

Adatbányászati feladatgyűjtemény tehetséges hallgatók számára

Buza Krisztián

Számítástudományi és Információelméleti Tanszék
Budapesti Műszaki és Gazdaságtudományi Egyetem

Tartalomjegyzék

- Modellek kiértékelése 3
- Mátrix faktorizáció 25
- Osztályozó algoritmusok 35
- További feladatok 61

Modellek kiértékelése

- Accuracy számítása
- Confusion Mátrix
- Precision, Recall, F-Measure számítása
- ROC görbe szerkesztése,
görbe alatti terület számítása

Háttér (1/3)

Egy hitelintézet különböző adatokat rögzít az ügyfelekről, ezek közé tartozik, hogy egy-egy ügyfélnek volt-e késedelme a korábban folyósított hitel törlesztő részleteinek visszafizetése során.

A hitelintézet két különböző elemző céget bíz meg azzal, hogy egy-egy modellt fejlesszenek, amely előrejelzi, hogy egy-egy adott ügyfél mekkora valószínűséggel lesz késedelmes a hitel visszafizetése során.

Tehát összesen két előrejelző modell van, az egyiket az egyik cég fejleszti, a másikat a másik.

Háttér (2/3)

A hitelintézet a „software as service” elv szerint azzal arányosan fizet az egyik illetve másik elemző cégnek, hogy az adott cég modelljét hány alkalommal használta a hitelintézet előrejelzésre.

Ezért a hitelintézet abban érdekelt, hogy eldöntse, melyik modell teljesít jobban, és csak a jobban teljesítő modellt használja, illetve – különlegesen kockázatos esetekben – mind a két modell szerinti előrejelzést figyelembe veszi.

A modellek kiértékeléséhez adott egy teszt adatbázis, amely független a modellek készítésekor (tanításakor) használt adatbázistól (ugyanazt az ügyfelet nem tartalmazza mind a két adatbázis). A modellek készítésekor használt adatbázishoz hasonlóan, a teszt adatbázis esetén is tudjuk, hogy melyik ügyfél volt késedelmes és melyik nem (a teszt adatbázis is múltbeli eseményeket ír le).

Háttér (3/3)

- Pozitív osztály: késedelmes ügyfelek
- Negatív osztály: jó ügyfelek

- Tesztadatok:

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28

Accuracy

- **1. Feladat: accuracy számítás**
- Válasszuk küszöbszámnak 0.5-t, azaz:
ha a modell jóslata szerint legalább 0.5 valószínűséggel tartozik egy példány a pozitív osztályba, akkor tekintsük úgy, hogy a modell az adott példányt a pozitív osztályba sorolta.
- Számoljuk ki mindkét modell accuracy-ját!
- Tesztadatok:

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28

Accuracy

- **Megoldás: accuracy számítás**
- accuracy = helyesen osztályozott példányok száma osztva az összes tesztpéldány számával
- Model 1: accuracy = $9 / 12 = 0.75$
- Model 2: accuracy = $8 / 12 = 0.66$

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	-	-	-	-	-	-	-	-	-	-	-	-
Model 2	+	+	+	-	+	-	+	-	-	-	-	-

Accuracy

2. feladat

- Tényleg jobb-e az első modell a másodiknál?
- Miért (mikor) lehet félrevezető az accuracy?

Accuracy

Megoldás

- Tényleg jobb-e az első modell a másodiknál?
 - Lásd következő feladatok
- Miért (mikor) lehet félrevezető az accuracy?
 - Ha kiegyensúlyozatlan az osztályok eloszlása
 - Minden hibát egyformának tekint
 - (pozitív ügyfél negatív osztályba sorolása illetve negatív ügyfél pozitív osztályba sorolása)
 - Nem független a választott küszöbszámtól

Confusion mátrix

- **3. Feladat: confusion mátrix számítás**
- Válasszuk küszöbszámnak 0.5-t, azaz:
ha a modell jóslata szerint legalább 0.5 valószínűséggel tartozik egy példány a pozitív osztályba, akkor tekintsük úgy, hogy a modell az adott példányt a pozitív osztályba sorolta.
- Adjuk meg mindkét modell confusion mátrix-át!
(true positiv-ok, true negative-ok, false positive-ok, false negative-ok számát)
- Tesztadatok:

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28

Confusion mátrix

- **Megoldás: confusion mátrix számítás**

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28



Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	-	-	-	-	-	-	-	-	-	-	-	-
Model 2	+	+	+	-	+	-	+	-	-	-	-	-

Model 1	osztályozó szerint		
	+	-	
Valóságban	+	0	3
	-	0	9

TP: 0, TN: 9, FP: 0, FN: 3

Model 2	osztályozó szerint		
	+	-	
Valóságban	+	2	1
	-	3	6

TP: 2, TN: 6, FP: 3, FN: 1

Súlyozott accuracy

- **4. Feladat:**

Az előbbi confusion mátrixok alapján (az előbbi küszöbszám mellett) számítsuk a modellek súlyozott accuracy-ját. Egy pozitív ügyfél negatív csoportba történő besorolásának súlya legyen kétszer akkora, mint az összes többi eset súlya.

Model 1		osztályozó szerint	
		+	-
Valóságban	+	0	3
	-	0	9

Model 2		osztályozó szerint	
		+	-
Valóságban	+	2	1
	-	3	6

Súlyozott accuracy

- Megoldás:**

Az előbbi confusion mátrixok alapján (az előbbi küszöbszám mellett) számítsuk a modellek súlyozott accuracy-ját. Egy pozitív ügyfél negatív csoportba történő besorolásának súlya legyen kétszer akkora, mint az összes többi eset súlya.

Model 1		osztályozó szerint	
		+	-
Valóságban	+	0	3
	-	0	9

Model 2		osztályozó szerint	
		+	-
Valóságban	+	2	1
	-	3	6

$$(0*w + 9*w) / (0*w + 3*2*w + 0*w + 9*w) = 9 / 15$$

$$(2*w + 6*w) / (2*w + 1*2*w + 3*w + 6*w) = 8 / 13$$

Precision, Recall, F-Measure

- **5. Feladat:**

Az előbbi confusion mátrixok alapján (az előbbi küszöbszám mellett) számítsuk a modell precision-ját, recall-ját, és F-measure-jét!

Model 1		osztályozó szerint	
		+	-
Valóságban	+	0	3
	-	0	9

Model 2		osztályozó szerint	
		+	-
Valóságban	+	2	1
	-	3	6

Precision, Recall, F-Measure

- **Megoldás:**

Az előbbi confusion mátrixok alapján (az előbbi küszöbszám mellett) számítsuk a modell precision-ját, recall-ját, és F-measure-jét!

Model 1		osztályozó szerint	
		+	-
Valóságban	+	0	3
	-	0	9

Model 2		osztályozó szerint	
		+	-
Valóságban	+	2	1
	-	3	6

P: $0 / 0 = ?$ (Java-ban: NaN)

R: $0 / 3 = 0$

F: $2 * 0 * (0/0) / (0 + (0/0)) = ?$ (NaN)

P: $2 / 5 = 0.4$

R: $2 / 3 = 0.67$

F: $2 * 0.4 * 0.67 / (0.4 + 0.67) = 0.5$

AUC (Area Under ROC curve)

- **6. feladat**

Szerkesszük meg az előbbi modellek ROC görbáját és számoljuk ki a görbe alatti területet!

