

INTELIGENȚĂ ARTIFICIALĂ



Sisteme inteligente

Sisteme care învață singure

– kNN si programare genetică –

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

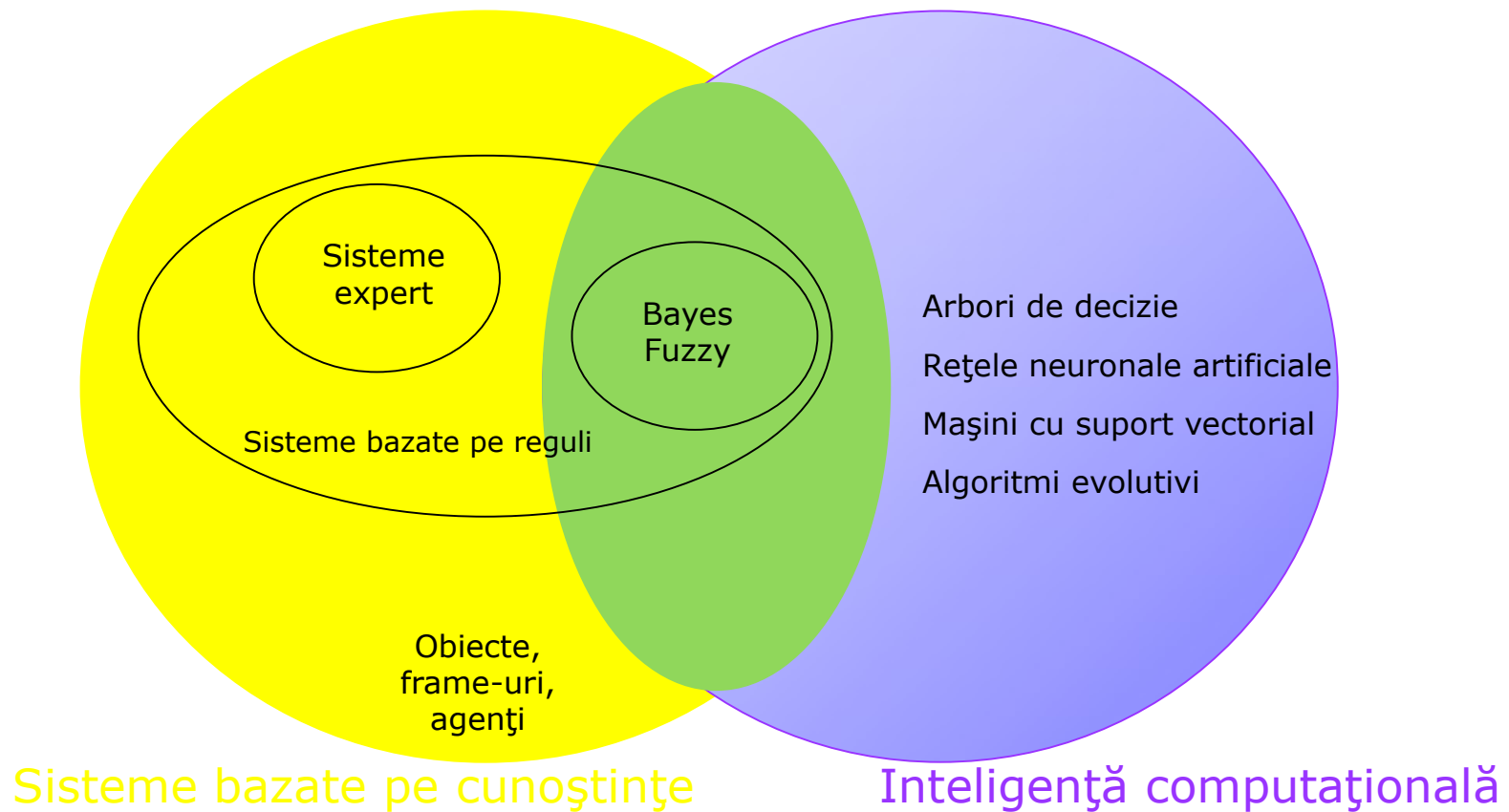
C. Sisteme inteligente

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

Materiale de citit și legături utile

- ❑ capitolul 15 din *C. Groșan, A. Abraham, Intelligent Systems: A Modern Approach, Springer, 2011*
- ❑ Capitolul 9 din *T. M. Mitchell, Machine Learning, McGraw-Hill Science, 1997*
- ❑ capitolul VI (18) din *S. Russell, P. Norvig, Artificial Intelligence: A Modern Approach, Prentice Hall, 1995*
- ❑ capitolul 10 și 11 din *C. Groșan, A. Abraham, Intelligent Systems: A Modern Approach, Springer, 2011*
- ❑ capitolul V din *D. J. C. MacKey, Information Theory, Inference and Learning Algorithms, Cambridge University Press, 2003*
- ❑ capitolul 3 din *T. M. Mitchell, Machine Learning, McGraw-Hill Science, 1997*

Sisteme inteligente



Sisteme inteligente – SIS – Învățare automată

□ Tipologie

■ În funcție de experiența acumulată în timpul învățării:

- SI cu învățare supervizată
- SI cu învățare nesupervizată
- SI cu învățare activă
- SI cu învățare cu întărire

■ În funcție de modelul învățat (algoritmul de învățare):

- Rețele neuronale artificiale
- **Mașini cu suport vectorial (MSV)**
- **Algoritmi evolutivi**
- **kNN**
- Arbori de decizie
- Modele Markov ascunse

Materiale de citit și legături utile

- ❑ capitolul VI (18) din *S. Russell, P. Norvig, Artificial Intelligence: A Modern Approach, Prentice Hall, 1995*
- ❑ capitolul 10 și 11 din *C. Groșan, A. Abraham, Intelligent Systems: A Modern Approach, Springer, 2011*
- ❑ capitolul V din *D. J. C. MacKey, Information Theory, Inference and Learning Algorithms, Cambridge University Press, 2003*
- ❑ capitolul 3 din *T. M. Mitchell, Machine Learning, McGraw-Hill Science, 1997*

Sisteme inteligente – SIS – Învățare automată

□ Tipologie

■ În funcție de experiența acumulată în timpul învățării:

- SI cu învățare supervizată
- SI cu învățare nesupervizată
- SI cu învățare activă
- SI cu învățare cu întărire

■ În funcție de modelul învățat (algoritmul de învățare):

- Arbori de decizie
- Rețele neuronale artificiale
- **Mașini cu suport vectorial (MSV)**
- Algoritmi evolutivi
- Modele Markov ascunse

Sisteme inteligente – SIS – MSV

- Mașini cu suport vectorial (MSV)
 - Definire
 - Tipuri de probleme rezolvabile
 - Avantaje
 - Dificultăți
 - Tool-uri

Sisteme inteligente – SIS – MSV

□ Definire

- Dezvoltate de Vapnik în 1970
- Popularizate după 1992
- Clasificatori liniari care identifică un hiperplan de separare a clasei pozitive de clasa negativă
- Au o fundamentare teoretică foarte riguroasă
- Funcționează foarte bine pentru date de volum mare (analiza textelor, analiza imaginilor)

■ Reamintim

- Problemă de învățare supervizată în care avem un set de date de forma:
 - (x^d, t^d) , cu:
 - $x^d \in \mathbf{R}^m \rightarrow x^d = (x^d_1, x^d_2, \dots, x^d_m)$
 - $t^d \in \mathbf{R} \rightarrow t^d \in \{1, -1\}$, $1 \rightarrow$ clasă pozitivă, $-1 \rightarrow$ clasă negativă
 - cu $d = 1, 2, \dots, n, n+1, n+2, \dots, N$
- Primele n date (se cunosc x^d și t^d) vor fi folosite drept bază de antrenament a MSV
- Ultimele $N-n$ date (se cunosc doar x^d , fără t^d) vor fi folosite drept bază de testare a MSV

Sisteme inteligente – SIS – MSV

□ Definire

- MSV găsește o funcție liniară de forma $f(\mathbf{x}) = \langle \mathbf{w} \cdot \mathbf{x} \rangle + b$, (\mathbf{w} -vector pondere) a.î.

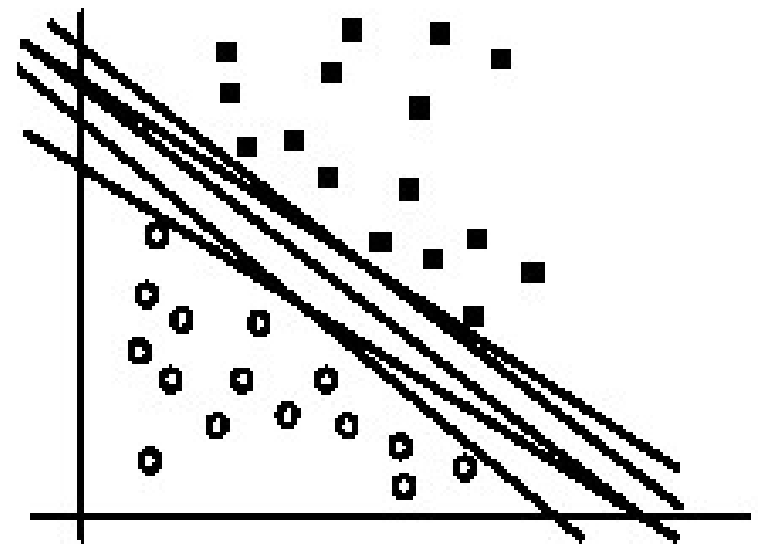
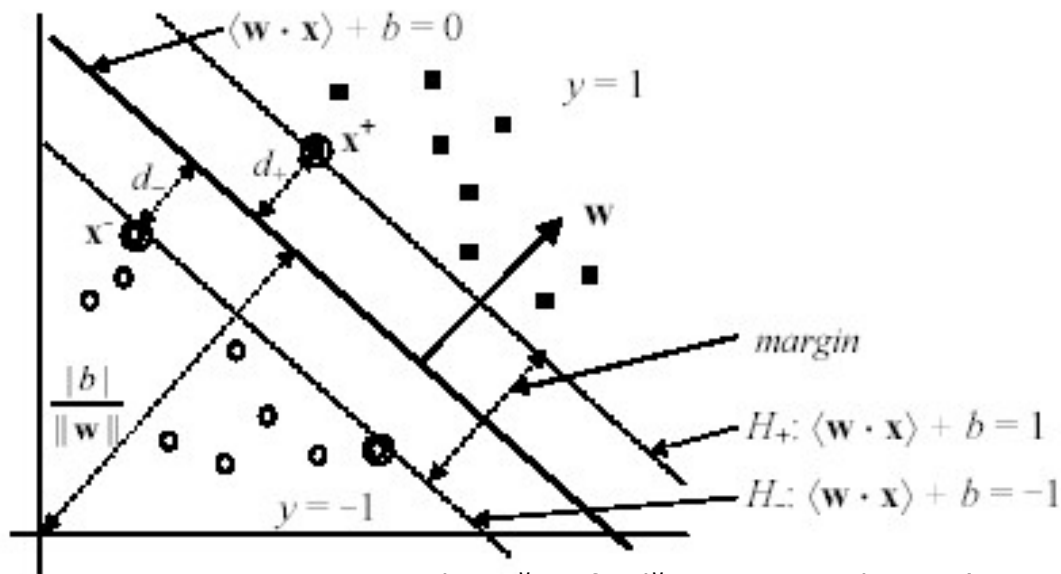
$$y_i = \begin{cases} 1 & \text{if } \langle \mathbf{w} \cdot \mathbf{x}_i \rangle + b \geq 0 \\ -1 & \text{if } \langle \mathbf{w} \cdot \mathbf{x}_i \rangle + b < 0 \end{cases}$$

- $\langle \mathbf{w} \cdot \mathbf{x} \rangle + b = 0 \rightarrow$ hiperplanul de decizie care separă cele 2 clase

