

Curs 4:

Tehnici de clasificare a datelor

Arbori de decizie și reguli de clasificare

Reminder

- Modele de clasificare:
 - Modele bazate pe reguli
 - Modele bazate pe instanțe
 - Modele probabiliste
 - Modele bazate pe funcții
- În construirea modelelor datele disponibile se împart în:
 - Set de antrenare – folosit în estimarea parametrilor modelului
 - Set de validare – folosit în alegerea caracteristicilor modelului (hiper-parametri)
 - Set de testare – folosit în evaluarea performanței
- Evaluarea performanței unui clasificator:
 - Construire matrice de confuzie
 - Calcul măsuri de performanță: acuratețe, sensibilitate / specificitate, precision / recall, F-score, AUC (area under receiver-operator-characteristic curve)

Structura

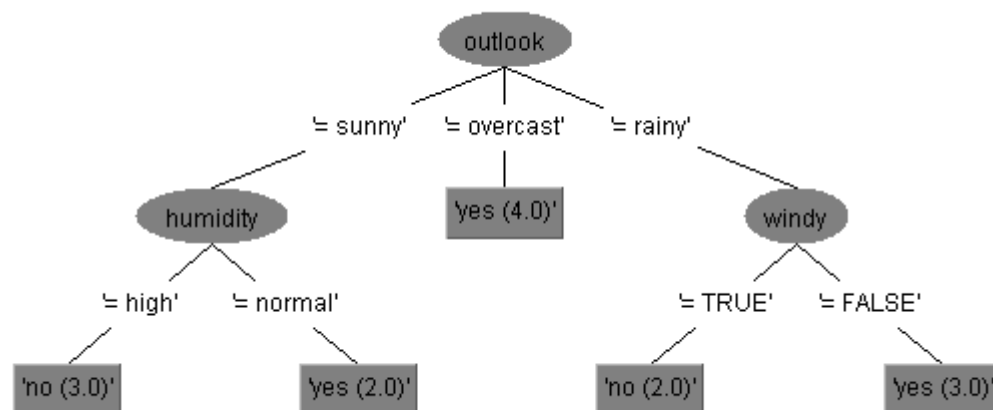
- Arbori de decizie
 - Criterii de ramificare
 - Algoritmi de construire
- Reguli de clasificare
 - Caracteristici ale unui set de reguli
 - Algoritmi de extragere a regulilor din date

Arbori de decizie

Set de date: weather/play

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

Arbore de decizie (construit folosind Weka)



Cum poate fi utilizat un arbore de decizie?

Ce clasă corespunde unei noi instanțe?

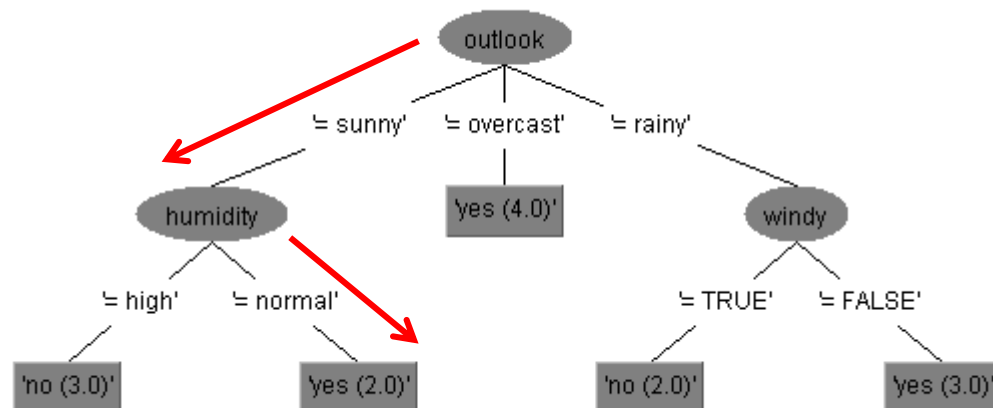
(outlook=sunny, temperature=mild, humidity=normal, windy=False)?

Arbori de decizie

Set de date: weather/play

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

Arbore de decizie (construit folosind Weka)



Cum poate fi utilizat un arbore de decizie?

Ce clasă corespunde unei noi instanțe?

(outlook=sunny, temperature=mild, humidity=normal, windy=False)?

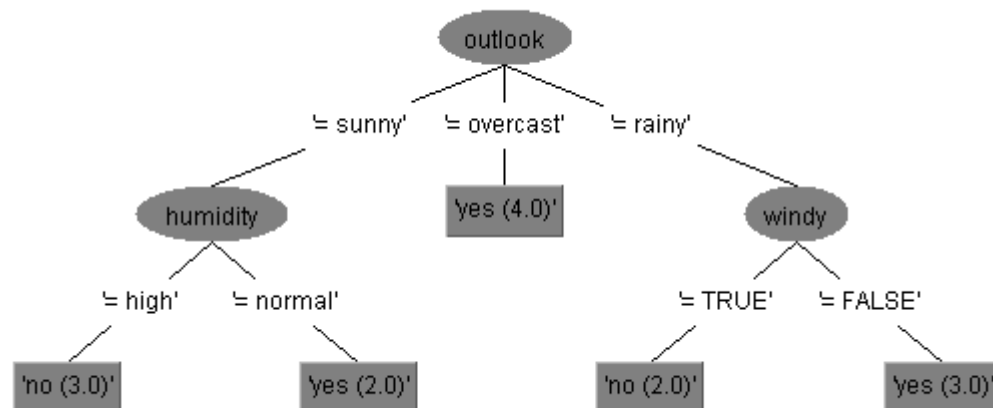
Clasa: Yes (clasa dominantă asociată nodului frunză)

Arbori de decizie

Set de date: weather/play

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

Arbore de decizie (construit folosind Weka)



Cum poate fi tradus într-un set de reguli de clasificare? Fiecare ramură conduce la o regulă

- Regula 1:** IF outlook=sunny and humidity=high THEN play=no
- Regula 2:** IF outlook=sunny and humidity=normal THEN play=yes
- Regula 3:** IF outlook=overcast THEN play=yes
- Regula 4:** IF outlook=rainy and windy=True THEN play=no
- Regula 5:** IF outlook=rainy and windy=False THEN play=yes

Arbori de decizie

Cum poate fi construit un arbore de decizie pornind de la date?