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28

AUC (Area Under ROC curve)

- **Megoldás**
- 1. lépés: rendezzük a példányokat a modell kimenete szerint

Valóság	-	-	+	-	-	-	+	-	-	-	+	-
Model 1	0.2	0.1	0.49	0.26	0.3	0.12	0.31	0.2	0.1	0.32	0.4	0.2
Model 2	0.6	0.51	0.8	0.12	0.54	0.39	0.53	0.46	0.41	0.37	0.49	0.28



Valóság	-	-	-	-	-	-	-	-	+	-	+	+
Model 1	0.1	0.1	0.12	0.2	0.2	0.2	0.26	0.3	0.31	0.32	0.4	0.49

Valóság	-	-	-	-	-	-	+	-	+	-	-	+
Model 2	0.12	0.28	0.37	0.39	0.41	0.46	0.49	0.51	0.53	0.54	0.6	0.8

AUC (Area Under ROC curve)

- 2. lépés: meghatározzuk TP-k és FP-k arányát minden értelmes küszöbszámra

Valóság	-	-	-	-	-	-	-	-	+	-	+	+	+
Model 1	0.1	0.1	0.12	0.2	0.2	0.2	0.26	0.3	0.31	0.32	0.4	0.49	
TP	3		3	3			3	3	3	2	2	1	0
FP	9	7	6				3	2	1	1	0	0	0
TN	0	2	3			6	7	8	8	9	9	9	9
FN	0	0	0			0	0	0	0	1	1	2	3
TPR	1		1	1			1	1	1	2/3	2/3	1/3	0
FPR	1	7/9	6/9				3/9	2/9	1/9	1/9	0	0	0

AUC (Area Under ROC curve)

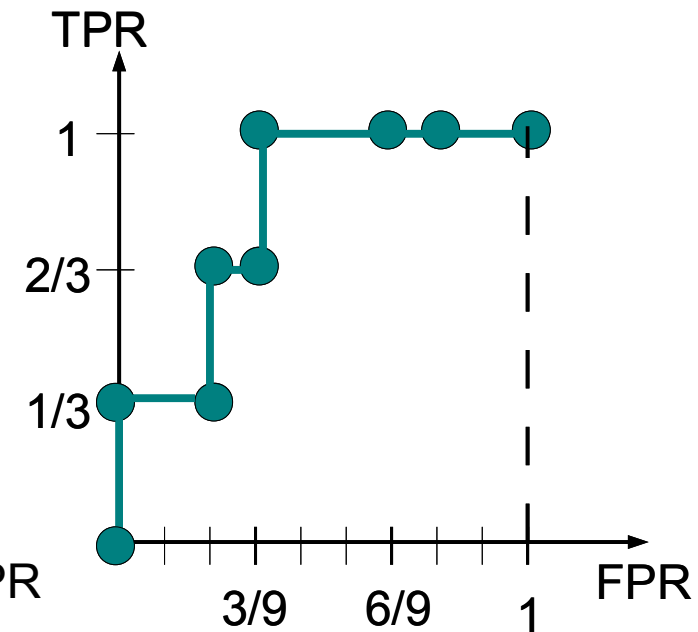
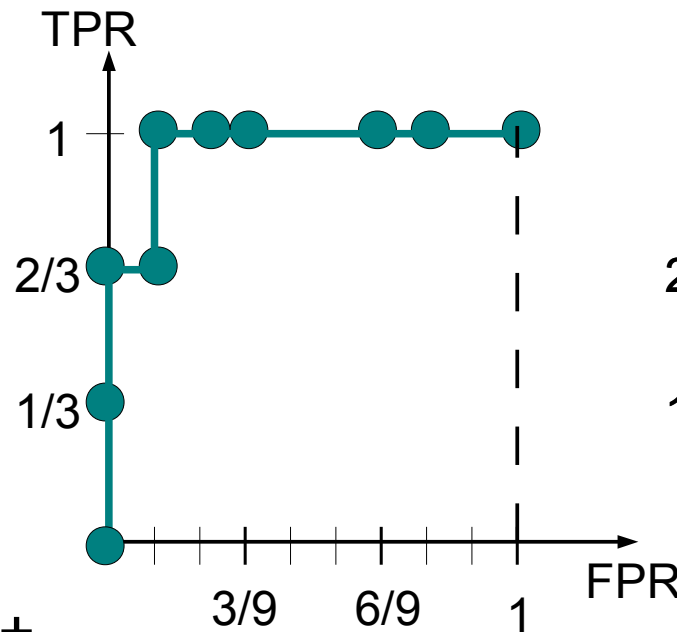
- 2. lépés: meghatározzuk TP-k és FP-k arányát minden értelmes küszöbszámra

Valóság	-	-	-	-	-	-	+	-	+	-	-	+	
Model 2	0.12	0.28	0.37	0.39	0.41	0.46	0.49	0.51	0.53	0.54	0.6	0.8	
TP	3	3	3	3	3	3	3	2	2	1	1	1	0
FP	9	8	7	6	5	4	3	3	2	2	1	0	0
TN	0	1	2	3	4	5	6	6	7	7	8	9	9
FN	0	0	0	0	0	0	0	1	1	2	2	2	3
TPR	1	1	1	1	1	1	1	2/3	2/3	1/3	1/3	1/3	0
FPR	1	8/9	7/9	6/9	5/9	4/9	3/9	3/9	2/9	2/9	1/9	0	0

AUC (Area Under ROC curve)

- 3. lépés:
ábrázoljuk
TP-k és FP-k
arányát

Görbe alatti terület:
 $(1/9) \cdot (2/3) + (8/9) \cdot 1 =$
 $= \mathbf{26/27}$
 $(1/3) \cdot (2/9) + (2/3) \cdot (1/9) +$
 $+ 1 \cdot (6/9) = \mathbf{22/27}$



...
TPR	1	1	1	1	1	1	1	1	2/3	2/3	1/3	0	0
FPR	1	7/9	6/9	3/9	2/9	1/9	1/9	0	0	0	0	0	0

...
TPR	1	1	1	1	1	1	1	2/3	2/3	1/3	1/3	1/3	0
FPR	1	8/9	7/9	6/9	5/9	4/9	3/9	3/9	2/9	2/9	1/9	0	0

ROC

- **7. feladat**

a) Hogy néz ki a „tökéletes” modell ROC görbéje? Mennyi a görbe alatti terület?

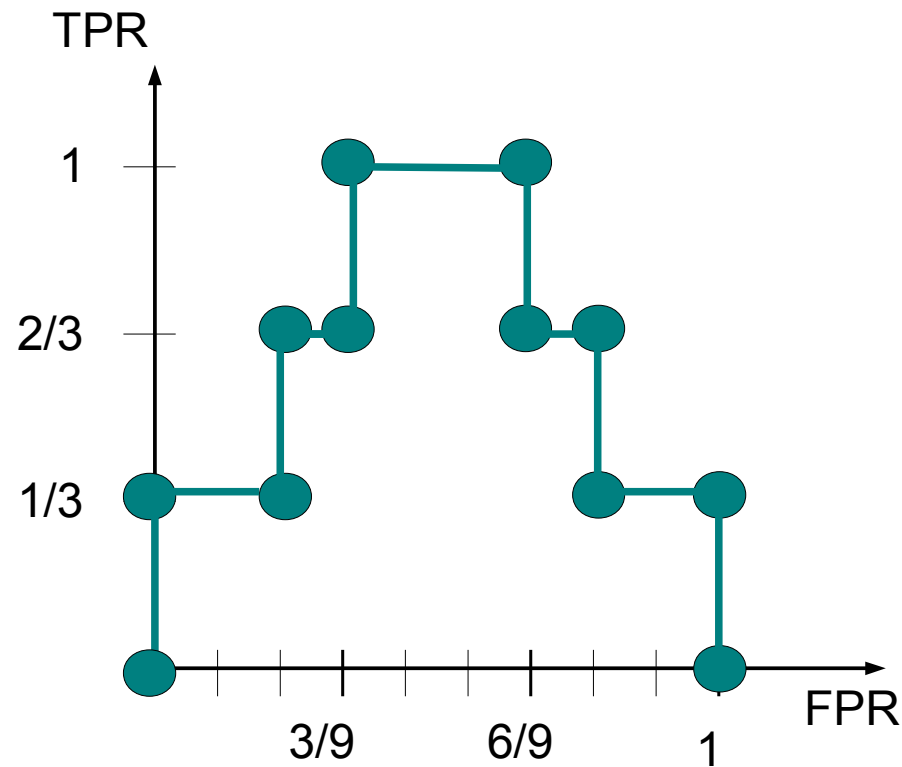
b) Hogy néz ki egy véletlenszerű előrejelzést adó modell ROC görbéje? Mennyi a görbe alatti területe?

- (Véletlenszerű előrejelzést adó modell:
a modell által előrejelzett osztálycímke független a példányok valóságos osztálycímkejétől)

c) Mi a ROC alatti terület előnye accuracy-val szemben?

