

Adatbányászat

Bevezetés

Tikk Domonkos

Áttekintés

- Mi fán terem az adatbányászat
- Módszertan
- Tipikus feladatok
- Eszközök
- Esettanulmány



Forrás: <http://www.adatmentes-adatvissza.hu/hu/adatmentes-adatbanyaszat-data-mining.html>

Mottó

Megfulladunk az adatban és tudásra éhezünk
(We are drowning in information but starved for
knowledge)



Adatbányászat – motiváció

- Gyorsuló ütemben növő adatmennyiség
- Üzleti igény az adatokban rejlő információk kinyerésére
- Definíció: döntéstámogatási **follyamat**, amely **érvényes, hasznos, rejtett** (korábban nem ismert) információt állít elő nagy mennyiségű – jellemzően adatbázisokban tárolt – adatból
(forrás: Abonyi (szerk): Adatbányászat a hatékonyság eszköze)
- Automatizálható folyamat
 - emberi erőforrás igénye alacsony
 - gyorsan generálhatóak az információk

A definíció elemei

- folyamat
 - nem dobozos termék, hanem átfogó tudást igényel az alkalmazása is
- érvényes
 - pontosság, statisztikai szignifikancia, teljesség
- hasznos
 - adjon új, értékes ismereteket
 - gyakran nehéz üzleti értéket meghatározni
- rejtett (előzőleg nem ismert)
 - hipotézis megerősítése vs. új felfedezése
 - prediktív vs. leíró adatbányászat

Trend

The production of data is expanding at an astonishing pace. Experts now point to a 4300% increase in annual data generation by 2020. Drivers include the switch from analog to digital technologies and the rapid increase in data generation by individuals and corporations alike.

■ Size of Total Data ■ Enterprise Created Data
■ Enterprise Managed Data

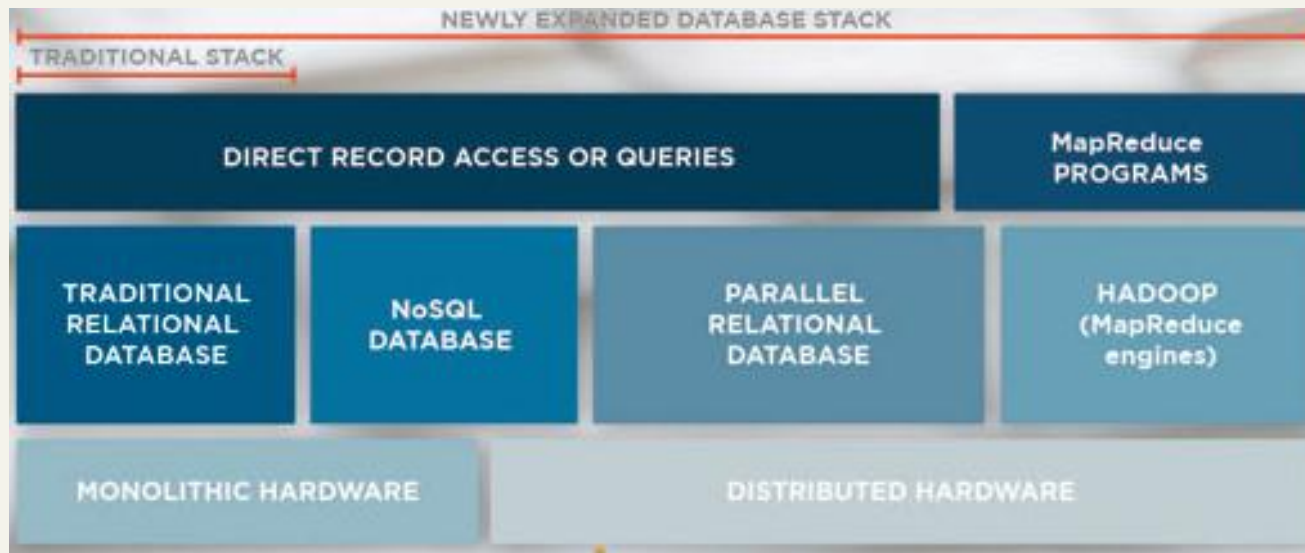
2020: MORE THAN 1/3 OF THE DATA PRODUCED WILL LIVE IN OR PASS THROUGH THE CLOUD.

2012: CUSTOMERS WILL START STORING 1 EB OF INFORMATION.



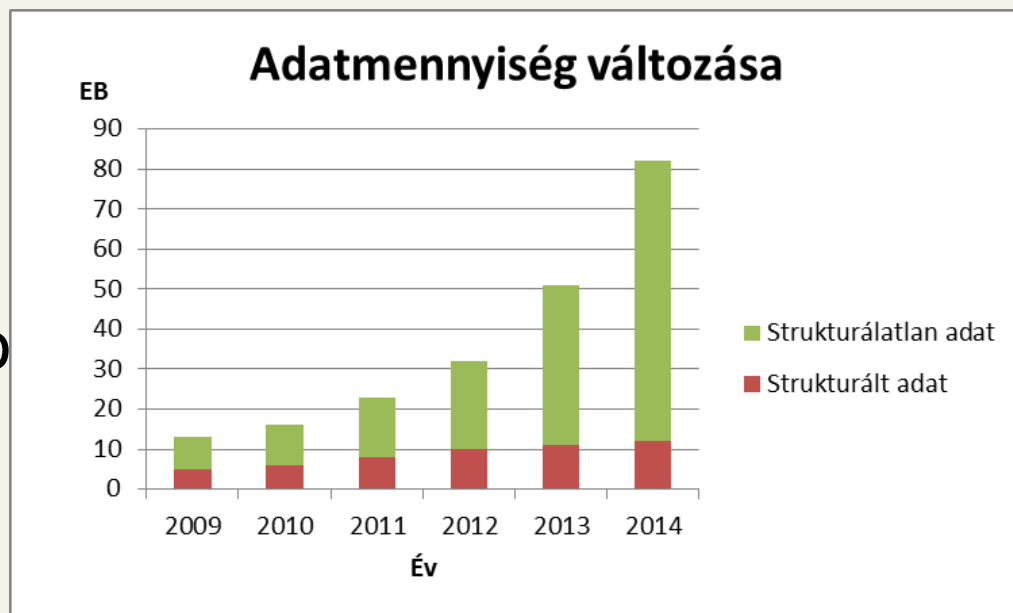
Adatbázis dimenzió

- Tradicionális relációs adatbázisban már csak az adatok ~20%-át tárolják
- Új technológiák: No-SQL DB, párhuzamos RelDB, elosztott rendszerek (Hadoop)
- Hardver szerint: monolitikus vs. elosztott

