interacțiune cu mediul

# Învățare prin întărire (reinforcement learning)

---

- Scopul unui algoritm de învățare prin întărire
  - Maparea unor stări în alte stări prin intermediul execuției unor acțiuni de interacțiune cu mediul cu scopul maximizării unei recompense
  
- Elemente
  - Mediul de interacțiune
  - Funcția de întărire (de recompensare)
  - Funcția de evaluare (a unei stări)

# Învățare prin întărire (reinforcement learning)

---

## □ Elemente → Mediul

### ■ Dinamic

- Acțiuni de nivel primar
  - Ex. Voltajul motoarelor
- Acțiuni de nivel înalt
  - Ex. Acceptarea unui job

### ■ Observabil

- Citiri de senzori
- Descrieri simbolice

# Învățare prin întărire (reinforcement learning)

---

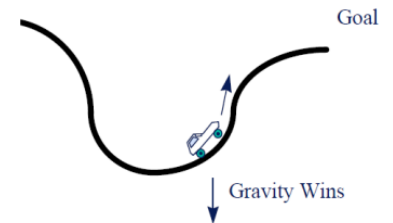
- Elemente → Funcția de întărire
  - Are scopul evaluării unei stări în care a ajuns agentul ca urmare a unei acțiuni
  - Măsoară recompensa primită ca urmare a efectuării unei acțiuni
  - Agentul va învăța să execute acele acțiuni care îi
    - maximizează recompensa sau
    - minimizează penalizarea

# Învățare prin întărire (reinforcement learning)

## □ Elemente → Funcția de întărire

### ■ Tipologie

- Funcții de tip recompensă întârziată pură ("probleme de evitat")
  - Orice stare este evaluată cu 0, mai puțin
  - starea finală care este evaluată cu
    - +1, dacă este o stare obiectiv
    - -1, dacă este o stare ne-obiectiv (stare "de evitat")
  - Ex.
    - jocul de table
- Funcții care minimizează timpul necesar atingerii scopului
  - Efectuarea acelor acțiuni care conduc la starea finală pe "drumul cel mai scurt"
  - Evaluare stări
    - 0, pentru starea finală
    - -1, pentru toate celelalte stări
  - Ex
    - Mașina care urcă un deal
- Funcții specifice jocurilor
  - Doi sau mai mulți jucători cu obiective opuse
  - Întărire de tip maximin, minimax
  - Ex
    - Rachete și avioane țintă



# Învățare prin întărire (reinforcement learning)

---

## □ Elemente → Funcția valoare

### ■ Utilitate

- Cum se aleg cele mai bune acțiuni?
- Cum se poate măsura utilitatea unei acțiuni?

### ■ Concepte de bază

- Politică  $\pi(s)$  → ce acțiuni trebuie executate din fiecare stare  $s$
- Valoarea unei stări  $V(s)$  = suma recompenselor primite ca urmare a execuției unei politici (set de acțiuni) din acea stare  $s$  până la o stare finală

### ■ Definiție

- o mapare de la valoarea unei stări la valoarea unei alte stări
- poate fi aproximată

# Învățare prin întărire (reinforcement learning)

---

## □ Elemente → Funcția valoare

- Valoarea aproximativă a unei stări  $V(s_t)$
- Valoarea optimă a unei stări  $V^*(s_t)$

$$V(s_t) = V^*(s_t) + e(s_t)$$

$$V(s_t) = r(s_t) + \gamma V(s_{t+1})$$

□ unde:

- $e(s_t)$  – eroarea de aproximare
  - $r(s_t)$  – recompensa primită în starea  $s_t$
- 
- Cum se aproximează valoarea unei stări (algoritmi)?
    - Q-learning
    - SARSA (State-Action-Reward-State-Action)

# Învățare prin întărire (reinforcement learning)

---

## □ Elemente → Funcția valoare → Q-learning

### ■ Ideea de bază

- $Q(s, a) = r(s, a) + \gamma \max_{a'} Q(s', a')$
- dintre toate acțiunile posibile  $a_i$  se alege acea acțiune care produce cea mai mare valoare  $Q(s, a_i)$

