

BME-VIK Méréstechnikai és Információs Rendszerek Tanszék
Intelligens rendszerek szakirány

EGÉSZSÉGÜGYI KÓDOLÁSTÁMOGATÓ RENDSZER FEJLESZTÉSE

BSc szakdolgozat

Készítette: Sáránci István
Külső konzulens: Héja Gergely (GYEMSZI)
Tanszéki konzulens: Strausz György

2012. január 5.

Tartalom

- Bevezető
- Osztályozó módszerek
- Implementáció
- Eredmények, további lehetőségek

BEVEZETŐ

A problémáról

- Egészségügyi kódolás
 - Betegségek formalizált ábrázolására
 - PI. BNO-rendszer
 - Diagnózisok, zárójelentések kódolása
 - Statisztika, finanszírozás
- Automatizálás igénye
 - Manuálisan lassú, drága
 - Legalább részben automatizálni

Mi automatizálható?

- Teljes automatizálás nehéz
 - Szakértők is tévednek
- Kevesebb kódnak kelljen manuálisan utánanézni
 1. Felhasználó begépezi a diagnózist
 2. A rendszer visszaad egy tipplistát
 3. A felhasználó megnézi a talált kódok leírását és dönt
 - **Ehhez szakértelem szükséges!**
- Felhasználó–szolgáltatás kommunikáció
 - Webes felület

Példa



- Bizonyosságértékek (relevancia)
 - Mivel érdemes kezdeni
 - Mennyire biztos magában a rendszer
 - Nem feltétlenül valószínűség (%)

Kód	Bizonyosság
I48H0	0,762
I4710	0,144
I4900	0,048
I5130	0,027
I4990	0,019

Megközelítési módok

- Szabályalapú
 - Elkészítéséhez tárgyterületi tudás kell
- Természetesnyelv-feldolgozás (NLP)
 - Morfológiai elemzés
 - Előfeldolgozás
 - Szintaktikai elemzés itt nem releváns
 - Szemantika
 - Ontológiákkal
- **Gépi tanulás, statisztikai osztályozás**
 - Nem kell tárgyterületi szaktudás
 - Minták kellene
 - Számításigényes
 - De egyre gyorsabb gépek

OSZTÁLYOZÓ MÓDSZEREK

Gépi tanulás, statisztika

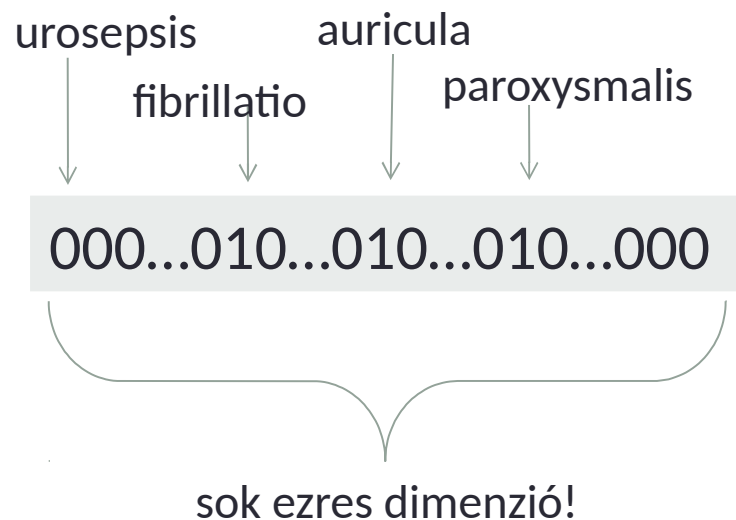
- Vektortér
- Naiv Bayes-háló
- Neurális hálózatok
 - Többrétegű perceptron
- SVM
- Kevert módszerek

Vektorizálás

- Vektorbemenet kell
- Szóhalmaz-modell (*bag of words*)
 - Elvész információ
 - Sorrend
 - Diagnózisoknál nem lényeges
 - Szóhasonlóságok (pl. morfológia, szinonimák)
 - Előfeldolgozás

Példa

„fibrillatio aricula paroxysmalis”



Vektortér

- Folytonossági hipotézis
- Összehasonlítás sokdimenziós vektortérben
 - Koszinuszos hasonlóság mérték
 - Koszinuszos hasonlóság mérték
 - $$\text{cosSim}(v, w) = \frac{v \cdot w}{\|v\| \cdot \|w\|}$$
- Gondok
 - Sok összehasonlítás, nem hatékony
 - Sok összehasonlítás, nem hatékony
 - Invertált index (96%-kal kevesebb összehasonlítás)
 - Invertált index (96%-kal kevesebb összehasonlítás)
 - Azonosan kezeli a dimenziókat (szavakat)
 - Azonosan kezeli a dimenziókat (szavakat)
 - Súlyozás kell
 - Súlyozás kell

