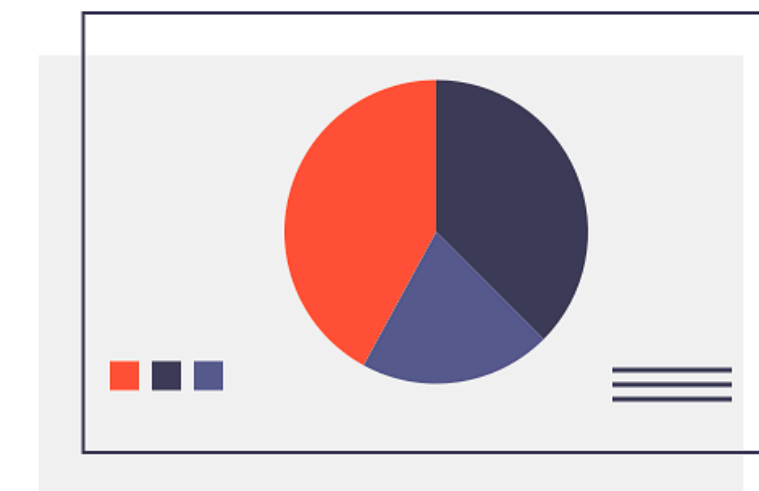


Futóverseny eredményeinek kiértékelését támogató eszköztár és alkalmazás fejlesztése

Témavezető: dr. Leitold Dániel

Ultramarathon

- Népszerű
- Nagy kihívás
- 4 napos esemény
- Többnyire publikus adathalmaz
- Data Science



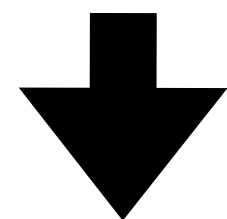
Felépítés

- Adatösszegyűjtés / tisztítás
- Irodalomkutatás
- Kutatások alkalmazása az adathalmazon
- Modell-építés
- Kiértékelés
- Alkalmazás fejlesztése



Adathalmaz

- 2015 - 2020 közötti eredmények
- Számos tulajdonság
- Összesítve is alacsony részvételi arány



- ~ 1000 sornyi adat

Nominális	Folytonos	Diszkrét	Ordinális
Név	Helyezés	Rajtszám	Kategória
Ország	Születési év	Esemény éve	
Klub	Napi időadatok		
Város	Eredmény		
Nem	Futott kilométer		
Rajtszám			

Első iteráció

Cél megfogalmazása

Teljesíthető?