### ■ Algoritm

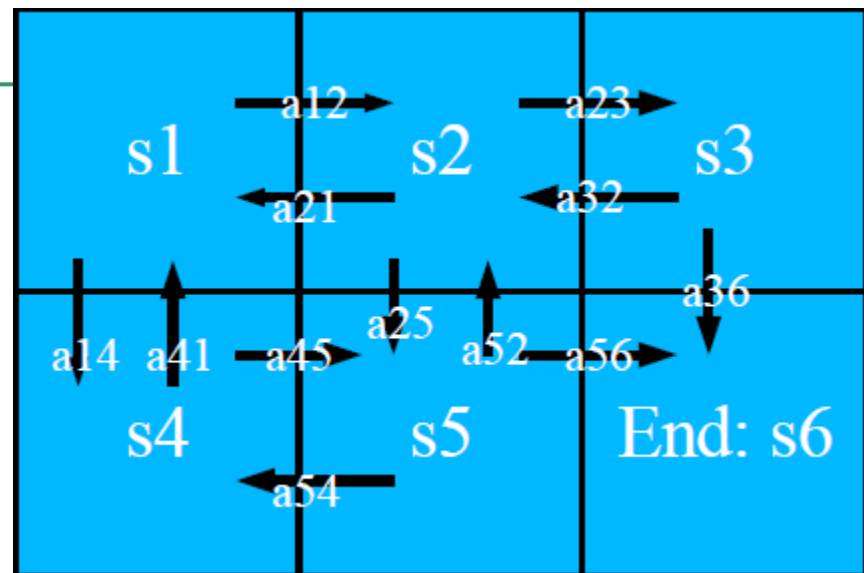
- Pentru fiecare pereche  $(s, a)$  se reține valoarea 0 (într-un tablou  $Q$ )
- Fie starea curentă  $s$
- Repetă
  - Se alege o acțiune  $a$  și se execută
  - Se înregistrează recompensa imediată  $r$
  - Se observă noua stare  $s'$
  - Se modifică valoarea  $Q(s, a) = r + \gamma \max_{a'} Q(s', a')$
  - $s = s'$
- Până când se ajunge la un număr de repetări prefixat



# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	0
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



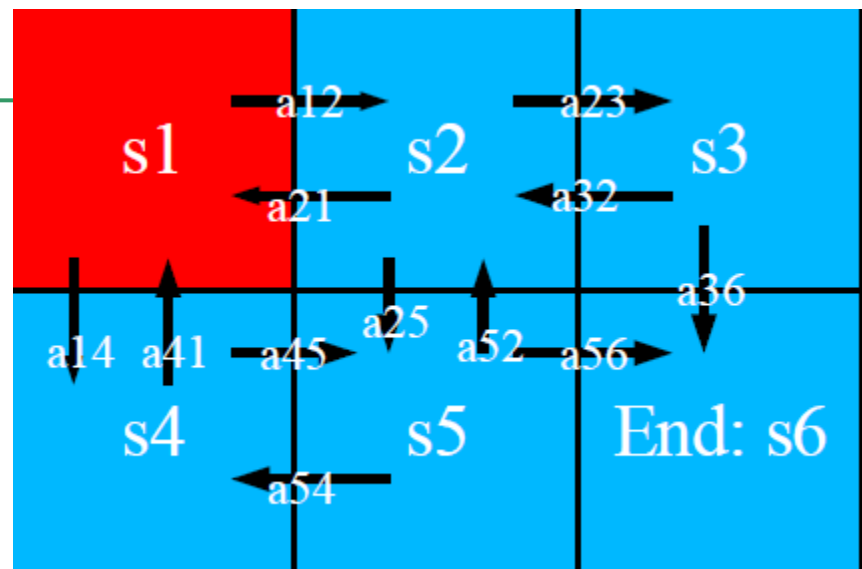
Pp:

- $r = 100$ , dacă se ajunge în  $s_6$   
0, altfel
- $\gamma = 0.5$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	0
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



poziție curentă:  $s_1$

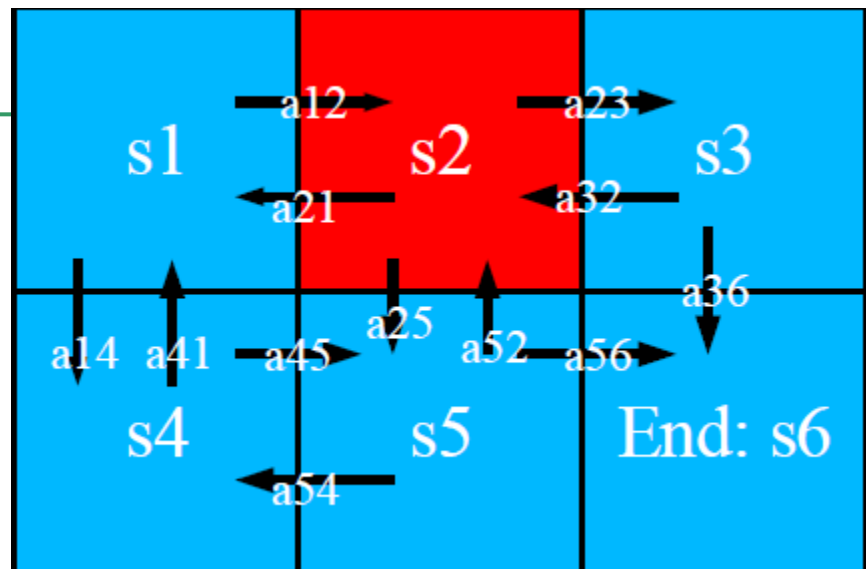
acțiuni posibile:  $a_{12}, a_{14}$

acțiune selectată:  $a_{12}$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	0
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



