

- dacă o persoană este bolnavă testul este 99% pozitiv
- dacă o persoană este sănătoasă testul este 0.1% pozitiv

Cifrele par a fi excelente, dar pentru comercializare suntem interesați de probabilitatea ca o persoana să fie bolnavă dacă testul este pozitiv

B evenimentul “persoana este bolnavă”

T evenimentul “testul este pozitiv”

$$P(B) = 0.0001; P(\bar{B}) = 0.9999; P(T|B) = 0.99; P(T|\bar{B}) = 0.001$$

$$P(B|T) = \frac{P(T|B) \times P(B)}{P(T|B) \times P(B) + P(T|\bar{B}) \times P(\bar{B})} = 0.09$$

Avem 9% șanse ca o persoana cu test pozitiv să fie într-adevăr bolnavă

Clasificarea Bayesiană naivă

- Realizați profilul persoanei care ia masa în oraș într-un restaurant sau un fast-food, folosind ca atribute vârsta, dacă este sau nu supraponderal, timpul disponibil (puțin/suficient) și venitul lunar, folosind următoarea mulțime de antrenament

	Vârsta	Supraponderal (sau cu tendințe)	Timp disponibil	Venit lunar (lei)	Unde mănâncă
1	39	nu	suficient	1800	Restaurant
2	26	nu	puțin	900	Fast-food
3	30	da	puțin	1100	Fast-food
4	60	nu	suficient	2400	Restaurant
5	28	nu	suficient	1700	Restaurant
6	24	nu	suficient	800	Fast-food
7	45	da	puțin	1900	Restaurant
8	56	nu	puțin	2600	Restaurant
9	27	da	suficient	2100	Restaurant
10	35	nu	puțin	1000	Fast-food
11	32	nu	suficient	800	Fast-food
12	59	da	suficient	1500	Restaurant
13	42	nu	puțin	3600	Restaurant
14	23	nu	suficient	700	Fast-food
15	38	da	suficient	2700	Restaurant
16	51	nu	puțin	1300	Restaurant
17	64	nu	suficient	2100	Restaurant
18	27	da	suficient	2500	Restaurant
19	23	da	puțin	800	Fast-food
20	43	nu	puțin	3700	Restaurant

Unde va mânca o persoană de 25 ani, care nu este supraponderală și nu are tendințe, care are timp suficient și un venit lunar de 1900 RON?

$$P(A_k | \Omega_i) = \frac{1}{\sqrt{2\pi} \cdot \sigma_{ki}} \cdot \exp\left(-\frac{(A_k - \mu_{ki})^2}{2 \cdot \sigma_{ki}^2}\right)$$

- În domeniul bancar avem problema estimării riscului acordării unui credit unei anumite persoane care are următoarele atribute: (datorii puține, garanții adecvate, venit lunar 2000 lei), utilizând următoarea mulțime de antrenament

client	RISC	datorii	garanții	venit lunar (lei)
1	înalt	multe	nu există	850
2	înalt	multe	nu există	1000
3	înalt	puține	nu există	600
4	înalt	puține	nu există	500
5	scăzut	puține	nu există	1800
6	înalt	puține	adecvate	500
7	înalt	puține	nu există	700
8	scăzut	puține	nu există	1600
9	scăzut	puține	nu există	2800
10	scăzut	multe	adecvate	1100
11	înalt	multe	nu există	500
12	înalt	multe	nu există	600
13	scăzut	multe	nu există	1600
14	înalt	multe	nu există	1400
15	înalt	multe	adecvate	450
16	înalt	puține	nu există	700
17	scăzut	puține	adecvate	1200
18	scăzut	puține	adecvate	3200
19	scăzut	puține	adecvate	1100
20	înalt	multe	nu există	400

- O companie dorește efectuarea unui studiu asupra timpului petrecut de telespectatori privind reclamele televizate la diverse produse și servicii. Se urmărește realizarea profilului telespectatorului care urmărește cel puțin jumătate din timpul acordat publicității, în cadrul unui program. Desenați arborele de clasificare și decizie corespunzător folosind trei măsuri studiate. Avem la dispoziție următoarea bază de date.

