

Üzleti Elemzések Módszertana

2. Előadás: Osztályozás

Kuknyó Dániel
Budapesti Gazdasági Egyetem

2023/24
2.félév

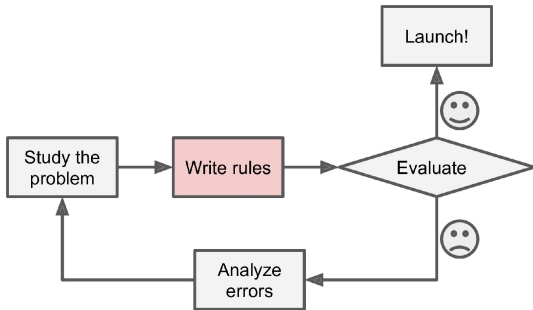
- 1 Bevezetés
- 2 Osztályozás
- 3 Osztályozás vagy regresszió?
- 4 Osztályozás jósága

- 1 Bevezetés
- 2 Osztályozás
- 3 Osztályozás vagy regresszió?
- 4 Osztályozás jósága

A determinisztikus szemléletmód

A hagyományos szoftverfejlesztési folyamatmodell eljárása:

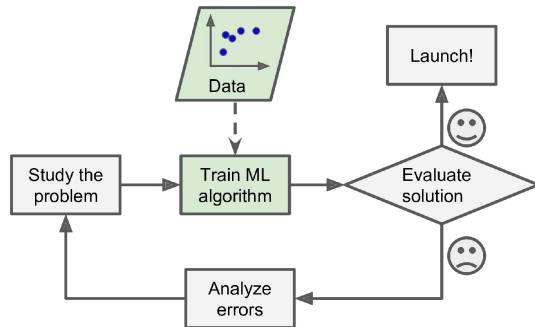
- 1 Az adott jelenség megfigyelése és adatok rögzítése
- 2 A megfigyelésekre olyan szabályok kidolgozása, amelyek jól leírják azt
- 3 A létrejött szabályrendszer kiértékelése
- 4 Rendszer fejlesztése a hibák alapján
- 5 Iteráció



A gépi tanulás szemléletmód

A gépi tanulás szemléletének
folyamatmodellje:

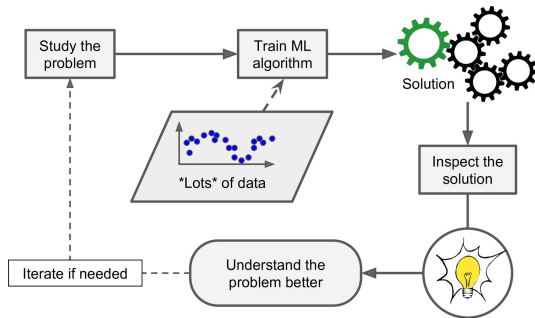
- 1 Adott jelenség megfigyelése és adatok rögzítése
- 2 Gépi tanulási modell tanítása az adatokon a szakterületi tudás segítségével
- 3 Modell kiértékelése
- 4 Hibák elemzése és kiértékelése
- 5 Iteráció



Tanítás automatizálása adatalapúan

Az gépi tanuló modellek tanítása és kiértékelése hosszú távon egy iteratív folyamat már létező keretrendszerrel, mint az MLOps. Ennek számos területen vannak előnyei:

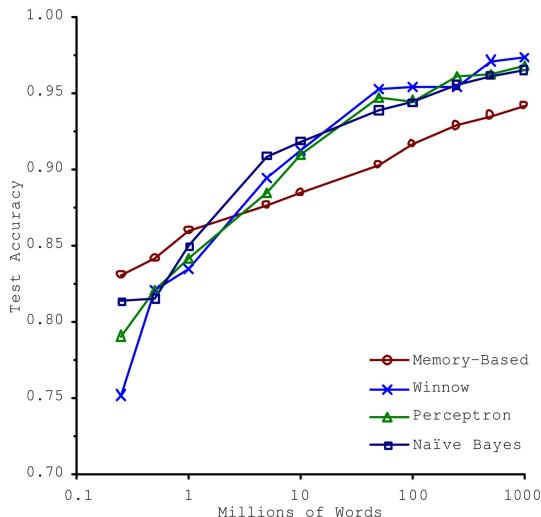
- Adaptáció az új adatokhoz
- Javuló modell teljesítmény
- Hibák és problémák azonosítása
- Új technológiai fejlődés integrálása
- Skálázhatóság és rugalmasság
- Szakterületi következtetések az elemzések által