acțiuni posibile din  $s_2$ :  $a_{21}, a_{23}, a_{25}$

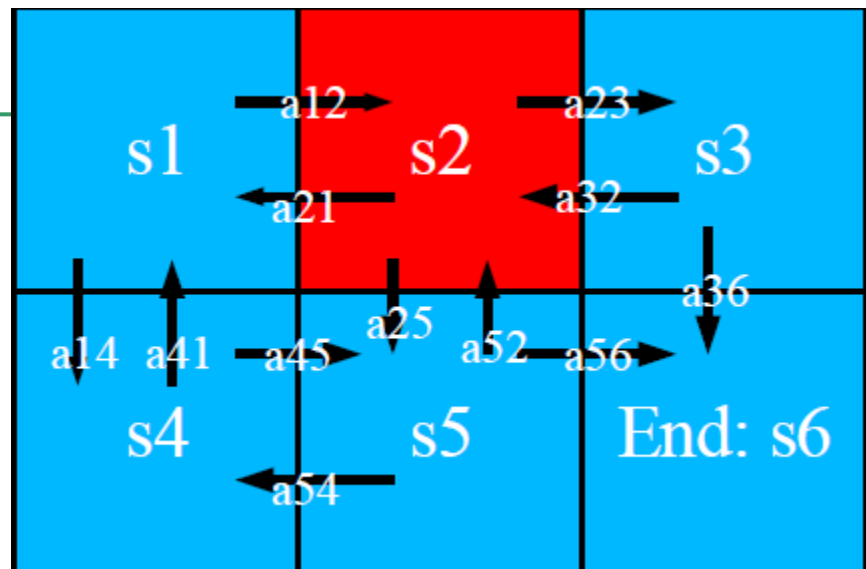
$$Q(s_1, a_{12}) = r + \gamma * \max\{Q(s_2, a_{21}), Q(s_2, a_{23}), Q(s_2, a_{25})\}$$

$$Q(s_1, a_{12}) = 0$$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	0
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



Poziție curentă:  $s_2$

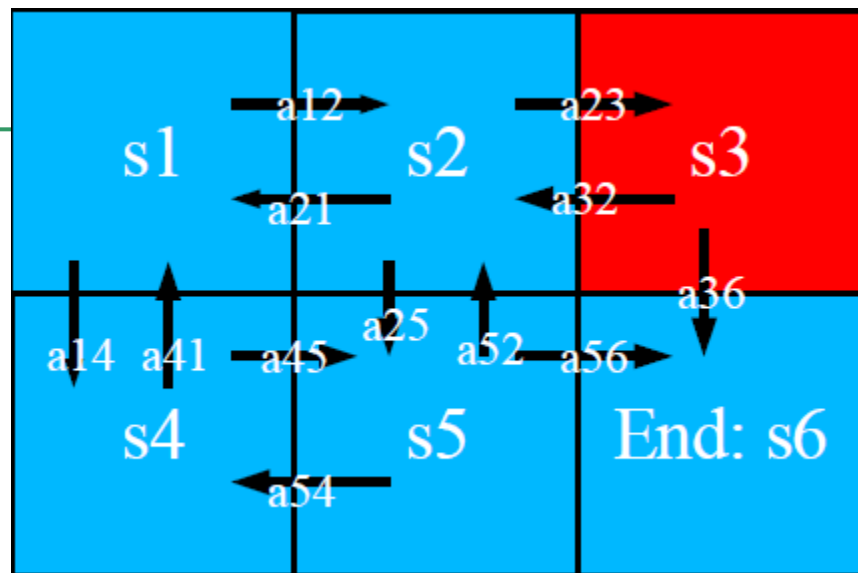
Acțiuni posibile:  $a_{21}, a_{23}, a_{25}$

Acțiune selectată:  $a_{23}$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	0
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



acțiuni posibile din  $s_3$ :  $a_{32}, a_{36}$

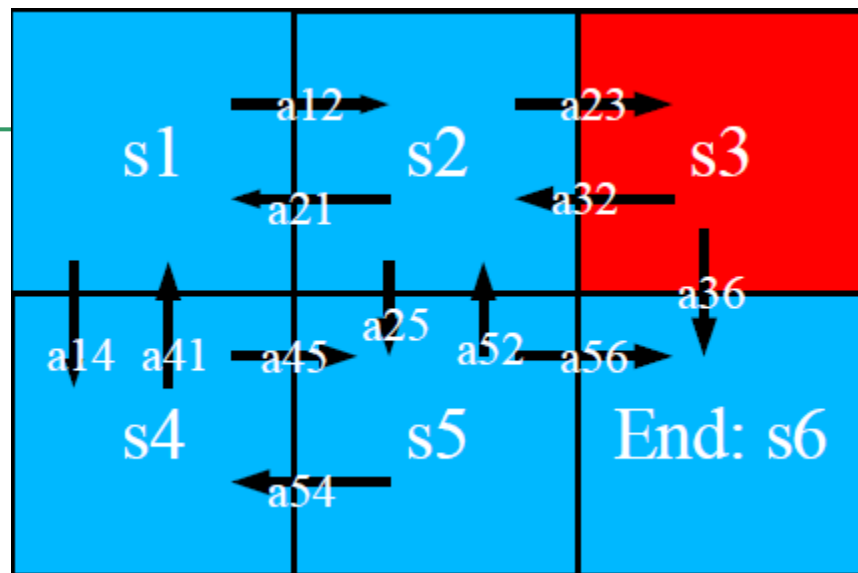
$$Q(s_2, a_{23}) = r + \gamma * \max\{Q(s_3, a_{32}), Q(s_3, a_{36})\}$$

$$Q(s_2, a_{23}) = 0$$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	0
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