- Se alege un atribut și se plasează în rădăcina arborelului
- Pt fiecare valoare posibilă a atributului (dintre cele prezente în setul de date) se construiește o ramură
- Se partiționează setul de date în subseturi corespunzătoare fiecărei ramuri
 - Dacă un subset conține date ce aparțin unei singure clase atunci el va corespunde unui nod frunză (nu se mai ramifică) – **nod pur**
 - Dacă subsetul conține date din mai multe clase atunci se continuă procesul de partiționare până când
 - se ajunge la un nod pur
 - pe ramura respectivă au fost deja analizate toate attributele
 - subsetul corespunzător ramurii este vid

Weather/play dataset

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

Problemă: în ce ordine ar trebui analizate attributele? Ce condiție de testare ar trebui asociată cu fiecare nod?

Arbori de decizie

In ce ordine ar trebui analizate attributele?

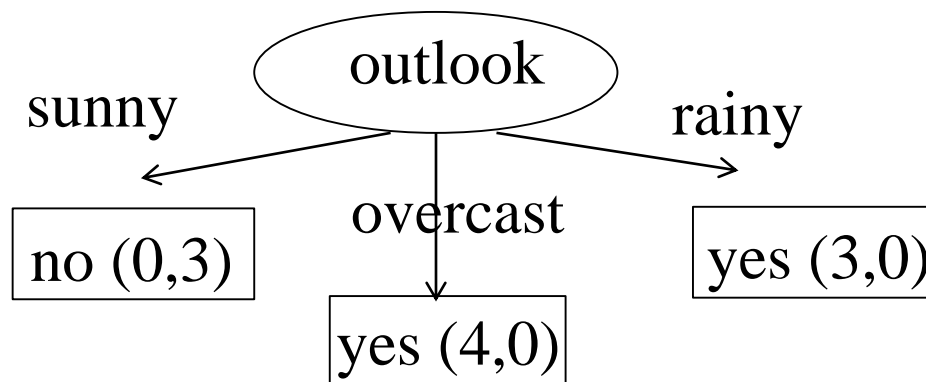
Set de date: weather/play (date selectate)

Ideea principală:

- Se selectează atributul care conduce la un arbore cât mai simplu adică un atribut cu grad de puritate cât mai mare (ideal ar fi ca pentru fiecare valoare a atributului datele corespunzătoare să aparțină aceleiași clase)

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	overcast	cool	normal	TRUE	yes
7	sunny	mild	high	FALSE	no
8	rainy	mild	normal	FALSE	yes
9	overcast	mild	high	TRUE	yes
10	overcast	hot	normal	FALSE	yes

Exemplu:



Obs:

- Toate nodurile frunză sunt “pure” (conțin date ce aparțin aceleiași clase)
- Conduce la reguli cu un singur atribut în membrul stâng
- O astfel de situație se întâmplă rar pentru date reale

Arbori de decizie

Principalele probleme ce trebuie soluționate la construirea unui arbore de decizie

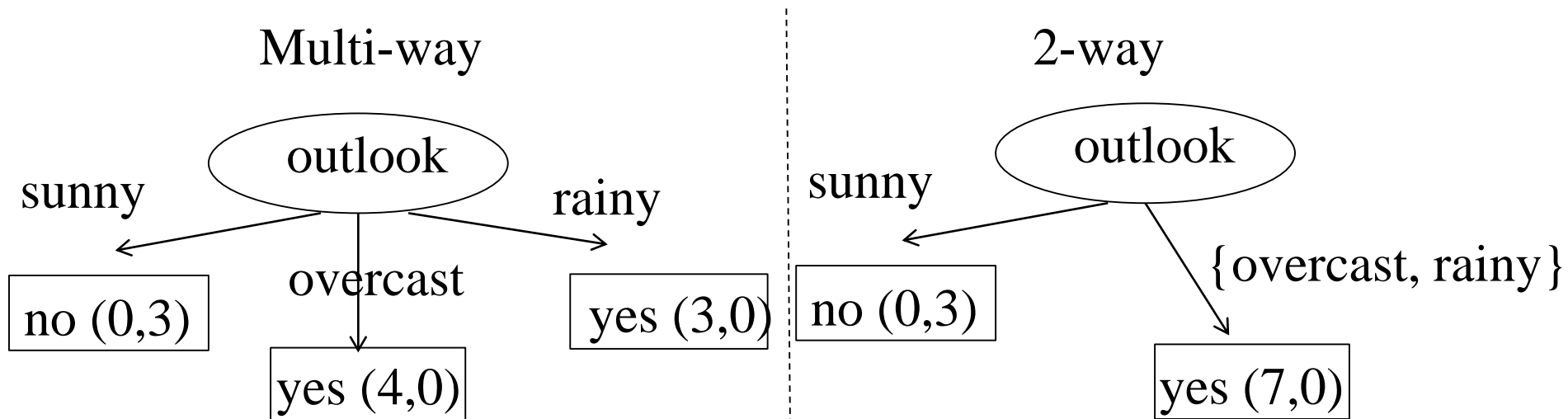
- **Ce atribut ar trebui selectat pentru partiționare?**
 - Cel cu puterea cea mai mare de discriminare = cel ce asigură partiționarea setului curent în subseturi cu **grad mare de puritate**
 - Criterii ce pot fi utilizate:
 - Bazate pe entropie (ex: câștig informațional)
 - Index Gini
 - Măsură a erorii de clasificare
- **Ce condiții de test trebuie asignate ramurilor corespunzătoare unui nod ?**
 - Depinde de tipul atributului
 - Nominal, ordinal: **atribut = valoare**
 - Continuu: **atribut < valoare** sau **valoare1 < atribut < valoare2**
 - Depinde de gradul de ramificare dorit:
 - Ramificare binară (setul curent de date este împărțit în două subseturi)
 - Ramificare multiplă (setul curent de date este împărțit în mai multe subseturi)

Arbori de decizie

- Ce condiții de test trebuie asigurate ramurilor corespunzătoare unui nod ?