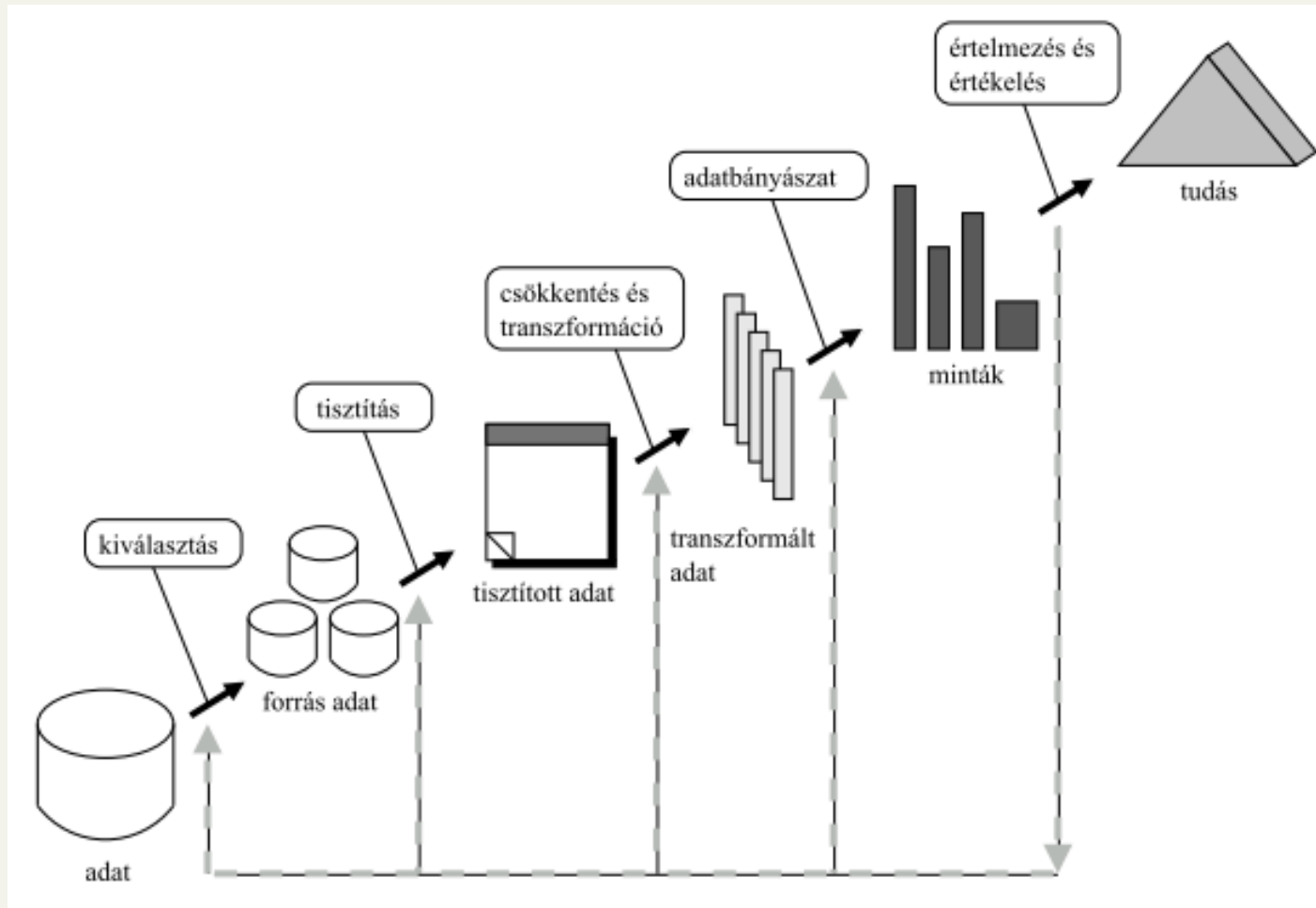


Méret és tartalom

- Az emberi tudás összmenyiségét 300 EB-ra becsülték (2007)
- Avatar számítógépes grafikával készített effektjeinek adatmennyisége: 1 PB
- Strukturált adat:
relációs DB
- Strukturálatlan adat:
szöveg, audio, video



A tudásfeltárás folyamata

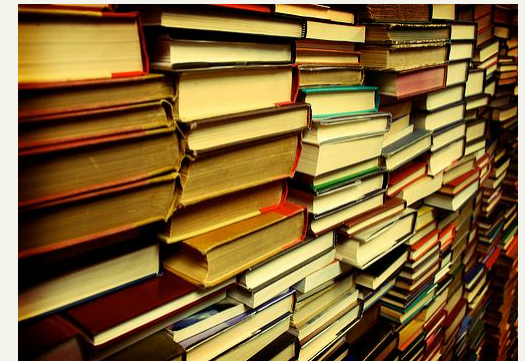


forrás: Bodon F.: Adatbányászat (elektronikus jegyzet)

Adattípusok

- Adatbázis-tartalom
- tranzakciós adatbázisok
- adattárházak
- térinformatikai adatok
- idősor és temporális adatok
- szöveges és multimédia adatok
- WWW
- heterogén adatbázisok

Company Name	CustomerTy	Order Date	Ship Date	PO Number
A. Datum Corporation	Commercial	3/23/2007	3/29/2007	10
A. Datum Corporation	Commercial	3/23/2007	3/29/2007	11
Contoso, Ltd	Government	10/16/2007	10/19/2007	26
Contoso, Ltd	Government	12/10/2007	12/11/2007	12A948
Trey Research	Government	6/21/2007	6/28/2007	30
Trey Research	Government	11/15/2007	11/18/2007	PO359194
Litware, Inc	Non-Profit	11/12/2007	11/14/2007	32
Litware, Inc	Non-Profit	1/12/2008	1/20/2008	PO1294204
New Customer	Non-Profit	12/20/2007	12/21/2007	PO128A438



Alkalmazási területek



pénzügyi szektor



tudomány



gyártástechnológia



közlekedés



jog



telekommunikáció



energiaipar

Pénzügyi szektor

- Bankkártya bűncselekmények
- Hitelképesség-elemzés
- Ügyfélszegmentáció
- Ügyfélérték számítás
- Lojalitás vizsgálat
- Keresztértékesítés
- Kampánymenedzsment
- Vásárlói kosár elemzés

Tudomány/egészségügy

- Kutatási eredmények kiértékelése
- Képek osztályozása
- Új kapcsolatok keresése tényadatokból
- Korreláció elemzés (hipotézis és tényleges mérések között)
- Gyógyszerforgalmi előrejelzések
- Betegségek és fizikai megfigyelések korrelációvizsgálata
- Kórházi monitorozó rendszerek

Telekommunikáció/energiaszektor

- Lemorzsolódás-előrejelzés
- Ügyfél-szegmentáció és termék targetálás
- Véleményalkotók azonosítása – hívási gráf elemzések
- Hálózati hiba előrejelzése
- Túl- és alulszámlázások azonosítása
- Csalás-felderítés