Poziție curentă:  $s_3$

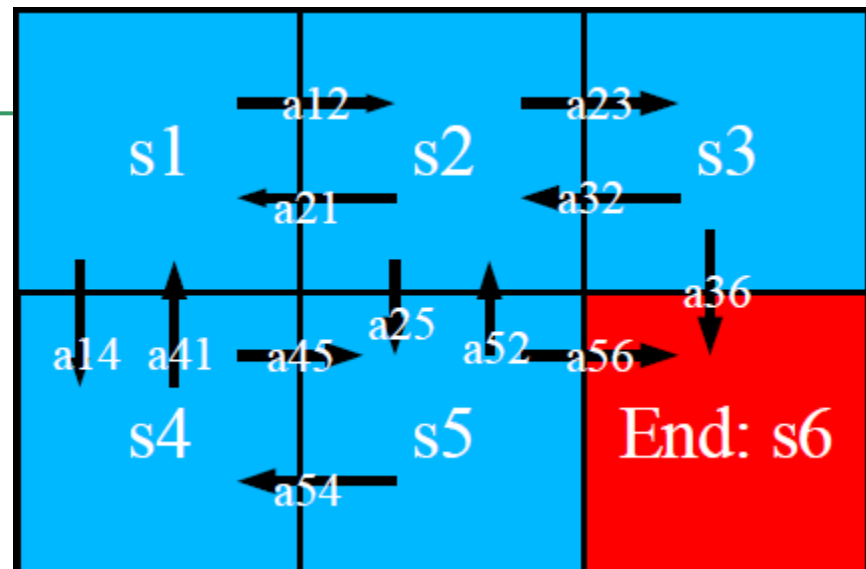
Acțiuni posibile:  $a_{32}, a_{36}$

Acțiune selectată:  $a_{36}$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s, a$	$Q(s, a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
<b><math>s_3, a_{36}</math></b>	<b>100</b>
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



$S_6$  – stare finală

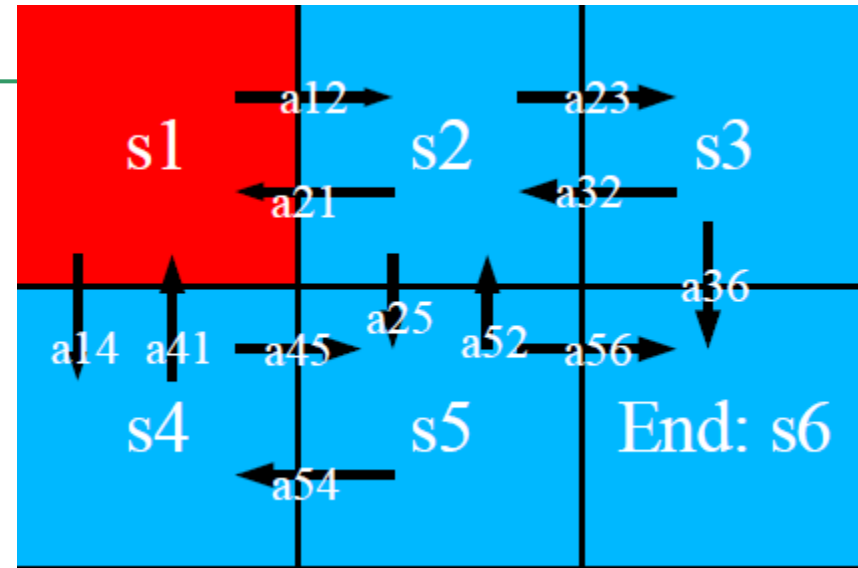
$$Q(s_3, a_{36}) = r$$

$$Q(s_3, a_{36}) = 100$$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	100
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