ROC

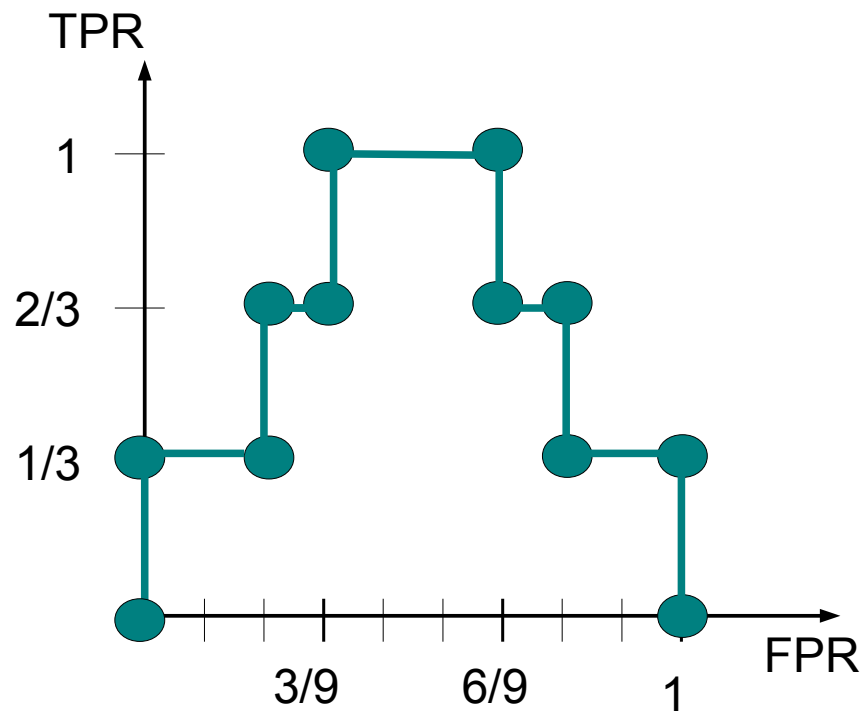
- **8. Feladat:** Lehet-e a ROC görbe ilyen?



ROC

- **Megoldás:**

Nem: a TP-k aránya nem csökkenhet az FP-k arányának növekedése mellett



Mátrix faktorizáció

Háttér



5	?	4	?	...
?	4	?	?	...
?	5	4	?	...
4	?	4	5	...
...

Online kereskedelmi rendszerben nyilván tartjuk, hogy melyik felhasználó milyen terméket (filmet) vásárolt és arról milyen visszajelzést adott (mennyire volt elégedett a termékkel). Ezek alapján kívánunk személyre szabott reklámokat küldeni a felhasználóknak.

Mátrix faktorizáció algoritmus

9. Feladat:

Hajtsa végre a tanult, ritka mátrixokat faktorizáló iteratív algoritmus egy javítási lépést az ábrában jelölt elemmel! Adott: $\varepsilon = 0.05$

$$\begin{array}{|c|c|} \hline 2 & 1 \\ \hline 2 & 2 \\ \hline 3 & 2 \\ \hline 1 & 1 \\ \hline 3 & 1 \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline 2 & 2 & 1 & 3 & 1 \\ \hline 1 & 0 & 3 & 3 & -1 \\ \hline \end{array} \approx \begin{array}{|c|c|c|c|c|} \hline 5 & ? & 4 & ? & 1 \\ \hline ? & 4 & ? & ? & ? \\ \hline ? & 4 & 4 & ? & 3 \\ \hline 4 & ? & 4 & 5 & ? \\ \hline 4 & ? & ? & 4 & 2 \\ \hline \end{array}$$

U **V** **M**

Mátrix faktorizáció algoritmus

Megoldás:

2	1
2	2
3	2
1	1
3	1

 \times

2	2	1	3	1
1	0	3	3	-1

 \approx

5	?	4	?	1
?	4	?	?	?
?	4	4	?	3
4	?	4	5	?
4	?	?	4	2

U V M

$$(3 \times 2) + (2 \times 0) = 6$$

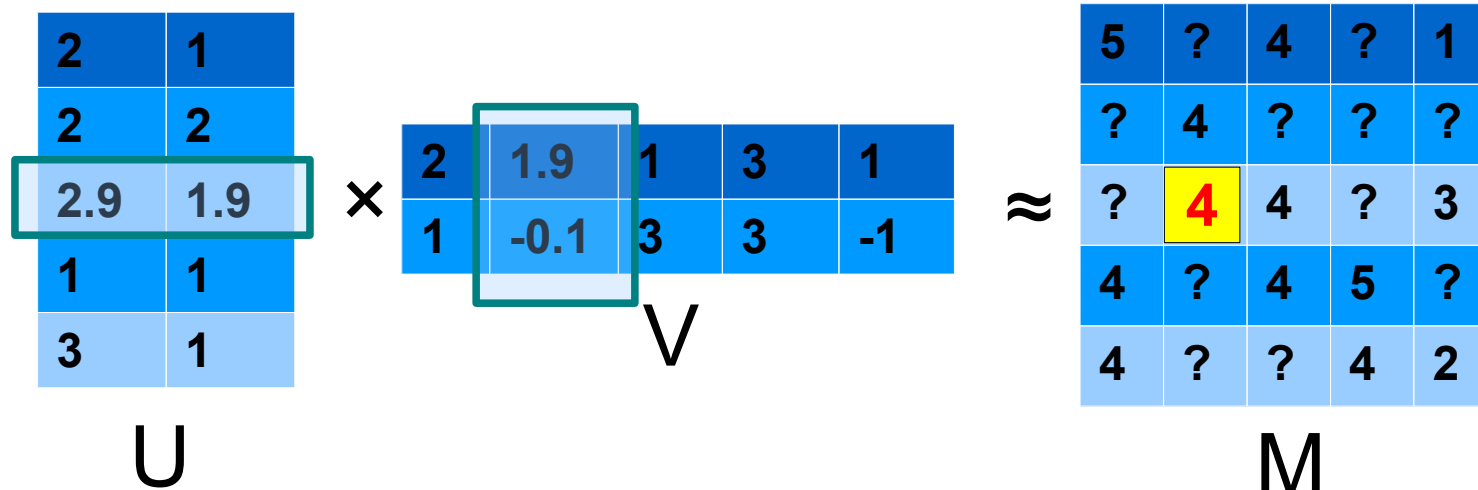
nagyobb, mint a választott elem \rightarrow csökkentjük a U és V aktuális sorában ill. oszlopában lévő számokat

$6 - 4 = 2$, $\epsilon = 0.05 \rightarrow 0.05 * 2 = 0.1$ -gyel csökkentjük a két jelölt vektor elemeit

- Megjegyzés:** valós algoritmusok ennél valamivel összetettebb frissítési lépést használnak, lásd pl. Yehuda Koren, Robert Bell, Chris Volinsky: *Matrix factorization techniques for Recommender Systems*, IEEE Computer, Aug. 2009

Mátrix faktorizáció algoritmus

Megoldás:



$$(3 \times 2) + (2 \times 0) = 6$$

nagyobb, mint a választott elem \rightarrow csökkentjük a U és V aktuális sorában ill. oszlopában lévő számokat

$6 - 4 = 2$, $\epsilon = 0.05 \rightarrow 0.05 * 2 = 0.1$ -gyel csökkentjük a két jelölt vektor elemeit

- Megjegyzés:** valós algoritmusok ennél valamivel összetettebb frissítési lépést használnak, lásd pl. Yehuda Koren, Robert Bell, Chris Volinsky: *Matrix factorization techniques for Recommender Systems*, IEEE Computer, Aug. 2009

Mátrix faktorizáció algoritmus

10. Feladat:

Kellőképpen sok iteráció után az algoritmus az alábbi U és V mátrixokat találja meg. Ezen U és V mátrixok alapján becsüljük meg M első sorbeli hiányzó értékeit!