Atribute nominale și ordinale:

- Ramificare multiplă (multi-way): atâtea ramuri câte valori posibile are atributul
- Ramificare binară (2-way): două ramuri



Arbori de decizie

Ce condiții de test trebuie asignate ramurilor corespunzătoare unui nod?

Atribute numerice:

- Trebuie discretizate în prealabil, după care se aplică strategia specifică atributelor nominale sau ordinale

Ce atribut se selectează pentru partiționare?

- Acel atribut care conduce la reducerea maximă în conținutul de informație necesar pentru a lua decizia corectă

Exemplu: information gain

Câștig informațional =
Entropia(distribuția datelor înainte de partiționare) –
EntropiaMedie (distribuția datelor după partiționare)

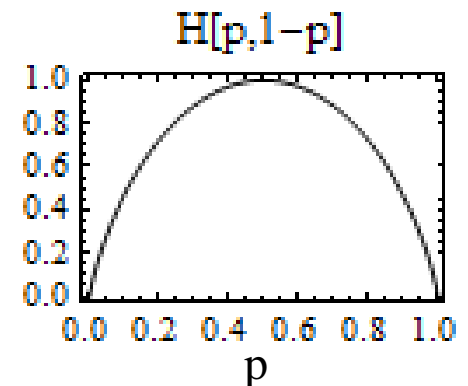
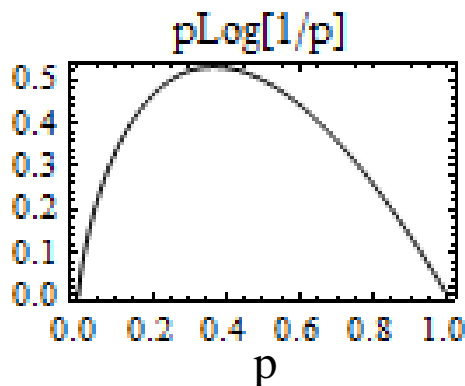
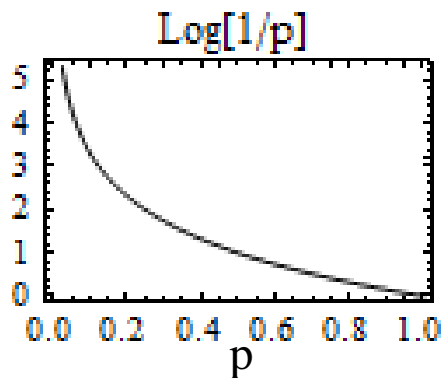
Reminder: entropie

Fie $D=(p_1, p_2, \dots, p_k)$ o distribuție de probabilitate. Entropia asociată acestei distribuții de probabilitate este caracterizată de:

$$H(D) = H(p_1, p_2, \dots, p_k) = -\sum_{i=1}^k p_i \log p_i$$

și poate fi interpretată ca o **măsură medie a incertitudinii** (sau surprizei) când se generează/selectează o valoare pe baza acestei distribuții

Caz particular: $k=2 \Rightarrow p_1=p, p_2=1-p$



Obs: Interpretare $\text{Log}[1/p]$: surpriza de a observa un eveniment caracterizat de o probabilitate mică (eveniment neașteptat) este mai mare decât cea corespunzătoare unui eveniment de probabilitate mai mare (eveniment așteptat)

Reminder: entropie

În contextul rezolvării problemelor de clasificare:

- $D = \{C_1, C_2, \dots, C_k\}$ (set de date distribuit în k clase)
- Distribuția de probabilitate (p_1, p_2, \dots, p_k) , $p_i = \text{card}(C_i) / \text{card}(D)$
- Fie A un atribut și v_1, v_2, \dots, v_{m_A} valorile posibile ale acestui atribut
- Fie D_j = setul de date din D pt care atributul A are valoarea v_j și P_j distribuția datelor din D_j în cele k clase (C_{ji} = set de date din clasa C_i care au valoarea v_j pt atributul A)
- **Câștigul informațional** obținut prin partiționarea setului de date folosind atributul A este:

$$IG(D, A) = H(D) - \sum_{j=1}^{m_A} P(D_j | A = v_j) H(D_j | A = v_j), \quad H(D) = - \sum_{i=1}^k p_i \log p_i$$

$$H(D_j | A = v_j) = - \sum_{i=1}^k p_{ij} \log p_{ij}, \quad p_{ij} = \frac{\text{card}(C_{ji})}{\text{card}(C_i)}$$

$$P(D_j | A = v_j) = \frac{\text{card}(D_j)}{\text{card}(D)}$$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

- Distribuția claselor (C_1 ="yes", C_2 ="no"):
- $p_1=9/14$, $p_2=5/14$
- $H(p_1, p_2)=0.94$

Outlook

	C1 (yes)	C2(no)	Frequency
Sunny	2/5	3/5	5/14
Overcast	4/4	0/4	4/14
Rainy	3/5	2/5	5/14

$$H(\text{sunny}) = -2/5 \cdot \log(2/5) - 3/5 \cdot \log(3/5) = 0.97$$

$$H(\text{overcast}) = -1 \cdot \log(1) - 0 = 0$$

$$H(\text{rainy}) = -3/5 \cdot \log(3/5) - 2/5 \cdot \log(2/5) = 0.97$$

$$IG(\text{outlook}) = 0.94 - 5/14 \cdot 0.97 - 4/14 \cdot 0 - 5/14 \cdot 0.97 = 0.94 - 0.69 = 0.25$$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

- Distribuția claselor(C_1 ="yes", C_2 ="no"):
- $p_1=9/14$, $p_2=5/14$
- $H(p_1, p_2)=0.94$

Temperature

	C1 (yes)	C2(no)	Frequency
Hot	2/4	2/4	4/14
Mild	4/6	2/6	6/14
Cool	3/4	1/4	4/14

$$H(\text{hot}) = -2/4 \cdot \log(2/4) - 2/4 \cdot \log(2/4)$$

$$H(\text{mild}) = -4/6 \cdot \log(4/6) - 2/6 \cdot \log(2/6)$$

$$H(\text{cool}) = -3/4 \cdot \log(3/4) - 1/4 \cdot \log(1/4)$$

$$IG(\text{temperature}) = 0.03$$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