JOC NOU

poziție curentă:  $s_1$

acțiuni posibile:  $a_{12}, a_{14}$

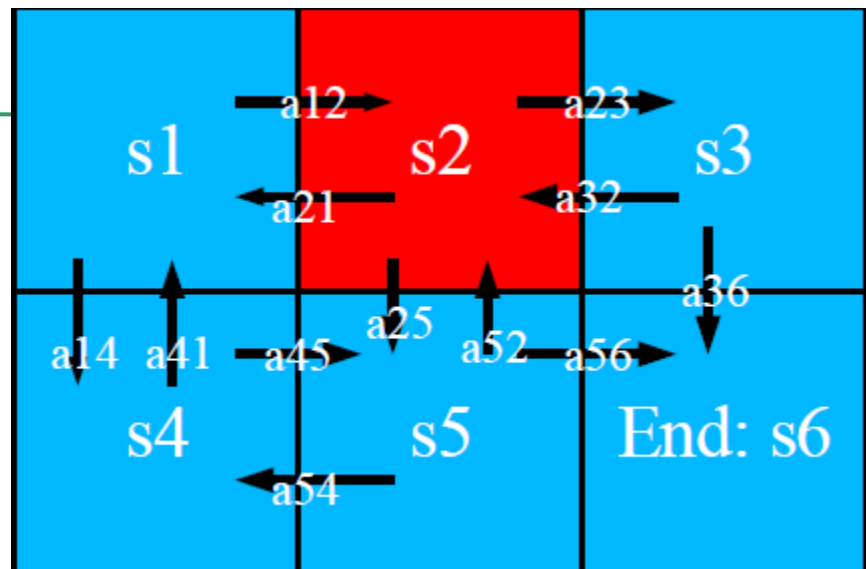
acțiune selectată:  $a_{12}$



# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	100
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



acțiuni posibile din  $s_2$ :  $a_{21}, a_{23}, a_{25}$

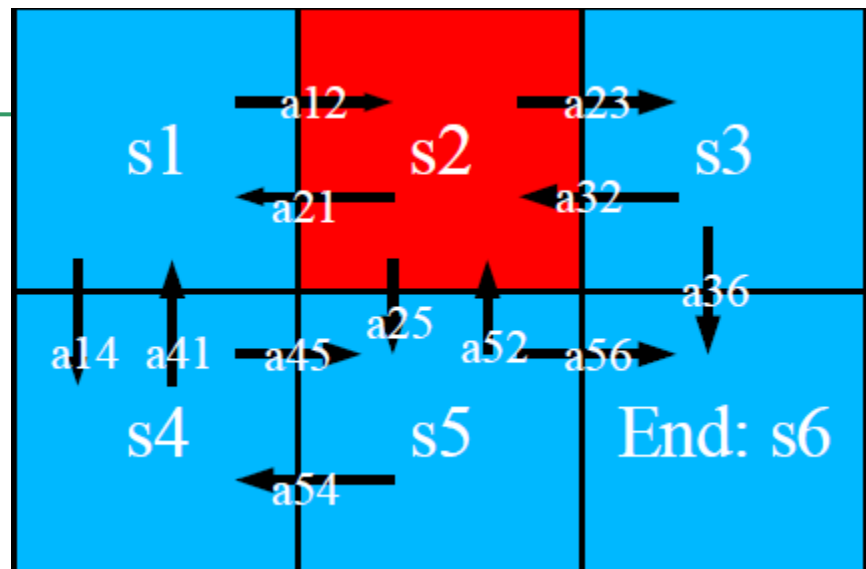
$$Q(s_1, a_{12}) = r + \gamma * \max\{Q(s_2, a_{21}), Q(s_2, a_{23}), Q(s_2, a_{25})\}$$

$$Q(s_1, a_{12}) = 0$$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	100
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



Poziție curentă:  $s_2$

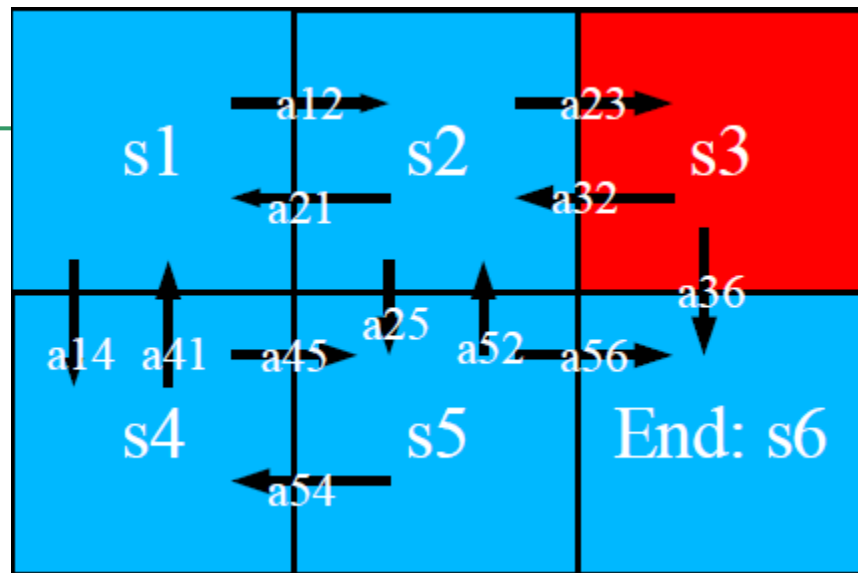
Acțiuni posibile:  $a_{21}, a_{23}, a_{25}$