Az adatok észszerűtlen hatékonysága

2001-es kutatásukban Michele Blanko és Eric Brill kimutatták, hogy a különböző ML algoritmusok **hasonlóan jól teljesítenek a természetes nyelvfelismerés területén mint a hagyományos algoritmusok**, ha elég sok adaton tanítják a modelleket. Ahogy ők fogalmaztak:

„Az eredmények azt mutatják, hogy újra kell gondolnunk, mire fordítjuk a pénzünket és erőforrásainkat: algoritmusok fejlesztésére, vagy adatgyűjtésre.”

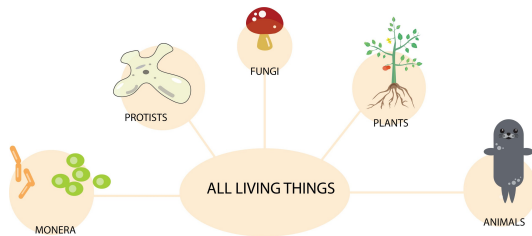


- 1 Bevezetés
- 2 **Osztályozás**
- 3 Osztályozás vagy regresszió?
- 4 Osztályozás jósága

Osztályozás

Osztályozás

Az osztályozás a felügyelt gépi tanulás egyik alapvető feladata, amelynek célja, hogy megtanuljon egy modellt vagy szabályrendszert egy adott bemeneti adat alapján annak **besorolására előre meghatározott kategóriákba vagy csoportokba.**

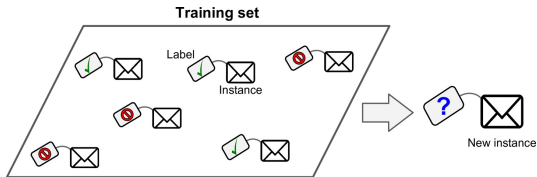


Five Kingdom system classification

Modellalapú osztályozás

Az osztályozó modell feladata, hogy a tanító adathalmaz alapján olyan szabályrendszert hozzon létre, ami **képes elszeparálni egymástól az egyedeket**.

Amennyiben érkezik egy új adatpont, a modell a saját szabályrendszere segítségével már **képes lesz becslést adni annak osztályára vonatkozóan**.

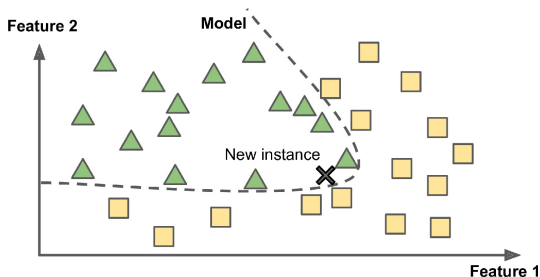


Modellalapú osztályozás

Döntési határ

Olyan **határérték**, amelyet a **modell állít be** az adatpontok különböző osztályokba való besorolásához.

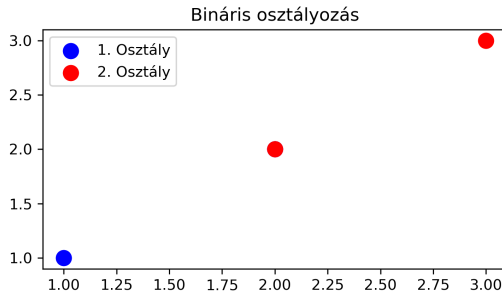
A határ **lehet egy vonal, egy sík vagy akár egy sokdimenziós felület**, attól függően, hogy milyen típusú osztályozó modellt használunk és milyen a bemeneti adatok dimenzionalitása.



