

Modele grafice. Rețele Bayesiene

Cursul 2

Programare și modelare probabilistă - anul III

Facultatea de Informatică, UAIC

e-mail: adrian.zalinescu@uaic.ro

web: <https://sites.google.com/view/fiicoursepmp/home>

16 Octombrie 2023

- 1 Modelarea dependențelor
- 2 Utilizarea rețelelor Bayesiene
 - Modelarea BN cu pgmpy

Modele grafice

Cum putem face inferențe într-un mediu aleator?

Putem folosi *modelele grafice direcționate*:

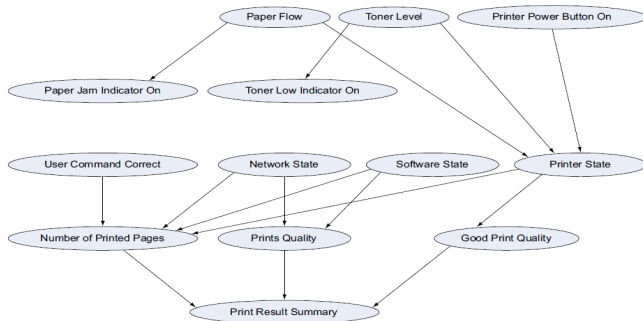
- o reprezentare a distribuției multivariate

$$P(X_1, X_2, \dots, X_n)$$

factorizată de dependențe condiționale.

- reprezentăm variabilele prin **noduri** și dependențele prin **muchii**.

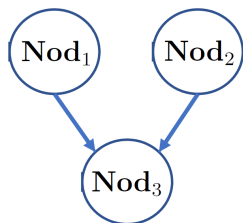
Exemplu:



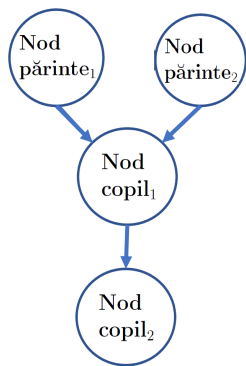
De ce folosim modelele grafice direcționate?

Modelele grafice sunt generative:

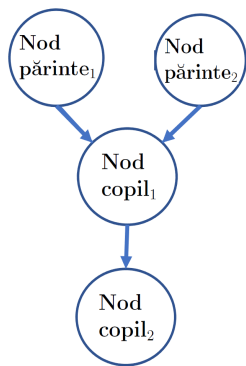
- pot defini distribuții de probabilitate;
- sunt expresive în termeni de structură a dependențelor;
- se pot genera realizări (eșantioane) din respectivele distribuții.



- *Nodurile* conțin distribuții condiționale (CPD);
- *Muchiile neorientate* definesc conexiunea între noduri (sau scheletul grafului);
- Informația sau influența curge de-a lungul *muchiilor orientate*;
- Nodurile conectate prin muchii se numesc *vecini*.

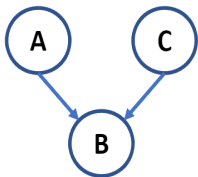


- *Nodurile părinți* influențează prin muchii orientate *nodurile copii*;
- Nodurile copii sunt *descendenți* ai nodurilor părinți;
- *Ascendenții* sunt părinții, bunicii, etc;
- Ex.: ascendenții nodului NC_2 :
 $A(NC_2) = \{NP_1, NP_2, NC_1\}$;
- *Descendenții* unui nod sunt influențați de acesta;
- Ex.: descendenții nodului NP_1 : $D(NP_1) = \{NC_1, NC_2\}$.

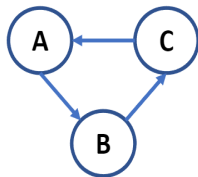


- *Gradul unui nod*: numărul de vecini;
- *Gradul interior al unui nod*: numărul de părinți;
- *Gradul exterior al unui nod*: numărul de copii;
- Exemplu: $IN(NC_1) = 2$;
- Exemplu: $OUT(NC_1) = 1$;
- Un *nod rădăcină* nu are ascendenți: mulțimea acestora este $\{NP_1, NP_2\}$;
- Un *nod frunză* nu are descendenți: mulțimea acestora este $\{NC_2\}$.