Sisteme inteligente – SIS – MSV

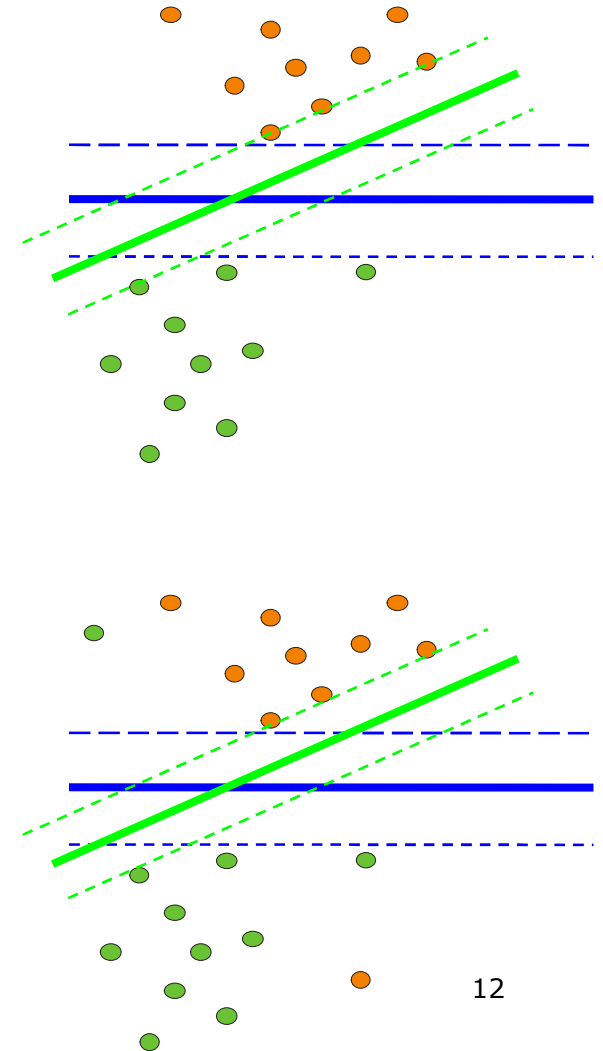
□ Definiere

- Pot exista mai multe hiperplane
 - Care este cel mai bun hiperplan?
- MSV caută hiperplanul cu cea mai largă margine (cel care micșorează eroarea de generalizare)
 - Algoritmul SMO (*Sequential minimal optimization*)



Sisteme inteligente – SIS – MSV

- Tipuri de probleme rezolvabile
 - Probleme de clasificare → Cazuri de date
 - Liniar separabile
 - Separabile
 - Eroarea = 0
 - Ne-separabile
 - Se relaxează constrângerile → se permit unele erori
 - C – coeficient de penalizare

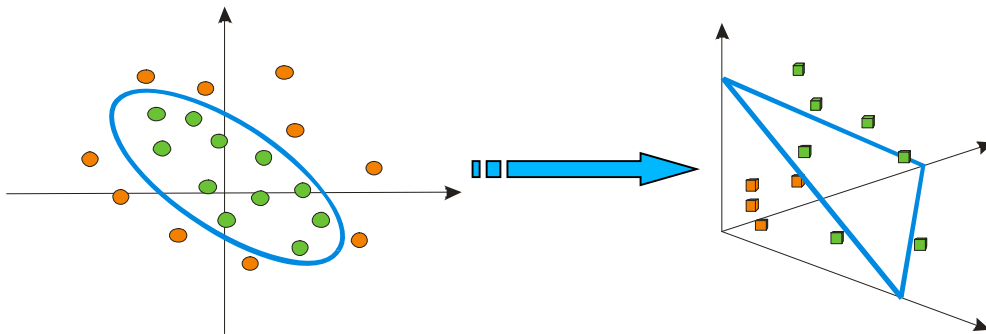


Sisteme inteligente – SIS – MSV

Cazuri de date

□ Non-liniar separabile

- Spațiul de intrare se transformă într-un spațiu cu mai multe dimensiuni (*feature space*), cu ajutorul unei funcții kernel, unde datele devin liniar separabile



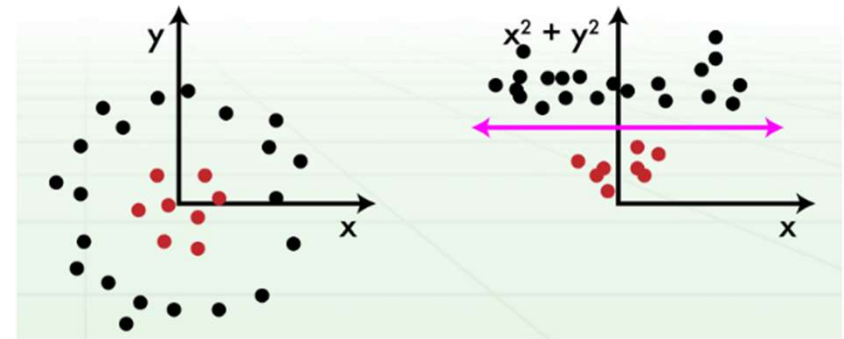
■ Kernele posibile

□ Clasice

- Polynomial kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = (\langle \mathbf{x}_1, \mathbf{x}_2 \rangle + 1)^d$
- RBF kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\sigma \|\mathbf{x}_1 - \mathbf{x}_2\|^2)$

□ Kernele multiple

- Liniare: $K(\mathbf{x}_1, \mathbf{x}_2) = \sum w_i K_i$
- Ne-liniare
 - Fără coeficienți: $K(\mathbf{x}_1, \mathbf{x}_2) = K_1 + K_2 * \exp(K_3)$
 - Cu coeficienți: $K(\mathbf{x}_1, \mathbf{x}_2) = K_1 + c_1 * K_2 * \exp(c_2 + K_3)$



Sisteme inteligente – SIS – MSV

□ Configurarea MSV

■ Parametrii unei MSV

□ Coeficientul de penalizare C

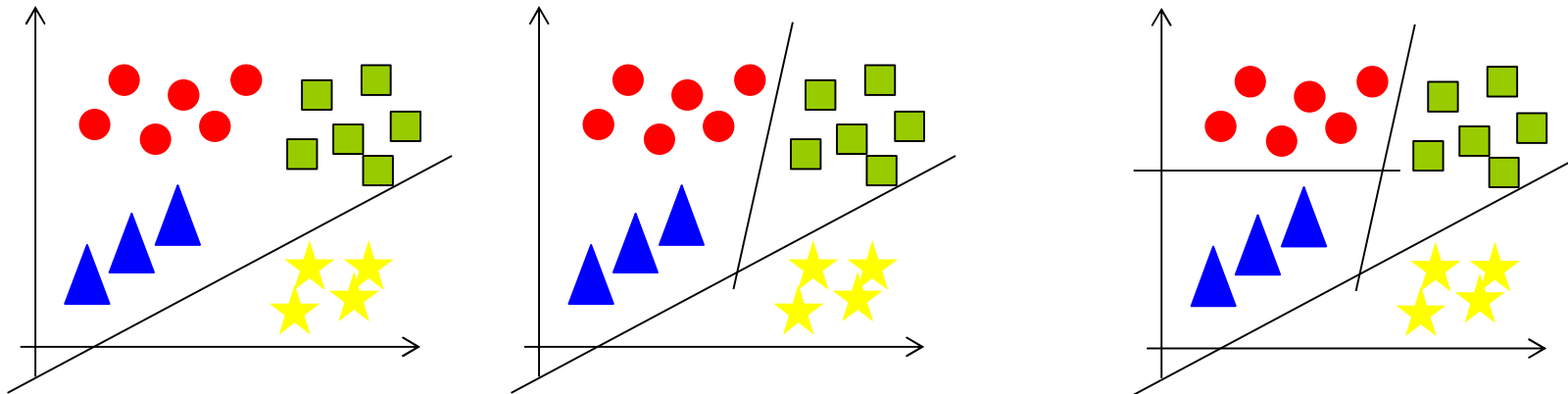
- C – mic → convergență lentă
- C – mare → convergență rapidă

□ Parametrii funcției kernel (care kernel și cu ce parametri)

- Dacă m (nr de attribute) este mult mai mare decât n (nr de instanțe)
 - MSV cu kernel liniar (MSV fără kernel) → $K(\mathbf{x}^{d1}, \mathbf{x}^{d2}) = \mathbf{x}^{d1} \cdot \mathbf{x}^{d2}$
- Dacă m (nr de attribute) este mare, iar n (nr de instanțe) este mediu
 - MSV cu kernel Gaussian $K(\mathbf{x}^{d1}, \mathbf{x}^{d2}) = \exp(-||\mathbf{x}^{d1} - \mathbf{x}^{d2}||^2 / 2\sigma^2)$
 - σ – dispersia datelor de antrenament
 - Attributele instanțelor trebuie normalizate (scalate la (0,1))
- m (nr de attribute) este mic, iar n (nr de instanțe) este mare
 - Se adaugă noi attribute, iar apoi
 - MSV cu kernel liniar

Sisteme inteligente – SIS – MSV

- MSV pentru probleme de clasificare supervizate cu mai mult de 2 clase
 - Una vs. restul (one vs. all)



Sisteme inteligente – SIS – MSV

□ MSV structurate

■ Învăţare automată

□ Normală $f: \mathcal{X} \rightarrow \mathbf{R}$

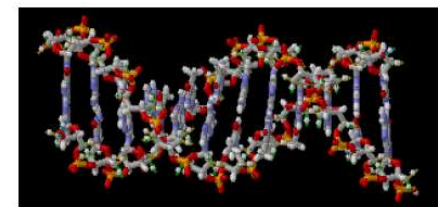
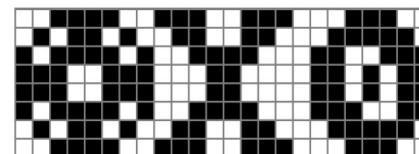
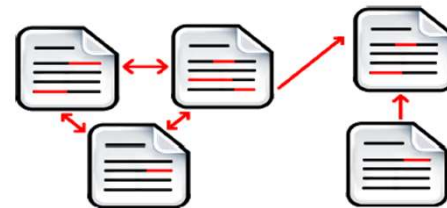
- Intrări de orice fel
- Ieşiri numerice (naturale, întregi, reale)

□ Structurată: $\mathcal{X} \rightarrow \mathcal{Y}$

- Intrări de orice fel
- Ieşiri de orice fel (simple sau structurate)

■ Informaţii structurate

- Texte şi hiper-texte
- Molecule şi structuri moleculare
- Imagini



Sisteme inteligente – SIS – MSV

□ MSV structurate

■ Aplicații

□ Procesarea limbajului natural

- Traduceri automate (ieșiri → propoziții)
- Analiza sintactică și/sau morfologică a propozițiilor (ieșiri → arborele sintactic și/sau morfologic)

□ Bioinformatică

- Predicția unor structuri secundare (ieșirile → grafe bi-partite)
- Predicția funcționării unor enzime (ieșirile → *path*-uri în arbori)

□ Procesarea vorbirii

- Transcrieri automate (ieșiri → propoziții)
- Transformarea textelor în voce (ieșiri → semnale audio)

□ Robotică

- Planificare (ieșirile → secvențe de acțiuni)

Sisteme inteligente – SIS – MSV

□ Avantaje

- Pot lucra cu orice fel de date (liniar separabile sau nu, distribuit uniform sau nu, cu distribuție cunoscută sau nu)
 - Funcțiile kernel care crează noi atribute (features) → straturile ascunse dintr-o RNA
- Dacă problema e convexă oferă o soluție unică → optimul global
 - RNA pot asocia mai multe soluții → optime locale
- Selectează automat mărimea modelului învățat (prin vectorii suport)
 - În RNA straturile ascunse trebuie configurate de către utilizator *a priori*
- Nu învață pe derost datele (overfitting)
 - RNA se confruntă cu problema overfitting-ului chiar și cand modelul se învață prin validare încrucișată