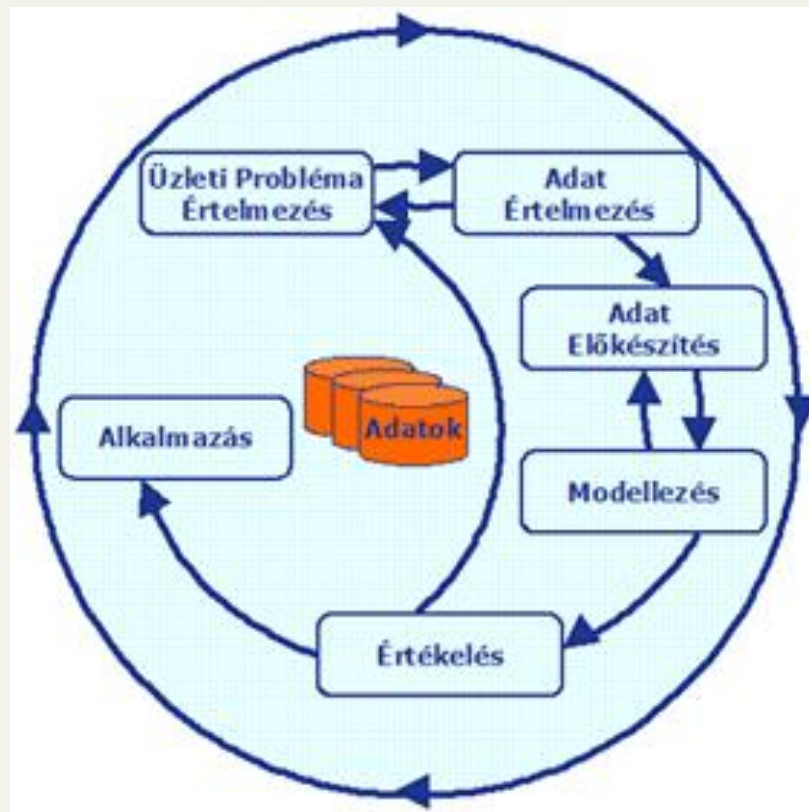
AZ ADATBÁNYÁSZAT MÓDSZERTANA

Az adatbányászat folyamata



CRISP-DM

- Cross-Industry Standard Process for Data Mining



Üzleti probléma értelmezése

- Üzleti célok megfogalmazása
 - üzleti háttér, üzleti cél és siker tényezők
- Helyzetfelmérés
 - erőforrások, követelmények, források, feltételezések
 - kockázatfelmérés, haszon és költségek
 - terminológia
- Adatbányászati célok definiálása
 - adatbányászati célok és sikerkritériumok
- Projektterv elkészítése
 - Eszközök és technikák értékelése

Adatértelmezés

- Kiindulási adatok gyűjtése
 - hozzáférés biztosítása, adatintegráció
- Ráttekintés az adatokra
 - főbb jellemzők (típusok, értéktartományok)
- Alap statisztikai jellemzők feltárása
 - lekérdezés, vizualizálás, értelmezés
 - célparaméter eloszlása, főbb dimenziók mentén való szegmentálás
- Adatminőségi vizsgálat
 - feltöltöttség, lefedettség, adathelyesség
- Minderről beszámoló készül

Adatok előkészítése

- Adatkiválasztás
 - a célok eléréséhez mely adatok hasznosak
- Adattisztítás
 - adatkitöltés, inkonzisztencia megszüntetése
- Új paraméterek bevezetése
 - Származtatott adatok, generált rekordok
- Adatintegráció
 - több forrás esetén
- Adatformátum módosítása
 - adatbányászati modellhez igazítás

Modellezés

- Modellező technika kiválasztása
 - eszközt találni a célhoz, adatfeltáró elemzés
- Modell tesztelésének meghatározása
 - Kiértékelési módszer, vizualizálás
- Modellalkotás
 - Paraméterbeállítás, modellek, dokumentálás
- Modell kiértékelése és megjelenítése
 - Fontos a jól vizualizálható eredmény
 - beállítások felülvizsgálata

Üzleti értékelés

- A modell üzleti célú értékelése
 - üzleti elvárásoknak megfelel?
 - éles környezetben tesztelhető
- A teljes elemzési folyamat felülvizsgálata
 - pl. adatok hosszú távú elérhetősége
- Következő lépések
 - döntés a felhasználhatóságról, üzleti bevezetésről

Üzleti alkalmazás

- Alkalmazás megtervezése
 - beépítés az üzleti folyamatokba
- Alkalmazás fenntartás és monitoring
 - tesztesetek, ellenőrzések beállítása
- Projekt tanulmány elkészítése
 - Beszámoló, prezentáció
- A projekt felülvizsgálata
 - éles eredmények kiértékelése (ROI!)
 - pozitívumok vs. negatívumok
 - elvárttól való eltérések elemzése

Mai óra

- Mi fán terem az adatbányászat
- Módszertan
- Tipikus feladatok
- Eszközök
- Esettanulmány



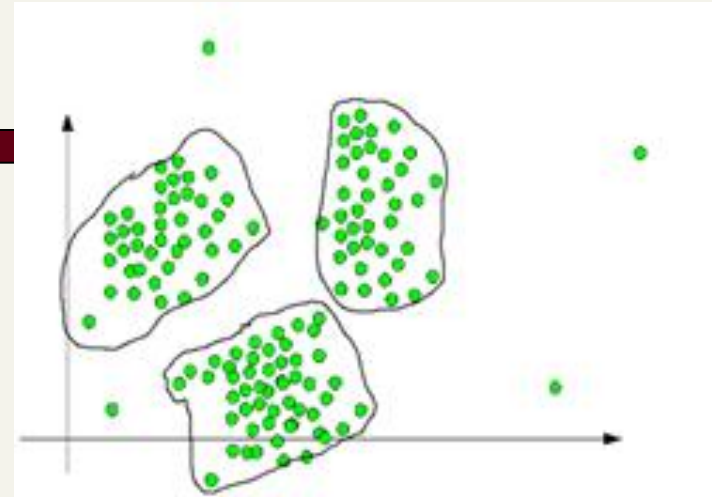
Forrás: <http://www.adatmentes-adatvissza.hu/hu/adatmentes-adatbanyaszat-data-mining.html>