Model építés

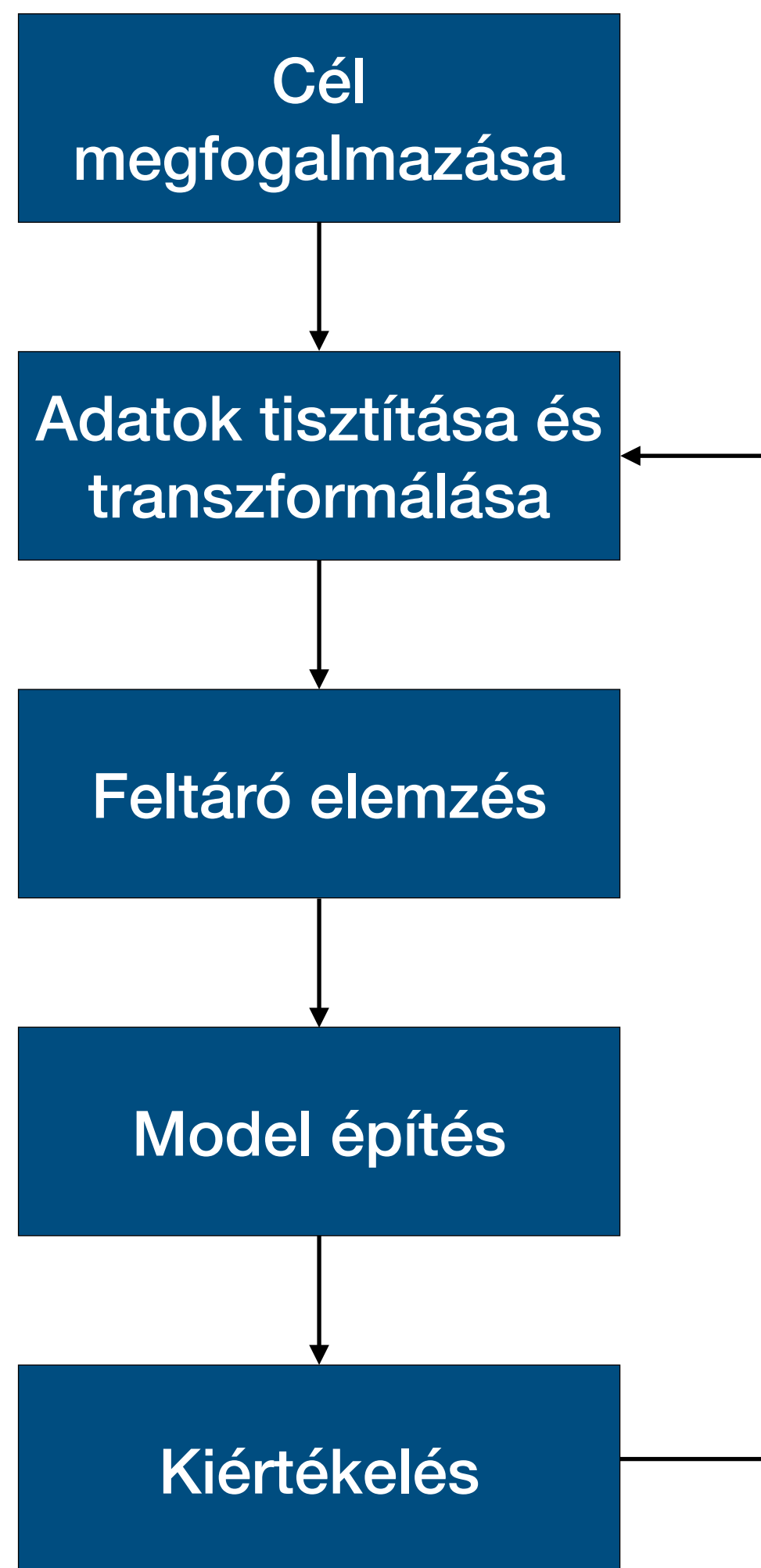
- Döntési fa osztályozó
- Random forest osztályozó

Adat előkészítés

- Oszlopok törlése / hozzáadása (Név, Város, stb), (Befejezte, Tempo, stb)
- Értelmetlen adatok törlése
- Időadatok átalakítása
- Adat diszkretizáció (vödrözés)