IDF-súlyozás

- „inverse document frequency”
- Gyakran szereplő szó kevésbé fontos
- Információelméletileg is alátámasztható

$$v'_j = v_j \cdot \log \left(\frac{|D|}{|\{w \in D \mid w_j = 1\}|} \right)$$

Bayes

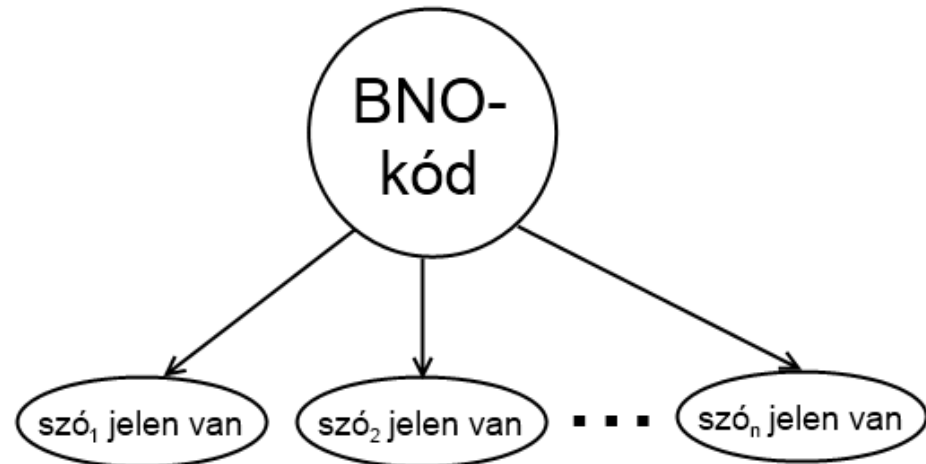
- **Relevancia: feltételes valószínűség**

$$p(C_i | \mathbf{v}) = p(\mathbf{v} | C_i) \cdot \frac{p(C_i)}{p(\mathbf{v})}$$

- feltételes együttes eloszlás becslése mintákból
- $p(\mathbf{v} | C_i)$ feltételes együttes eloszlás becslése mintákból
 - **Kezdeti, sok dimenziós egyengetegre, sok dimenziós paraméter**
- **Bayes-tétel** árolandó paraméter
- **Bayes-tétel** több leírás feltételes függetlenségek fennállásakor
- **Naïv Bayes** több leírás feltételes függetlenségek fennállásakor
- **Naïv Bayes** Minden attribútum (szó) páronként feltételesen független, ismerve a célváltozót (BNO-kód)
 - Minden attribútum (szó) páronként feltételesen független, ismerve a célváltozót (BNO-kód)

Naiv Bayes

$$p(\mathbf{v}|C_i) = \prod_j p(v_j|C_i)$$

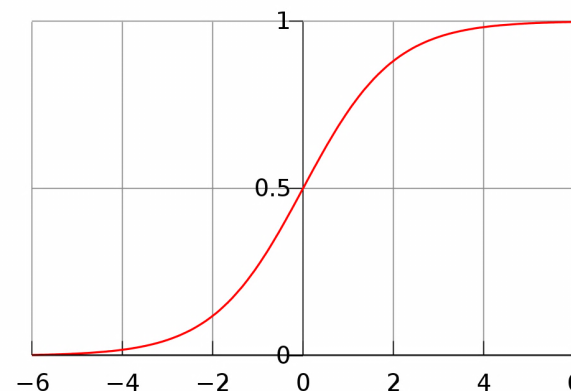
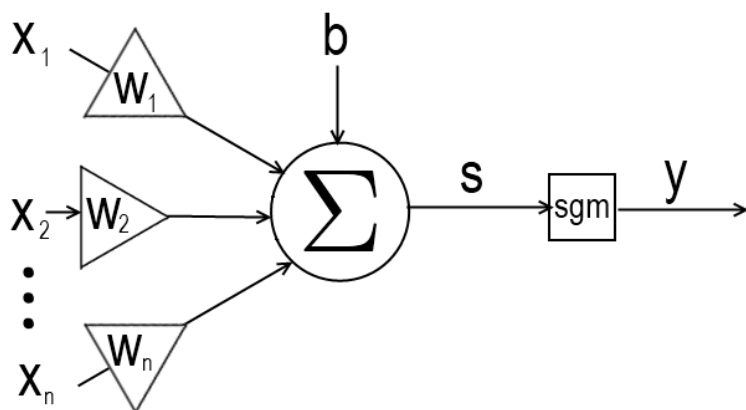


- Ha a mintából becsülve 0, elvesz információ
 - Laplace-simítással nem lehet 0
- Ha $p(v_j|C_i)$ a mintából becsülve 0, elvesz információ
 - Laplace-simítással nem lehet 0

$$\hat{p}(v_j|C_i) = \frac{\theta + |\{\mathbf{w} \mid w_j = v_j \wedge \mathcal{D}(\mathbf{w}) = C_i\}|}{2 \cdot \theta + |\{\mathbf{w} \mid \mathcal{D}(\mathbf{w}) = C_i\}|}$$