FELADATTÍPUSOK

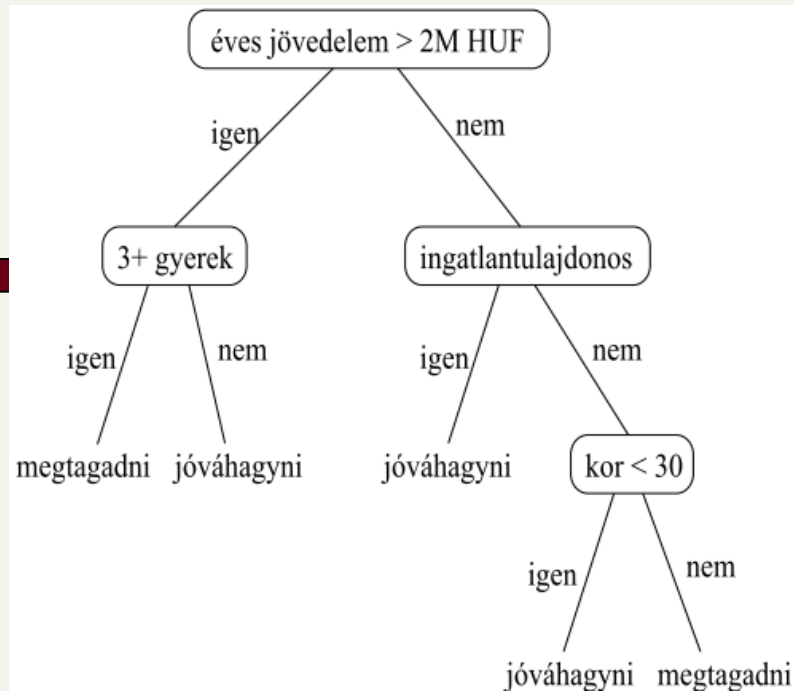
Csoportosítás

- Szegmentálás, klaszterezés
- *Felügyelet nélküli tanulás*
- Adatok felosztása csoportokra
- Csoporton belül hasonlóak
- Csoportok viszont különbözőek
- Hány csoport legyen?
- Milyen legyen a felosztás struktúrája (egyszintű, hierarchikus)
- Példa: piacszegmentálás



Osztályozás

- Klasszifikáció, kategorizálás
- *Felügyelt tanulás*
- Tanuló és tesztadatok
- Modellépítés (-generálás)
- Alkalmazás (előrejelzés)
- Függvény: bemenetekből kimenetet állít elő (osztálycímke)
- Példa:
 - hitelelbírálás
 - égitestek besorolása (galaxis, közeli csillag, egyéb)



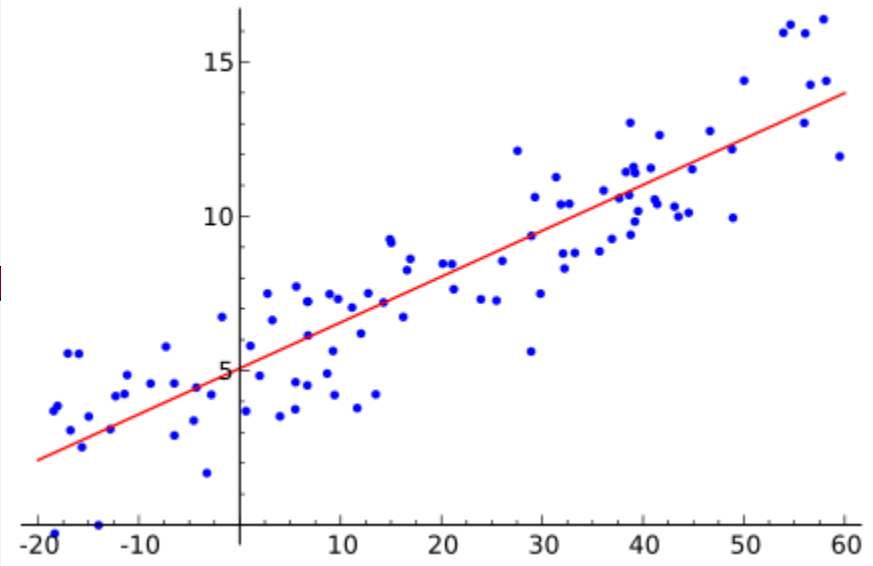
Asszociáció

- Gyakori elemhalmazok
- Objektumok közti összefüggés
- Kapcsolat erőssége
- Példa: vásárlóikosár-elemzés
- Ha valaki vesz A és B terméket akkor C-t is vesz
- Konfidencia, támogatottság
- Más feladat: gyakori sorozatok (adatszekvenciák), gyakori epizódok (részben rendezett)

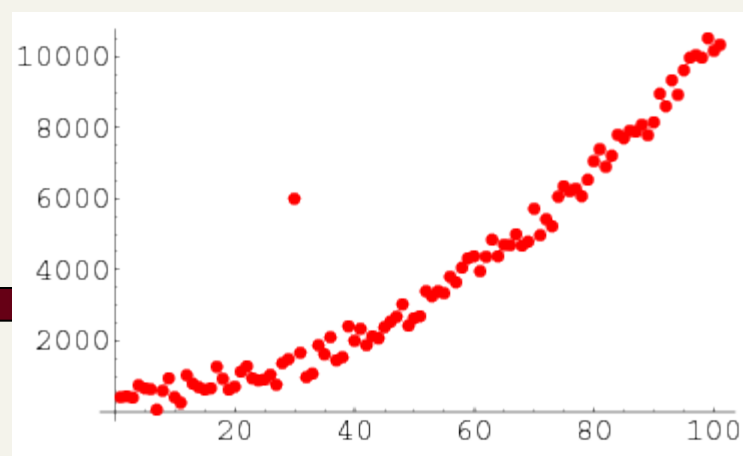


Regresszió

- Függvény illesztés
- osztályozáshoz hasonló
- az adatban rejlő sajátosságok modellezése,
- kimenet numerikus értéke, nem kategorikus adat (osztályozás)
- Adatbányász: modellkiválasztás (lineáris, polinom, logaritmikus, hiperfelület)
- Példa:
 - időbeni előrejelzés: BUX index alakulása
 - statikus: betegség valószínűségi orvosi adatok alapján