□ Dificultăți

- Doar atribute reale
- Doar clasificare binară
- Background matematic dificil

□ Tool-uri

- LibSVM → <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- Weka → SMO
- SVMLight → <http://svmlight.joachims.org/>
- SVM Torch → <http://www.torch.ch/>
- <http://www.support-vector-machines.org/>

Sisteme inteligente – SIS – Învățare automată

□ Programare genetică

- Definire
- Proiectare
- Avantaje
- Limite
- Versiuni

Sisteme inteligente – SIS – PG

Reamintim

- Învățare supervizată → problemă de regresie (Studiul legăturii între variabile)
 - Se dă un set de n date (exemple, instanțe, cazuri)
 - date de antrenament – sub forma unor perechi ($attribute_data_i, ieșire_i$), unde
 - $i = 1, n$ (n = nr datelor de antrenament)
 - **$attribute_data_i = (atr_{i1}, atr_{i2}, \dots, atr_{im})$** , m – nr atributelor (caracteristicilor, proprietăților) unei date
 - **$ieșire_i$ – un număr real**
 - date de test
 - sub forma (**$attribute_data_i$**), $i = n+1, N$ ($N-n$ = nr datelor de test)
 - Să se determine
 - o funcție (necunoscută) care realizează corespondența atribut – ieșire pe datele de antrenament
 - Ieșirea (valoarea) asociată unei date (noi) de test folosind funcția învățată pe datele de antrenament
- Cum găsim forma (expresia) funcției?
 - Algoritmi evolutivi → Programare genetică

Sisteme inteligente – SIS – PG

Reamintim

- Algoritmi evolutivi
 - Inspirați din natură (biologie)
 - Iterativi
 - Bazați pe
 - populații de potențiale soluții
 - căutare aleatoare ghidată de
 - Operații de selecție naturală
 - Operații de încrucișare și mutație
 - Care procesează în paralel mai multe soluții
- Metafora evolutivă

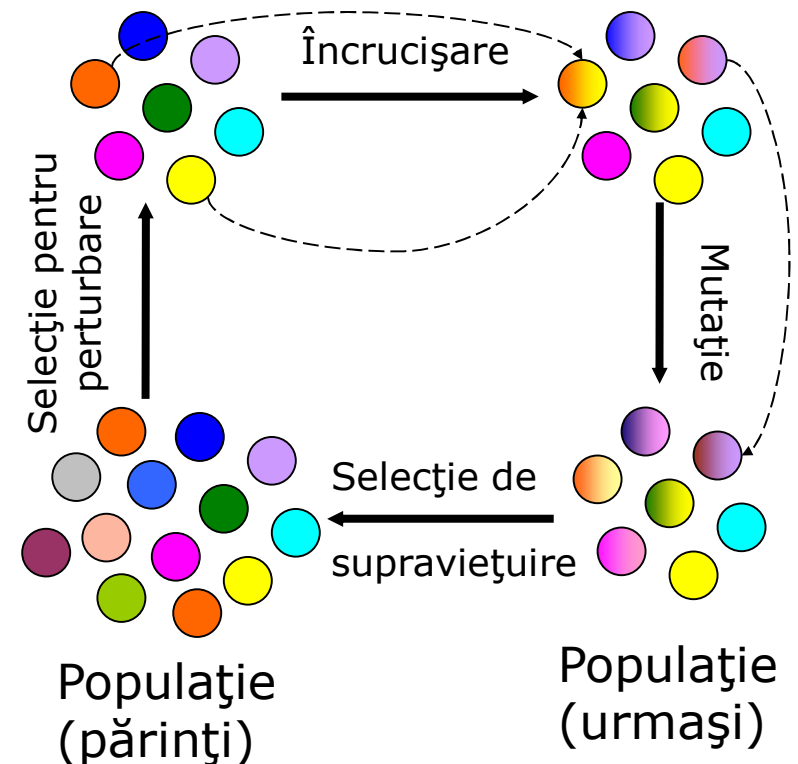
Evoluție naturală	Rezolvarea problemelor
Individ	Soluție potențială (candidat)
Populație	Mulțime de soluții
Cromozom	Codarea (reprezentarea) unei soluții
Genă	Parte a reprezentării
Fitness (măsură de adaptare)	Calitate
Încrucișare și mutație	Operații de căutare
Inteligință artificială - Sisteme inteligente (SVM, GP, kNN, arbori de decizie)	
Mediu	Spațiul de căutare al problemei

Sisteme inteligente – SIS – PG

Reamintim

□ Algoritmi evolutivi

```
Initializare populație P(0)
Evaluare P(0)
g := 0; //generația
CâtTimp (not condiție_stop) execută
    Repetă
        Selectează 2 părinți p1 și p2 din P(g)
        Încrucișare(p1,p2) => o1 și o2
        Mutație(o1) => o1*
        Mutație(o2) => o2*
        Evaluare(o1*)
        Evaluare(o2*)
        adăugare o1* și o2* în P(g+1)
    Până când P(g+1) este completă
    g := g + 1
Sf CâtTimp
```



Sisteme inteligente – SIS – PG

□ Definire

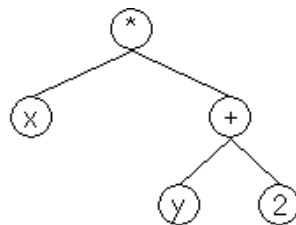
- Propusă de John Koza în 1988
- <http://www.genetic-programming.org/>
- Un tip particular de algoritmi evolutivi
- Cromozomi
 - sub formă de arbore care codează mici programe
- Fitness-ul unui cromozom
 - Performanța programului codat în el
- Scopul PG
 - Evoluarea de programe de calculator
 - AG evoluează doar soluții pentru probleme particulare

Sisteme inteligente – SIS – PG

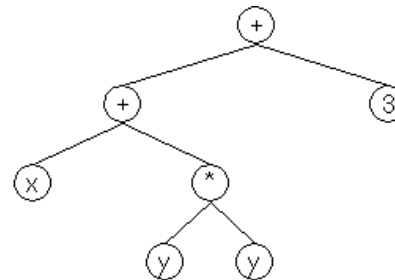
□ Proiectare

■ Reprezentarea cromozomilor

- Foarte importantă, dar este o sarcină dificilă
- Cromozomul = un arbore cu noduri de tip
 - Funcție → operatori matematici (+, -, *, /, sin, log, if, ...)
 - Terminal → attribute ale datelor problemei sau constante (x, y, z, a, b, c, ...)
- care codează expresia matematică a unui program (problema regresiei → a unei funcții)



$x(y+2)$



$x+y^2+3$

Sisteme inteligente – SIS – PG

□ Proiectare

■ Fitness

- Eroarea de predicție – diferența între ceea ce dorim să obținem și ceea ce obținem de fapt
- pp o problemă de regresie cu următoarele date de intrare (2 atribute și o ieșire) și 2 cromozomi:
 - $c_1 = 3x_1 - x_2 + 5 = (((3, x_1, *), x_2, -), 5, +)$
 - $c_2 = 3x_1 + 2x_2 + 2 \quad f^*(x_1, x_2) = 3x_1 + 2x_2 + 1$ – necunoscută

x_1	x_2	$f^*(x_1, x_2)$	$f_1(x_1, x_2)$	$f_2(x_1, x_2)$	$ f^* - f_1 $	$ f^* - f_2 $
1	1	6	7	7	1	1
0	1	3	4	4	1	1
1	0	4	8	5	4	1
-1	1	0	1	1	1	1
					$\Sigma = 7$	$\Sigma = 4$

→ c_2 e mai bun
ca c_1

Sisteme inteligente – SIS – PG

□ Proiectare

■ Fitness

- Eroarea de predicție – diferența între ceea ce dorim să obținem și ceea ce obținem de fapt

- pp o problemă de clasificare cu următoarele date de intrare (2 atribute și o ieșire) și 2 cromozomi:

- $c_1 = 3x_1 - x_2 + 5$ $f_1(1,1) = 7$ -: $\text{sigm}(7) - (0,1)$ -: *Theta*

- $c_2 = 3x_1 + 2x_2 + 2$

x_1	x_2	$f^*(x_1, x_2)$	$f_1(x_1, x_2)$	$f_2(x_1, x_2)$	$ f^* - f_1 $	$ f^* - f_2 $
1	1	Yes	Yes	Yes	0	0
0	1	No	Yes	No	1	0
1	0	Yes	No	No	1	1
-1	1	Yes	No	yes	1	0
					$\Sigma=3$	$\Sigma=1$

→ c_2 e mai bun
ca c_1

Sisteme inteligente – SIS – PG

□ Proiectare

■ Inițializarea cromozomilor

- Generare aleatoare de arbori corecți → programe valide (expresii matematice valide)
- Se stabilește o adâncime maximă a arborilor D_{\max}

□ 3 metode de inițializare

- *Full* → fiecare ramură a rădăcinii are adâncimea D_{\max}
 - Nodurile aflate la o adâncime $d < D_{\max}$ se inițializează cu una dintre funcțiile din F
 - Nodurile aflate la o adâncime $d = D_{\max}$ se inițializează cu unul dintre terminalele din T
- *Grow* → fiecare ramură a rădăcinii are o adâncime $< D_{\max}$
 - Nodurile aflate la o adâncime $d < D_{\max}$ se inițializează cu un element din $F \cup T$
 - Nodurile aflate la o adâncime $d = D_{\max}$ se inițializează cu unul dintre terminalele din T
- *Ramped half and half* → $\frac{1}{2}$ din populația de cromozomi se inițializează folosind metoda *full*, $\frac{1}{2}$ din populația de cromozomi se inițializează folosind metoda *grow*

Sisteme inteligente – SIS – PG

□ Proiectare

- Operatori genetici → Selecția pentru recombinare
 - similar oricărui algoritm evolutiv
 - recomandare → selecție proporțională
 - over-selection → pentru populații f mari
 - Se ordonează populația pe baza fitness-ului și se împarte în 2 grupuri:
 - Grupul 1: cei mai buni x% cromozomi din populație
 - Grupul 2: restul de (100-x)% cromozomi din populație
 - Pentru populații cu 1000, 2000, 4000, 8000 de cromozomi, x este stabilit la 32%, 16%, 8%, respectiv 4%
 - 80% din operațiile de selecție vor alege cromozomi din grupul 1, 20% din grupul 2

Sisteme inteligente – SIS – PG

□ Proiectare

■ Operatori genetici → Selecția de supraviețuire

□ Scheme

- Generațională
- steady-state

□ Probleme

- *Bloat* → supraviețuirea celui mai “gras” individ (dimensiunea cromozomilor crește de-a lungul evoluției)
- Soluții
 - Interzicerea operatorilor de variație care produc descendenți prea mari
 - Presiunea economiei (zgârceniei) – penalizarea cromozomilor prea mari

Sisteme inteligente – SIS – PG

□ Proiectare

■ Operatori genetici → Încrucișare și mutație

□ Parametri

- O probabilitate p de alegere între încrucișare și mutație
 - $p = 0$ (cf. Koza) sau $p = 0.05$ (cf. Banzhaf)
- O probabilitate p_c și respectiv p_m de stabilire a nodului care urmează a fi supus modificării