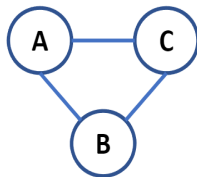
Tipuri de Grafuri



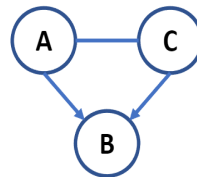
Graf aciclic
orientat



Graf ciclic
orientat



Graf
neorientat

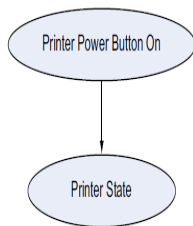


Graf parțial
orientat

Tipuri de dependență în modele grafice

- *dependențe direcționate* de la o variabilă la alta.
De obicei modelează o relație cauză-efect.
- *dependențe nedirecționate*, care modelează relații dintre variabile ce nu se influențează direct una pe cealaltă.

Exemplu:



Starea imprimantei depinde de starea butonului de pornire.

Atenție! Direcția de inferență nu e neapărat cea de dependență.

Relații cauză-efect

- *temporale*. Exemplu: se apasă butonul de oprire mai întâi și apoi imprimanta se va opri.
- *cauză-efect de stare*. Exemplu: o variabilă reprezintă dacă butonul de oprire este în starea “off”, iar cealaltă dacă imprimanta este oprită.
- *starea de adevăr a unei măsurări*. Când o variabilă *măsoară* valoarea unei alte variabile. Exemplu: o variabilă care reprezintă dacă becul ce indică conectarea imprimantei la sursa de curent electric este aprins. Măsurătorile sunt de obicei observate și, în acest caz, raționamentul merge în sens invers, către valoare.
- *de la parametru la o variabilă ce utilizează parametrul*. Exemplu: presupunerea despre o monedă, care reprezintă probabilitatea de a obține “stemă”, și o aruncare a monezii.

Relații asimetrice adiționale

- *parte-întreg*. Exemplu: o imprimantă cu toner și suport pentru alimentarea cu hârtie. Defectarea unei componente duce la defectarea întregii imprimante. Alteori, proprietăți ale întregului pot influența proprietăți ale componentelor.
- *specific-general*. Exemplu: un utilizator care a avut de a face cu multe probleme specifice de tip “paper jam” sau “poor printing quality” va ști cum să rezolve o problemă mai generală de tip “poor printing” (specific → general). La crearea unui obiect, se generează mai întâi proprietățile generale și apoi cele specifice (general → specific).
- *concret/detaliat - abstract/sumar*. Exemplu: relația dintre rezultatul unui test și nota finală. Mai multe rezultate la teste influențează nota finală.

Dependențe nedirecționate

- *două efecte ale aceeași cauze.* Exemplu: două măsurători ale aceleleași valori, dar nu există o variabilă pentru acea valoare.
- *două cauze ale aceluiași efect cunoscut.* Exemplu: De regulă, suportul de alimentare cu hârtie și tonerul sunt independente. În cazul când starea generală a imprimantei este un efect (cunoscut), atunci cele două cauze devin depedente prin efect. Aceasta se numește *dependență indusă*.

Dependențe în modele grafice orientate

Factorizarea unei distribuții reduce mult complexitatea computațională:

- O distribuție bivariată poate fi factorizată ca o distribuție condițională și una necondițională:

$$P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X).$$

- Observăm că factorizarea nu este unică.

În cazul multivariat, formula înmulțirii sau formula înlănțuirii probabilităților stă la baza factorizării distribuțiilor:

- Mai întâi,

$$P(X_1, X_2, \dots, X_n) = P(X_1 | X_2, \dots, X_n) P(X_2, \dots, X_n).$$

- Continuând,

$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_1 | X_2, \dots, X_n) \cdot P(X_2 | X_3, \dots, X_n) \cdot \dots \\ &\quad \dots \cdot P(X_{n-1} | X_n) \cdot P(X_n). \end{aligned}$$

- Factorizarea poate fi efectuată în orice altă ordine.
- Din păcate, a obține cea mai bună factorizare este o problemă NP-dificilă.

O *rețea Bayesiană* este o reprezentare a unui model probabilist cu trei componente:

- 1 o mulțime de **variabile** împreună cu domeniile valorilor corespunzătoare;
- 2 un **graf aciclic orientat** în care fiecare variabilă este un nod;
- 3 pentru fiecare variabilă, o **distribuție de probabilitate condițională (CPD)** peste variabilele care definesc părinții.