1.57	0.90
1.77	1.13
1.13	1.76
1.04	1.47
1.44	0.86

U

\times

1.86	1.37	1.31	1.67	0.47
1.64	1.38	1.70	2.08	1.19

V

\approx

5	?	4	?	1
?	4	?	?	?
?	4	4	?	3
4	?	4	5	?
4	?	?	4	2

M

Mátrix faktorizáció algoritmus

Megoldás:

M hiányzó elemeinek becslése az U és V mátrix szorzataként adódik.

1.57	0.90
1.77	1.13
1.13	1.76
1.04	1.47
1.44	0.86

U

\times

1.86	1.37	1.31	1.67	0.47
1.64	1.38	1.70	2.08	1.19

V

\approx

5	A	4	B	1
?	4	?	?	?
?	4	4	?	3
4	?	4	5	?
4	?	?	4	2

M

$$A = (1.57 * 1.37) + (0.9 * 1.38) = 3.39$$

$$B = (1.57 * 1.67) + (0.9 * 2.08) = 4.49$$

Mátrix faktorizáció algoritmus

- **11. feladat:**

a) Milyen probléma adódhat ha ε -t túl kicsinek vagy túl nagyoknak választjuk?

b) Hogyan kerülhetjük el a túltanulást (overfitting-et) ?

Mátrix faktorizáció algoritmus

- **Megoldás**

a) Milyen probléma adódhat ha ε -t túl kicsinek vagy túl nagyoknak választjuk?

Ha ε túl kicsi \rightarrow túl keveset csökkentünk \rightarrow
lassan találjuk meg a legjobb közelítést

Ha ε túl nagy \rightarrow túl sokat csökkentünk \rightarrow
túl nagy léptekben haladunk, lehet, hogy nem
jutunk az optimumhoz, hanem „átlépünk” felette

Mátrix faktorizáció algoritmus

- **Megoldás:**

b) Hogyan kerülhetjük el a túltanulást (overfitting-et) ?

- Hold-out módszer:

Az ismert elemek (tanító adatok) közül kiválasztunk néhányat, amelyeket „félre teszünk”, azaz amelyek sohasem lesznek a „kiválasztott” elemek a javítási lépések során. Egy-egy javítás végrehajtása után megnézzük nem romlott-e (sokat) a „félretett elemek” becslése.

- Regularizáció

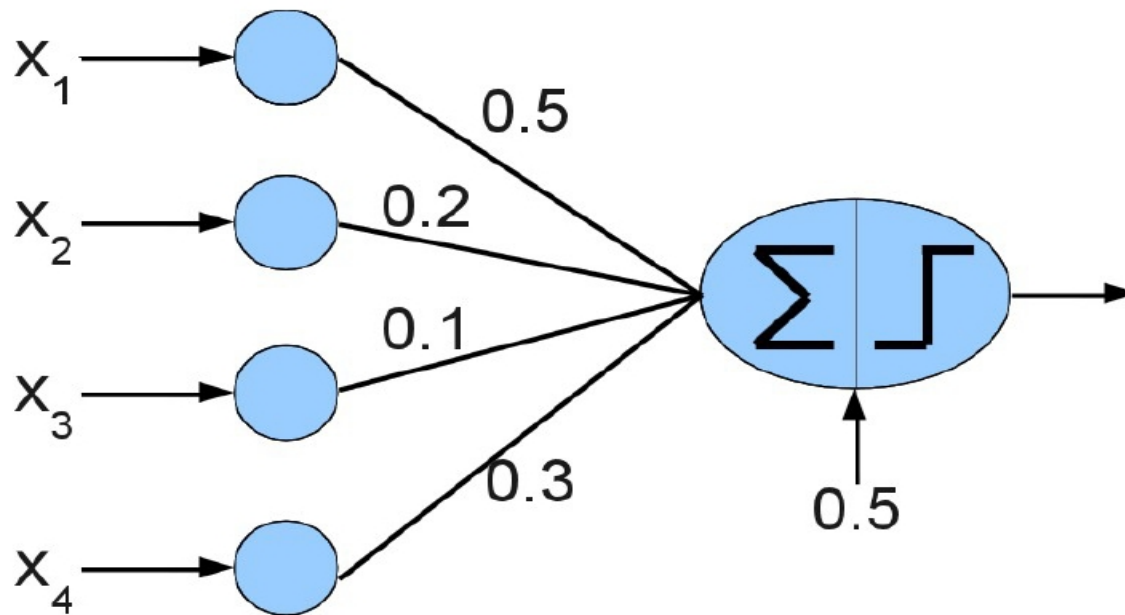
Pl. U és V oszlopainak ill. sorainak számát viszonylag kicsire választjuk

Osztályozó algoritmusok

- Perceptron algoritmus
- Naive Bayes
- Döntési fák

Perceptron algoritmus

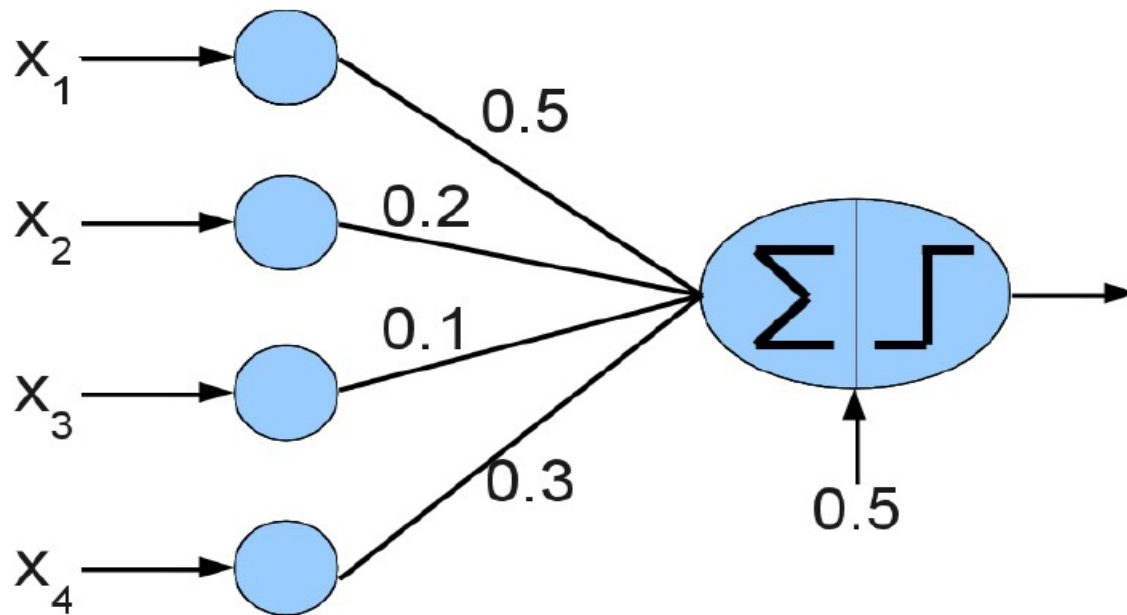
- **12. Feladat:** Adott az alábbi perceptron. A két osztálycímekét +1-gyel és -1-gyel jelöljük.
a) Hogyan osztályozza ez a perceptron az $(x_1=1, x_2=-0.8, x_3=-0.3, x_4=1.5)$ példányt?