	Timp vizionare	Venit	Domiciliu(urban/rural)	Vârsta	Salariat/șomer/pensionar
1	<1/2	1000	urban	25	salariat
2	<1/2	1800	urban	30	salariat
3	>1/2	800	rural	40	șomer
4	<1/2	2500	urban	35	salariat
5	<1/2	400	rural	45	salariat
6	>1/2	700	urban	65	pensionar
7	<1/2	3500	urban	50	salariat
8	>1/2	300	rural	28	șomer
9	<1/2	500	urban	51	șomer
10	>1/2	1100	rural	62	pensionar
11	>1/2	1500	rural	48	salariat
12	>1/2	900	urban	69	pensionar
13	<1/2	2000	rural	41	salariat
14	>1/2	1000	urban	37	salariat
15	<1/2	450	rural	60	pensionar
16	<1/2	1700	rural	41	salariat
17	>1/2	1300	urban	34	salariat
18	<1/2	400	urban	32	șomer
19	<1/2	1500	rural	46	salariat

20	>1/2	900	urban	71	pensionar
21	>1/2	600	urban	52	șomer

- O firmă de asigurări prezintă următoarele date referitoare la asigurații Casco; pentru o nouă politică în sensul aplicării unei ponderi clienților predispuși la accidente, construiți un arbore de clasificare și decizie pe baza acestor date și pe baza acestuia deduceți regulile de clasificare.

risc	vârsta	tipul mașinii	vechime	a avut accident
înalt	peste 30	sport	sub 2 ani	da
moderat	peste 30	familial	sub 2 ani	nu
scăzut	peste 30	familial	peste 2 ani	nu
înalt	sub 30	sport	peste 2 ani	da
înalt	sub 30	familial	peste 2 ani	da
scăzut	peste 30	familial	peste 2 ani	nu
moderat	sub 30	sport	peste 2 ani	nu
moderat	peste 30	familial	sub 2 ani	da
moderat	sub 30	sport	peste 2 ani	nu
înalt	peste 30	sport	sub 2 ani	da

- Într-un supermarket, pe parcursul unei zile au loc 100 000 de tranzacții; dintre acestea 7000 conțin cafea, 1000 conțin lapte condensat și 10000 conțin biscuiți. Cafea și lapte condensat se regăsesc în 800 tranzacții, biscuiți și lapte condensat în 300 tranzacții, biscuiți și cafea în 2000 tranzacții în timp ce cafea, lapte condensat și bicuiți se găsesc în 100 de tranzacții.

Calculați suportul, încrederea și diferența de nivel pentru regulile cu trei articole, respectiv două articole, deduse din datele prezentate. Deduceți care este regula cea mai bună.

$$\text{Confidence} = \frac{\text{frecventa aparitiilor}}{\text{frecventa antecedent}}$$

$$\text{Diferența de nivel} = \frac{\text{confidence}}{\text{frecventa consecinta}}$$

- Aplicați metoda k-nearest neighbor luând ca mulțime de antrenament aceasta bază de date restrânsă de iriși, în stabilirea tipului de iris ce are următoarele caracteristici (4,15,25,49)

Tipul de Iris: 0 Setosa; 1 Virginica; 2 Versicolor

PW lățimea petalei; PL lungimea petalei;

SW lățimea sepalei; SL lungimea sepalei

Type	PW	PL	SW	SL
0	2	14	33	50
1	24	56	31	67
1	23	51	31	69
0	2	10	36	46
1	20	52	30	65
1	19	51	27	58
2	13	45	28	57
2	16	47	33	63
1	17	45	25	49

2	14	47	32	70
0	2	16	31	48
1	19	50	25	63
0	1	14	36	49
0	2	13	32	44
2	12	40	26	58
1	18	49	27	63
2	10	33	23	50
0	2	16	38	51
0	2	16	30	50
1	21	56	28	64
0	4	19	38	51
0	2	14	30	49
2	10	41	27	58
2	15	45	29	60
0	2	14	36	50
1	19	51	27	58
0	4	15	34	54
1	18	55	31	64
2	10	33	24	49
0	2	14	42	55
2	14	39	27	52
2	12	39	27	58
1	23	57	32	69
2	15	42	30	59
2	13	44	23	63

- Considerăm un lot de 14 pacienți dintre care 6 au cancer hepatic (HCC). Pentru a decide dacă un pacient are sau nu cancer hepatic, studiile clinice arată că se iau în considerare anumite enzime serice (de obicei un număr de 15). Pentru simplificarea calculului, le vom lua în considerare pe cele mai importante din punct de vedere clinic: alkaline phosphatase (FA), g-glutamyl transferase (gGT), leucine amino peptidase (LAP) și colesterol (C)

	FA	gCT	LAP	C	diagnostic
P1	162	255	74	258	non HCC
P2	210	324	93	220	non HCC
P3	259	208	115	266	non HCC
P4	120	114	89	171	non HCC
P5	246	173	98	210	non HCC
P6	138	189	76	165	non HCC
P7	132	177	48	138	non HCC
P8	152	183	105	178	non HCC
P9	422	488	183	292	HCC
P10	607	259	65	275	HCC
P11	446	283	176	309	HCC
P12	460	1053	145	221	HCC
P13	680	381	280	275	HCC
P14	561	450	180	164	HCC

Utilizând metoda k-nearest neighbor stabiliți diagnosticul următorilor pacienți:
P' (180,119,114,171) și P'' (236,270,88,150)