Acțiune selectată:  $a_{23}$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	50
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	100
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



acțiuni posibile din  $s_3$ :  $a_{32}, a_{36}$

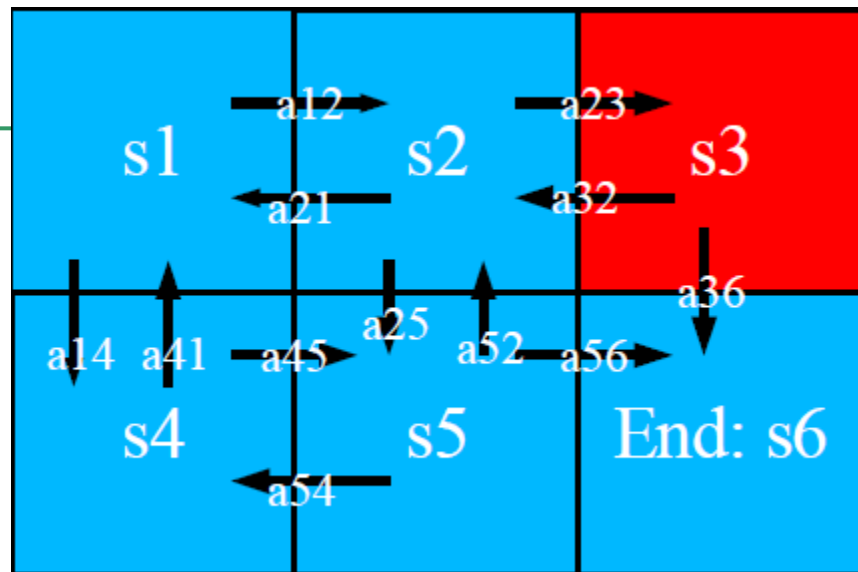
$$Q(s_2, a_{23}) = r + \gamma * \max\{Q(s_3, a_{32}), Q(s_3, a_{36})\}$$

$$Q(s_2, a_{23}) = 50$$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	50
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	100
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



Poziție curentă:  $s_3$

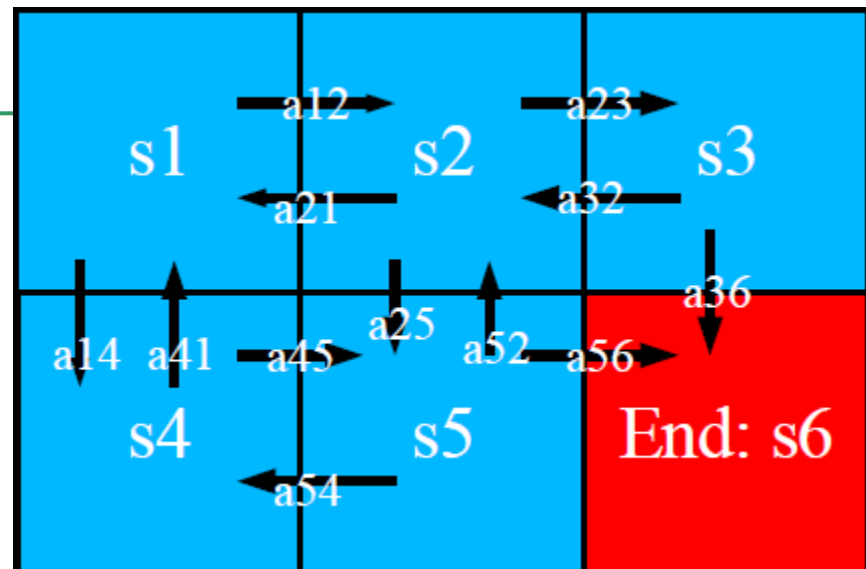
Acțiuni posibile:  $a_{32}, a_{36}$

Acțiune selectată:  $a_{36}$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s, a$	$Q(s, a)$
$s_1, a_{12}$	0
$s_1, a_{14}$	0
$s_2, a_{21}$	0
$s_2, a_{23}$	0
$s_2, a_{25}$	0
$s_3, a_{32}$	0
$s_3, a_{36}$	100
$s_4, a_{41}$	0
$s_4, a_{45}$	0
$s_5, a_{54}$	0
$s_5, a_{52}$	0
$s_5, a_{56}$	0