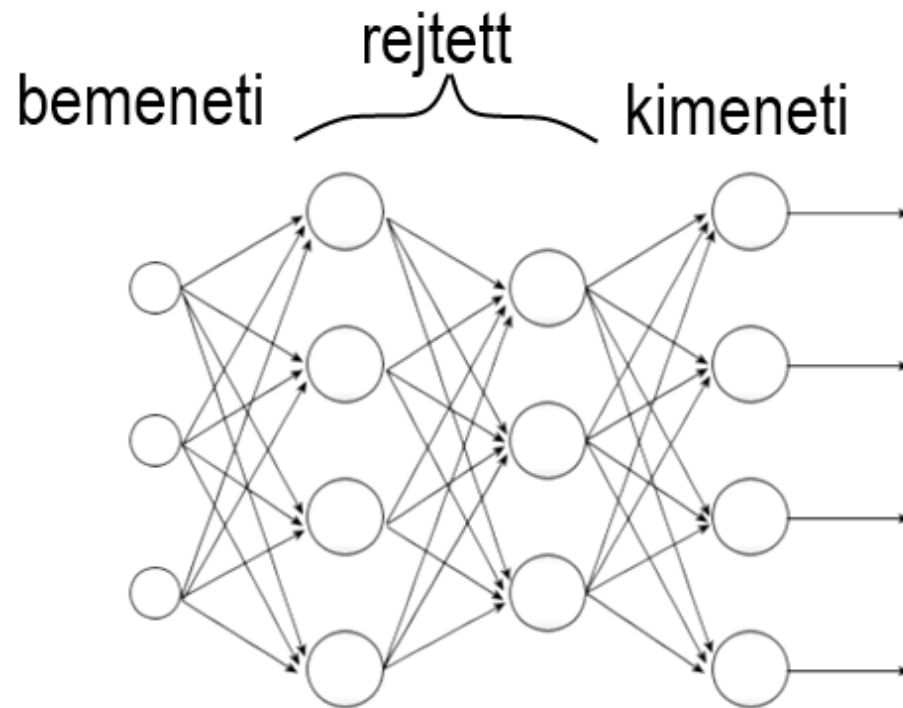
Neurális hálózatok

- Biológiai mechanizmusok alapján konstruált számító rendszerek
- Felügyelt tanulásra egyik legnépszerűbb a többrétegű perceptron (MLP)
 - Elemi neuronja: perceptron, kimenetén szigmoid nemlinearitással



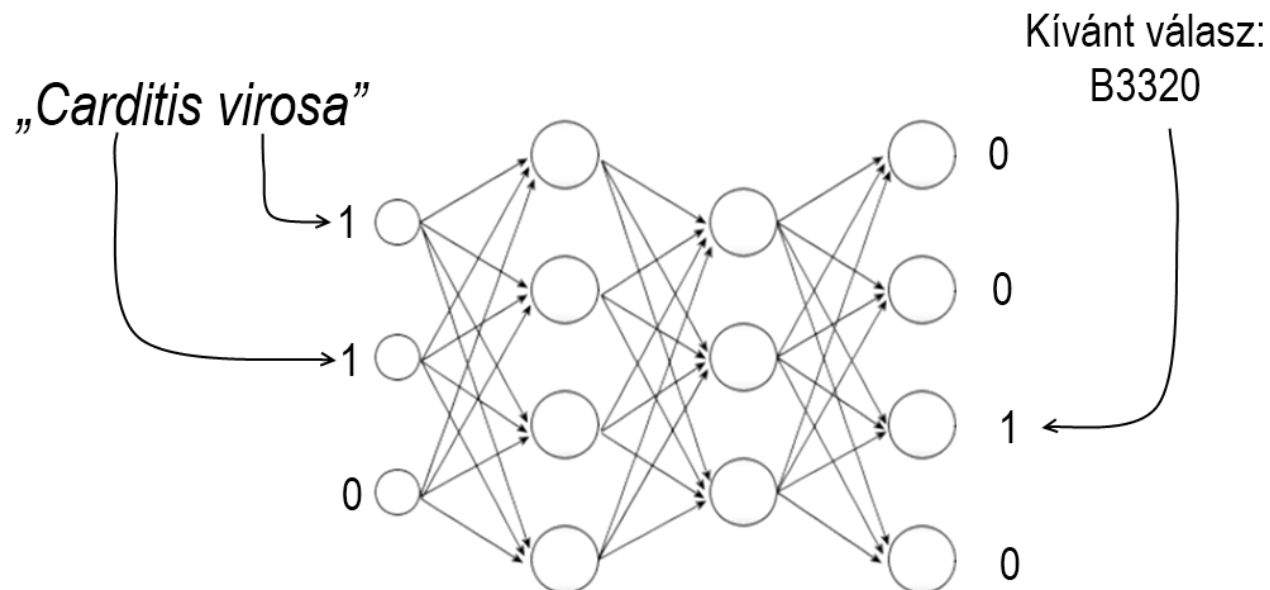
$$sgm(s) = \frac{1}{1 + \exp(-s)}$$

Többrétegű perceptron



Többrétegű perceptron

- Kimeneti neuronok a kódok relevanciáját becslik
 - Több ezer kimeneti neuron



Többrétegű perceptron

- Tanítása hibakritérium alapján

- Várható mégyzetes hiba

- Maximum likelihood (más néven keresztentrópia)