Eltéréselemzés



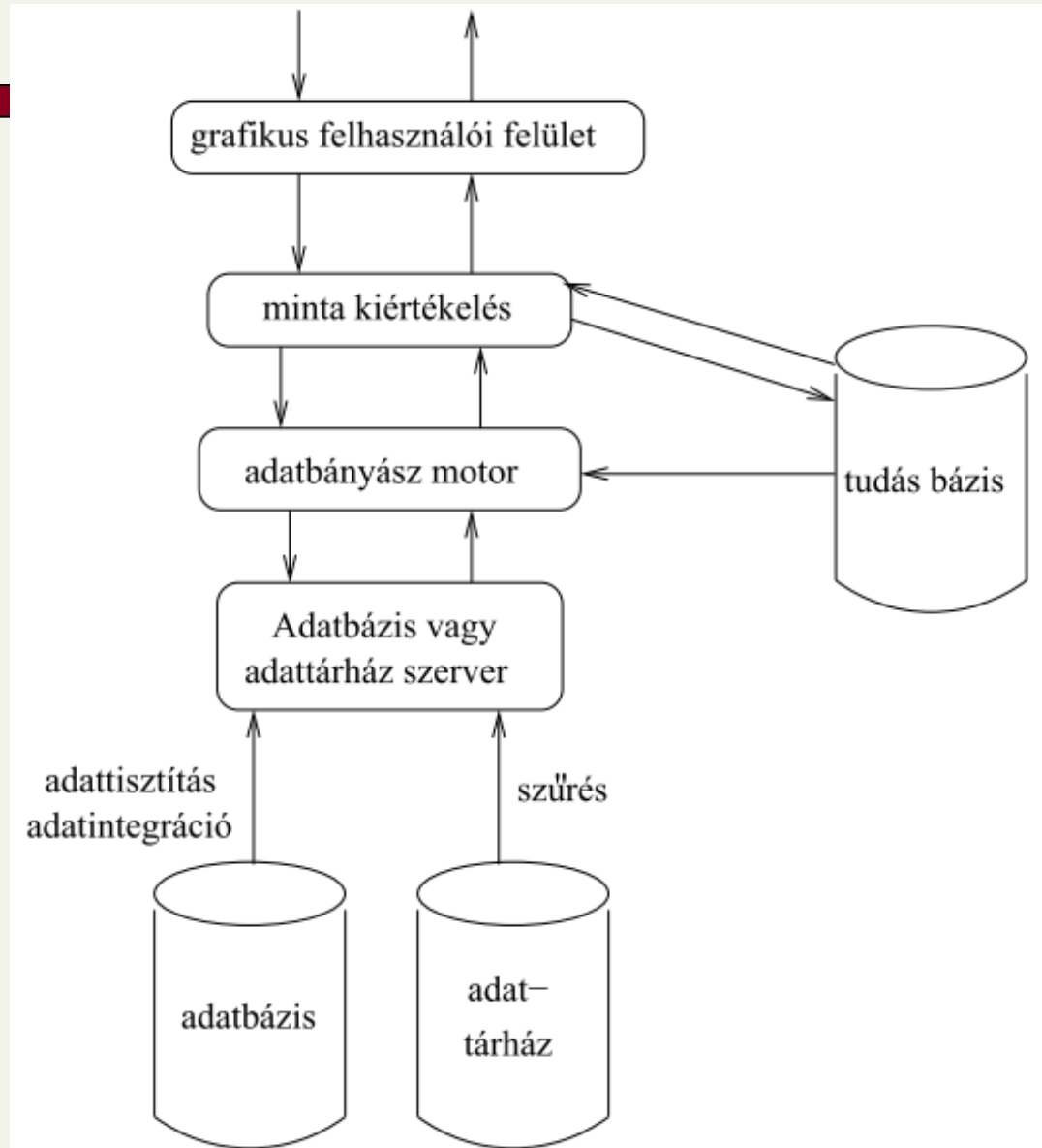
- outlier analízis
- olyan adatok azonosítása, amelyek eltérnek az elvárttól
- Lehet zaj, mérési hiba, kivétel (ekkor szűrni lehet)
- Alkalmas csalások kiderítésére
- Példa:
 - hitelkártya-visszaélések
 - áramlopás
 - biztonsági elemzés





ESZKÖZTÁRAK

Alapelemek



Eszközök

- általános adatkezelők
 - Excel
- programozási keretrendszer:
 - Matlab
- piaci szoftverek
 - IBM/SPSS Clementine, Statsoft Statistica, SAS Data Miner
- adatbázis-kezelők adatbányász kiegészítései
 - Oracle, MySQL, IBM
- Ingyenes rendszerek
 - WEKA, Rapidminer

IBM/SPSS Clementine

The screenshot displays the IBM/SPSS Clementine software interface. The main workspace shows a workflow diagram with the following components and connections:

- Input:** A document icon labeled "snapshotrainN.db" is connected to a "table" icon.
- Processing:** The "table" icon connects to a "type" icon (a hexagon with four colored dots).
- Modeling:** The "type" icon connects to a "pep" icon (a diamond with a person and a gear).
- Outputs:** The "type" icon also connects to three triangular icons labeled "children", "children x mortgage ..", and "children v. income".
- Deployment:** The "pep" icon connects to a "Publisher" icon (a square with a diamond).

The right-hand pane is divided into two sections:

- Streams:** Contains a tree view with "Stream1" and "mailshot2".
- CRISP-DM Classes:** Shows a hierarchical structure for an "(unsaved project)":
 - Business Understanding
 - Data Understanding
 - Data Preparation
 - Modeling
 - Evaluation
 - Deployment

The bottom of the interface features a toolbar with various icons for different operations, including "Table", "Matrix", "Analysis", "Data Audit", "Statistics", "Quality", "Report", "Set Globals", "Publisher", "Database", "Flat File", "SPSS Export", "SAS Export", "Excel", and "SPSS Procedure". The status bar at the very bottom indicates "Server: Local Server" and "25Mb / 43Mb".