$S_6$  – stare finală

$$Q(s_3, a_{36}) = r$$

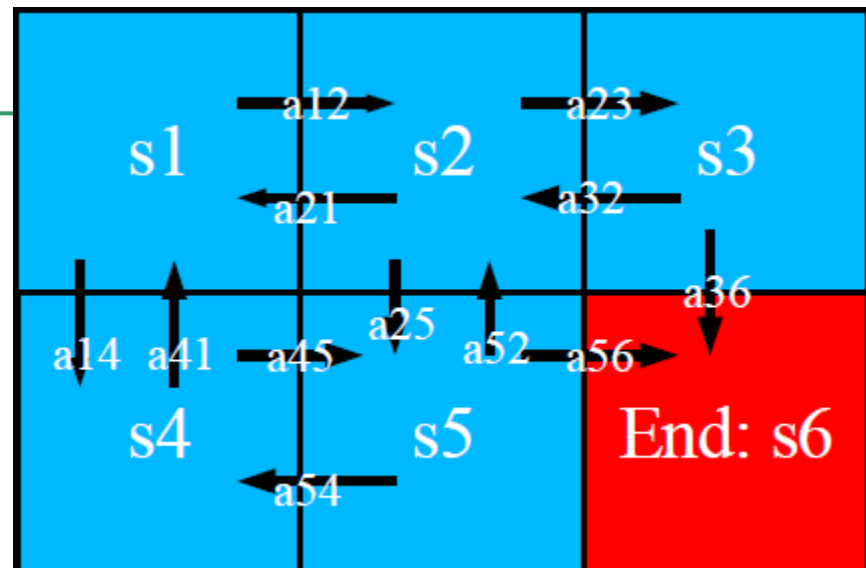
$$Q(s_3, a_{36}) = 100$$

# Învățare prin întărire (reinforcement learning)

## Q-learning (exemplu)

$s,a$	$Q(s,a)$
$s_1, a_{12}$	25
$s_1, a_{14}$	25
$s_2, a_{21}$	12.5
$s_2, a_{23}$	50
$s_2, a_{25}$	25
$s_3, a_{32}$	25
<b><math>s_3, a_{36}</math></b>	<b>100</b>
$s_4, a_{41}$	12.5
$s_4, a_{45}$	50
$s_5, a_{54}$	25
$s_5, a_{52}$	25
$s_5, a_{56}$	100

.....



# Învățare prin întărire (reinforcement learning)

---

## □ Caracteristici

- Învățare din recompense
- Interacțiuni cu sarcinile
  - Secvențe de stări, acțiuni și recompense
- Lumi incerte și nedeterministe
- Consecințe întârziate
- Învățare direcționată către țintă
- Echilibru între explorare și exploatare