Exemplu de rețea Bayesiană

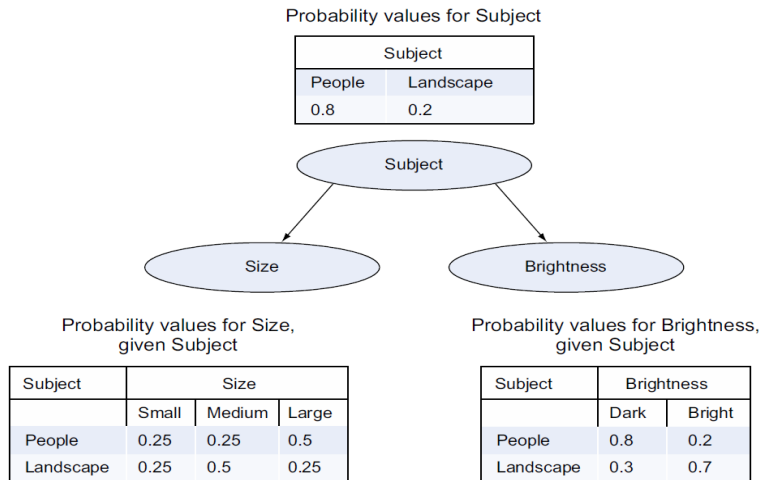
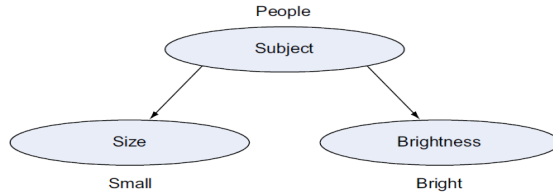


Figure 5.4 A three-node Bayesian network

$P(\text{Subject} + \text{People})$

Subject	
People	Landscape
0.8	0.2



$P(\text{Size} = \text{Small} / \text{Subject} = \text{People})$

Subject	Size		
	Small	Medium	Large
People	0.25	0.25	0.5
Landscape	0.25	0.5	0.25

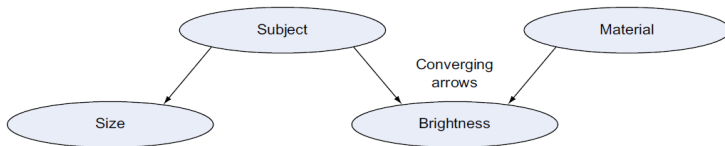
$P(\text{Brightness} = \text{Bright} / \text{Subject} = \text{People})$

Subject	Brightness	
	Dark	Bright
People	0.8	0.2
Landscape	0.3	0.7

$$P(\text{Subject} = \text{People}, \text{Size} = \text{Small}, \text{Brightness} = \text{Bright}) = 0.8 \times 0.25 \times 0.2 = 0.04$$

Convergența arcelor și dependența indusă

- 1 “Subject” și “Material” sunt independente când nimic nu e observat.
- 2 “Subject” și “Material” NU sunt independente condițional când “Brightness” este observată.



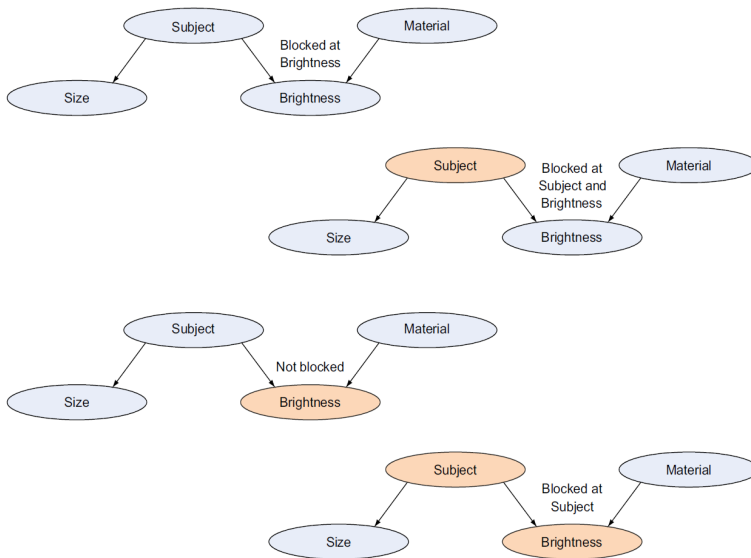
Raționamentul într-o rețea (d-separare)

- Într-o rețea, raționamentul poate evolua pe un drum atât timp cât acesta este *neblocat* la o anumită variabilă.
- În multe cazuri un drum este *blocat* la o variabilă dacă variabila este observată.

Exemplu: dacă “Subject” este observat, atunci drumul “Size-Subject-Brightness” este blocat. Dacă se observă “Size”, atunci nu se poate schimba încrederea despre “Brightness” dacă “Subject” este și ea observată.