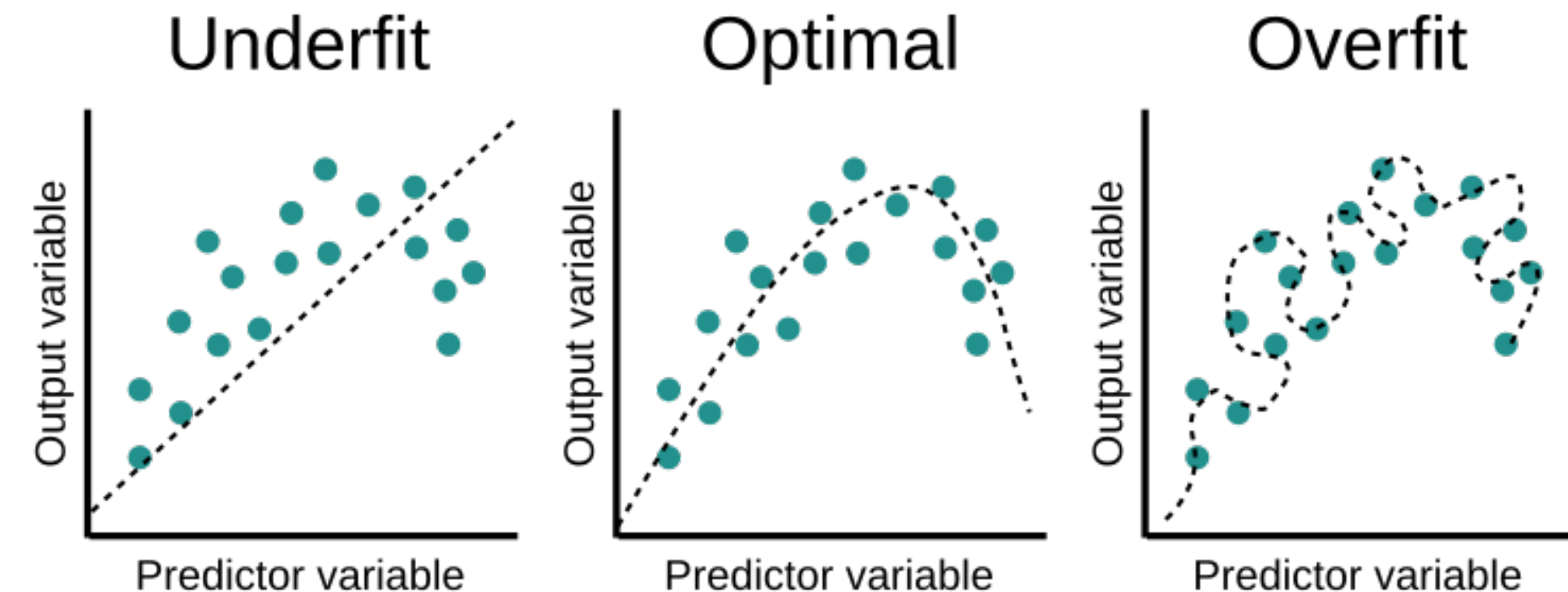
Kiértékelés

- Pontatlan modell (60-65%)
- Több adatra van szükség
- Szervezők bevétele, 2015 előtti eredmények feldolgozása
- Új iteráció



Problémák

- Kis adathalmaz, bővítve is kicsi
2511 sor
- Imbalanced adathalmaz
~ 20% adta fel a versenyt
- Zajos adatok, nagyban eltérő információk
Verseny hossza, stb
- Használhatatlan adatok
2008 - 2011 közötti adatok
- 2/3 vágás nem elegendő
Más megközelítésre van szükség
- Túlillesztés gyakori kicsi adathalmazok esetén
- High bias, mivel egyenlőtlen az adathalmaz
~ 95% pontos modellek, a bővített adathalmazokon
- Tisztítás során sok hamis információt generálhat



Megoldások

- Más gépi tanuló algoritmusok alkalmazása
- Irodalom kutatás
- Feature engineering
- Tulajdonságok bővítése
- Mintavételezési eljárások alkalmazása
- Keresztvalidálási eljárások
- Importance analysis

Nominális	Folytonos	Diszkrét	Ordinális
Név	Helyezés	Rajtszám	Kategória
Ország	Születési év	Esemény éve	Hőmérséklet
Klub	Napi időadatok		Tempó kategóriák
Város	Eredmény		Felhősség
Nem	Futott kilométer		Légnyomás
Rajtszám	Napi szakaszok hossza		Eső mértéke
	Napi szakasz tempók		Szél
	Napi tempók		Széllökések
	Időjárás adatok		

Keresztvalidáció

- Olyan esetekben nagyban ajánlott az alkalmazása, amikor egy model pontossága bizonytalan (sokat ugrál a pontossága)
- Általános megközelítés a 2/3 -os vágás, ami általában működik alkalmas adathalmazokon
- Imbalanced adathalmazokon erősen high bias eredményt ad, ami sok pontatlan előrejelzést eredményez.
- Több fajta keresztvalidálási eljárás is létezik pl.: (Leave One Out CV), de a mi esetünkben kettő vizsgálatán lesz a hangsúly.

K-fold Cross Validation

- Pro
 - Teljesítményt növelheti
 - Ismeretlen adathalmazokat segít jobban megérteni
 - Kis adathalmazon optimális
- Cons
 - Imbalanced adathalmazon high bias
 - Nagy adathalmazon lassú lehet

Stratified K-fold CV

- Pro
 - Valós pontosságot tovább növeli speciális adathalmazokon (pl. Imbalanced).
 - Minden iterációban azonos eloszlású az osztálycímke
 - Kis adathalmazon optimális
- Cons
 - Lassabb