WEKA

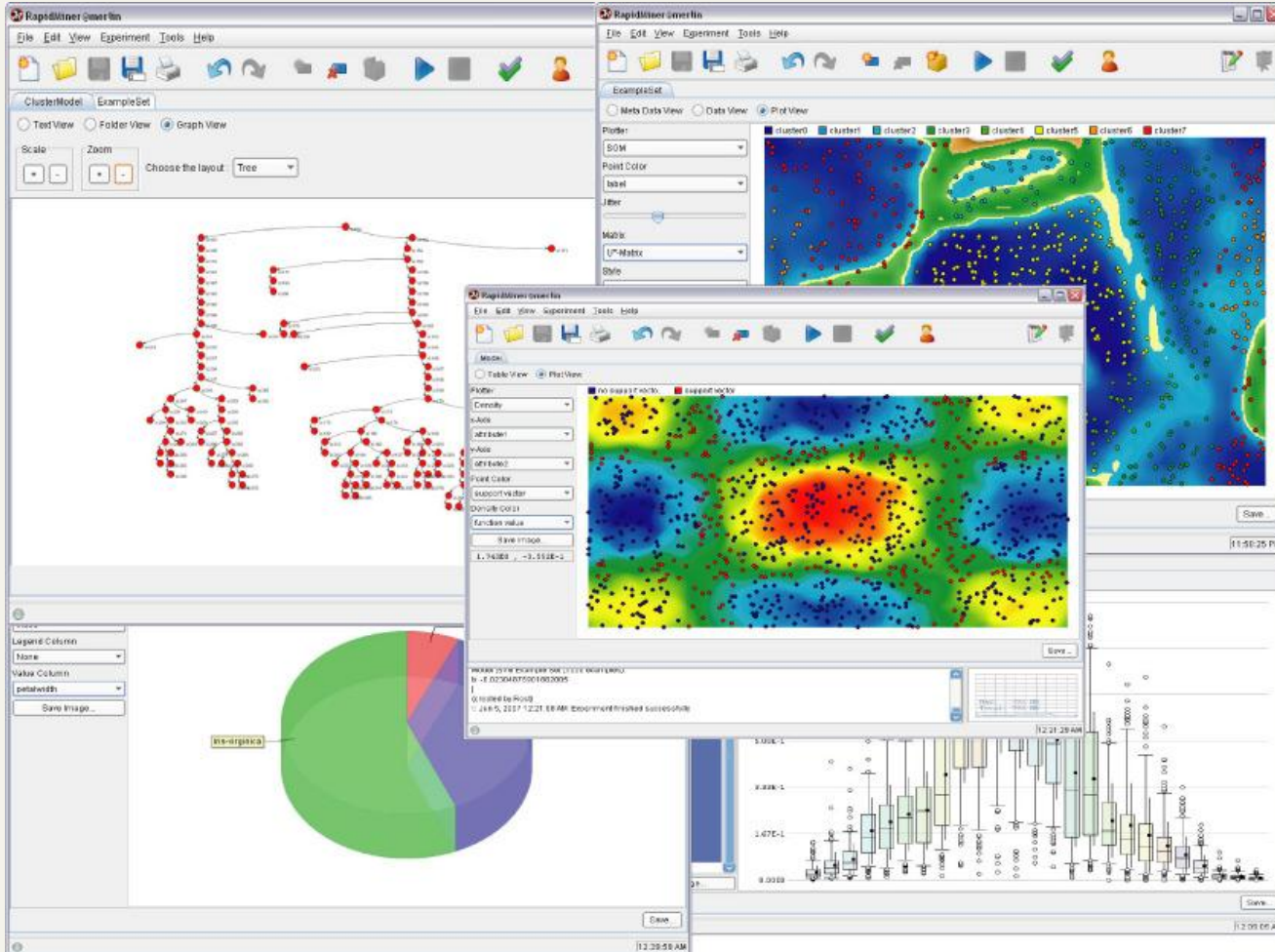
The screenshot displays the Weka Explorer application window, which is divided into several functional areas:

- Classifier:** Shows the selected classifier as "J48 -C 0.25 -M 2".
- Test options:** Includes radio buttons for "Use training set" (selected), "Supplied test set", "Cross-validation" (with 10 folds), and "Percentage split" (with 66%).
- Classifier output:** Displays the results of stratified cross-validation:

```
=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      144           96  %
Incorrectly Classified Instances     6             4  %
Kappa statistic                     0.94
Mean absolute error                  0.035
Root mean square error
```
- Tree View:** A decision tree visualization is shown, detailing splits based on petalwidth, petallength, and petalwidth attributes. The tree structure is as follows:
 - Root node: petalwidth (split at 0.6)
 - Left branch (<= 0.6): Iris-setosa (50.0)
 - Right branch (> 0.6): petalwidth (split at 1.7)
 - Left branch (<= 1.7): petallength (split at 4.9)
 - Left branch (<= 4.9): Iris-versicolor (48.0/1.0)
 - Right branch (> 4.9): petalwidth (split at 1.5)
 - Left branch (<= 1.5): Iris-virginica (3.0)
 - Right branch (> 1.5): Iris-versicolor (3.0/1.0)
 - Right branch (> 1.7): Iris-virginica (46.0/1.0)

- Visualize:** A scatter plot of Iris data is shown with axes for "Y: petalwidth (Num)" and "X: petalwidth". The plot includes a legend for Iris-versicolor and Iris-virginica, and a "Jitter" slider.

Rapidminer





ESETTANULMÁNY

Tüdőembólia detektálása

- Adatok:
 - numerikus gyanús régiókról (3D pixel – 116 jellemző); első 3: (x,y,z) az adatok szemantikája nem ismert; [-1; 1] be normált
 - beteg azonosító (egy beteghez több mérés)
 - pozitív és negatív minták címkével (beteg/nem)
- Feladatok:
 - új minták osztályozása
 - beteg emberek azonosítása
 - egészséges emberek azonosítása 100%-kkal

Feladat felépítése

- tanítóadatok:
 - 3303 adatsor; 46 beteg és 20 egészséges eset
- tesztadatok:
 - 1391 adatsor, 33 eset
- Felügyelt tanulás (osztályozás)
 - olyan modell építése, amely egy adott mintáról el tudja dönteni, hogy beteg-e vagy sem

Kényszerfeltételek

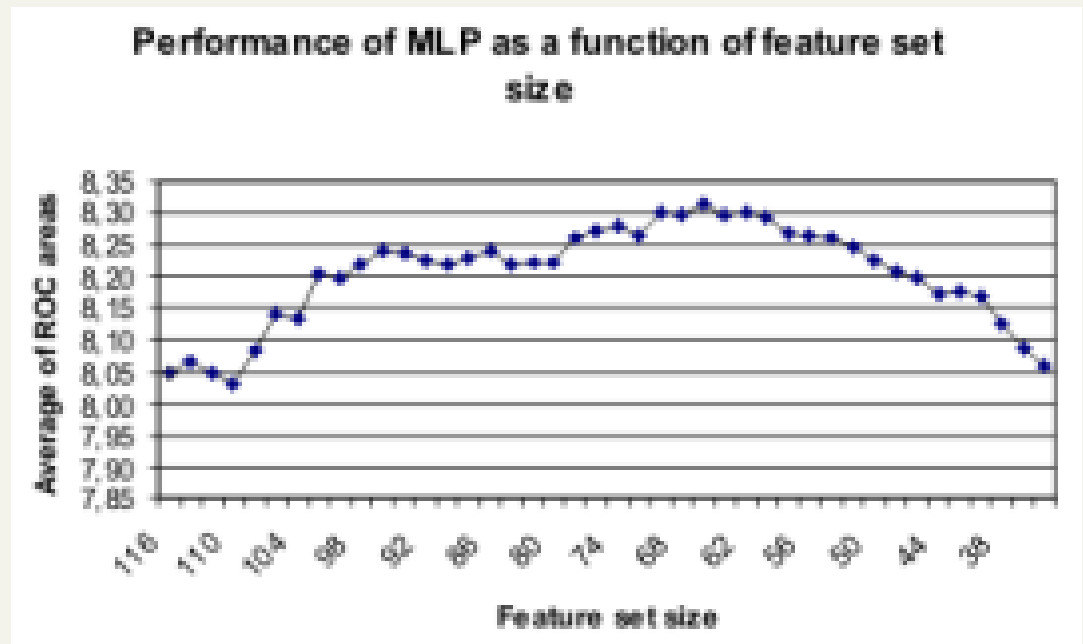
- hibás pozitív (FP) esetek minimalizálása
 - „farkast kiáltani” faktor csökkentése
- 3 küszöbérték adott, hogy páciensenként mennyi lehet a hibás esetek aránya (FP rates: 2; 4; 10)
- mérések:
 - helyes pozitív (TP) adatok azonosítása adott FP kényszerfeltételek mellett (#TP-PE)
 - hány beteg páciens ismer fel adott FP kényszerfeltételek mellett a rendszer (#TP-P)
 - helyesen azonosított egészséges páciensek száma