$$J(W) = \mathbb{E}_x \left\{ \sum_{j=1}^n [d_j(x) - y_j(x, W)]^2 \right\}$$

- Maximum likelihood (más néven keresztentrópia)

- Súlytérben a hiba negatív gradiense felé mozdulunk

$$J(W) = \mathbb{E}_x \left\{ - \sum_j (d_j(x) \cdot \log[y_j(x, W)] + (1 - d_j(x)) \cdot \log[1 - y_j(x, W)]) \right\}$$

- Tanítás ciklusokban

- Hiba-visszaterjesztés (láncszabály alapján)

- Amíg a hiba csökken (csúszóablak)

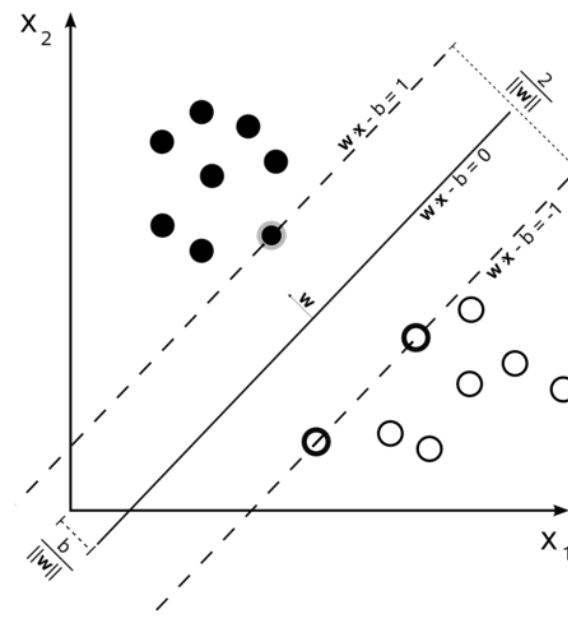
- Tanítás ciklusokban

- Hiba-visszaterjesztés (láncszabály alapján)

- Amíg a hiba csökken (csúszóablak)

Szupport vektor gép (SVM)

- Lineáris változat
 - Nagy margó (biztonsági sáv) az elválasztó hipersík körül
- Lágú lineáris
 - „Rossz helyen” lévő minta megengedett, de büntetjük
 - Együttes optimalizálás (súlyozva)
- Nemlineáris
 - Nagydimenziójú jellemzőtérben lineáris
 - Szövegosztályozásnál általában nem szükséges



Többosztályos SVM

- Dekomponálás bináris feladatokra
- Megkülönböztethetünk
 - Egy osztályt a többitől (őket egybevonva)
 - Páronként az osztályokat
- SVM feladat átfogalmazása (több feltétel)

Keveréses módszerek

Kód	Bizonyosság
B	0,27
C	0,26
A	0,25
D	0,11
E	0,11

Kód	Bizonyosság
D	0,28
A	0,25
B	0,17
E	0,16
C	0,14



Kód	Bizonyosság
A	0,25
B	0,22
C	0,20
D	0,195
E	0,135



Kiértékelési mérce

- **Találási arány (teljesség, recall)**
 - A lista megengedett hosszától függ
 - $$R(h) = \frac{\# \text{ helyes kód a listán}}{\# \text{ tesztbemenetek}}$$
 - Súlyozás lehetséges a listán elfoglalt hely szerint (nem csak, hogy rajta van-e), relevancia szerint stb.
 - Súlyozás lehetséges a listán elfoglalt hely szerint (nem csak, hogy rajta van-e), relevancia szerint stb.