K-fold Cross Validation

Pl.: 1000 sor adatra 5 fold

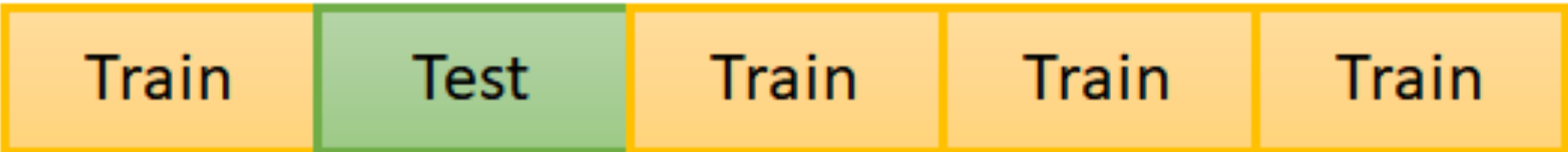
200 test, 800 train
180/20

1st Iteration



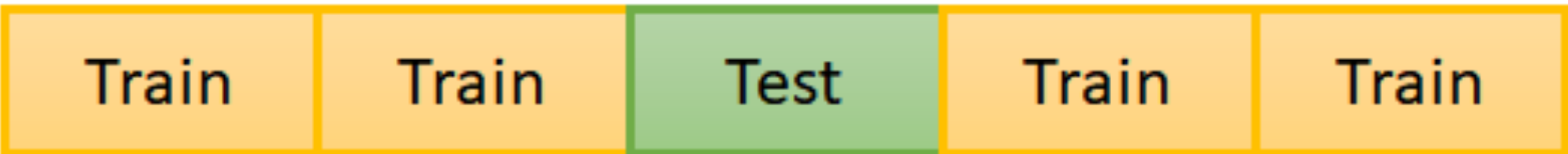
200 test, 800 train
198/2

2nd Iteration



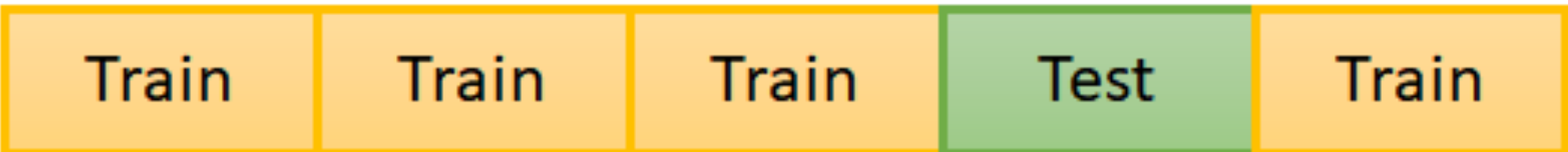
200 test, 800 train
120/80

3rd Iteration



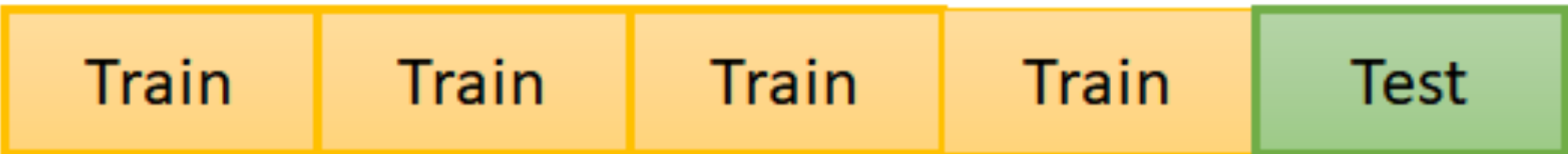
200 test, 800 train
110/90

4th Iteration



200 test, 800 train
200/0

5th Iteration



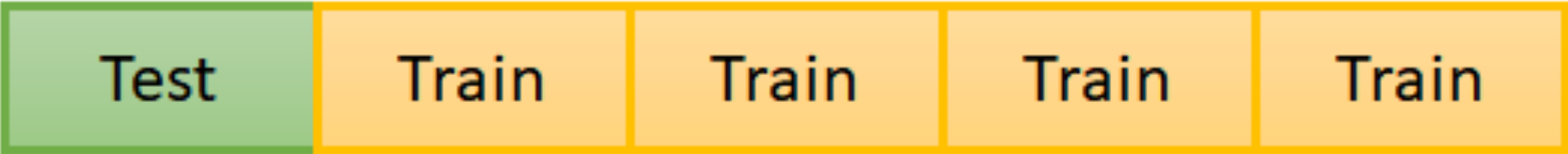
Mean