- Distribuția claselor (C_1 ="yes", C_2 ="no"):
- $p_1=9/14$, $p_2=5/14$
- $H(p_1, p_2)=0.94$

Humidity

	C1 (yes)	C2(no)	Frequency
High	3/7	4/7	7/14
Normal	6/7	1/7	7/14

$$H(\text{high}) = -3/7 \cdot \log(3/7) - 4/7 \cdot \log(4/7)$$

$$H(\text{normal}) = -6/7 \cdot \log(6/7) - 1/7 \cdot \log(1/7)$$

$$IG(\text{humidity}) = 0.15$$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

- Distribuția claselor(C_1 ="yes", C_2 ="no"):
- $p_1=9/14$, $p_2=5/14$
- $H(p_1, p_2)=0.94$

Windy

	C1 (yes)	C2(no)	Frequency
False	6/8	2/8	8/14
True	3/6	3/6	6/14

$$H(\text{false}) = -6/8 \cdot \log(6/8) - 2/8 \cdot \log(2/8)$$

$$H(\text{true}) = -3/6 \cdot \log(3/6) - 3/6 \cdot \log(3/6)$$

$$IG(\text{windy}) = 0.05$$

Alegerea atributului de ramificare

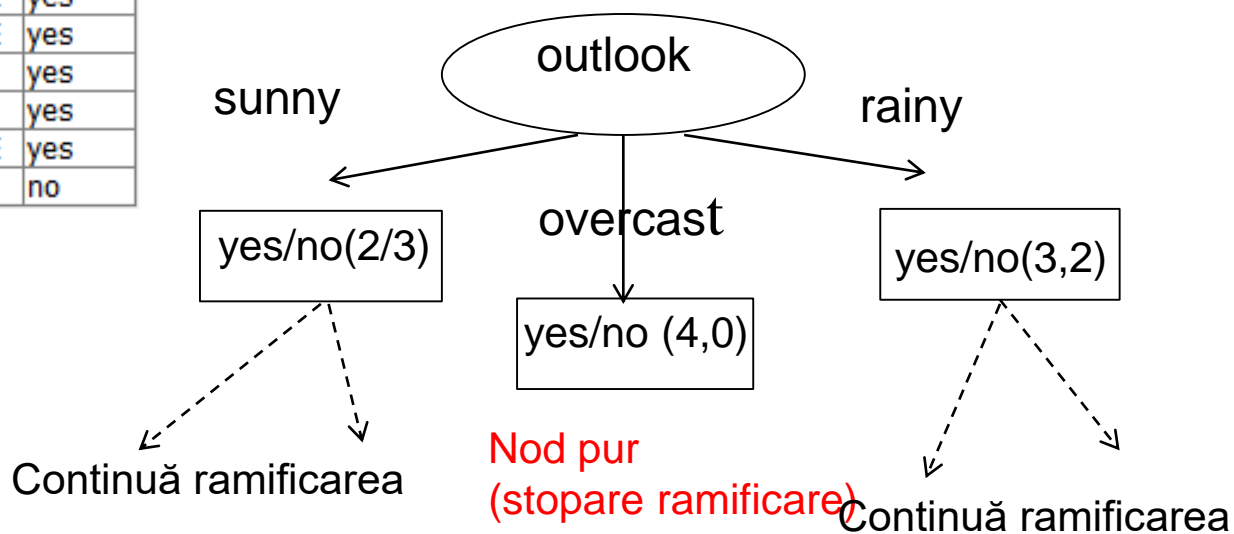
Exemplu

Relation: weather.symbolic					
No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no

Câștigul informațional al fiecărui atribut:

- $IG(\text{outlook})=0.25$
- $IG(\text{temperature})=0.03$
- $IG(\text{humidity})=0.15$
- $IG(\text{windy})=0.05$

Primul atribut selectat: **outlook**

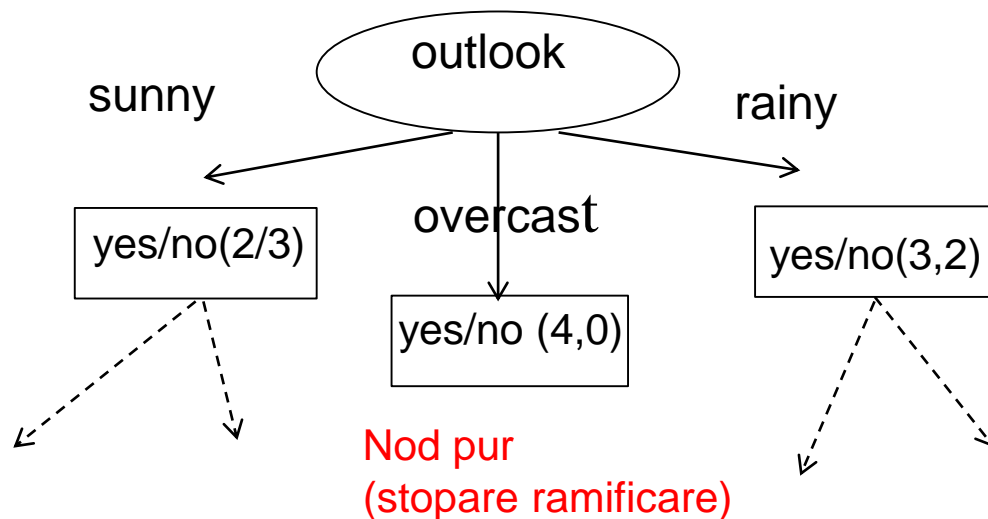


Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no



Temperature

	C1 (yes)	C2 (no)	Freq.
Hot	0/2	2/2	2/5
Mild	1/2	1/2	2/5
Cool	1/1	0/1	1/5

Câștig informațional pt attributele rămase:

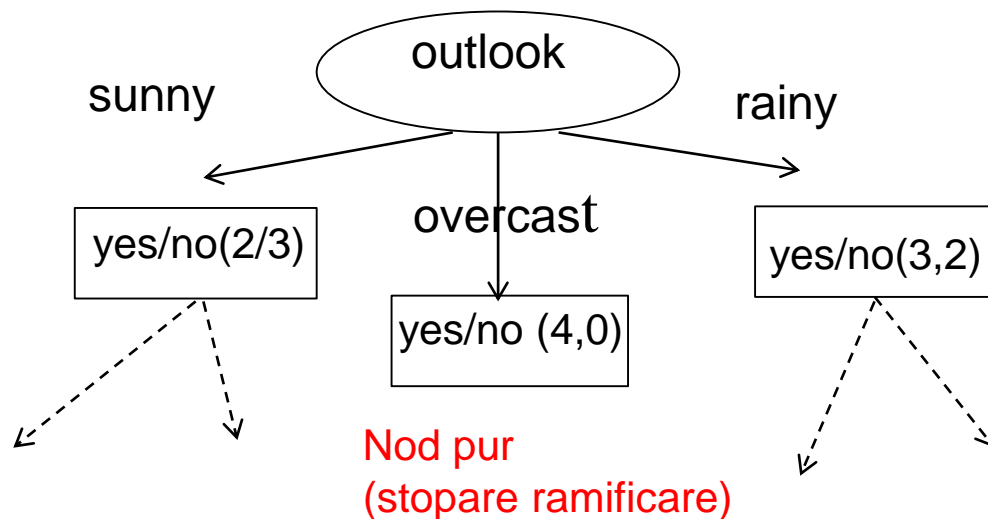
- Entropia pt subsetul "sunny" :
 $H(D(\text{sunny})) = -2/5 \cdot \log(2/5) - 3/5 \cdot \log(3/5) = 0.97$
- $H(\text{hot})=0$, $H(\text{mild})=1$, $H(\text{cool})=0$
- $IG(\text{temperature}) = 0.97 - 2/5 = 0.57$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no



Humidity

	C1 (yes)	C2 (no)	Freq.
High	0/3	3/3	3/5
Nor mal	2/2	0/2	2/5

Câștig informațional pt attributele rămase:

▪ Entropia pt subsetul "sunny" :

$$H(D(\text{sunny})) = -2/5 \cdot \log(2/5) - 3/5 \cdot \log(3/5) = 0.97$$

▪ $H(\text{high})=0$, $H(\text{normal})=0$

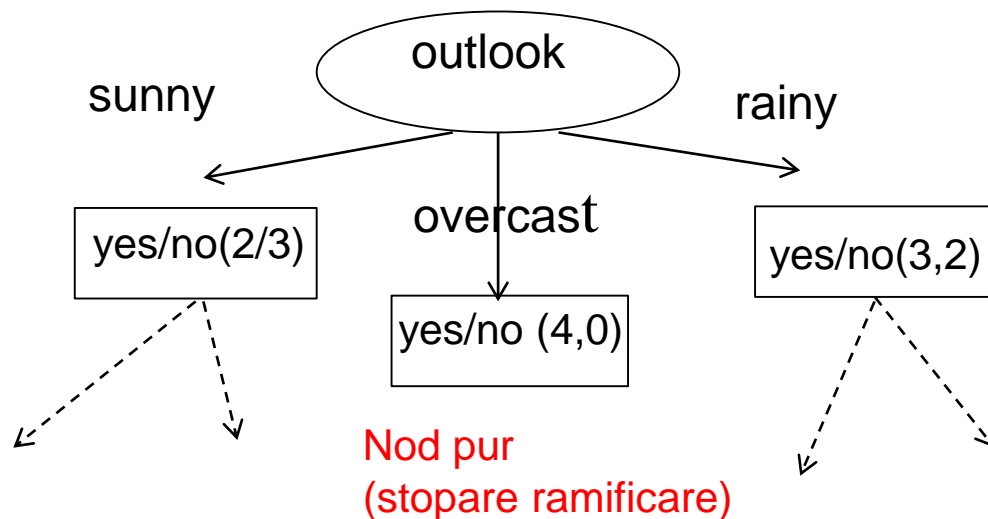
▪ $IG(\text{humidity}) = 0.97 - 0 = 0.97$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no



Windy

	C1 (yes)	C2 (no)	Freq.
false	1/3	2/3	3/5
true	1/2	1/2	2/5

Câștig informațional pt attributele rămase:

▪ Entropia pt subsetul "sunny" :

$$H(D(\text{sunny})) = -2/5 \cdot \log(2/5) - 3/5 \cdot \log(3/5) = 0.97$$

▪ $H(\text{false})=0$, $H(\text{true})=1$

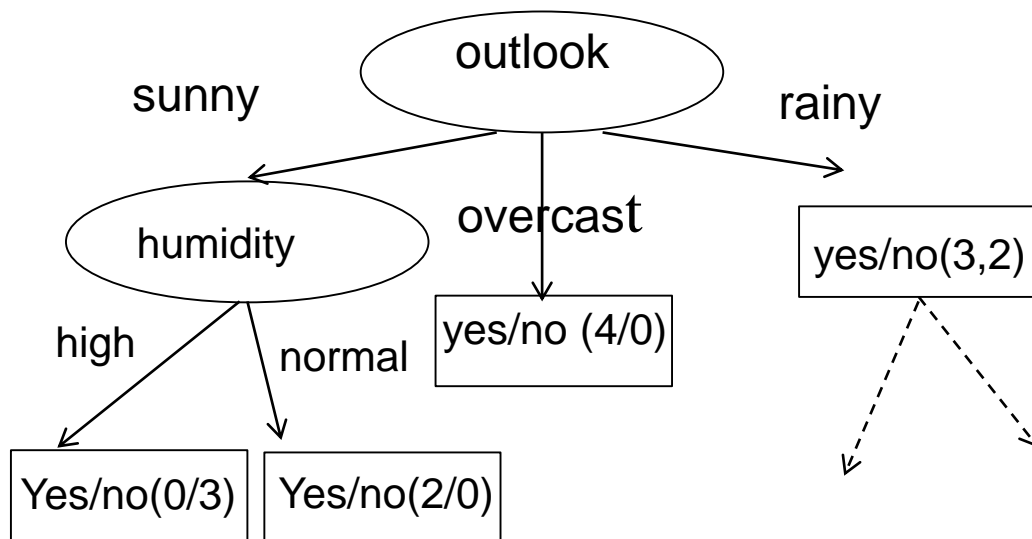
▪ $IG(\text{windy}) = 0.97 - 0.95 = 0.02$