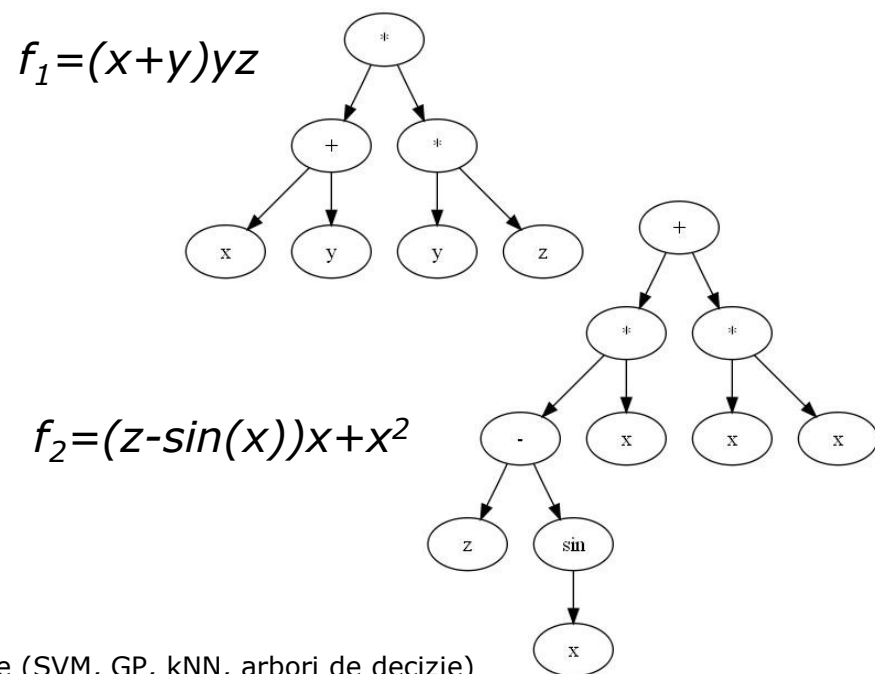
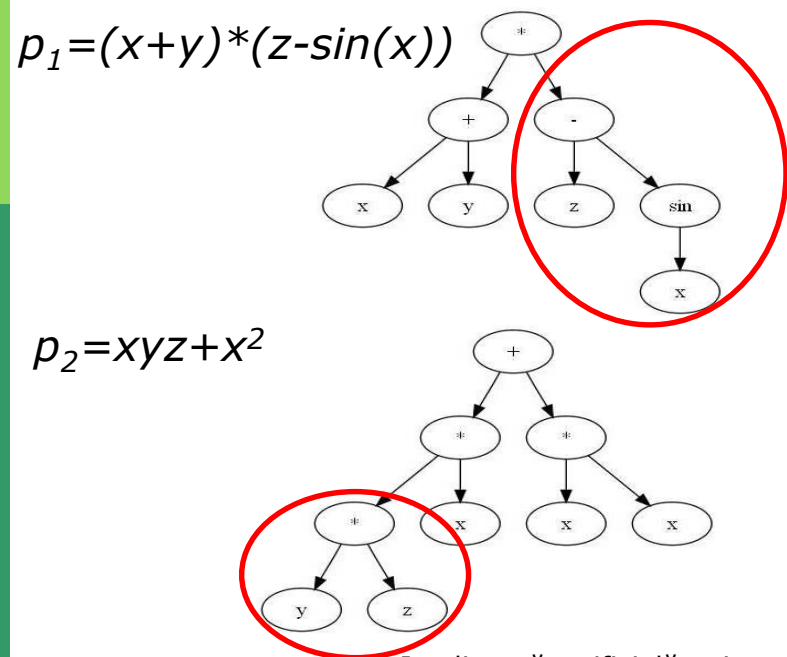
□ Dimensiunea descendenților diferă de dimensiune părinților

Sisteme inteligente – SIS – PG

Proiectare

Operatori genetici → Încrucișare

- Cu punct de tăietură – se interchimbă doi sub-arbori
- Punctul de tăietură se generează aleator



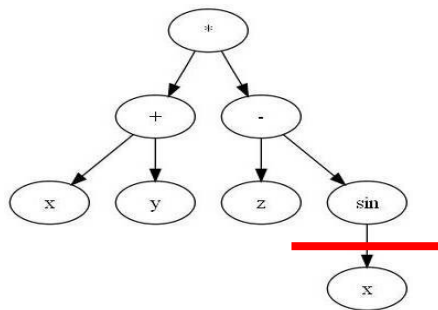
Sisteme inteligente – SIS – PG

□ Proiectare

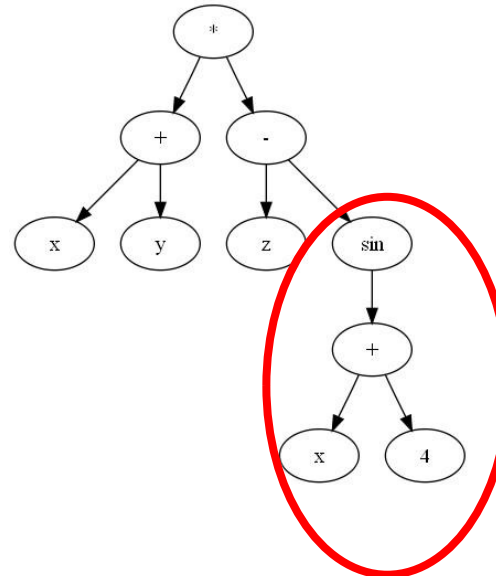
■ Operatori genetici → Mutație

- Mutație de tip *grow* → Înlocuirea unei frunze cu un nou sub-arbore

$$p = (x + y) * (z - \sin(x))$$



$$f = (x + y) * (z - \sin(x + 4))$$



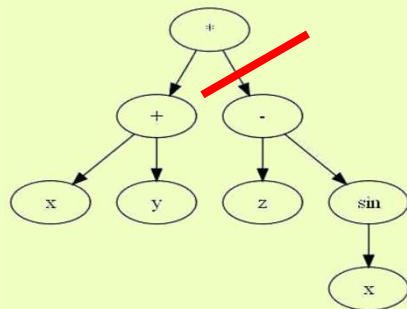
Sisteme inteligente – SIS – PG

□ Proiectare

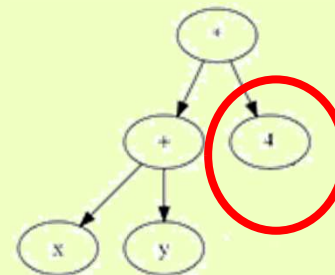
■ Operatori genetici → Mutație

- Mutație de tip *shrink* → Înlocuirea unui sub-arbore cu o frunză

$$p = (x + y) * (z - \sin(x))$$



$$f = (x + y) * 4$$



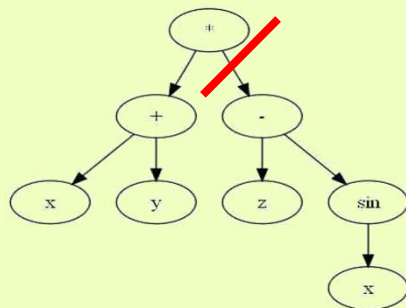
Sisteme inteligente – SIS – PG

□ Proiectare

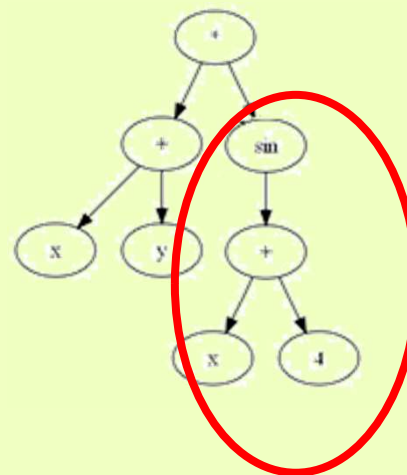
■ Operatori genetici → Mutație

- Mutație de tip *Koza* → Înlocuirea unui nod (intern sau frunză) cu un nou sub-arbore

$$p = (x + y) * (z - \sin(x))$$



$$f = (x + y) * \sin(x + 4)$$



Sisteme inteligente – SIS – PG

□ Proiectare

■ Operatori genetici → Mutație

□ Mutație de tip *switch*

- selectarea unui nod intern și re-ordonarea sub-arborilor săi

□ Mutație de tip *cycle*

- selectarea unui nod și înlocuirea lui cu unul nou de același tip (intern – cu o funcție – sau frunză – cu un terminal)

Sisteme inteligente – SIS – PG

□ comparație AG și PG

■ Forma cromozomilor

- AG – cromozomi liniari
- PG – cromozomi ne-liniari

■ Dimensiunea cromozomilor

- AG – fixă
- PG – variabilă (în adâncime sau lățime)

■ Schema de creare a descendenților

- AG – încrucișare și mutație
- PG – încrucișare sau mutație

Sisteme inteligente – SIS – PG

□ Avantaje

- PG găsește soluții problemelor care nu au o soluție optimă
 - Un program pentru conducerea mașinii → nu există o singură soluție
 - Unele soluții implică un condus sigur, dar lent
 - Alte soluții implică o viteză mare, dar un risc ridicat de accidente
 - Conducerea mașinii \leftrightarrow compromis între viteză mare și siguranță
- PG este utilă în problemele a căror variabile se modifică frecvent
 - Conducerea mașinii pe autostradă
 - Conducerea mașinii pe un drum forestier

□ Limite

- Timpul mare necesar evoluției pentru identificarea soluției

Sisteme inteligente – SIS – PG

□ Versiuni ale PG

- PG liniară (Cramer, Nordin)
- Gene Expression Programming (Ferreira)
- Multi Expression Programing (Oltean)
- Gramatical Evolution (Ryan, O'Neill)
- Cartesian Genetic Programming (Miller)

Sisteme inteligente – SIS – PG

□ PG liniară

- Evoluarea de programe scrise într-un limbaj imperativ (calculul fitness-ului nu necesită interpretare) → viteză mare de lucru
- Reprezentare
 - Vector de instrucțiuni, fiecare instrucțiune fiind de forma (pp ca aritatea maximală a unei funcții din F este n):
 - $\text{Index_op, registru_out, registru_in}_1, \text{registru_in}_2, \dots, \text{registru_in}_n$

- $v_i = v_j * v_k$ // instruction operating on two registers
- $v_i = v_j * c$ // instruction operating on one register and one constant
- $v_i = \sin(v_j)$ // instruction operating on one register

```
void LGP_program (double v[11])
{
    ...
    v[8] = v[0] - 10;
    v[6] = v[2] * v[0];
    v[5] = v[8] * 7;
    v[4] = v[2] - v[0];
    v[10] = v[1]/v[4];
    v[3] = sin(v[1]);
    v[1] = v[8] - v[6];
    v[7] = v[10] * v[3];
    v[9] = v[0] + v[7];
    v[2] = v[7] + 3;
    ...
}
```

```
void LGP_effective_program (double v[11])
{
    ...
    v[4] = v[2] - v[0];
    v[10] = v[1]/v[4];
    v[3] = sin(v[1]);
    v[7] = v[10] * v[3];
    v[9] = v[0] + v[7];
    ...
}
```

Sisteme inteligente – SIS – PG

□ PG liniară

■ Inițializare

- Aleatoare
- Restricții
 - Lungimea inițială a cromozomului (nr de instrucțiuni)

■ Operatori genetici de variație

- Încrucișare – cu 2 puncte de tăietură
- Mutație
 - Micro mutație → schimbarea unui operand sau operator (nu se modifică dimensiunea cromozomului)
 - Macro mutație → inserarea sau eliminarea unei instrucțiuni (se modifică dimensiunea cromozomului)

Sisteme inteligente – SIS – PG

□ PG liniară

■ Avantaje

- Evoluare într-un limbaj de nivel redus (low-level)

■ Dezavantaje

- Numărul de regiștri necesari (numărul de attribute ale problemei)

■ Resurse

- Register Machine Learning Technologies <http://www.aimlearning.com>
- Peter Nordin's home page <http://fy.chalmers.se/~pnordin>
- Wolfgang Banzhaf's home page <http://www.cs.mun.ca/~banzhaf>
- Markus Brameier's home page <http://www.daimi.au.dk/~brameier>

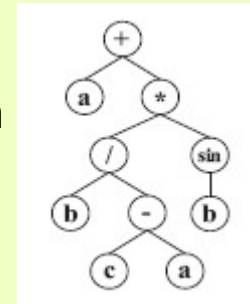
Sisteme inteligente – SIS – PG

■ Gene Expression Programming (GEP)

■ Ideea de bază

- Reprezentarea liniară a expresiilor codabile în arbori (prin parcurgerea în lățime a acestora – breadth-first)

$$C = +a * /Sb - bcacabbc$$



■ Reprezentare

- Un cromozom este format din mai multe gene
 - Legate între ele prin + sau *
- Fiecare genă este formată din:
 - Cap
 - conține h elemente → funcții și terminale
 - Coadă
 - conține doar terminale, în număr de $t = (n-1)*h+1$, unde n – aritatea maximă a unei funcții din F

Sisteme inteligente – SIS – PG

□ GEP

■ Inițializare

- Aleatoare, cu elemente din F și T conform regulilor precizate anterior

■ Operatori de variație

□ Încrucișare

- La nivel de alelă
 - Cu un punct de tăietură
 - Cu două puncte de tăietură
- Încrucișare la nivel de genă
 - Cromozomii schimbă între ei anumite gene (plasate pe aceeași poziție)

□ Mutație

- La nivel de alelă
 - se modifică un element din cap sau coadă, respectând regulile de la inițializare

□ Transpoziții

Sisteme inteligente – SIS – PG

□ GEP

■ Avantaje

- Codarea în cromozomi a unor programe corecte datorită separării unei gene în cap și coadă

■ Dezavantaje

- Cromozomii multi-genă
 - Câte gene?
 - Cum se leagă genele între ele?