Keveréses módszerek

- Több osztályozó eredménylistájának összevonása
- Jobb eredmény, mint a legjobb különálló
 - Ha a hibás kódok „véletlenszerűek” vagy a sorrendjük véletlenszerű, de a jó kód nagyrészkön szerepel
 - Optimális konstans súlyok keresése (**konstans súlyozó**)
 - Ha a bemeneti tér különböző részein jók
 - Bemenetfüggő optimális súly tanulása (**súlybecslő**)
 - Tanuljuk meg, melyik hol milyen jó (**jóságbecslő**)

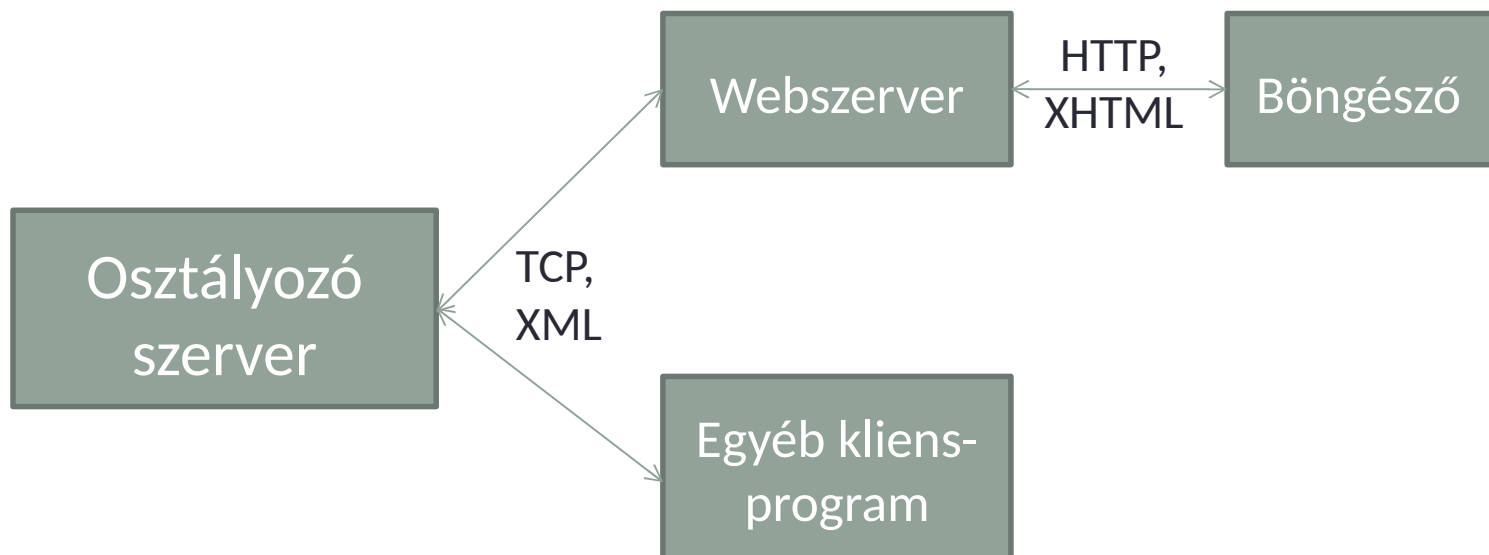
Keveréses módszerek

- Keverendő osztályozók előállítása
 - Eltérő modell
 - Eltérő tanítóminták
 - BNO hierarchikus
 - Főcsoportonként külön osztályozó
(**szakértőegyüttes**, mixture of experts)
 - Külön kapuzó osztályozó tanulja: melyik bemenet melyik főcsoporthoz

IMPLEMENTÁCIÓ

Implementáció

- Alapvető felépítés



Osztályozóprogram

- Különböző módok
 - Tanítás és a kapott osztályozó fájlba szerializálása
 - Fájlból beolvasott osztályozóval TCP-n figyelés, érkező kérések (diagnózis) kiszolgálása (kódlista)
 - Keresztkiértékelés
- Szükséges
 - Osztályozó felépítésének megadása (XML)
 - Üzenetek (kérdés, válasz) formátuma (XML)

XML példák: Lekérés szervertől

```
<?xml version="1.0" encoding="UTF-8"  
standalone="no"?>  
<request>  
  <inputs>  
    <input>ruptura corpus lutea haemorrhagia</input>  
  </inputs>  
  <resultcount>5</resultcount>  
</request>
```