- În alte cazuri, un drum este *blocat* la o variabilă dacă variabila este neobservată și devine *neblocat* când variabila devine observată (exemplu pe slide-ul următor).

Exemple de drumuri blocate/neblocate



- Când se raționează de la un efect X la o cauză Y , X nu este independent de Y , dar devine independent condițional când e observat Z , dacă Z blochează drumul de la X la Y .
- Același lucru se întâmplă când raționamentul merge de la o cauză la un efect indirect sau între două efecte ale aceleiași cauze.
- Pentru două cauze X și Y ale aceleiași efect Z , lucrurile merg în sens contrar. X și Y sunt independente, dar nu și independente condițional când e observat Z (din cauza dependenței induse).

Dependențe în rețele Bayesiene

- Într-un graf orientat, alegerea părinților determină semantica.
- Factorizarea întregii distribuții este definită de *semantica globală* a unui graf aciclic orientat.
- Semantica globală este specificată de

$$P(X) = \prod_{i=1:d} P(X_i | \pi(X_i)),$$

unde X_i , $1 \leq i \leq d$ reprezintă nodurile grafului, $X = (X_1, \dots, X_d)$, iar $\pi(X_i)$ mulțimea părinților lui X_i .

Exemplu de factorizare a unei distribuții

- Un student caută un loc de muncă la o companie IT;
- Angajatorul ia decizia de angajare bazat pe un test de aptitudini și pe calitatea scrisorii de recomandare de la profesorul de ML;
- Scorul la testul de aptitudini depinde doar de inteligența studentului;
- Nota la ML depinde de inteligența studentului, dar și de gradul de dificultate al materiei;
- Din păcate, profesorul nu-și mai amintește de student și își bazează scrisoarea de recomandare doar pe nota obținută.

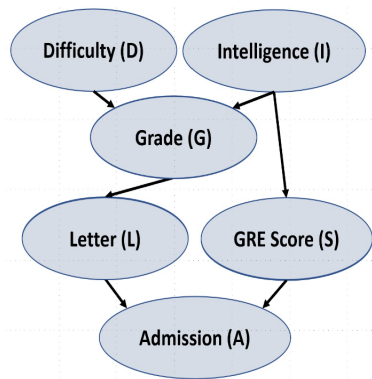
Distribuția comună pentru acest exemplu este dată de:

$$P(D, I, S, G, L, A),$$

unde

- $D \in \{0, 1\}$: dificultatea cursului de ML;
- $I \in \{0, 1\}$: inteligența studentului;
- $S \in \{0, 1\}$: scorul la testul de aptitudini;
- $G \in \{0, 1, 2\}$: nota (grade) la cursul de ML;
- $L \in \{0, 1\}$: calitatea scrisorii de recomandare;
- $A \in \{0, 1\}$: decizia de angajare.

Astfel, tabelul complet al distribuției are $2 \times 2 \times 2 \times 3 \times 2 \times 2 = 96$ intrări.



- $D \perp I$;
- $L \perp D, I \mid G$
- $S \perp D, G, L \mid I$;
- $A \perp D, I, G \mid L, S$.

Astfel,

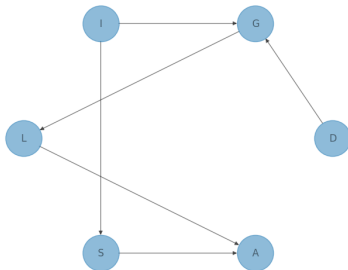
$$P(D, I, G, L, S, A) = P(D)P(I)P(G|I, D)P(L|G)P(S|I)P(A|L, S).$$

Modelarea BN cu pgmpy

```
from pgmpy.models import BayesianNetwork
from pgmpy.factors.discrete import TabularCPD
import networkx as nx

student_model = BayesianNetwork([('D', 'G'), ('I', 'G'), ('G', 'L'),
                                ('I', 'S'), ('L', 'A'), ('S', 'A')])

pos = nx.circular_layout(student_model)
nx.draw(student_model, with_labels=True, pos=pos, alpha=0.5, node_size=2000)
```



Definirea variabilelor rădăcină:

```
CPD_D = TabularCPD(variable='D', variable_card=2, values=[[0.3], [0.7]])
```

```
print(CPD_D)
```

```
CPD_I = TabularCPD(variable='I', variable_card=2, values=[[0.2], [0.8]])
```

```
print(CPD_I)
```

```
▶ +-----+-----+
  | D(0) | 0.3 |
  +-----+-----+
  | D(1) | 0.7 |
  +-----+-----+
  | I(0) | 0.2 |
  +-----+-----+
  | I(1) | 0.8 |
  +-----+-----+
```