Az osztályozás fajtái

Bináris osztályozás

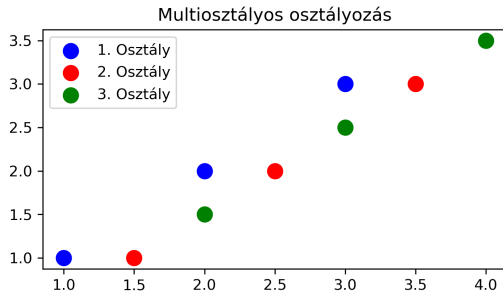
A modell két lehetséges osztály közül valamelyikbe sorolja be az egyedeket. Minden egyedhez csakis 1 osztály tartozhat.



Az osztályozás fajtái

Multiosztályos osztályozás

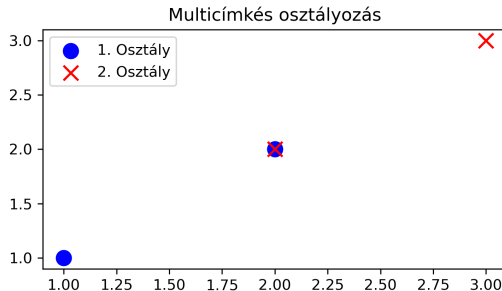
Több, mint két lehetséges kategória létezik, amibe az egyedek besorolhatók, ezek közül az egyikbe fog sorolódni az egyed. Minden egyedhez legalább és legfeljebb 1 osztály tartozik.



Az osztályozás fajtái

Multicímkes osztályozás

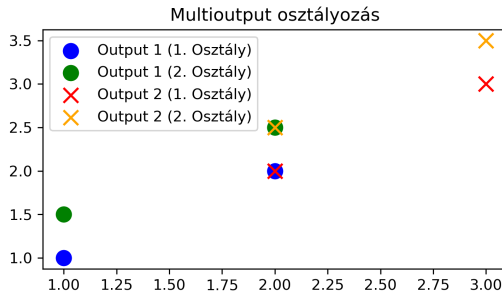
Minden mintaegyedhez több bináris vagy multicímkes címke kategóriából tartozhat osztály.



Az osztályozás fajtái

Multioutput osztályozás

A multicímkes osztályozás generalizált változata. Egy egyedhez egy multicímkes halmazból több elem is tartozhat.



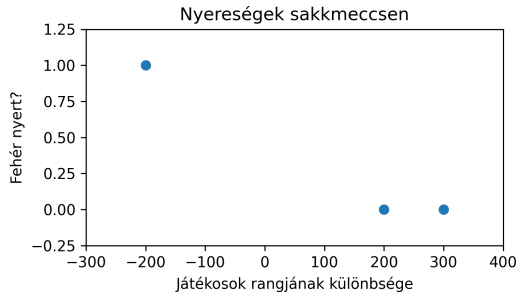
- 1 Bevezetés
- 2 Osztályozás
- 3 Osztályozás vagy regresszió?
- 4 Osztályozás jósága

Példa: a probléma bemutatása

A következő kis adathalmaz három sakkjátszmának rögzítette az eredményét. Minden meccs esetén rögzítésre kerültek a következő rekordok:

Különbség	Nyertes
200	0
-200	1
300	0

Ebben az esetben az x változó, a **két játékos rangjának különbsége** a fehér és fekete játékos különbségét jelzi, az y célváltozó pedig egy azt a valószínűséget jelenti, hogy **a fehér nyert-e**.



Példa: lineáris predikció

Az adathalmazra egy lineáris regresszor modellt illesztve az eredmény a következő:

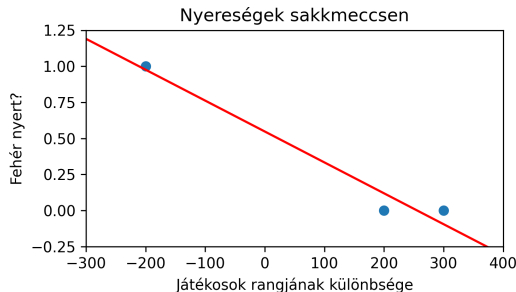
Különbség	Nyertes	Predikció
200	0	0.11
-200	1	0.97
300	0	-0.1

Ebben az esetben a lineáris modell:

$$\hat{y} = \theta_0 + \theta_1 \cdot x$$

Ahol \hat{y} a modell predikciója a nyertesre vonatkozóan, θ_0 a konstans torzítás, θ_1 a függvény meredeksége és x a két játékos rangjának különbsége.