Perceptron algoritmus

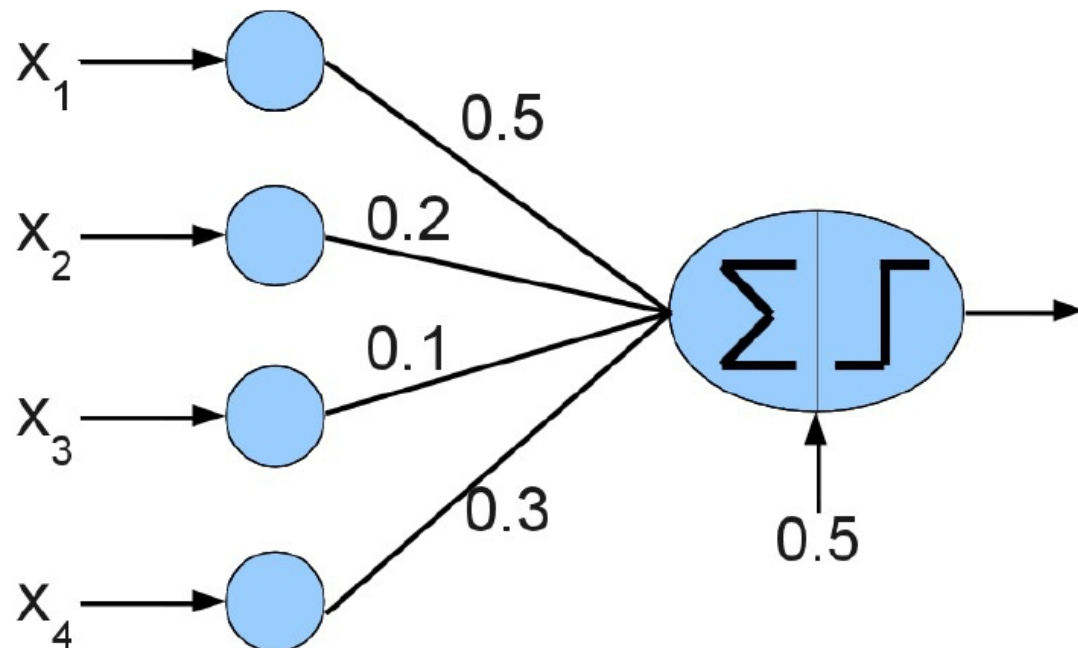
- **Megoldás:**

$1*0.5+(-0.8)*0.2+(-0.3)*0.1+1.5*0.3 - 0.5 = 0.26$, ami nagyobb 0-nál \rightarrow +1 osztályba sorolja



Perceptron algoritmus

- **12. Feladat:** b) Az ($x_1=1, x_2=-0.8, x_3=-0.3, x_4=1.5$) valódi osztálycímkeje -1 , és ezt a példányt a perceptron tanításához használjuk, $\lambda = 0.1$. Hajtsa végre a tanítási algoritmus egy javító lépését! Változnak-e az élek súlyai, ha igen, hogyan?



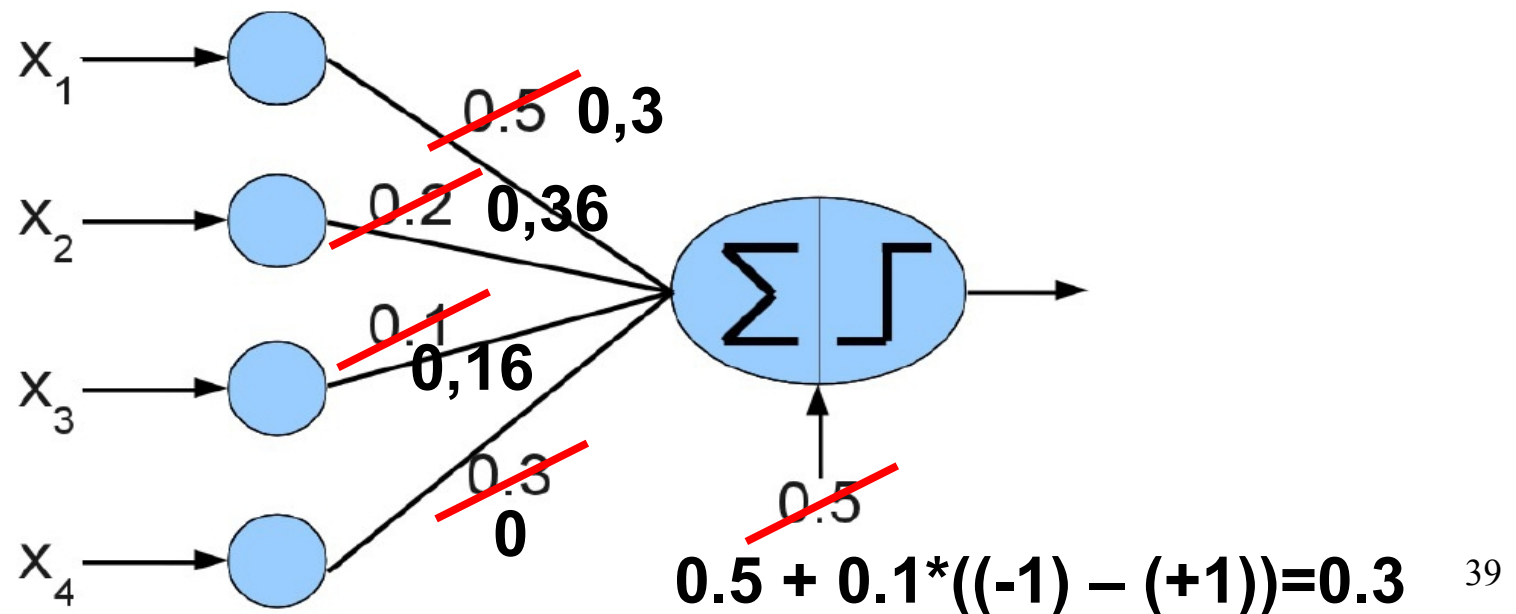
Perceptron algoritmus

- **Megoldás:**

b) i-dik új súly = i-dik régi súly +
+ λ *(valódi osztálycímke – felismert osztálycímke)* x_i

Tehát: $w_1^{\text{új}} = 0.5 + 0.1 * ((-1) - (+1)) * 1 = 0.3$

$$w_2^{\text{új}} = 0.2 + 0.1 * ((-1) - (+1)) * (-0.8) = 0.36$$



Naive Bayes

- **13. Feladat**

Étteremláncunk TelePocak néven új akciót tervez indítani. Néhány embert megkérdeztünk arról, hogy érdekelné őket az akció. Ez alapján készítendő egy olyan osztályozó, amely eldönti, hogy kiket érdekelhet még az akció.

Adott megkérdezett emberek adatait tartalmazó adatbázis (következő lap), készítsen egy Naive Bayes osztályozót és ezzel osztályozza a további potenciális ügyfeleket.

Naive Bayes

A megkérdezés eredménye

Sorsz.	Életkor	Testsúly	Sportol	Érdekli-e az akció
1	fiatal	alacsony	igen	igen
2	idős	közepes	nem	nem
3	középkorú	magas	nem	igen
4	idős	közepes	igen	nem
5	fiatal	magas	nem	igen
6	középkorú	alacsony	nem	nem
7	idős	alacsony	nem	nem
8	fiatal	közepes	nem	igen
9	középkorú	magas	igen	igen
10	idős	közepes	igen	nem

Naive Bayes

Osztályozandó potenciális ügyfelek

Sorsz.	Életkor	Testsúly	Sportol	Érdekli-e az akció
11	fiatal	közepes	igen	?
12	idős	alacsony	igen	?
13	középkorú	közepes	igen	?
14	középkorú	közepes	nem	?

Naive Bayes

- **Megoldás**

- Jelölés:

Életkor $\rightarrow X$, Testsúly $\rightarrow Y$, sportol $\rightarrow Z$, érdekli-e $\rightarrow C$

- $P(C | X, Y, Z)$ -t szeretnénk becsülni.

- Bayes-tétel szerint $P(C | X, Y, Z)$ egyenlő

$$\frac{P(X, Y, Z | C) * P(C)}{P(X, Y, Z)}$$

- Elég a számlálót kiszámolnunk (lásd később)

Naive Bayes

- Naive Bayes feltételezi, hogy az attribútumok adott osztályváltozó („érdekli-e”) esetén feltételesen függetlenek:

$$P(X, Y, Z | C) = P(X | C) * P(Y | C) * P(Z | C)$$

- Ezért tehát $P(X | C)$ -t, $P(Y | C)$ -t, $P(Z | C)$ -t és $P(C)$ -t kell becsülnünk a megadott tanító adatok alapján.