Nehézségek

- Zajos és kevés adat
 - hogyan generálták az adatokat
 - különböző gépek, szakértői címkézés
- adatsorok szemantikája ismeretlen
- térbeli összefüggések az adatok között nem azonosíthatóak
- atipikus a célfüggvény és a kényszerfeltétel

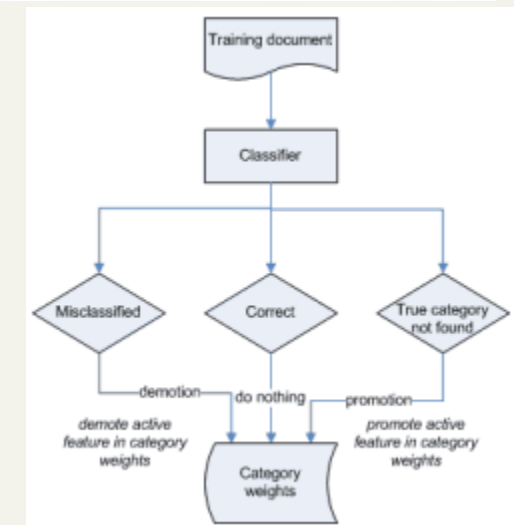
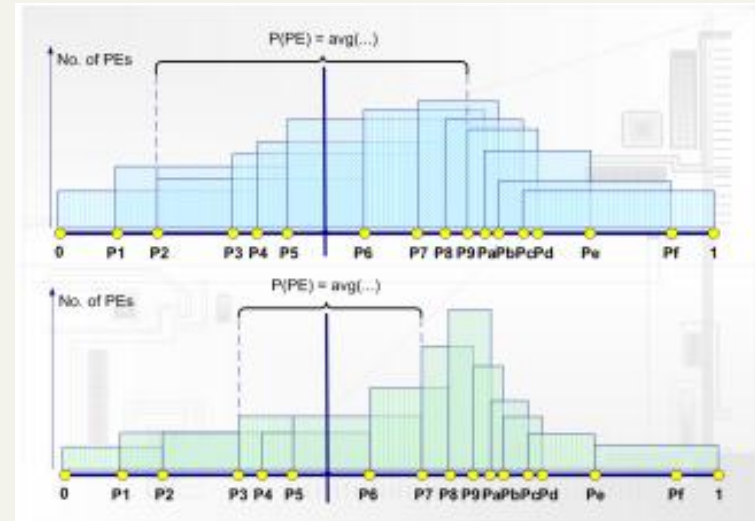
Adattisztítás

- A 116 adatelemből melyik használható
 - pontosít
 - gyorsít



Adatbányászati módszerek

- Osztályozók
 - statisztikai alapú
 - neurális háló A
 - neurális háló B
- Kombináció
 - osztályozó bizottság
 - feladat specifikusan
 - konfidencia értékek
 - osztályozó szinten
 - predikció szinten



Kombináció

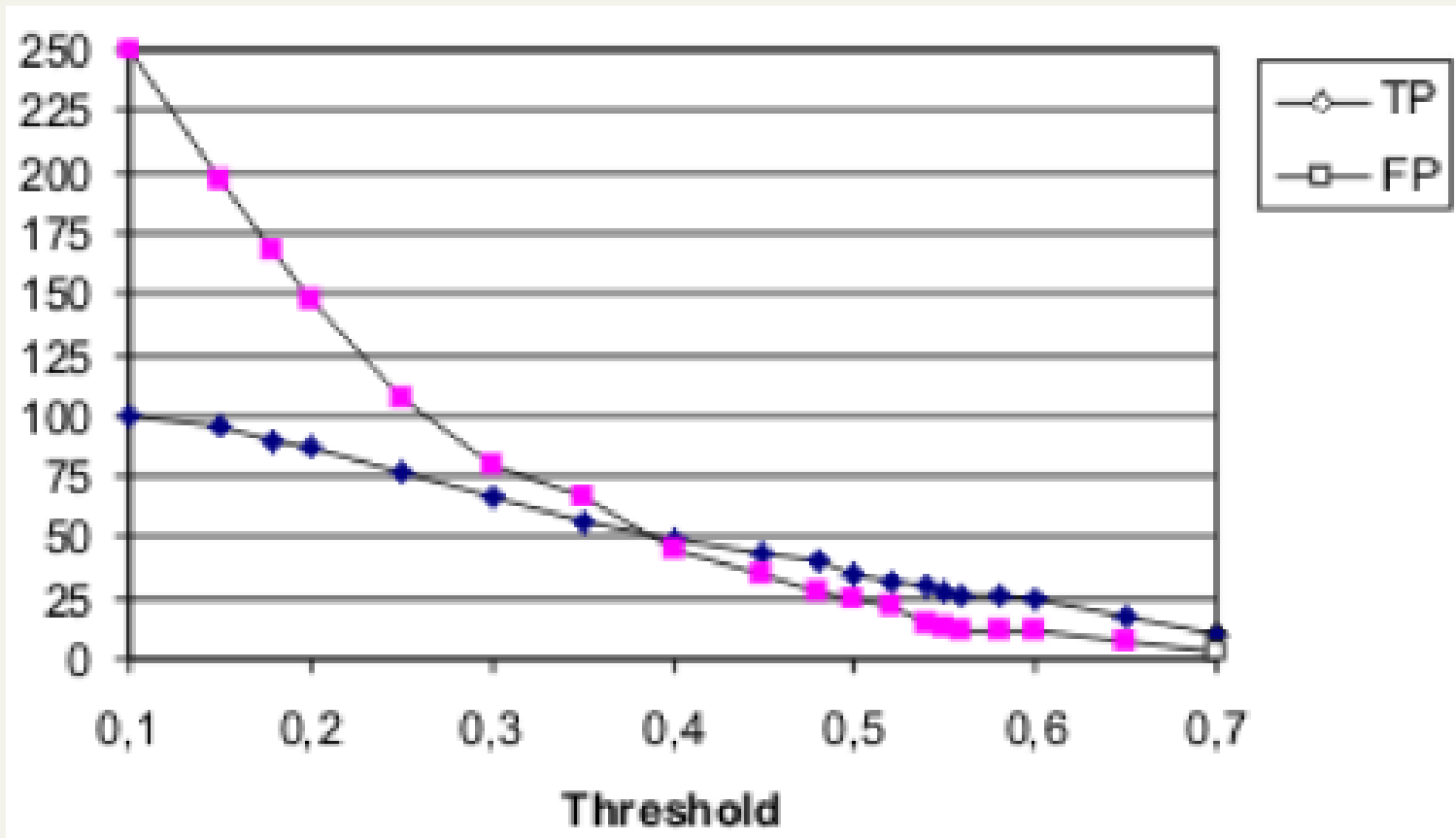
- Parametrizált vétó stratégia
- Biztossgái tagok:
 - 3 módszer, (1-1-2) beállítással: 4 adatsor
- Osztályozók kimenetének súlyozási szabályai
 - Egyöntetű pozitív döntés esetén: pozitív
 - 2-3 pozitív és nincs vétó: pozitív
 - 1 pozitív és nincs gyenge vétó: pozitív
 - különben negatív