■ Resurse

- Gene Expression Programming website, <http://www.gepsoft.com>
- Heitor Lopes's home page <http://www.cpgei.cefetpr.br/~hslopes/index-english.html>
- Xin Li's home page <http://www.cs.uic.edu/~xli1>
- GEP in C# <http://www.c-sharpcorner.com/Code/2002/Nov/GEPAlgorithm.asp>

Sisteme inteligente – SIS – PG

□ Multi Expression Programming (MEP)

■ Ideea de bază

- Cromozomul este format din mai multe gene, fiecare genă → cod cu 3 adrese
 - Similar PG liniară, dar mai rapid

■ Reprezentare

- Liniară
- O genă conține o funcție (unară sau binară) și pointeri spre argumentele sale
- Cromozomul codează mai multe potențiale soluții → fiecare soluție corespunde unei gene
 - Calitatea unei soluții (gene) = suma (peste datele de antrenament) între ceea ce trebuia obținut și ceea ce se obține
 - Calitatea unui cromozom = fitness-ul celei mai bune gene

Sisteme inteligente – SIS – PG

□ MEP

■ Inițializare

- Prima genă trebuie să fie un terminal
- Restul genelor pot conține
 - un terminal sau
 - o funcție (unară sau binară) și pointeri spre argumentele sale
 - Argumentele unei funcții poziționate în a i-a genă trebuie să fie poziționate în cromozom la indici mai mici decât i

■ Operatori de variație

- Încrucișare → schimbarea unor gene între părinți
 - Cu un punct de tăietură
 - Cu două puncte de tăietură
 - Uniformă
- Mutație → modificarea unei gene
 - Prima genă → generarea unui nou terminal
 - Restul genelor → generarea unui terminal sau a unei funcții (simbolul funcției și argumentele funcției)
 - Generarea are loc la fel ca la inițializare

Sisteme inteligente – SIS – PG

□ MEP

■ Avantaje

- Ieșire dinamică corespunzătoare unui cromozom
 - Complexitatea programului (expresiei) căutat(e)
 - Programe (expresii) de lungime variabilă obținute fără operatori speciali
 - Programe de lungime exponențială codate în cromozomi de lungime polinomială

■ Dezavantaje

- Complexitatea decodării pt date de antrenament necunoscute → evoluarea strategiilor de joc

■ Resurse

- Mihai Oltean's home page <http://www.cs.ubbcluj.ro/~>
- Crina Groșan's home page <http://www.cs.ubbcluj.ro/~cgrosan>
- MEP web page <http://www.mep.cs.ubbcluj.ro>
- MEP in C# <http://www.c-sharpcorner.com>

Sisteme inteligente – SIS – PG

□ Grammatical Evolution (GE)

■ Ideea de bază

- Evoluarea de programe în forma Backus-Naur (program exprimat sub forma unei gramatici cu simboluri terminale și non-terminale, simbol de start, reguli/producții)

■ Reprezentare

- String binar de codons (grupuri de 8 biți) → care regulă a gramaticii tb aplicată
- Exemplu

- $G = \{N, T, S, P\}$, $N = \{+, -, *, /, \sin, (,)\}$, $T = \{\text{expr}, \text{op2}, \text{op1}\}$, $S = \langle \text{expr} \rangle$, iar P este:

- $\langle \text{expr} \rangle ::= a|b|c| \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle | (\langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle) | \langle \text{op1} \rangle \langle \text{expr} \rangle$

- $\langle \text{op2} \rangle ::= +|-|*|/$,

- $\langle \text{op1} \rangle ::= \sin$

- $C_{GE}^* = (9\ 12\ 12\ 3\ 15\ 7\ 11\ 4\ 2\ 5\ 0\ 6\ 11\ 0\ 1\ 7\ 12)$

- $S = \langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle \rightarrow a \langle \text{op2} \rangle \langle \text{expr} \rangle \rightarrow a + \langle \text{expr} \rangle \rightarrow a + \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle \rightarrow a + \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle \rightarrow a + b \langle \text{op2} \rangle \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle$

$C_{GE} = (00001001\ 00001100\ 00001100\ 00000011\ 00001111\ 00000111\ 00001011\ 00000100\ 00000010\ 00000101\ 00000000\ 00000110\ 00001011\ 00000000\ 00000001\ 00000111\ 00001100)$

$a + b / \langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle$
 $a + b / (\langle \text{expr} \rangle \langle \text{op2} \rangle \langle \text{expr} \rangle) \langle \text{op2} \rangle \langle \text{expr} \rangle$
 $a + b / (c \langle \text{op2} \rangle \langle \text{expr} \rangle) \langle \text{op2} \rangle \langle \text{expr} \rangle$
 $a + b / (c - \langle \text{expr} \rangle) \langle \text{op2} \rangle \langle \text{expr} \rangle$
 $a + b / (c - a) \langle \text{op2} \rangle \langle \text{expr} \rangle$
 $a + b / (c - a) * \langle \text{expr} \rangle$
 $a + b / (c - a) * \langle \text{op1} \rangle \langle \text{expr} \rangle$
 $a + b / (c - a) * \sin \langle \text{expr} \rangle$

$$E = a + b / (c - a) * \sin(b)$$

Sisteme inteligente – SIS – PG

□ GE

■ Inițializare

- Stringul binar este inițializat aleator cu 0 și 1 fără restricții → programe valide
- Decodarea se termină când s-a obținut un program complet
 - Dacă s-au terminat codons și încă nu s-a format tot programul, se reiau codons de la primul element → wrapping

■ Operatori de variație

- Încrucișare
 - Cu un punct de tăietură
- Mutație
 - Schimbarea probabilistică a unui bit în opusul său
- Duplicare
 - O secvență de gene este copiată la sfârșitul cromozomului
- Pruning
 - Eliminarea genelor ne-folosite în procesul de transformare (decodare) a cromozomului

Sisteme inteligente – SIS – PG

□ GE

■ Avantaje

- Evoluarea de programe scrise în limbaje a căror instrucțiuni pot fi exprimate ca reguli de tip BNF
- Reprezentarea poate fi schimbată prin modificarea gramaticii

■ Dezavantaje

- Wrapping-ul la infinit → limitarea repetărilor și penalizarea cromozomilor care depășesc un anumit prag de repetări

■ Resurse

- Grammatical Evolution web page, <http://www.grammatical-evolution.org>
- Conor Ryan's home page, <http://www.csis.ul.ie/staff/conorryan>
- Michael O'Neill's home page, <http://ncra.ucd.ie/members/oneillm.html>
- John James Collins's home page, <http://www.csis.ul.ie/staff/jjcollins>
- Maarten Keijzer's home page, <http://www.cs.vu.nl/~mkeijzer>
- Anthony Brabazon's home page <http://ncra.ucd.ie/members/brabazont.html>

Sisteme inteligente – SIS – PG

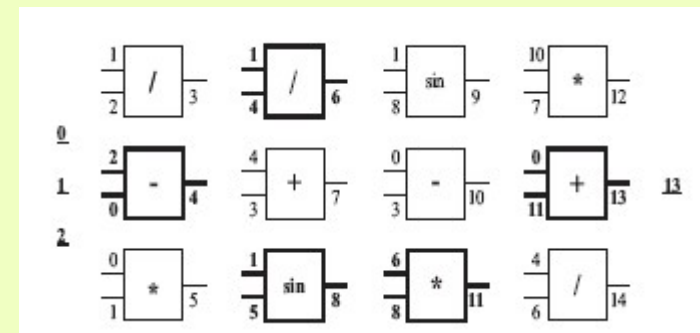
□ Cartesian Genetic Programming (CGP)

■ Ideea de bază

- Cromozomi sub formă de graf (matrice) → programe mai complexe decât cele din arbori

■ Reprezentare

- În sistem cartezian (matrice de noduri)
- Un nod are asociate
 - O funcție
 - Intrări
 - Ieșiri
- Ouputul cromozomului
 - Outputul oricărui nod



$C = (1, 2, 3, 2, 0, 1, 0, 1, 2, 1, 4, 3, 4, 3, 0, 1, 5, 4, 1, 8, 4, 0, 3, 1, 6, 8, 2, 10, 7, 2, 0, 11, 0, 4, 6, 3, 13)$

Sisteme inteligente – SIS – PG

□ CGP

■ Inițializare

- Aleatoare
- Intrările oricărui nod trebuie să fie noduri de pe coloanele anterioare
 - Nodurile de pe prima coloană au ca intrări caracteristicile datelor de antrenament

■ Operatori de variație

- Încrucișare
 - Nu se aplică
- Mutație
 - Modificarea elementelor unui nod

Sisteme inteligente – SIS – PG

□ CGP

■ Avantaje

- Evoluarea indicelui nodului care furnizează ieșirea programului codat în cromozom
- Programul evoluat poate avea una sau mai multe ieșiri

■ Dezavantaje

- Stabilirea numărului de coloane influențează rezultatele obținute

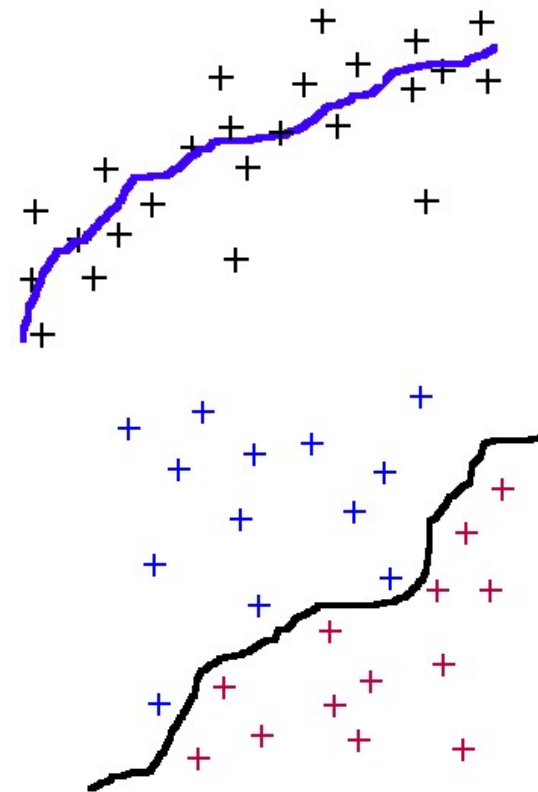
■ Resurse

- Julian. F. Miller's home page
<http://www.elec.york.ac.uk/intsys/users/jfm7>
- Lukás Sekanina's home page <http://www.fit.vutbr.cz/~sekanina/>

Sisteme inteligente – SIS – PG

□ Aplicații

- Probleme în care există o relație între intrări și ieșiri
- Probleme de regresie
- Probleme de clasificare