Stratified K-fold Cross Validation

Pl.: 1000 sor adatra 5 fold

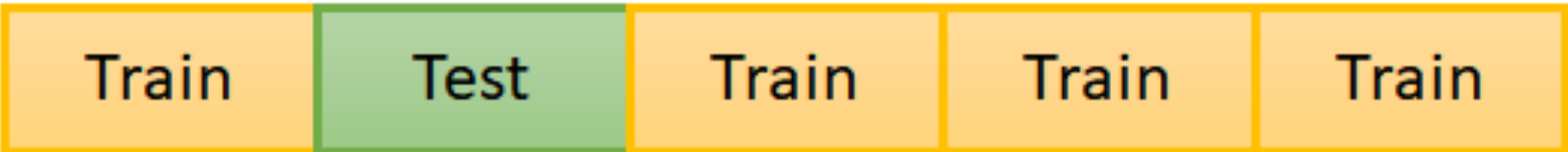
200 test, 800 train
120/80

1st Iteration



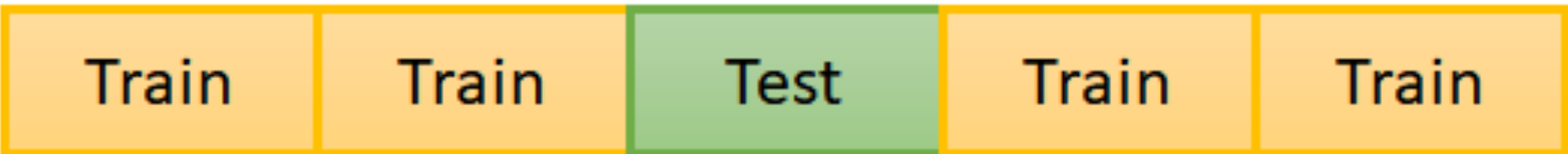
200 test, 800 train
120/80

2nd Iteration



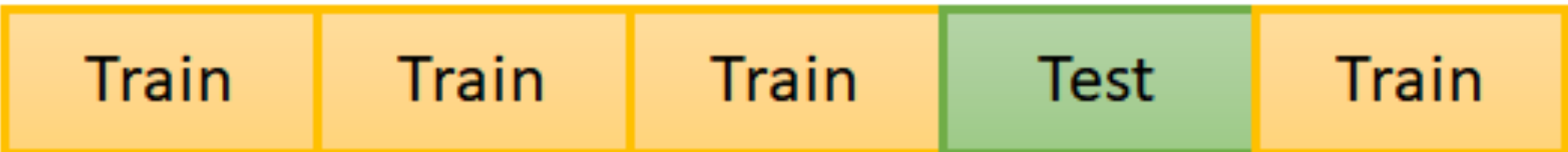
200 test, 800 train
120/80

3rd Iteration



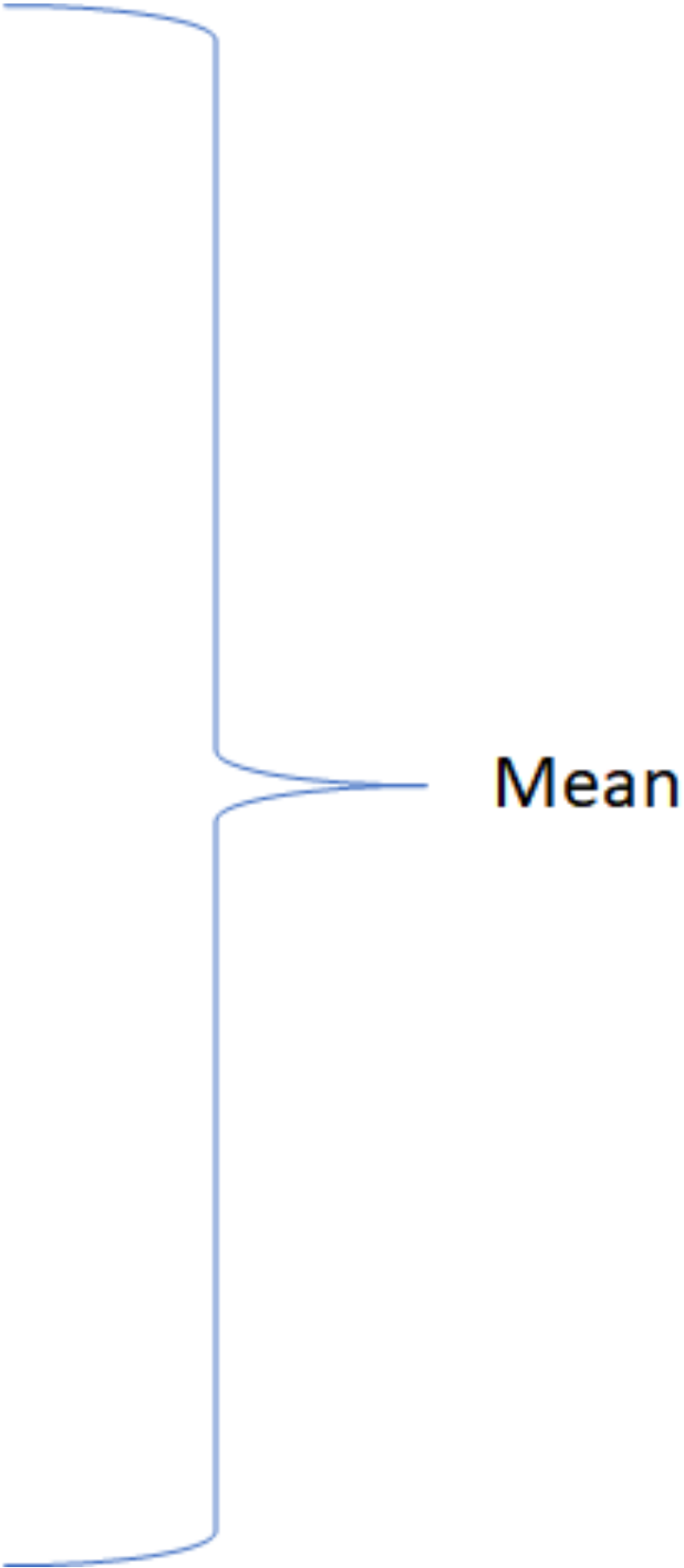
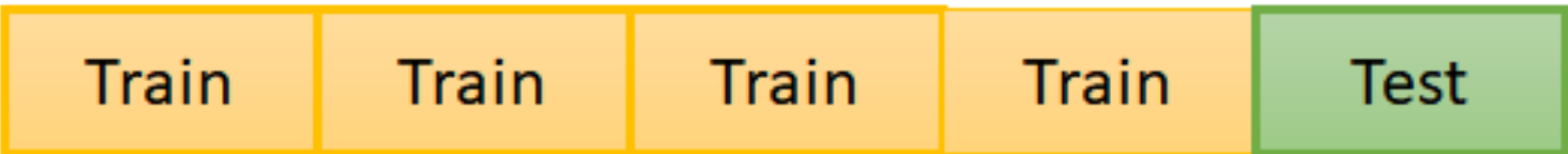
200 test, 800 train
120/80

4th Iteration



200 test, 800 train
120/80

5th Iteration



Második iteráció

Adatok előkészítése

- Összes adathalmaz migrálása
- Normalizáció
- Teljesebb adattisztítás
- Outlier detection
- Feature Engineering

Feltáró elemzés

- Correlation analysis
- DataPrep alkalmazása
- Vizualizáció
- Adathalmaz részletes értelmezése

Model építés

- Random forest osztályozó
- SVM osztályozó (kis adathalmazokra alkalmas)
- Naive Bayesian osztályozó
- Neurális háló osztályozó
- Keresztvalidációs eljárások alkalmazása
- Over / Under sampling eljárások alkalmazása