XML példák: Szerver válasza

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<resultsets>
  <resultset>
    <result>
      <class>N8310</class>
      <confidence>0.38768336305190293</confidence>
    </result>
    <result>
      <class>D6990</class>
      <confidence>0.17026160892877706</confidence>
    </result>
    <result>
      <class>H4310</class>
      <confidence>0.15751956314998392</confidence>
    </result>
    <result>
      <class>R58H0</class>
      <confidence>0.15043689961141468</confidence>
    </result>
    <result>
      <class>K2280</class>
      <confidence>0.13409856525792138</confidence>
    </result>
  </resultset>
</resultsets>
```

Összetett osztályozó megadása (részlet)

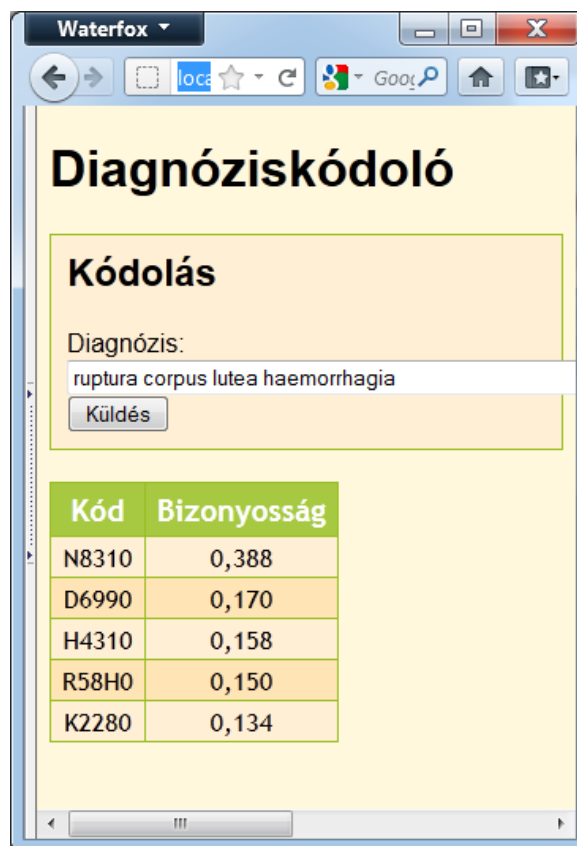
```
<classifier id="team">  
  <type>ConstantWeightTeam</type>  
  <param name="c1">  
    <classifier ref="moe" />  
  </param>  
  <param name="c2">  
    <classifier ref="idf" />  
  </param>  
  <param name="ratio">0.4</param>  
</classifier>
```

Keretrendszer

- Java nyelven készítettem el
- Osztályozási feladatok kezelése általánosan
 - Tanítás
 - Kiértékelés
 - Bemeneti és kimeneti transzformációk
 - Mintahalmazok manipulációja (keverés, vágás, iteráció)
 - Eredményhalmazok kezelése
 - Párhuzamosítás egyszerűen
- Egyszerű bővíthetőség
 - Új osztályozók
 - Új kiértékelési mércék
 - Új transzformációk (pl. nyelvi feldolgozás)

Webes felület

- JSP prototípus oldal

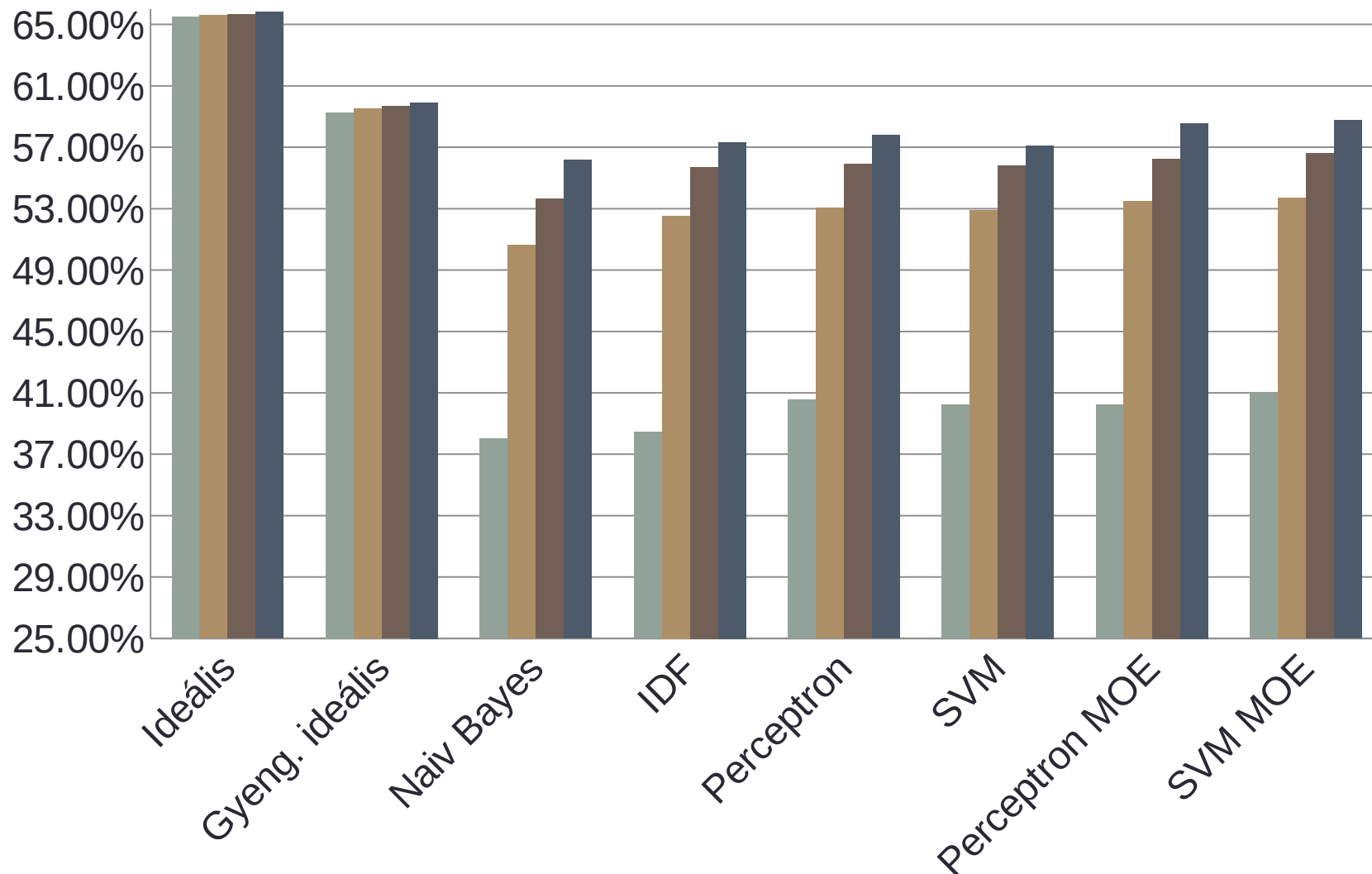