Alegerea atributului de ramificare

Exemplu

Relation: weather.symbolic

No.	outlook Nominal	temperature Nominal	humidity Nominal	windy Nominal	play Nominal
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	yes
14	rainy	mild	high	TRUE	no



Câștig informațional pt attributele rămase:

- $IG(\text{temperature}) = 0.97 - 2/5 = 0.57$
- $IG(\text{humidity}) = 0.97 - 0 = 0.97$
- $IG(\text{windy}) = 0.97 - 0.95 = 0.02$

Alegerea atributului de ramificare

Obs:

- Câștigul informațional favorizează attributele caracterizate printr-un număr mare de valori
- Pentru a limita această influență se poate utiliza raportul (**Gain Ratio**):

$$GainRatio(D, A) = \frac{IG(D, A)}{H(p_1^A, p_2^A, \dots, p_{m_A}^A)}$$

$$p_j^A = \frac{card(D, A = v_j)}{card(D)}$$

(proportia de date care au valoarea v_j pt atributul A)

Alegerea atributului de ramificare

- Atributul de ramificare poate fi determinat folosind **indexul Gini** = măsură a gradului mediu de impuritate a subseturilor de date obținute prin ramificarea bazată pe atributul A (cu cât mai mică valoarea cu atât mai bună)

index Gini pt atributul A_i

$$G(A_i) = \frac{1}{N} \sum_{j=1}^{r_i} n_{ij} G(v_{ij}), \quad G(v_{ij}) = 1 - \sum_{k=1}^K p_{ijk}^2$$

n_{ij} = numarul de instante pt care A_i are valoarea v_{ij}

$$p_{ijk} = \frac{\text{numar de instante in } C_k \text{ cu } A_i = v_{ij}}{\text{numar de instante cu } A_i = v_{ij}}$$

Algoritmi pentru construirea arborilor de decizie

ID3 (Iterative Dichotomiser):

- Intrare: set de date D
- Ieșire: arbore de decizie (noduri interne etichetate cu atribute, noduri frunză etichetate cu clase, muchii etichetate cu valori ale atributelor)

```
DTinduction (D, DT, N) /* D=set date, DT=arbore de decizie, N=nod */  
  find the best splitting attribute A  
  label node N with A  
  construct the splitting predicates (branches) for N  
  FOR each branch i from N DO  
    construct the corresponding data set  $D_i$   
    create a new child node  $N_i$   
    IF <stopping condition>  
      THEN label  $N_i$  with the dominant class in  $D_i$  ( $N_i$  is a leaf node)  
    ELSE DTinduction( $D_i$ , DT,  $N_i$ )
```

Algoritmi pentru construirea arborilor de decizie

C4.5 = îmbunătățire a algoritmului ID3 pt a trata:

- **Atribute continue:**
 - Incorporează procedură de discretizare
- **Valori absente:**
 - Datele ce conțin valori absente sunt ignorate sau
 - Valorile absente sunt imputate
- **Atribut de ramificare:**
 - Utilizează Gain Ratio pt selecția atributului
- **Simplificare sau fasonare (Pruning):**
 - Anumiți subarbori sunt înlocuiți cu noduri frunză (dacă eroarea de clasificare nu crește semnificativ) – abordare bottom-up

Obs:

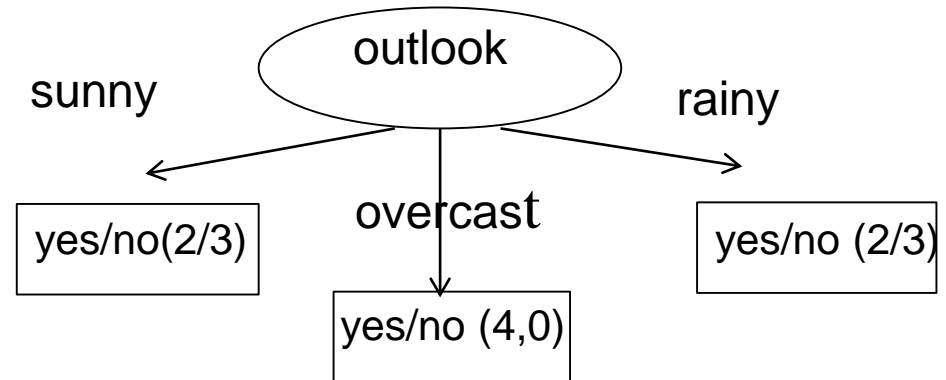
C5.0 – varianta comerciala a algoritmului C4.5 – accesibil pt analiza din 2011

J48 – implementarea din Weka a algoritmului C4.5

Algoritmi pentru construirea arborilor de decizie

Simplificare (pruning):

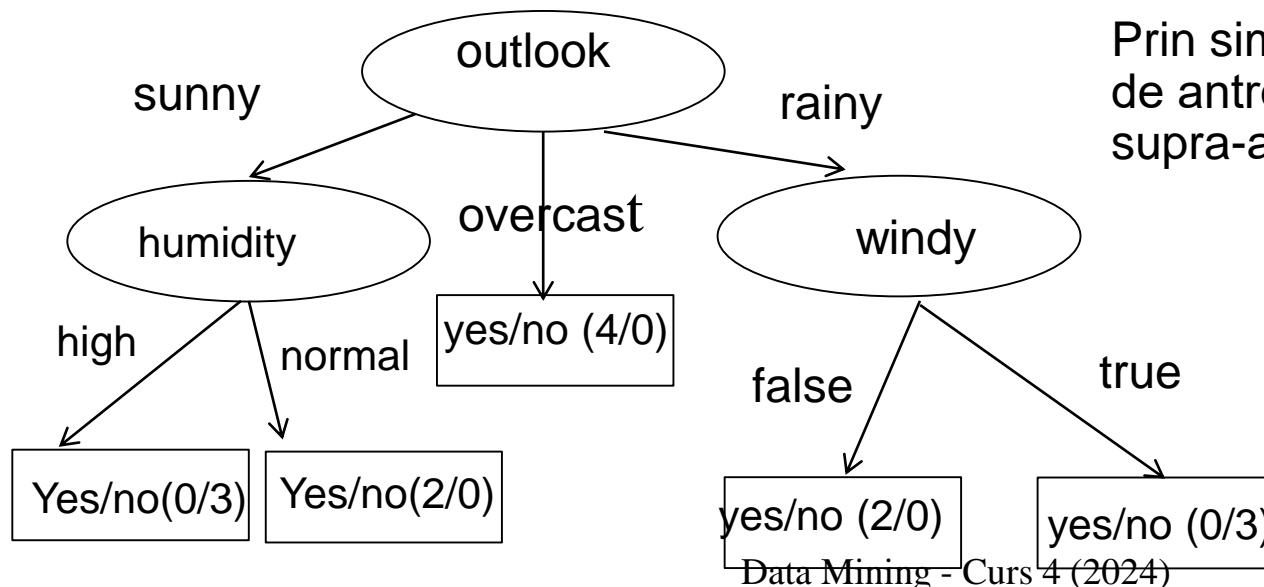
- Anumiți subarbori sunt înlocuiți cu noduri frunză (dacă eroarea de clasificare nu crește semnificativ) – abordare bottom-up