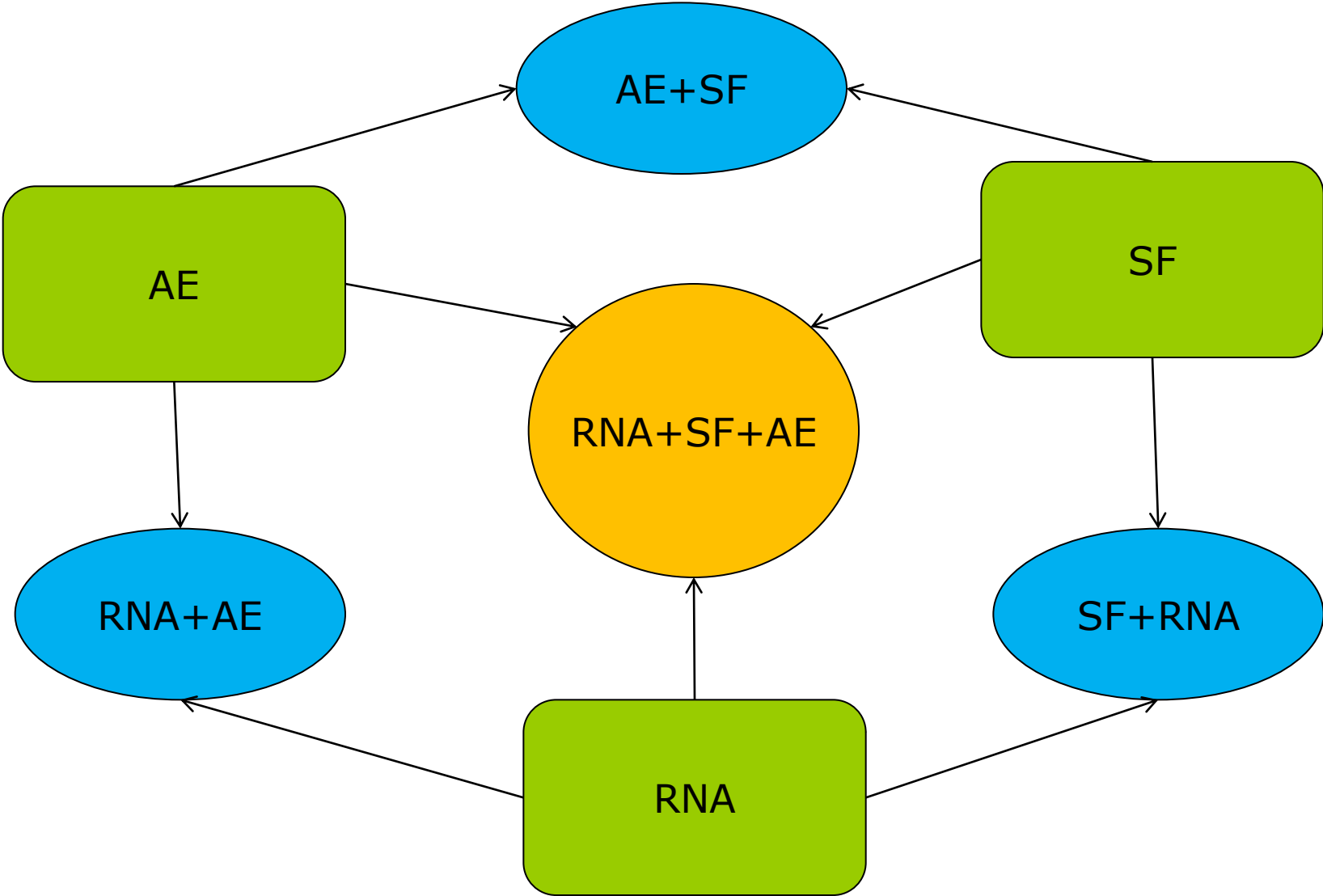
# Sisteme inteligente – sisteme hibride

---

- O combinație între 2 sau mai multe sisteme inteligente
  - Sisteme bazate pe reguli
    - Sisteme fuzzy (SF)
  - Sisteme bazate pe învățare automată
    - Rețele neuronale artificiale (RNA)
    - Algoritmi evolutivi (AE)



# Sisteme inteligente – sisteme hibride



# Sisteme inteligente – sisteme hibride

---

## □ Sisteme neuronale fuzzy (SF+RNA)

### ■ Modele cooperative

- RNA determină parametrii SF
  - funcțiile de apartenență → aproximare
- și/sau
  - regulile fuzzy → clustering (SOM – Self Organised Map)
- pe baza datelor de antrenament

### ■ Modele concurente

- RNA asistă SF în determinarea parametrilor

# Sisteme inteligente – sisteme hibride

---

## □ Sisteme evolutive fuzzy (AE+SF)

- Adaptarea evolutivă a funcțiilor de apartenență
  - Cromozomii codează parametrii diferitelor funcții de apartenență trapezoidale, triunghiulare, logistice, Laplace, Gaussian, etc)
  - Necesită existența unei baze de reguli
  - Se optimizează performanța unui motor de inferență deja existent
- Învățarea evolutivă a regulilor *if-then*
  - Cromozomii codează una sau toate regulile din baza de cunoștințe → se construiește un nou motor de inferență
  - Complexitate mărită

# Sisteme inteligente – sisteme hibride

---

- Rețele neuronale evolutive (RNA+AE)
  - Adaptarea evolutivă a design-ului RNA
    - Adaptarea evolutivă a ponderilor RNA
    - Adaptarea evolutivă a structurii RNA și a funcțiilor de activare a nodurilor RNA
    - Adaptarea evolutivă a regulilor de învățare

# Recapitulare



## □ Sisteme care învață singure (SIS)

### ■ Instruire (învățare) automata (Machine Learning - ML)

- Învățare supervizată → datele de antrenament (stări ale problemei) sunt deja etichetate cu elemente din  $E$ , iar datele de test trebuie etichetate cu una dintre etichetele din  $E$  pe baza unui model (învățat pe datele de antrenament) care face corespondența date-etichete. Nu există interacțiune cu mediul.
- Învățare nesupervizată → datele de antrenament (stări ale problemei) NU sunt etichetate, trebuie învățat un model de etichetare, iar apoi datele de test trebuie etichetate cu una dintre etichetele identificate de model. Nu există interacțiune cu mediul.
- Învățare cu întărire → datele de antrenament (stări ale problemei și acțiuni posibile) sunt etichetate cu măsuri de recompensă astfel încât să se poată optimiza recompensa primită prin efectuarea anumitor acțiuni. În plus, există interacțiune cu mediul prin intermediul acestor acțiuni.

### ■ Sisteme hibride

- O combinație de 2 sau mai multe sisteme inteligente
- Care se "ajută" între ele
  - În optimizarea diferitelor elemente sau algoritmi

# Cursul următor

---

## A. Scurtă introducere în Inteligența Artificială (IA)

## B. Rezolvarea problemelor prin căutare

- Definirea problemelor de căutare
- Strategii de căutare
  - Strategii de căutare neinformate
  - Strategii de căutare informate
  - Strategii de căutare locale (Hill Climbing, Simulated Annealing, Tabu Search, Algoritmi evolutivi, PSO, ACO)
  - Strategii de căutare adversială

## C. Sisteme inteligente

- Sisteme care învață singure
  - Arbori de decizie
  - Rețele neuronale artificiale
  - Mașini cu suport vectorial
  - Algoritmi evolutivi
  - Q-learning
- Sisteme hibride
- Sisteme bazate pe reguli în medii certe
- Sisteme bazate pe reguli în medii incerte (Bayes, factori de certitudine, Fuzzy)

---

□ Informațiile prezentate au fost colectate din diferite surse de pe internet, precum și din cursurile de inteligență artificială ținute în anii anteriori de către:

- Conf. Dr. Mihai Oltean –  
[www.cs.ubbcluj.ro/~moltean](http://www.cs.ubbcluj.ro/~moltean)
- Lect. Dr. Crina Groșan -  
[www.cs.ubbcluj.ro/~cgrosan](http://www.cs.ubbcluj.ro/~cgrosan)
- Prof. Dr. Horia F. Pop -  
[www.cs.ubbcluj.ro/~hfpop](http://www.cs.ubbcluj.ro/~hfpop)