EREDMÉNYEK

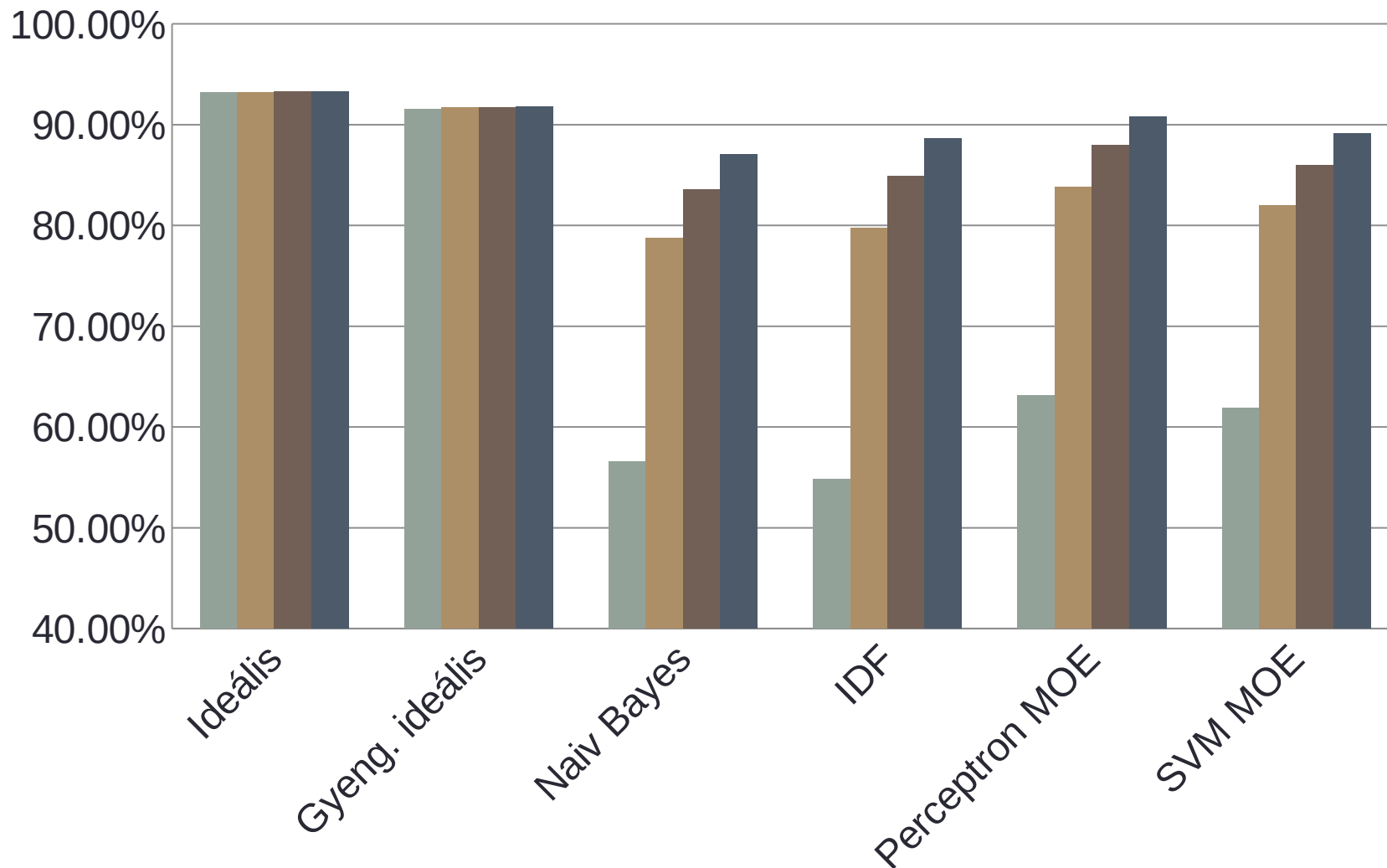
Mintahalmazok

- Magyar (3081 minta)
 - Tisztítatlan
 - Tisztított
- Német (93863 minta)
 - Feldolgozatlan
 - Szótövezett
 - Morfológiailag elemzett

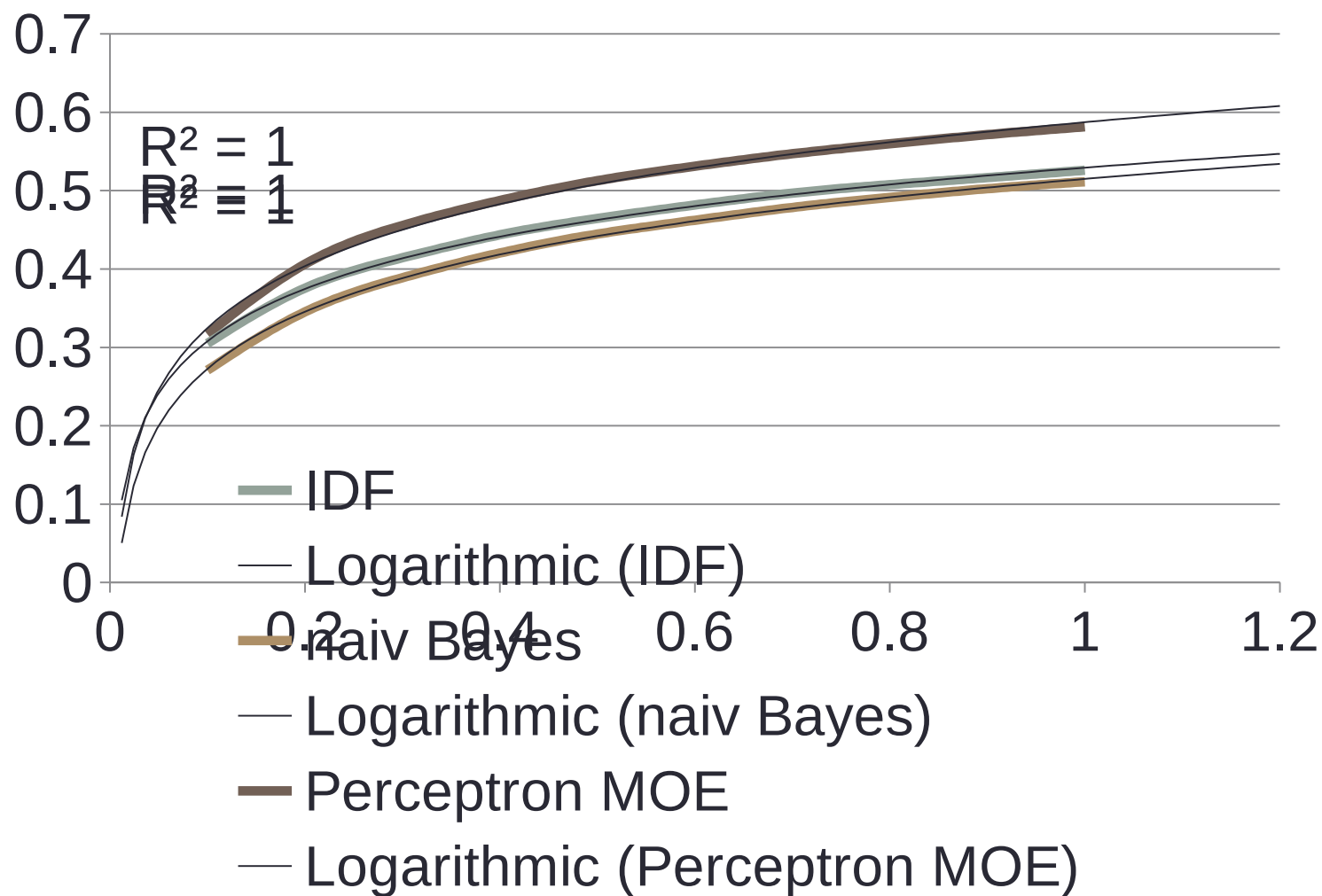
Magyar minta (tisztított)



Német minta (morfológiailag elemzett)



Tanulási görbék (német feldolgozatlan)



Továbbfejlesztési lehetőségek

- A különböző algoritmusok hasonlóan teljesítenek
- Nagyobb különbséget okoz a minta jósága, mérete

- Előfeldolgozás finomítása
 - Morfológiai elemzés
- Ontológiák használata
 - Pl. speciális fogalmak általánosítása
 - Szinonimák felismerése

- Felhasználói felület bővítése, BNO-adatbázissal összekötés

KÖSZÖNÖM A
FIGYELMET!