Arbore : error = 0

Arbore simplificat: error = 4/14

Prin simplificare eroarea pe setul de antrenare crește dar riscul de supra-antrenare poate să scadă



Extragerea regulilor de clasificare

Reminder: regulile de clasificare sunt structuri de tip IF ... THEN care conțin:

- In partea de **antecedent** (membrul stâng): condiții privind valorile atributelor (pot fi expresii logice care implică mai multe attribute)
- In partea de **consecință** (membrul drept): eticheta clasei

Example:

IF outlook=sunny THEN play=no

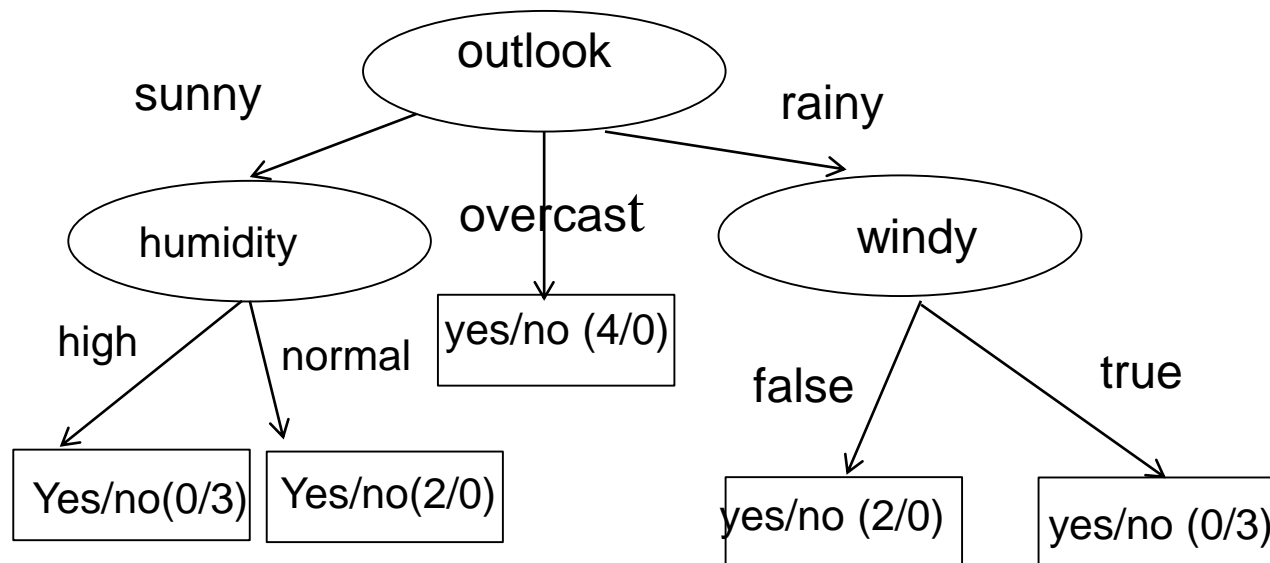
IF outlook=rainy THEN play=no

IF outlook=overcast THEN play=yes

Extragerea regulilor de clasificare

- Dacă regulile sunt extrase dintr-un arbore de decizie atunci fiecare ramură conduce la o regulă
- Condițiile referitoare la noduri aflate pe aceeași ramură se combină prin AND:
IF (outlook=sunny) and (humidity=high) THEN play=no
- Regulile corespunzând unor ramuri diferite dar conducând la aceeași consecință (aceeași etichetă de clasă) pot fi reunite prin disjuncție (OR) între părțile de antecedent:

IF (outlook=sunny and humidity=high) OR (outlook=rainy and windy =true) THEN play=no



Extragerea regulilor de clasificare

Alta variantă: Regulile de clasificare pot fi extrase direct din date printr-un proces de învățare utilizând algoritmi de acoperire (covering algorithms)

Noțiuni:

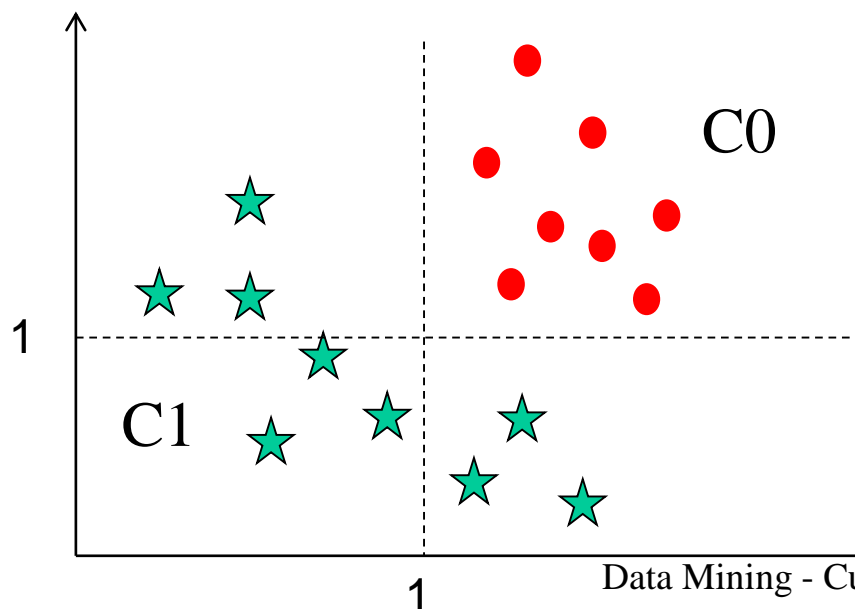
- O regulă **acoperă** o dată dacă valorile atributelor se potrivesc cu condițiile din antecedentul regulii
- Similar, despre o dată se spune că **activează** o regulă dacă valorile atributelor se potrivesc cu condițiile din antecedentul regulii
- **Suportul unei reguli (support)** = fracțiunea din setul total de date care este acoperită de către regulă și aparține aceleiași clase ca cea din regulă
$$= |\text{cover}(R) \cap \text{class}(R)| / |D|$$
- **Gradul de încredere în regulă (rule confidence)** = fracțiunea din datele acoperite de regulă care au aceeași clasă ca cea specificată de regulă
$$= |\text{cover}(R) \cap \text{class}(R)| / |\text{cover}(R)|$$
 - $\text{cover}(R)$ = subsetul de date acoperit de R
 - $\text{class}(R)$ = subsetul de date care au aceeași clasă cu cea specificată în R
 - D = setul de date

Extragerea regulilor de clasificare

Noțiuni:

- **Reguli mutual exclusive** = regiunile acoperite de reguli sunt disjuncte (o instanță activează o singură regulă)
- **Set complet de reguli** = fiecare instanță activează cel puțin o regulă

Obs: dacă setul de reguli e complet și regulile sunt mutual exclusive atunci decizia privind apartenența unei date la o clasă este simplu de luat



Exemplu:

R1: IF $x > 1$ and $y > 1$ THEN C0

R2: IF $x \leq 1$ THEN C1

R3: IF $x > 1$ and $y \leq 1$ THEN C1

Ce se întâmplă însă dacă aceste proprietăți nu sunt satisfăcute?

Extragerea regulilor de clasificare

Obs: dacă regulile nu sunt mutual exclusive atunci pot să apară **conflicte** (o instanță poate activa reguli care au asociate clase diferite)

Conflictele pot fi rezolvate în unul dintre următoarele moduri:

- **Ordonarea regulilor** (pe baza unui criteriu) și decizia se ia conform primei reguli activate (prima regulă care se potrivește cu instanța).
- **Criteriul de ordonare** poate fi corelat cu:
 - **calitatea regulii** (e.g. grad de încredere) – regulile cu grad mai mare de încredere sunt mai bune
 - **specificitatea regulii** – o regulă este considerată mai bună dacă este mai specifică (e.g. reguli care corespund claselor rare)
 - **complexitatea regulii** (e.g. numărul de condiții din partea de antecedent a regulii) – regulile mai simple sunt mai bune

Extragerea regulilor de clasificare

Obs: aceste criterii pot fi conflictuale (o regulă cu coeficient mare de încredere nu este neapărat o regulă simplă)

- Rezultatul se obține considerând **clasa dominantă** pe baza tuturor regulilor activate de către instanță

Extragerea regulilor de clasificare

Algoritm secvențial de acoperire:

Intrare: set de date

Ieșire: set ordonat de reguli

Pas 1: se **selectează una dintre clase** și se identifică cea mai “bună” regulă care acoperă datele din D care au clasa selectată. Se adaugă regula la lista de reguli.

Pas 2: Se elimină datele din D care activează regula adăugată. Dacă încă există clase netratate și date în D se reia de la Pas 1

Obs:

- Aceasta este structura generală a algoritmilor secvențiali de acoperire
- Algoritmii pot să difere în funcție de strategia de selecție a claselor

Extragerea regulilor de clasificare

Exemplu: algoritmul RIPPER (Repeated Incremental Pruning to Produce Error Reduction)

Particularități:

- Setul de date e divizat la început în **growing set** (folosit pentru construirea unui set de reguli care acoperă setul de date) și **pruning set** (folosit pt simplificarea regulilor, de ex. prin eliminarea unor attribute din membrul stâng al regulii - se alege varianta de simplificare care reduce cel mai mult eroarea pe **pruning set**)
- **Ordonare bazată pe clase:** clasele sunt selectate crescător după dimensiune (clasele rare sunt selectate prima dată)
- Regulile corespunzătoare unei clase sunt plasate consecutiv în lista de reguli

Extragerea regulilor de clasificare

Exemplu: algoritmul RIPPER (Repeated Incremental Pruning to Produce Error Reduction)

Particularități:

- Adăugarea unei noi reguli corespunzătoare unei clase este stopată:
 - Când regula devine prea complexă
 - Când ‘următoarea’ regulă are o eroare de clasificare (pe setul de validare) mai mare decât un prag prestabilit
- Dacă la sfârșit rămân date “neacoperite” atunci se poate defini o regulă de tipul “**catch all**” căreia i se asociază clasa dominantă din setul de date “neacoperite”

Sumar

- Arbori de decizie
 - Criterii de ramificare: Information Gain, Gain Ratio, Gini Index
 - Algoritmi de construire: ID3, C4.5
 - Simplificare (pruning)
 - Avantaje și dezavantaje
- Reguli de clasificare
 - Caracteristici: suport, încredere
 - Algoritmi de extragere a regulilor din date (RIPPER)
- Modelele bazate pe reguli sunt interpretabile
 - <https://christophm.github.io/interpretable-ml-book/tree.html>
 - interpretable models: <https://github.com/csinva/imodels>

Cursuri următoare

Curs 5:

- Clasificatori bazați pe prototipuri (k Nearest Neighbours)
- Modele probabiliste (Naive Bayes classifiers)

Curs 6:

- Rețele neuronale (Multilayer Perceptrons)
- Clasificatori bazați pe vectori suport (Support Vector Machines)