Mitől függ a vétó értéke

- Megengedett hibás minták aránya (FP rates: 2; 4; 10)
- Osztályozó pontossága, amit keresztvalidációval mértünk
- esetenként az osztályozó által adott konfidenciaérték (nem mindenütt adott)
- osztályozók belső küszöbértéke

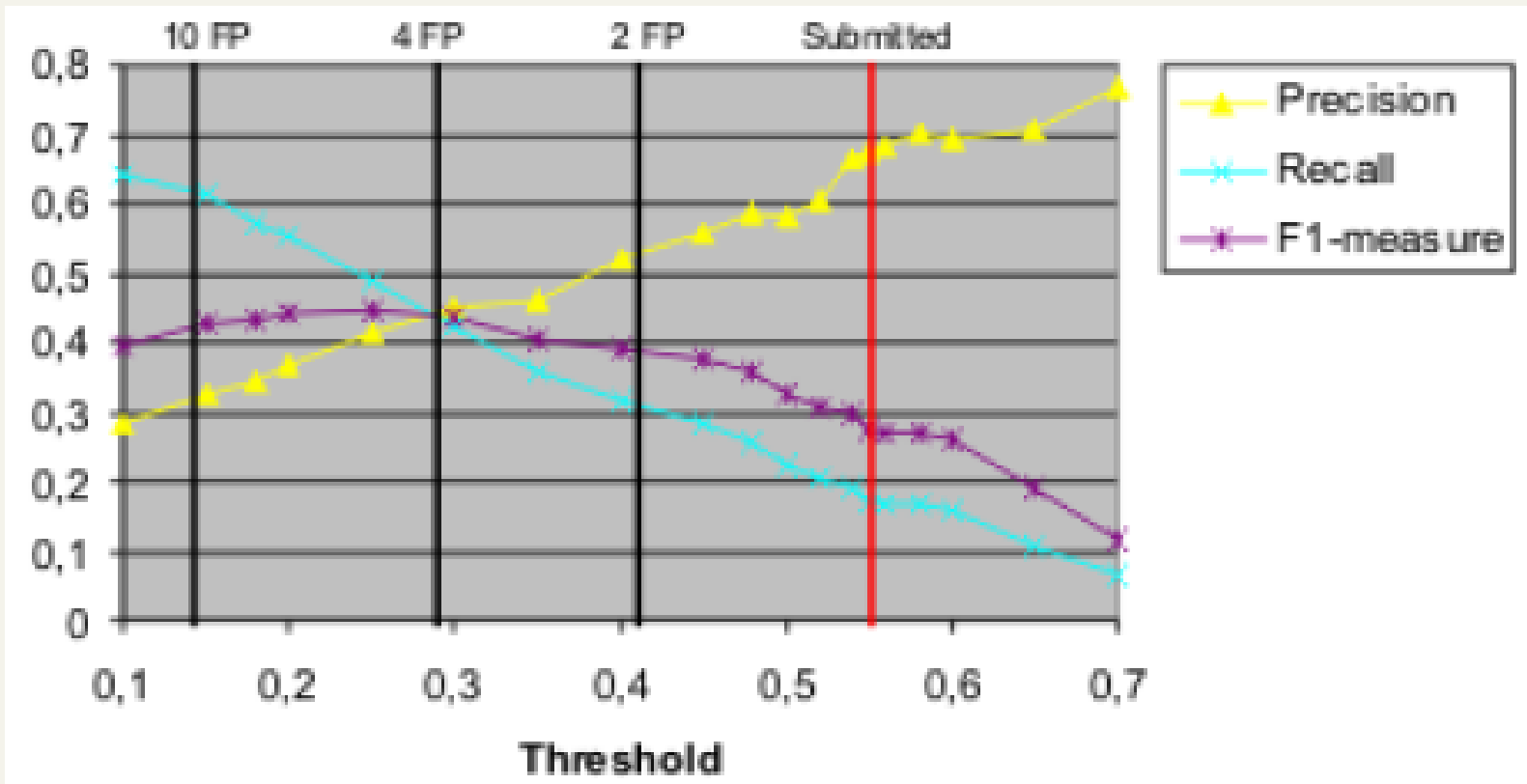
konzervatív vs. bátor

- a helyes és hibás találatok függvényében



Hagyományos IR mértékekkel

- pontosság (precision): a találatok közül mennyi helyes (PE)
- fedés (recall): hányat talál meg a tényleges PE-k közül



Eredmények

- Egyes osztályozókra

	Approaches			
	Statistical	ANN	Hitec-62	Hitec-116
Number of true PEs (23)	43	42	32	77
Number of false PEs (23)	62	12	18	91
Uniquely found PEs (23)	6	7	3	27

Kombináció

Combinations					
Stat & ANN	Stat & Hitec-62	Stat & Hitec-116	ANN & Hitec-62	ANN & Hitec-116	Hitec-62 & Hitec-116
22	22	36	24	33	27
8	5	19	3	5	8

Task1		
Voting for 2FP	Voting for 4FP	Voting for 10FP
51	80	90
21	57	109

Task2		
Voting for 2FP	Voting for 4FP	Voting for 10FP
49	64	98
18	58	144

Végeredmény

Method/submission	Precision (%)	Recall (%)	FP rate
Statistical	40.95	27.56	2.70
ANN	77.77	26.92	0.52
HITEC-116	45.83	49.35	3.96
HITEC-62	64.00	20.51	0.78
Task 1a	70.83	32.69	0.91
Task 1b	58.39	51.28	2.48
Task 1c	43.06	57.69	4.74
Task 2a	73.13	31.41	0.78
Task 2b	52.46	41.02	2.52
Task 2c	40.50	62.82	6.26

- pontosság: a találatok közül mennyi helyes (PE)
- fedés: hányat talál meg a tényleges PE-k közül