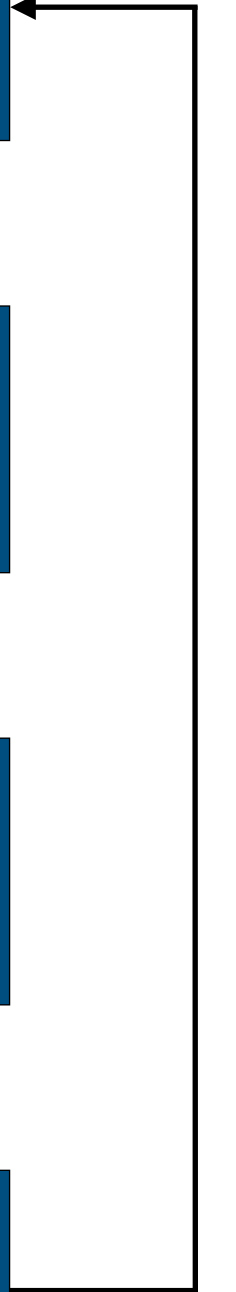
Cél
megfogalmazása

Adatok tisztítása és
transzformálása

Feltáró elemzés

Model építés

Kiértékelés



Konklúzió

60% -> 95%

Nem valós teljesítmény növekedés, az adathalmaz bővült, de ezzel az adathalmaz kiegyensúlyozatlansága is növekedett. Emiatt a modell high bias módon osztályozott, ezért szükséges volt különböző mintavételezési eljárást alkalmazni.

Modellek

Az munka folyamán váltakozó pontosságú modelleket lehetett építeni a tisztított adathalmazból. Ezek közül a legalkalmasabbak a Support Vector Machine alapú osztályozók, és a Naive Bayesian osztályozók.