Definirea variabilor cu un părinte:

```
CPD_L = TabularCPD(variable='L', variable_card=2,  
                    values=[[0.9, 0.6, 0.01],  
                             [0.1, 0.4, 0.99]],  
                    evidence=['G'],  
                    evidence_card=[3])  
CPD_S = TabularCPD(variable='S', variable_card=2,  
                    values=[[0.8, 0.1],  
                             [0.2, 0.9]],  
                    evidence=['I'],  
                    evidence_card=[2])  
print(CPD_L,CPD_S)
```

```
▶ +-----+-----+-----+-----+  
  | G      | G(0) | G(1) | G(2) |  
  +-----+-----+-----+-----+  
  | L(0)   | 0.9  | 0.6  | 0.01 |  
  +-----+-----+-----+-----+  
  | L(1)   | 0.1  | 0.4  | 0.99 |  
  +-----+-----+-----+-----+  
  +-----+-----+-----+  
  | I      | I(0) | I(1) |  
  +-----+-----+-----+  
  | S(0)   | 0.8  | 0.1  |  
  +-----+-----+-----+  
  | S(1)   | 0.2  | 0.9  |  
  +-----+-----+-----+
```

Definirea variabilor cu doi părinți:

```
CPD_G = TabularCPD(variable='G', variable_card=3,  
                    values=[[0.3, 0.7, 0.02, 0.2],  
                             [0.4, 0.25, 0.08, 0.3],  
                             [0.3, 0.05, 0.9, 0.5]],  
                    evidence=['I', 'D'],  
                    evidence_card=[2, 2])
```

```
print(CPD_G)
```

```
CPD_A = TabularCPD(variable='A', variable_card=2,  
                    values=[[0.9, 0.8, 0.7, 0.2],  
                             [0.1, 0.2, 0.3, 0.8]],  
                    evidence=['L', 'S'],  
                    evidence_card=[2, 2])
```

```
print(CPD_A)
```



I	I(0)	I(0)	I(1)	I(1)
D	D(0)	D(1)	D(0)	D(1)
G(0)	0.3	0.7	0.02	0.2
G(1)	0.4	0.25	0.08	0.3
G(2)	0.3	0.05	0.9	0.5
L	L(0)	L(0)	L(1)	L(1)
S	S(0)	S(1)	S(0)	S(1)
A(0)	0.9	0.8	0.7	0.2
A(1)	0.1	0.2	0.3	0.8

Adăgarea distribuțiilor condiționale la model:

```
student_model.add_cpds(CPD_D, CPD_I, CPD_S, CPD_G, CPD_L, CPD_A)  
student_model.get_cpds()
```

```
▶ [<TabularCPD representing P(D:2) at 0x17aeca4e5f0>,  
   <TabularCPD representing P(I:2) at 0x17ae87d5420>,  
   <TabularCPD representing P(S:2 | I:2) at 0x17aecb65570>,  
   <TabularCPD representing P(G:3 | I:2, D:2) at 0x17ae87d5120>,  
   <TabularCPD representing P(L:2 | G:3) at 0x17aecb66aa0>,  
   <TabularCPD representing P(A:2 | L:2, S:2) at 0x17ae87a6a70>]
```

Verificarea modelului:

```
student_model.check_model()
```

```
▶ True
```

Verificarea independențelor:

`student_model.local_independencies(['D', 'G', 'S', 'I', 'L', 'A'])`

▷ (D \perp S, I)
(G \perp S | D, I)
(S \perp G, L, D | I)
(I \perp D)
(L \perp S, D, I | G)
(A \perp G, D, I | S, L)

Inferența:

```
from pgmpy.inference import VariableElimination
```

```
infer = VariableElimination(student_model)
```

```
posterior_p = infer.query(["D", "I"], evidence={"A": 1})
```

```
print(posterior_p)
```

```
▶ +-----+-----+-----+
  | D      | I      | phi(D,I) |
  +=====+=====+=====+
  | D(0)   | I(0)   | 0.0298   |
  +-----+-----+-----+
  | D(0)   | I(1)   | 0.3291   |
  +-----+-----+-----+
  | D(1)   | I(0)   | 0.0492   |
  +-----+-----+-----+
  | D(1)   | I(1)   | 0.5918   |
  +-----+-----+-----+
```