Naive Bayes

Sor.	X	Y	Z	C
1	fiatal	alacsony	igen	igen
2	idős	közepes	nem	nem
3	középkorú	magas	nem	igen
4	idős	közepes	igen	nem
5	fiatal	magas	nem	igen
6	középkorú	alacsony	nem	nem
7	idős	alacsony	nem	nem
8	fiatal	közepes	nem	igen
9	középkorú	magas	igen	igen
10	idős	közepes	igen	nem

$$P(X=fiatal \mid C=igen) = 3/5$$

$$P(X=idős \mid C=igen) = 0$$

$$P(X=középk. \mid C=igen) = 2/5$$

$$P(X=fiatal \mid C=nem) = 0$$

$$P(X=idős \mid C=nem) = 4/5$$

$$P(X=középk. \mid C=nem) = 1/5$$

$$P(Y=alacsony \mid C=igen) = 1/5$$

$$P(Y=közepes \mid C=igen) = 1/5$$

$$P(Y=magas \mid C=igen) = 3/5$$

$$P(Y=alacsony \mid C=nem) = 2/5$$

$$P(Y=közepes \mid C=nem) = 3/5$$

$$P(Y=magas \mid C=nem) = 0$$

Naive Bayes

- **Megoldás (folyt.)**

Sor.	X	Y	Z	C
1	fiatal	alacsony	igen	igen
2	idős	közepes	nem	nem
3	középkorú	magas	nem	igen
4	idős	közepes	igen	nem
5	fiatal	magas	nem	igen
6	középkorú	alacsony	nem	nem
7	idős	alacsony	nem	nem
8	Fiatal	közepes	nem	igen
9	középkorú	magas	igen	igen
10	idős	közepes	igen	nem

$$P(Z=igen \mid C=igen) = 2/5$$

$$P(Z=nem \mid C=igen) = 3/5$$

$$P(Z=igen \mid C=nem) = 2/5$$

$$P(Z=nem \mid C=nem) = 3/5$$

$$P(C=igen) = 5/10$$

$$P(C=nem) = 5/10$$

Naive Bayes

Sorsz.	X	Y	Z	C
11	fiatal	közepes	igen	?

- A becsült valószínűségek alapján kiszámoljuk, hogy melyik osztálynak mekkora a valószínűsége az osztályozandó példányok esetében

$$P(C=\text{igen} \mid X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen}) =$$

$$= \frac{P(C=\text{igen}) * P(X=\text{fiatal} \mid C=\text{igen}) * P(Y=\text{közepes} \mid C=\text{igen}) * P(Z=\text{igen} \mid C=\text{igen})}{P(X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen})}$$

$$= \frac{5/10 * 3/5 * 1/5 * 2/5}{P(X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen})} = \frac{6 / 250}{P(X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen})}$$

Hasonlóképpen: $P(C=\text{nem} \mid X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen}) =$

$$= \frac{5/10 * 0 * 3/5 * 2/5}{P(X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen})} = \frac{0}{P(X=\text{fiatal}, Y=\text{közepes}, Z=\text{igen})}^{17}$$

Döntési fa

- **14. Feladat:**

A TelePocak akció adatbázisa alapján építendő egy döntési fa, amely képes eldönteni, hogy mely ügyfeleket érdekli potenciálisan az akció. A döntési fa építésekor a GINI-index szerinti javítás alapján válassza ki a legjobb vágást, csak multiway split-eket vizsgáljon. A vágás entrópiáját (SplitINFO) ne vegye figyelembe.

Döntési fa

- **Megoldás**

Sorsz.	Életkor	Testsúly	Sportol	Érdekl-e az akció
1	fiatal	alacsony	igen	igen
2	idős	közepes	nem	nem
3	középkorú	magas	nem	igen
4	idős	közepes	igen	nem
5	fiatal	magas	nem	igen
6	középkorú	alacsony	nem	nem
7	idős	alacsony	nem	nem
8	fiatal	közepes	nem	igen
9	középkorú	magas	igen	igen
10	idős	közepes	igen	nem

Kezdeti GINI: $1 - ((5/10)^2 + (5/10)^2) = 0.50$

Döntési fa

Sorsz.	Életkor	Érdekl-e az akció
1	fiatal			igen
2	idős			nem
3	középkorú			igen
4	idős			nem
5	fiatal			igen
6	középkorú			nem
7	idős			nem
8	fiatal			igen
9	középkorú			igen
10	idős			nem

Kezdeti GINI: 0.48

Életkor szerinti vágás

GINI-jének számítása:

$$\text{GINI}(\text{fiatal}) = 1 - ((3/3)^2 + (0/3)^2) = 0$$

$$\text{GINI}(\text{köz.}) = 1 - ((2/3)^2 + (1/3)^2) = 4/9$$

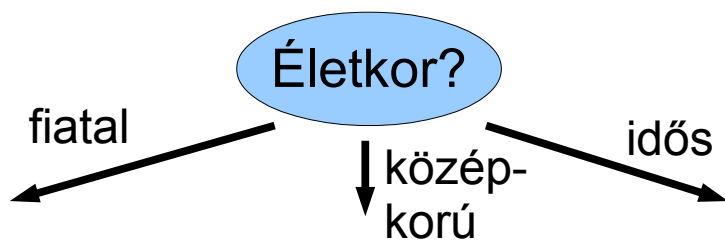
$$\text{GINI}(\text{idős}) = 1 - ((0/4)^2 + (4/4)^2) = 0$$

Életkor szerinti vágás GINI-je:

$$0.3 \cdot 0 + 0.3 \cdot (4/9) + 0.4 \cdot 0 = 0.13$$

GINI nyereség életkor szerinti

$$\text{vágás esetén: } 0.50 - 0.13 = 0.37$$



Döntési fa

Sorsz.	...	Testsúly	...	Érdekl-e az akció
1		alacsony		igen
2		közepes		nem
3		magas		igen
4		közepes		nem
5		magas		igen
6		alacsony		nem
7		alacsony		nem
8		közepes		igen
9		magas		igen
10		közepes		nem

Kezdeti GINI: 0.48

Testsúly szerinti vágás
GINI-jének számítása:

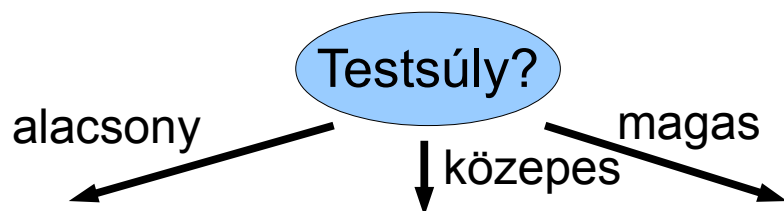
$$\begin{aligned} \text{GINI}(\text{alacsony}) &= \\ &= 1 - ((1/3)^2 + (2/3)^2) = 4/9 \end{aligned}$$

$$\begin{aligned} \text{GINI}(\text{közepes}) &= \\ &= 1 - ((1/4)^2 + (3/4)^2) = 6/16 \end{aligned}$$

$$\begin{aligned} \text{GINI}(\text{magas}) &= \\ &= 1 - ((3/3)^2 + (0/3)^2) = 0 \end{aligned}$$

Testsúly szerinti vágás GINI-je:
 $0.3 * (4/9) + 0.4 * (6/16) + 0.3 * 0 = 0.28$

GINI nyereség testsúly szerinti
vágás esetén: $0.50 - 0.28 = 0.22$



Döntési fa

Sorsz.	Sportol	Érdekl-e az akció
1			igen	igen
2			nem	nem
3			nem	igen
4			igen	nem
5			nem	igen
6			nem	nem
7			nem	nem
8			nem	igen
9			igen	igen
10			igen	nem

Kezdeti GINI: 0.48

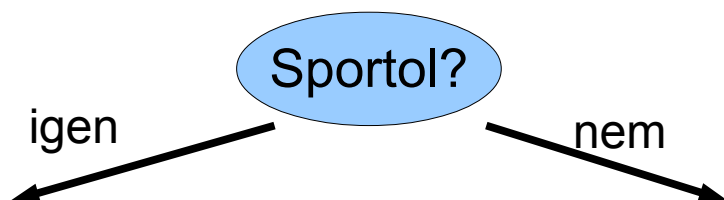
Sportol szerinti vágás
GINI-jének számítása:

$$\begin{aligned} \text{GINI}(\text{sportol}=\text{igen}) &= \\ &= 1 - ((2/4)^2 + (2/4)^2) = 0.5 \end{aligned}$$

$$\begin{aligned} \text{GINI}(\text{sportol}=\text{nem}) &= \\ &= 1 - ((3/6)^2 + (3/6)^2) = 0.5 \end{aligned}$$