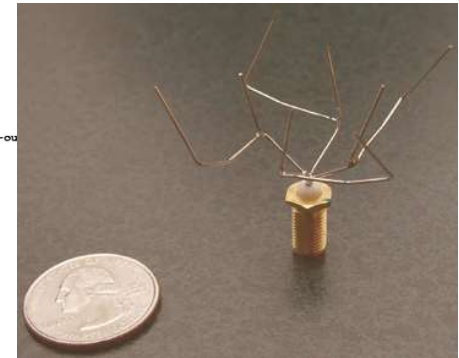
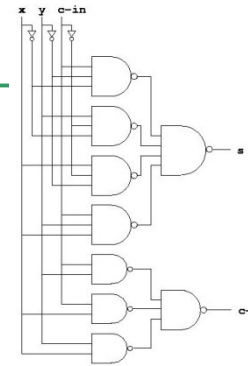


Sisteme inteligente – SIS – PG

□ Aplicații

■ Probleme de design

□ Evoluarea de circuite digitale



□ Evoluarea de antene

- http://idesign.ucsc.edu/projects/evo_antenna.html

□ Evoluarea de programe (scrise într-un anumit limbaj)

□ Evoluarea de picturi și muzică

- <http://www.cs.vu.nl/~gusz/>



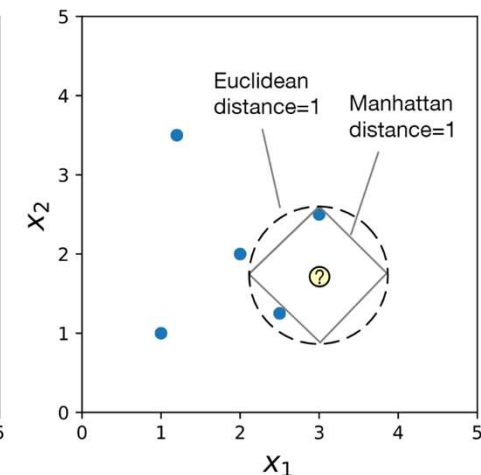
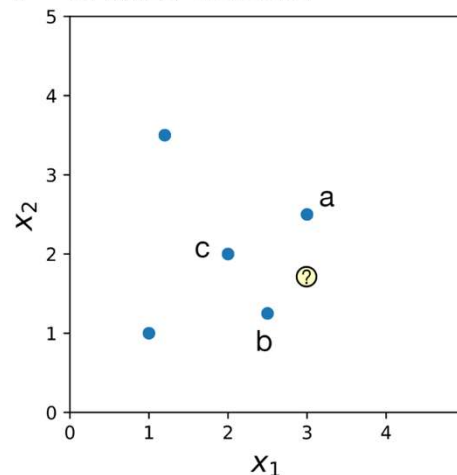
□ Altele

- <http://www.genetic-programming.com/humancompetitive.html>

Sisteme inteligente – SIS - kNN

Cei mai apropiați k vecini (*k-nearest neighbours - kNN*)

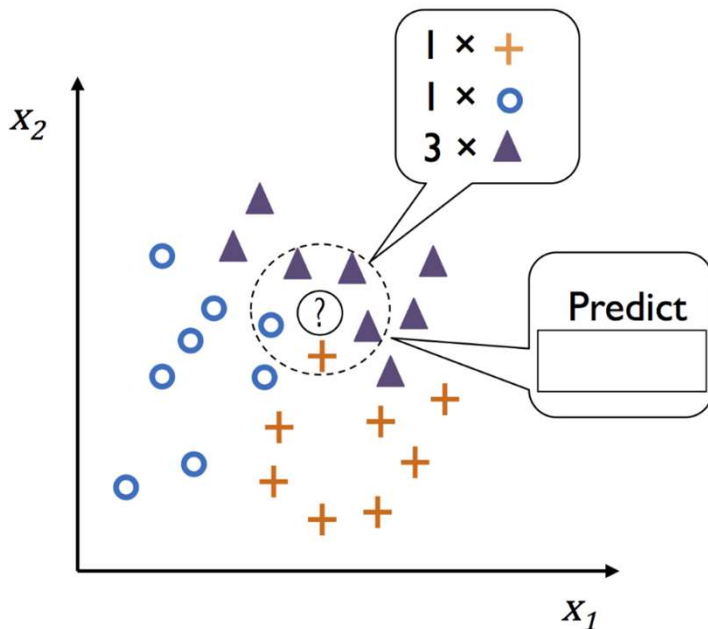
- ❑ Unul dintre cei mai simpli algoritmi de clasificare
- ❑ În etapa de antrenament, algoritmul doar citește datele de intrare (atributele și clasa fiecărei instanțe)
- ❑ În etapa de testare, pentru o nouă instanță (fără clasă):
 - se caută (printre instanțele de antrenament) **cei mai apropiați k vecini**
 - ❑ distanța Minkowski (Manhattan, Euclidiană) – attribute continue
 - ❑ distanța Hamming, Levensthein – analiza textelor
 - ❑ alte distanțe (funcții kernel)



- se preia clasa majoritară a acestor k vecini

Sisteme inteligente – SIS - kNN

- Unul dintre cei mai simpli algoritmi de clasificare
- În etapa de antrenament, algoritmul doar citește datele de intrare (atributele și clasa fiecărei instanțe)
- În etapa de testare, pentru o nouă instanță (fără clasă)
 - se caută (printre instanțele de antrenament) cei mai apropiați k vecini și
 - **se preia clasa majoritară a acestor k vecini**



o o o o + + + + + +

Majority vote: +

Plurality vote: +

o o o + + + ▲ ▲ ▲ ▲

Majority vote: None

Plurality vote: ▲

Sisteme inteligente – SIS - kNN

□ Tool-uri

■ Sklearn (python)

- <https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>

■ Weka (java)

- <https://weka.sourceforge.io/doc.dev/weka/classifiers/lazy/IBk.html>

□ Biblio

- https://github.com/rasbt/stat479-machine-learning-fs19/tree/master/02_knn

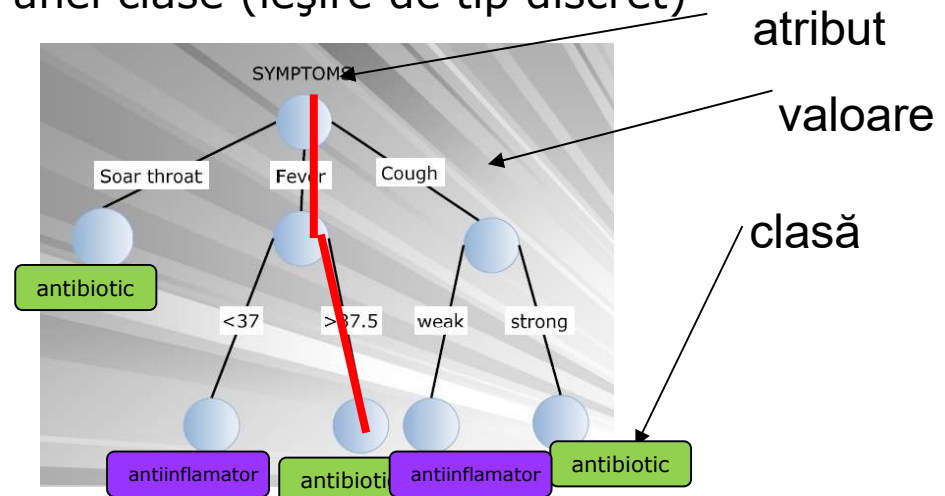
Sisteme inteligente – SIS – Arbori de decizie

□ Scop

- Divizarea unei colecții de articole în seturi mai mici prin aplicarea succesivă a unor reguli de decizie → adresarea mai multor întrebări
 - Fiecare întrebare este formulată în funcție de răspunsul primit la întrebarea precedentă
- Elementele se caracterizează prin informații non-metrice

□ Definire

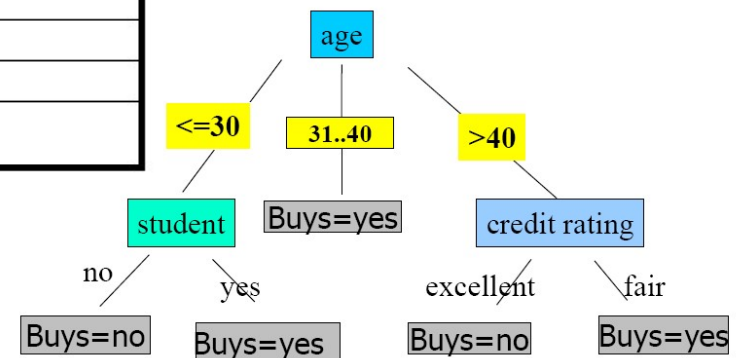
- Fiecare nod intern corespunde unui atribut
- Fiecare ramură de sub un nod (atribut) corespunde unei valori a atributului
- Fiecare frunză corespunde unei clase (ieșire de tip discret)



Sisteme inteligente – SIS – Arbori de decizie

□ Exemplu

rec	Age	Income	Student	Credit_rating	Buys_computer(CLASS)
r1	<=30	High	No	Fair	No
r2	<=30	High	No	Excellent	No
r3	31...40	High	No	Fair	Yes
r4	>40	Medium	No	Fair	Yes
r5	>40	Low	Yes	Fair	Yes
r6	>40	Low	Yes	Excellent	No
r7	31...40	Low	Yes	Excellent	Yes
r8	<=30	Medium	No	Fair	No
r9	<=30	Low	Yes	Fair	Yes
r10	>40	Medium	Yes	Fair	Yes
r11	<=30	Medium	Yes	Excellent	Yes
r12	31...40	Medium	No	Excellent	Yes
r13	31...40	High	Yes	Fair	Yes
r14	>40	Medium	No	Excellent	No



Sisteme inteligente – SIS – Arbori de decizie

□ Definire

■ Arborele de decizie

- Un graf special → arbore orientat bicolor
- Conține noduri de 3 tipuri:
 - Noduri de decizie → posibilitățile decidentului (ex. Diversele examinări sau tratamente la care este supus pacientul) și indică un test pe un atribut al articolului care trebuie clasificat
 - Noduri ale hazardului – evenimente aleatoare în afara controlului decidentului (rezultatul examinărilor, efectul terapiilor)
 - Noduri rezultat – situațiile finale cărora li se asociază o utilitate (apreciată aprioric de către un pacient generic) sau o etichetă
- Nodurile de decizie și cele ale hazardului alternează pe nivelele arborelui
- Nodurile rezultat – noduri terminale (frunze)
- Muchiile arborelui (arce orientate) → consecințele în timp (rezultate) ale deciziilor, respectiv ale realizării evenimentelor aleatoare (pot fi însoțite de probabilități)

■ Fiecare nod intern corespunde unui atribut

■ Fiecare ramură de sub un nod (atribut) corespunde unei valori a atributului

■ Fiecare frunză corespunde unei clase

Sisteme inteligente – SIS – Arbori de decizie

□ Tipuri de probleme

- Exemplele (instanțele) sunt reprezentate printr-un număr fix de atribute, fiecare atribut putând avea un număr limitat de valori
- Funcția obiectiv ia valori de tip discret
- AD reprezintă o disjuncție de mai multe conjuncții, fiecare conjuncție fiind de forma atributul a_i are valoarea v_j
- Datele de antrenament pot conține erori
- Datele de antrenament pot fi incomplete
 - Anumitor exemple le pot lipsi valorile pentru unele atribute
- **Probleme de clasificare**
 - Binară
 - exemple date sub forma $[(\text{atribut}_{ij}, \text{valoare}_{ij}), \text{clasă}_i, i=1,2,\dots,n, j=1,2,\dots,m, \text{clasă}_i \text{ putând lua doar 2 valori}]$
 - Multi-clasă
 - exemple date sub forma $[(\text{atribut}_{ij}, \text{valoare}_{ij}), \text{clasă}_i, i=1,2,\dots,n, j=1,2,\dots,m, \text{clasă}_i \text{ putând lua } k \text{ valori}]$
- **Probleme de regresie**
 - AD se construiesc similar cazului problemei de clasificare, dar în locul etichetării fiecărui nod cu eticheta unei clase se asociază nodului o valoare reală sau o funcție dependentă de intrările nodului respectiv
 - Spațiul de intrare se împarte în regiuni de decizie prin tăieturi paralele cu axele Ox și Oy
 - Are loc o transformare a ieșirilor discrete în funcții continue
 - Calitatea rezolvării problemei
 - Eroare (pătratică sau absolută) de predicție