Az adatpontokra egy lineáris regressziós függvényt illesztve az illesztett modell a következő lesz:

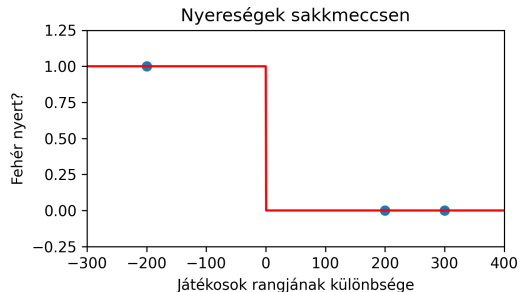


Példa: következtetések

A lineáris modell nem minden esetben ad racionális predikciót az adathalmazra vonatkozóan.

Negatív valószínűségeket nem értelmeztettek!

Éppen ezért ha a modellezés célváltozója egy valószínűség, szükség van arra, hogy az illesztett modell szélsőértéke 0 legyen ha a hely $-\infty$ és 1 ha a hely ∞ .

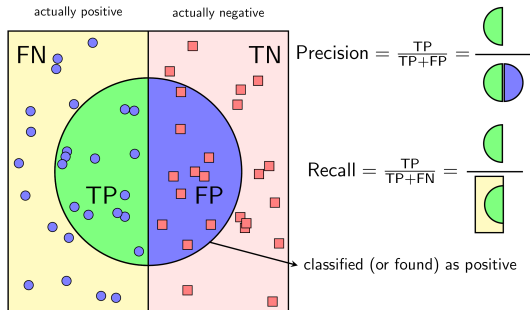


- 1 Bevezetés
- 2 Osztályozás
- 3 Osztályozás vagy regresszió?
- 4 Osztályozás jósága**

Az osztályozás teljesítményének mérése

- **Valós pozitív (TP):** Pozitív egyed, és annak is van osztályozva
- **Valós negatív (TN):** Negatív egyed, és annak is van osztályozva
- **Hamis pozitív (FP):** Negatív egyed, de pozitívnak van osztályozva
- **Hamis negatív (FN):** Pozitív egyed, de negatívnak van osztályozva

Ennek alapján két fő mutatószám áll elő, amellyel egy osztályozó modellt lehetséges értékelni:



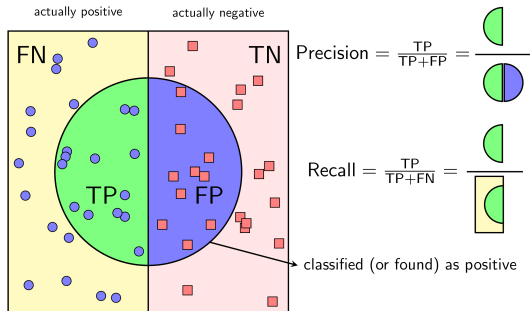
Az osztályozás teljesítményének mérése

Ennek alapján két fő mutatószám áll elő, amellyel egy osztályozó modellt lehetséges értékelni:

Pontosság

Megadja, hogy a pozitívnak osztályozott egyedek közül mekkora hányad volt ténylegesen pozitív:

$$P = \frac{TP}{TP + FP}$$



Az osztályozás teljesítményének mérése

Ennek alapján két fő mutatószám áll elő, amellyel egy osztályozó modellt lehetséges értékelni:

Visszahívás

Megadja, hogy az összes pozitív egyed mekkora hányadát osztályozta a modell pozitívnak:

$$R = \frac{TP}{TP + FN}$$