Sportol szerinti vágás GINI-je:
 $0.4 * 0.5 + 0.6 * 0.5 = 0.5$

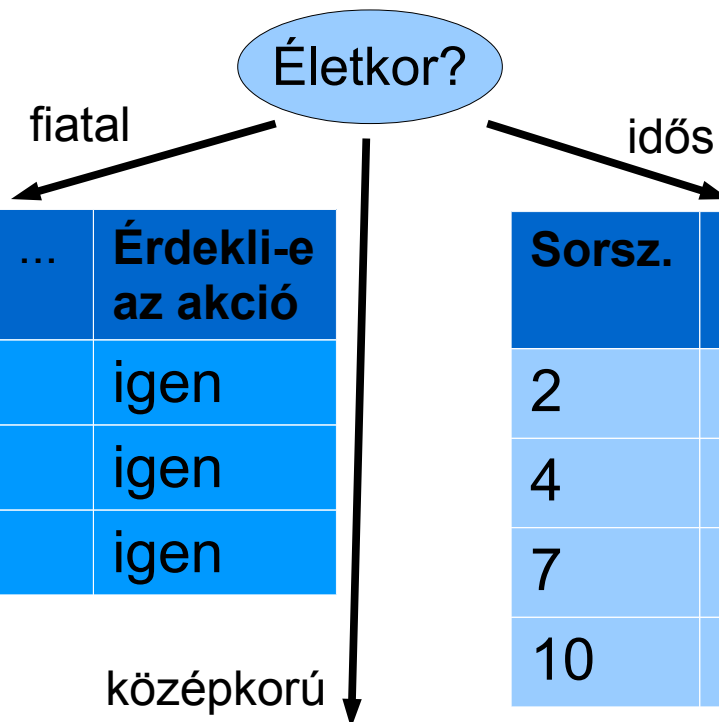
GINI nyereség „sportol” szerinti
vágás esetén: $0.50 - 0.50 = 0$



Döntési fa

- A nyereség a vizsgált vágások estén:
 - **Életkor: 0.37 (legjobb vágás)**
 - Testsúly: 0.22
 - Sportol: 0

Döntési fa

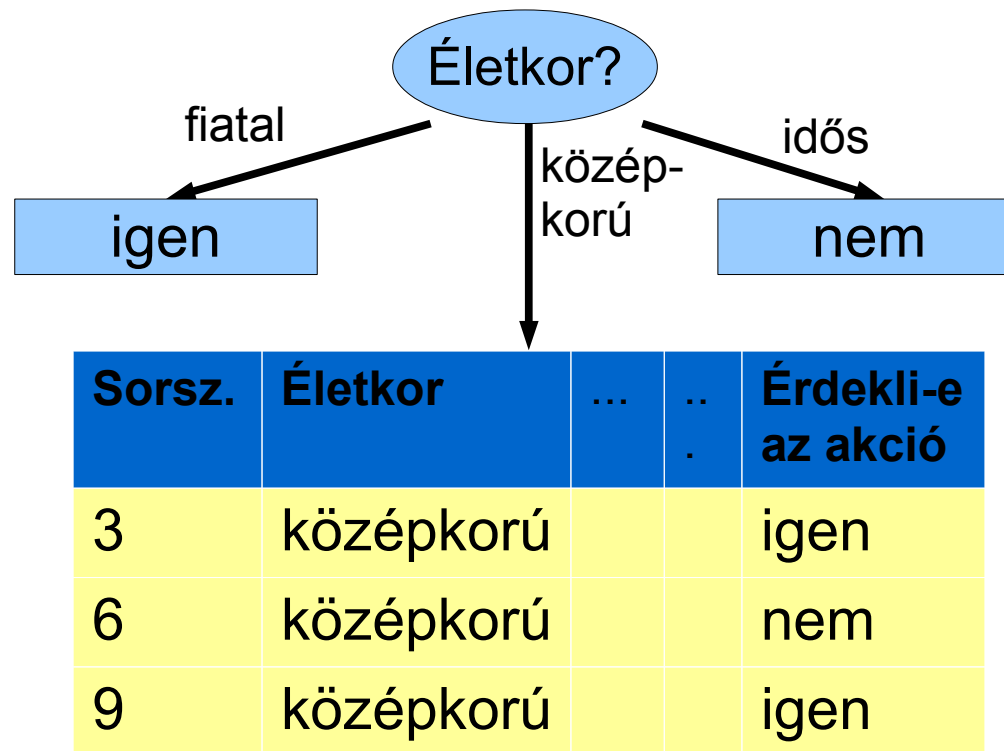


Sorsz.	Életkor	Érdekl-e az akció
1	fiatal			igen
5	fiatal			igen
8	fiatal			igen

Sorsz.	Életkor	Érdekl-e az akció
2	idős			nem
4	idős			nem
7	idős			nem
10	idős			nem

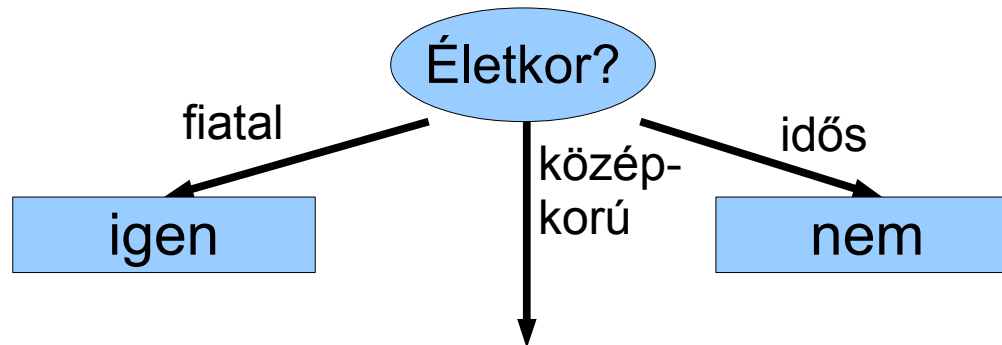
Sorsz.	Életkor	Érdekl-e az akció
3	középkorú			igen
6	középkorú			nem
9	középkorú			igen

Döntési fa



- A két szélső ágon minden példány ugyanazon osztályba tartozik → osztály nevével címkézett levél
- A középső ágon tovább vágunk (ha nincs más leállási feltétel, ami teljesül)

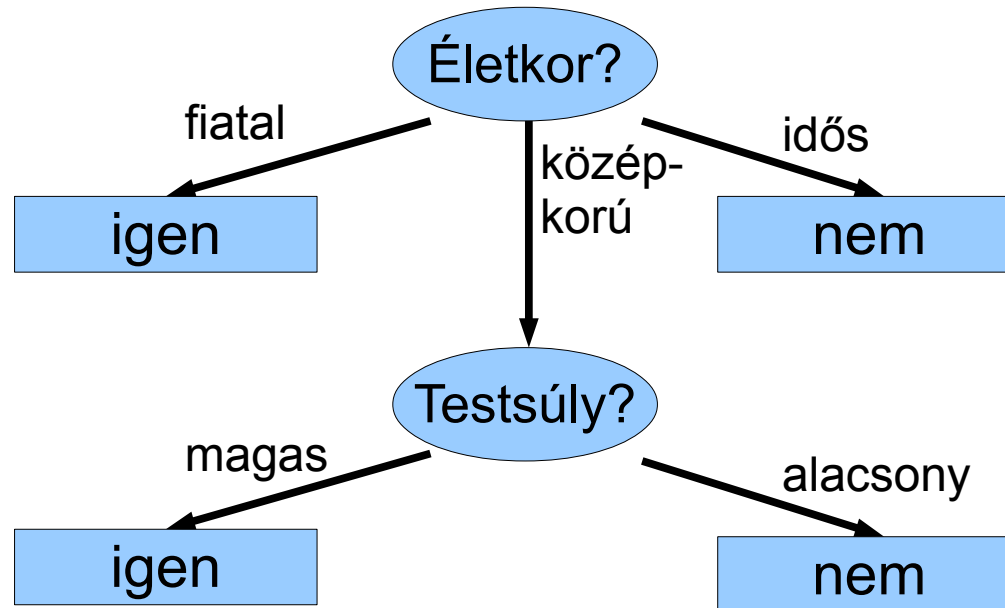
Döntési fa



Sorsz.	Életkor	Testsúly	Sportol	Érdekl-e az akció
3	középkorú	magas	nem	igen
6	középkorú	alacsony	nem	nem
9	középkorú	magas	igen	igen