Sisteme inteligente – SIS – Arbori de decizie

□ Proces

- Construirea (creșterea, inducția) arborelui
 - Se bazează pe un set de date de antrenament
 - Lucrează de jos în sus sau de sus în jos (prin divizare – *splitting*)
- Utilizarea arborelui ca model de rezolvare a problemelor
 - Ansamblul deciziilor efectuate de-a lungul unui drum de la rădăcină la o frunză formează o regulă
 - Regulile formate în AD sunt folosite pentru etichetarea unor noi date
- Tăierea (curățirea) arborelui (pruning)
 - Se identifică și se mută/elimină ramurile care reflectă zgomote sau excepții

Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD

■ Divizarea datelor de antrenament în subseturi pe baza caracteristicilor datelor

- Un nod → întrebare legată de o anumită proprietate a unui obiect dat
- Ramurile ce pleacă din nod → etichetate cu posibilele răspunsuri la întrebarea din nodul curent
- La început toate exemplele sunt plasate în rădăcină
 - La pornire, un atribut va fi rădăcina arborelui, iar valorile atributului vor deveni ramuri ale rădăcinii
- Pe următoarele nivele exemplele sunt partiționate în funcție de atribute → ordinea considerării atributelor
 - Pentru fiecare nod se alege în mod recursiv câte un atribut (cu valorile lui pe ramurile descendente din nodul curent)
- Divizarea → greedy în luarea deciziilor

■ Proces iterativ

- Reguli de oprire
 - toate exemplele aferente unui nod fac parte din aceeași clasă → nodul devine frunză și este etichetat cu Ci
 - Nu mai sunt exemple → nodul devine frunză și este etichetat cu clasa majoritară în setul de date de antrenament
 - nu mai pot fi considerate noi atribute

Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD

■ Exemplu

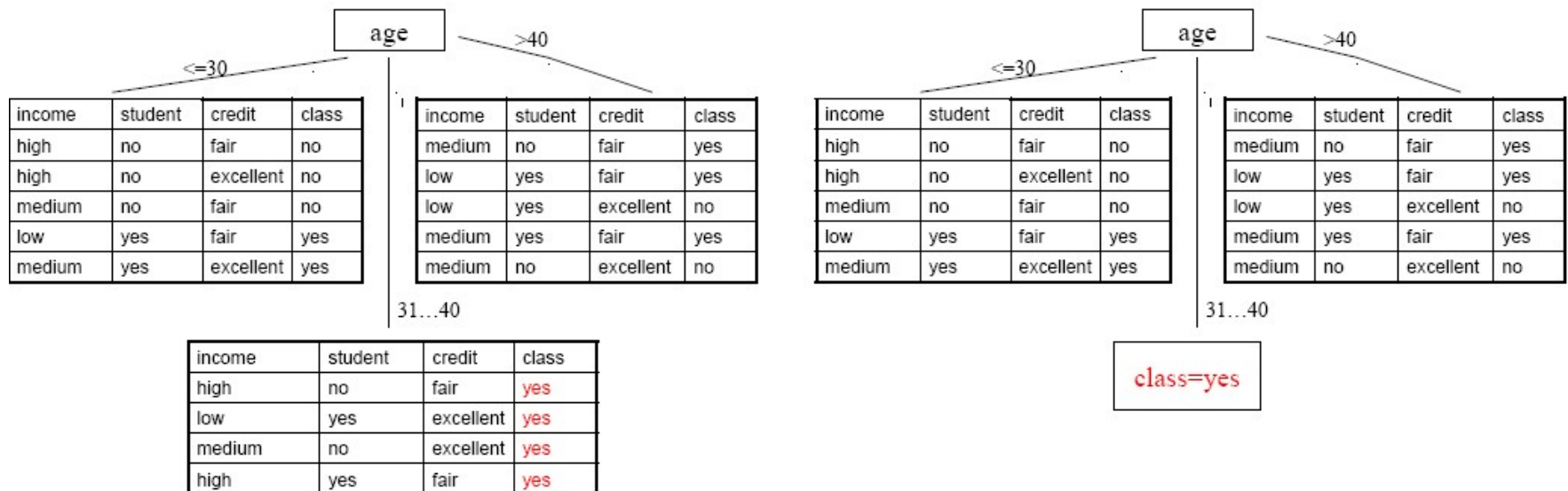
rec	Age	Income	Student	Credit_rating	Buys_computer(CLASS)
r1	<=30	High	No	Fair	No
r2	<=30	High	No	Excellent	No
r3	31...40	High	No	Fair	Yes
r4	>40	Medium	No	Fair	Yes
r5	>40	Low	Yes	Fair	Yes
r6	>40	Low	Yes	Excellent	No
r7	31...40	Low	Yes	Excellent	Yes
r8	<=30	Medium	No	Fair	No
r9	<=30	Low	Yes	Fair	Yes
r10	>40	Medium	Yes	Fair	Yes
r11	<=30	Medium	Yes	Excellent	Yes
r12	31...40	Medium	No	Excellent	Yes
r13	31...40	High	Yes	Fair	Yes
r14	>40	Medium	No	Excellent	No

Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD

■ Exemplu

- Pentru rădăcină se alege atributul *age*

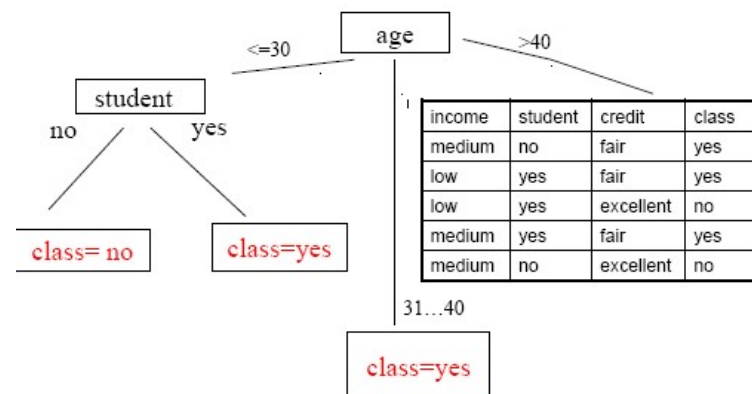
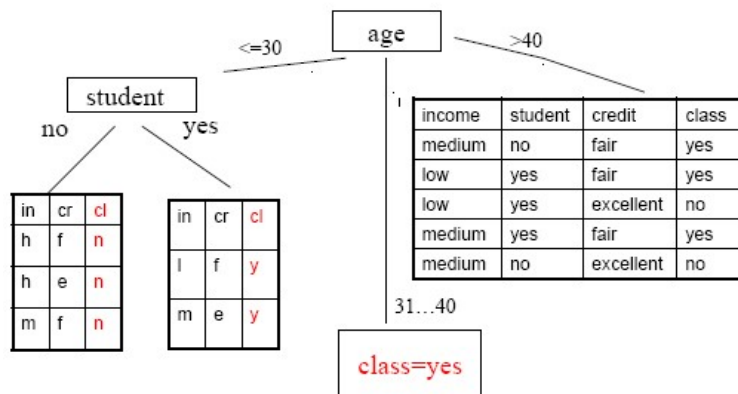


Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD

■ Exemplu

- Pentru rădăcină se alege atributul *age*
- Pe ramura ≤ 30 se alege atributul *student*

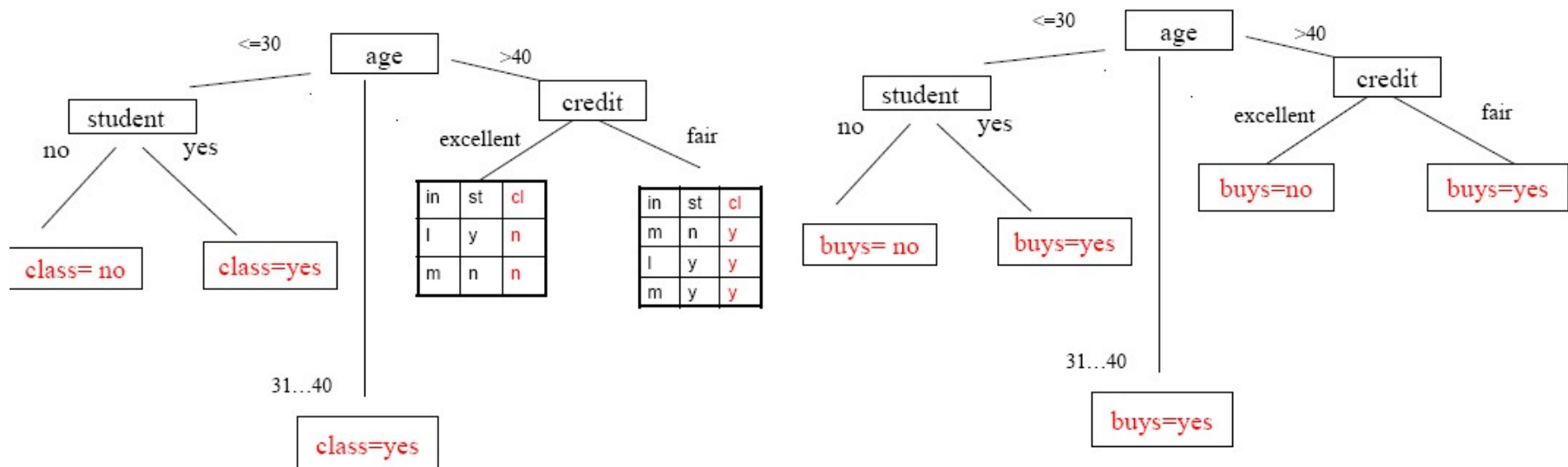


Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD

■ Exemplu

- Pentru rădăcină se alege atributul *age*
- Pe ramura ≤ 30 se alege atributul *student*
- Pe ramura > 40 se alege atributul *credit*



Sisteme inteligente – SIS – Arbori de decizie

- Proces → Construirea AD → Algoritmul ID3/C4.5
 - Greedy, recursiv, top-down, divide-and-conquer

```
generare(D, A){ //D – partiționare a exemplor de antrenament, A – lista de atribute
    Crearea unui nod nou N
    Dacă exemplele din D fac parte dintr-o singură clasă C atunci
        nodul N devine frunză și este etichetat cu C
        returnează nodul N
    Altfel
        Dacă A=∅ atunci
            nodul N devine frunză și este etichetat cu clasa majoritară în D
            returnează nodul N
        Altfel
            atribut_separare = Selectează_atribut(D, A)
            Etichetează nodul N cu atribut_separare
            Pentru fiecare valoare posibilă vj a lui atribut_separare
                Fie Dj mulțimea exemplor din D pentru care atribut_separare = vj
                Dacă Dj = ∅ atunci
                    Atașează nodului N o frunză etichetată cu clasa majoritară în D
                Altfel
                    Atașează nodului N un nod returnat de generare(Dj, A – atribut_separare)
            Returnează nodul N
}
```

Sisteme inteligente – SIS – Arbori de decizie

- Proces → Construirea AD → Algoritmul ID3/C4.5
 - Selectează_atribut(D, A) → Alegerea atributului aferent unui nod (rădăcină sau intern)
 - Aleatoare
 - Atributul cu cele mai puține/multe valori
 - Pe baza unei ordini prestabilite a atributelor
 - Câștigul de informație
 - Rata câștigului
 - Indicele Gini
 - Distanța între partițiile create de un atribut

Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD → Algoritmul ID3/C4.5 → Selectare atribut

■ Câștigul de informație

□ O măsură de impuritate

- 0 (minimă) dacă toate exemplele fac parte din aceeași clasă
- 1 (maximă) dacă avem număr egal de exemple din fiecare clasă

□ Se bazează pe entropia datelor

- măsoară impuritatea datelor
- numărul sperat (așteptat) de biți necesari pentru a coda clasa unui element oarecare din setul de date
- clasificare binară (cu 2 clase): $E(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$, unde
 - p_+ - proporția exemplelor pozitive în setul de date S
 - p_- - proporția exemplelor negative în setul de date S
- clasificare cu mai multe clase: $E(S) = \sum_{i=1, 2, \dots, k} -p_i \log_2 p_i$ - entropia datelor relativ la atributul țintă (atributul de ieșire), unde
 - p_i - proporția exemplelor din clasa i în setul de date S

□ câștigul de informație (*information gain*) al unei caracteristici a (al unui atribut a) datelor

- Reducerea entropiei setului de date ca urmare a eliminării atributului a
- $Gain(S, a) = E(S) - \sum_{v \in \text{valori}(a)} |S_v| / |S| E(S_v)$
- $\sum_{v \in \text{valori}(a)} |S_v| / |S| E(S_v)$ - informația scontată

Sisteme inteligente – SIS – Arbori de decizie

□ Proces → Construirea AD → Algoritmul ID3/C4.5 → Selectare atribut

□ Câștigul de informație

□ exemplu

	a1	a2	a3	Clasa
d1	mare	roșu	cerc	clasa 1
d2	mic	roșu	pătrat	clasa 2
d3	mic	roșu	cerc	clasa 1
d4	mare	albastru	cerc	clasa 2

$$S = \{d1, d2, d3, d4\} \rightarrow p_+ = 2/4, p_- = 2/4 \rightarrow E(S) = -p_+ \log_2 p_+ - p_- \log_2 p_- = 1$$

$$S_{v=\text{mare}} = \{d1, d4\} \rightarrow p_+ = 1/2, p_- = 1/2 \rightarrow E(S_{v=\text{mare}}) = 1$$

$$S_{v=\text{mic}} = \{d2, d3\} \rightarrow p_+ = 1/2, p_- = 1/2 \rightarrow E(S_{v=\text{mic}}) = 1$$

$$S_{v=\text{roșu}} = \{d1, d2, d3\} \rightarrow p_+ = 2/3, p_- = 1/3 \rightarrow E(S_{v=\text{roșu}}) = 0.923$$

$$S_{v=\text{albastru}} = \{d4\} \rightarrow p_+ = 0, p_- = 1 \rightarrow E(S_{v=\text{albastru}}) = 0$$

$$S_{v=\text{cerc}} = \{d1, d3, d4\} \rightarrow p_+ = 2/3, p_- = 1/3 \rightarrow E(S_{v=\text{cerc}}) = 0.923$$

$$S_{v=\text{patrat}} = \{d2\} \rightarrow p_+ = 0, p_- = 1 \rightarrow E(S_{v=\text{patrat}}) = 0$$

$$\text{Gain}(S, a) = E(S) - \sum_{v \in \text{valori}(a)} |S_v| / |S| E(S_v)$$

$$\text{Gain}(S, a_1) = 1 - (|S_{v=\text{mare}}| / |S| E(S_{v=\text{mare}})) + |S_{v=\text{mic}}| / |S| E(S_{v=\text{mic}})) = 1 - (2/4 * 1 + 2/4 * 1) = 0$$

$$\text{Gain}(S, a_2) = 1 - (|S_{v=\text{roșu}}| / |S| E(S_{v=\text{roșu}})) + |S_{v=\text{albastru}}| / |S| E(S_{v=\text{albastru}})) = 1 - (3/4 * 0.923 + 1/4 * 0) = 0.307$$

$$\text{Gain}(S, a_3) = 1 - (|S_{v=\text{cerc}}| / |S| E(S_{v=\text{cerc}})) + |S_{v=\text{patrat}}| / |S| E(S_{v=\text{patrat}})) = 1 - (3/4 * 0.923 + 1/4 * 0) = 0.307$$

Sisteme inteligente – SIS – Arbori de decizie

- Proces → Construirea AD → Algoritmul ID3/C4.5 → Selectare atribut

- Rata câștigului

- Penalizează un atribut prin încorporarea unui termen – *split information* – sensibil la gradul de împrăștiere și uniformitate în care atributul separă datele

- *Split information* – entropia relativ la valorile posibile ale atributului a
- S_v – proporția exemplelor din setul de date S care au atributul a evaluat cu valoarea v

- $splitInformation(S, a) = - \sum_{v=value(a)} \frac{|S_v|}{|S|} \log_2 \frac{|S_v|}{|S|}$

Sisteme inteligente – SIS – Arbori de decizie

▣ Studiu de caz - Problema vampirilor

Items	Shadow	Complexio n	Garlic	Accent	Vampir e
i1	?	Pale	Yes	None	No
i2	Yes	Ruddy	Yes	None	No
i3	?	Ruddy	No	None	Yes
i4	No	Average	No	Heavy	Yes
i5	?	Average	No	Odd	Yes
i6	Yes	Pale	No	Heavy	No
i7	Yes	Average	No	Heavy	No
i8	?	Ruddy	Yes	Odd	No

Sisteme inteligente – SIS – Arbori de decizie

□ Proces

- Construirea arborelui
- Utilizarea arborelui ca model de rezolvare a problemelor
 - Ideea de bază
 - Se extrag regulile formate în arborele anterior construit → Reguli extrase din arborele dat în exemplul anterior:
 - IF *age* = " ≤ 30 " AND *student* = "no" THEN *buys_computer* = "no"
 - IF *age* = " ≤ 30 " AND *student* = "yes" THEN *buys_computer* = "yes"
 - IF *age* = "31...40" THEN *buys_computer* = "yes"
 - IF *age* = " > 40 " AND *credit_rating* = "excellent" THEN *buys_computer* = "no"
 - IF *age* = " > 40 " AND *credit_rating* = "fair" THEN *buys_computer* = "yes"
 - Regulile sunt folosite pentru a clasifica datele de test (date noi). Fie *x* o dată pentru care nu se știe clasa de apartenență → Regulile se pot scrie sub forma unor predicate astfel:
 - IF *age*(*x*, ≤ 30) AND *student*(*x*, no) THEN *buys_computer* (*x*, no)
 - IF *age*(*x*, ≤ 30) AND *student* (*x*, yes) THEN *buys_computer* (*x*, yes)

Sisteme inteligente – SIS – Arbori de decizie

□ Proces

- Construirea arborelui
- Utilizarea arborelui ca model de rezolvare a problemelor
 - Dificultăți
 - *Underfitting* (sub-potrivire) → AD indus pe baza datelor de antrenament este prea simplu → eroare de clasificare mare atât în etapa de antrenare, cât și în cea de testare
 - *Overfitting* (supra-potrivire, învățare pe derost) → AD indus pe baza datelor de antrenament se potrivește prea accentuat cu datele de antrenament, nefiind capabil să generalizeze pentru date noi
 - Soluții:
 - fasonarea arborelui (pruning) → Îndepărtarea ramurilor nesemnificative, redundante → arbore mai puțin stufos
 - validare cu încrucișare

Sisteme inteligente – SIS – Arbori de decizie

□ Proces

- Construirea arborelui
- Utilizarea arborelui ca model de rezolvare a problemelor
- Tăierea (fasonarea) arborelui
 - Necesitate
 - Odată construit AD, se pot extrage reguli (de clasificare) din AD pentru a putea reprezenta cunoștințele sub forma regulilor *if-then* atât de ușor de înțeles de către oameni
 - O regulă este creată (extrasă) prin parcurgerea AD de la rădăcină până la o frunză
 - Fiecare pereche (atribut, valoare), adică (nod, muchie), formează o conjuncție în ipoteza regulii (partea dacă), mai puțin ultimul nod din drumul parcurs care este o frunză și reprezintă consecința (ieșirea, partea atunci) regulii
 - Tipologie
 - Prealabilă (*pre-pruning*)
 - Se oprește creșterea arborelui în timpul inducției prin sistarea divizării unor noduri care devin astfel frunze etichetate cu clasa majoritară a exemplurilor aferente nodului respectiv
 - Ulterioară (*post-pruning*)
 - După ce AD a fost creat (a crescut) se elimină ramurile unor noduri care devin astfel frunze → se reduce eroarea de clasificare (pe datele de test)

Sisteme inteligente – SIS – Arbori de decizie

□ Tool-uri

- <http://webdocs.cs.ualberta.ca/~aixplore/learning/DecisionTrees/Applet/DecisionTreeApplet.html>
- WEKA → J48
- <http://id3alg.altervista.org/>
- <http://www.rulequest.com/Personal/c4.5r8.tar.gz>

□ Biblio

- <http://www.public.asu.edu/~kirkwood/DASstuff/decisiontrees/index.html>

Sisteme inteligente – SIS – Arbori de decizie

□ Avantaje

- Ușor de înțeles și interpretat
- Permit utilizarea datelor nominale și categoriale
- Logica deciziei poate fi urmărită ușor, regulile fiind vizibile
- Lucrează bine cu seturi mari de date

□ Dezavantaje

- Instabilitate → modificarea datelor de antrenament
- Complexitate → reprezentare
- Greu de manevrat
- Costuri mari pt inducerea AD
- Inducerea AD necesită multă informație

Sisteme inteligente – SIS – Arbori de decizie

□ Dificultăți

- Existența mai multor arbori
 - Cât mai mici
 - Cu o acuratețe cât mai mare (ușor de “citit” și cu performanțe bune)
 - Găsirea celui mai bun arbore → problemă NP-dificilă
- Alegerea celui mai bun arbore
 - Algoritmi euristici
 - ID3 → cel mai mic arbore acceptabil
 - → teorema lui Occam: “always choose the simplest explanation”
- Attribute continue
 - Separarea în intervale
 - Câte intervale?
 - Cât de mari sunt intervalele?
- Arbori prea adânci sau prea stufoși
 - Fasonarea prealabilă (pre-pruning) → oprirea construirii arborelui mai devreme
 - Fasonarea ulterioară (post-pruning) → înlăturarea anumitor ramuri

Învățare supervizată – algoritmi

Arbori de decizie

□ Tool-uri

- <http://webdocs.cs.ualberta.ca/~aixplore/learning/DecisionTrees/Applet/DecisionTreeApplet.html>
- WEKA → J48
- <http://id3alg.altervista.org/>
- <http://www.rulequest.com/Personal/c4.5r8.tar.gz>
- <https://scikit-learn.org/stable/modules/tree.html>

□ Biblio

- <http://www.public.asu.edu/~kirkwood/DASTuff/decisiontrees/index.html>
- https://github.com/rasbt/stat479-machine-learning-fs19/tree/master/06_trees



Recapitulare

- Sisteme care învață singure (SIS)
 - Algoritmi de programare genetică (PG)
 - Algoritmi evolutivi cu cromozomi sub formă de arbore
 - Cromozomii
 - Arborescenți
 - Matriciali
 - Liniari
 - codează potențiale soluții de tipul
 - Expresiilor matematice → probleme de regresie/clasificare
 - Expressilor de tip Boolean → probleme de tip EvenParity / proiectare de circuite digitale
 - Programelor → evoluarea de cod sursă pentru rezolvarea unor probleme

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
- Sisteme bazate pe reguli
- Sisteme hibride

Cursul următor – Materiale de citit și legături utile

- ❑ capitolul 15 din *C. Groșan, A. Abraham, Intelligent Systems: A Modern Approach, Springer, 2011*
- ❑ Capitolul 9 din *T. M. Mitchell, Machine Learning, McGraw-Hill Science, 1997*

□ 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

■ Lect. Dr. Crina Groșan -
www.cs.ubbcluj.ro/~cgrosan

■ Prof. Dr. Horia F. Pop -
www.cs.ubbcluj.ro/~hfpop