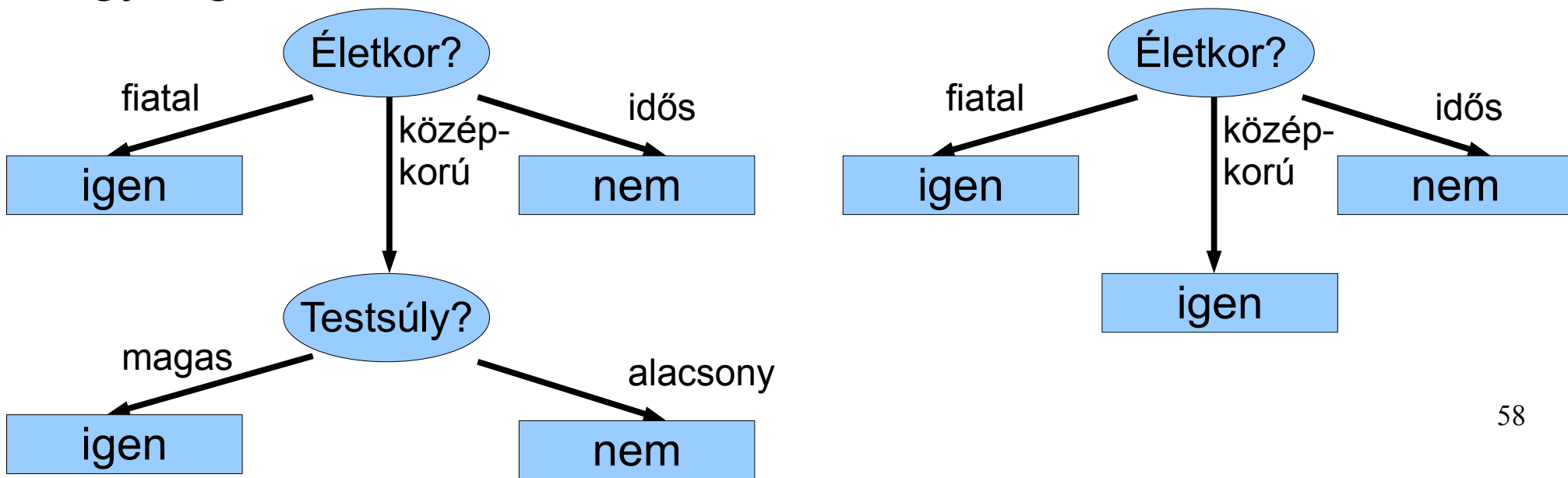
- Vágjunk a testsúly szerint! (Teljesen tiszta vágást eredményez → GINI-növekedést is maximalizálja)

Döntési fa



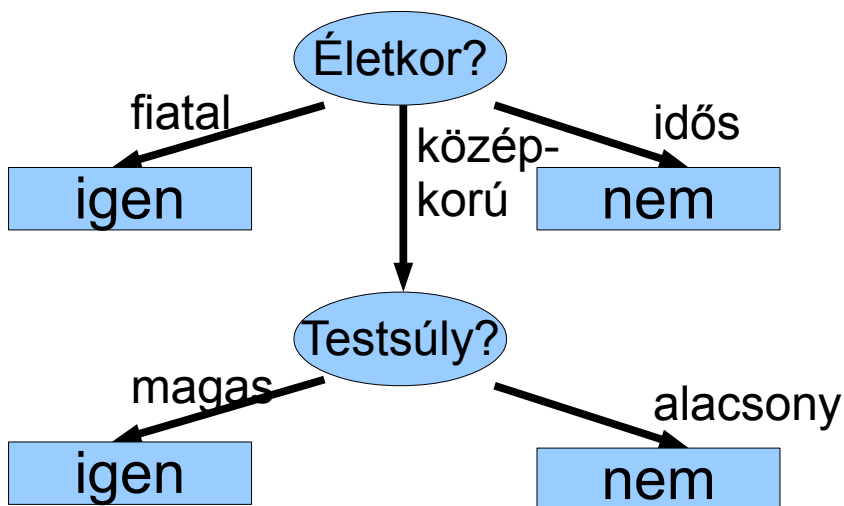
Döntési fa - MDL

- **15. Feladat:** Az előbbi adatbázishoz egy másik döntési fát is építettünk. A második döntési fa építésekor kikötöttük, hogy minden ágra legalább 3 példány jusson (lásd alább). Minimum Description Length (MDL) elv alapján a két fa közül várhatóan melyik fa fogja jobban osztályozni az ismeretlen osztályba tartozó példányokat? Tételezzük fel, hogy 1-1 osztályozási hiba description length-je 2 egység, a döntési fa minden nem-levél csomópontjának description length-je 3 egység, minden levél csomópont description length-je 1 egység.

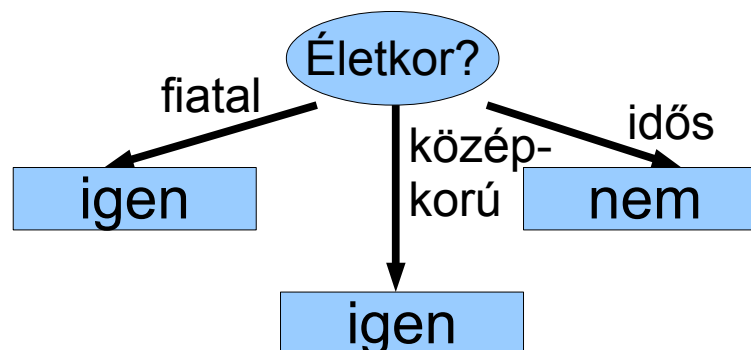


Döntési fa - MDL

- **Megoldás:** A description length-ek:



$$\begin{aligned} & 2 \cdot 3 \\ & + 4 \cdot 1 \\ & + 0 \cdot 2 \\ & = \mathbf{10} \end{aligned}$$



$$\begin{aligned} & 3 \quad (\text{nem-levél csomópontok}) \\ & + 3 \cdot 1 \quad (\text{levél csomópontok}) \\ & + 2 \cdot 1 \quad (\text{osztályozási hibák}) \\ & = \mathbf{8} \end{aligned}$$

- MDL szerint a második (jobb oldali) fát preferáljuk, mert annak description length-e kisebb.

Döntési fa

- **16. Feladat:**

- a) Milyen esetekben fontos a vágás entrópiájának (SplitINFO) figyelembe vétele?
- b) Mit jelent a ReducedErrorPruning?
Mikor használjuk?
- c) Miért szokás kikötni, hogy csak olyan vágást engedünk, amely során minden ágra legalább N darab példány kerül? (Az N egy felhasználó által állítható paraméter.)
- d) Használhatunk-e a GINI helyett más mértéket annak eldöntésére, hogy mi a legjobb vágás?

További feladatok

Waveletek

- **17. feladat:**
Számolja ki az alábbi jel Haar-Wavelet transzformáltját!

4	8	12	20	18	16	15	13	12	8	6	4	3	2	1	0
---	---	----	----	----	----	----	----	----	---	---	---	---	---	---	---

Waveletek

- **Megoldás:**

(Megjegyzés: az egyszerűség kedvéért a számítás során a skálázástól eltekintettünk.)

4	8	12	20	18	16	15	13	12	8	6	4	3	2	1	0
6	16	17	14	10	5	2.5	0.5	-4	-8	2	2	4	2	1	1
11	15.5	7.5	1.5	-10	3	5	2	-4	-8	2	2	4	2	1	1
13.25	4.5	-4.5	6	-10	3	5	2	-4	-8	2	2	4	2	1	1
8.875	8.75	-4.5	6	-10	3	5	2	-4	-8	2	2	4	2	1	1