Cross Validation

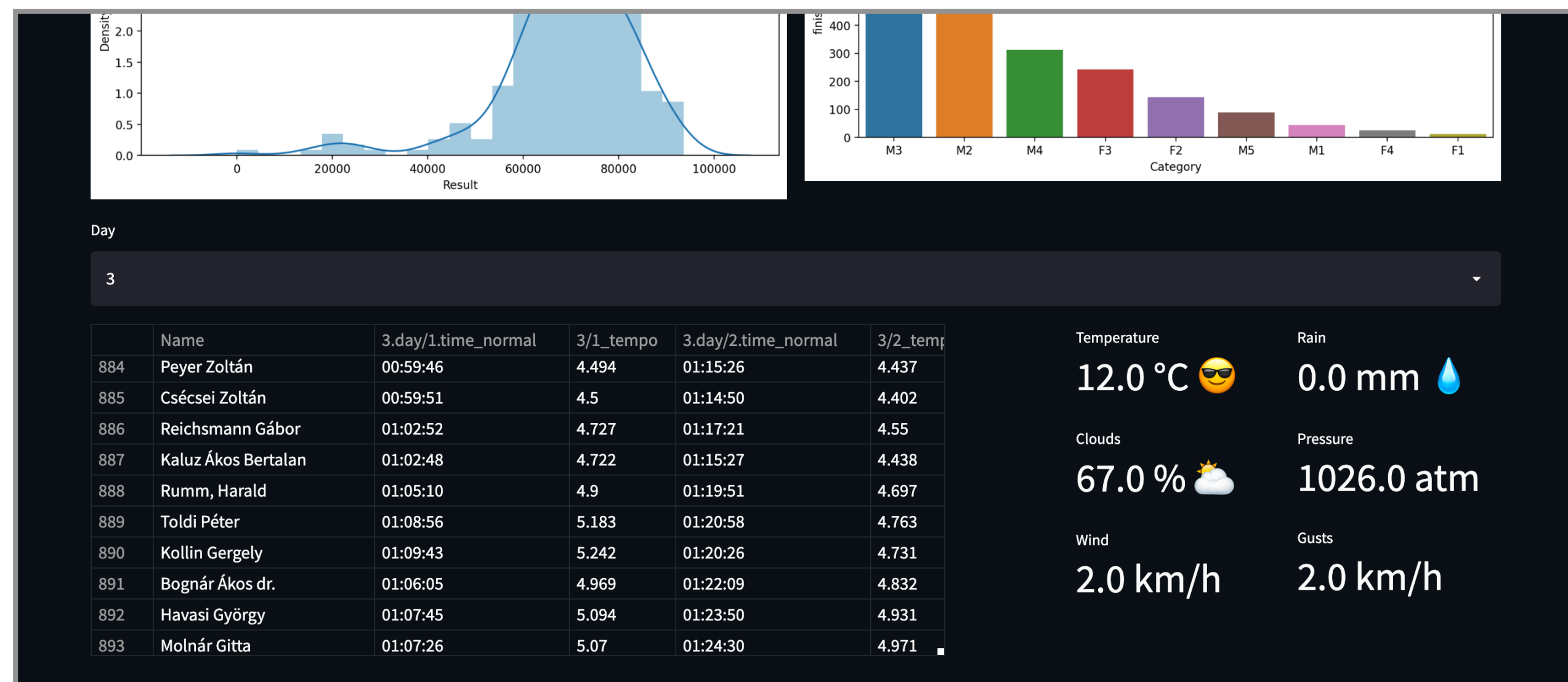
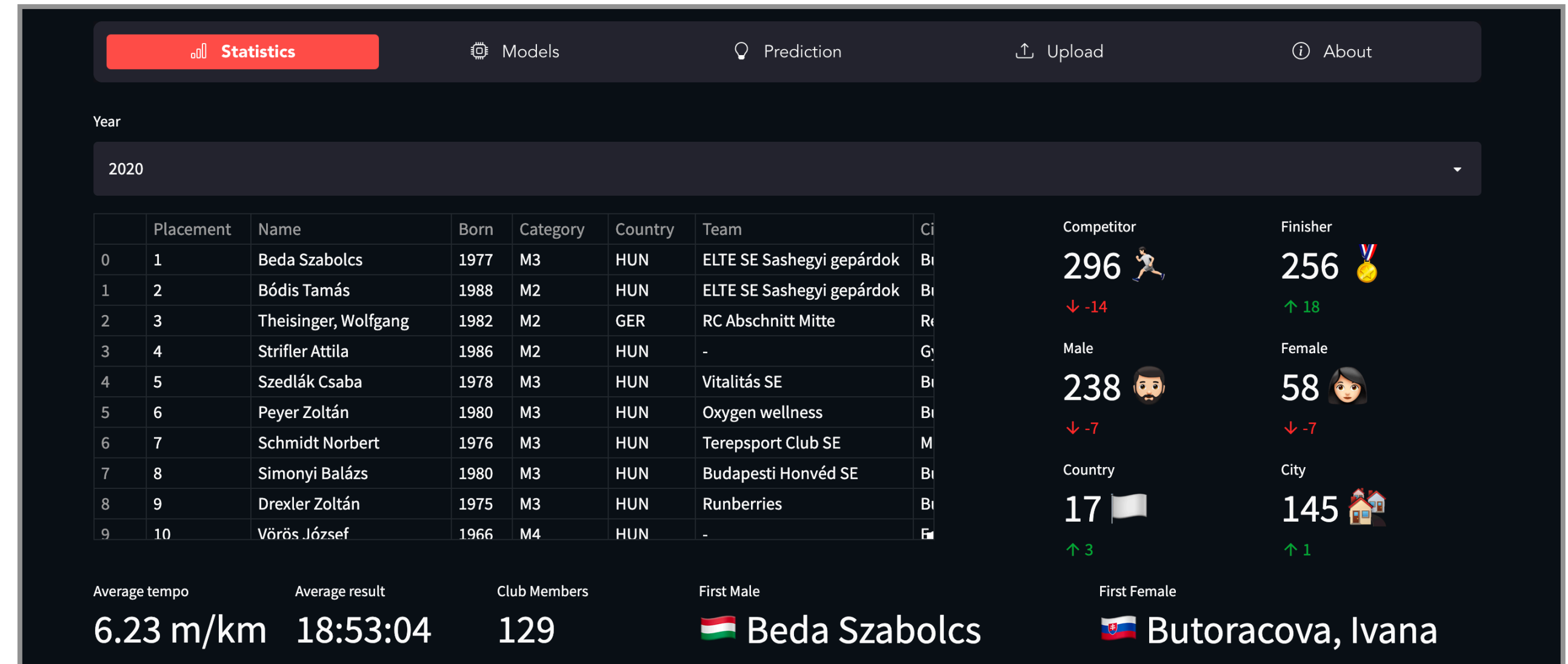
A keresztvalidálási eljárások alkalmazásával ténylegesen sikerült az épített modellek teljesítményét javítani. A legjobban 10 k-nál látszódott a javulás és a Naive Bayesian osztályozóval sikerült a legpontosabb modellt építeni, ami 78,23%

Eredmény

Az első iterációhoz képest nagyobb teljesítményű és pontosabb modellt sikerült építeni. Az adathalmazra és a felhasznált tulajdonságra a legalkalmasabban a Naive Bayesian osztályozó alkalmazható.

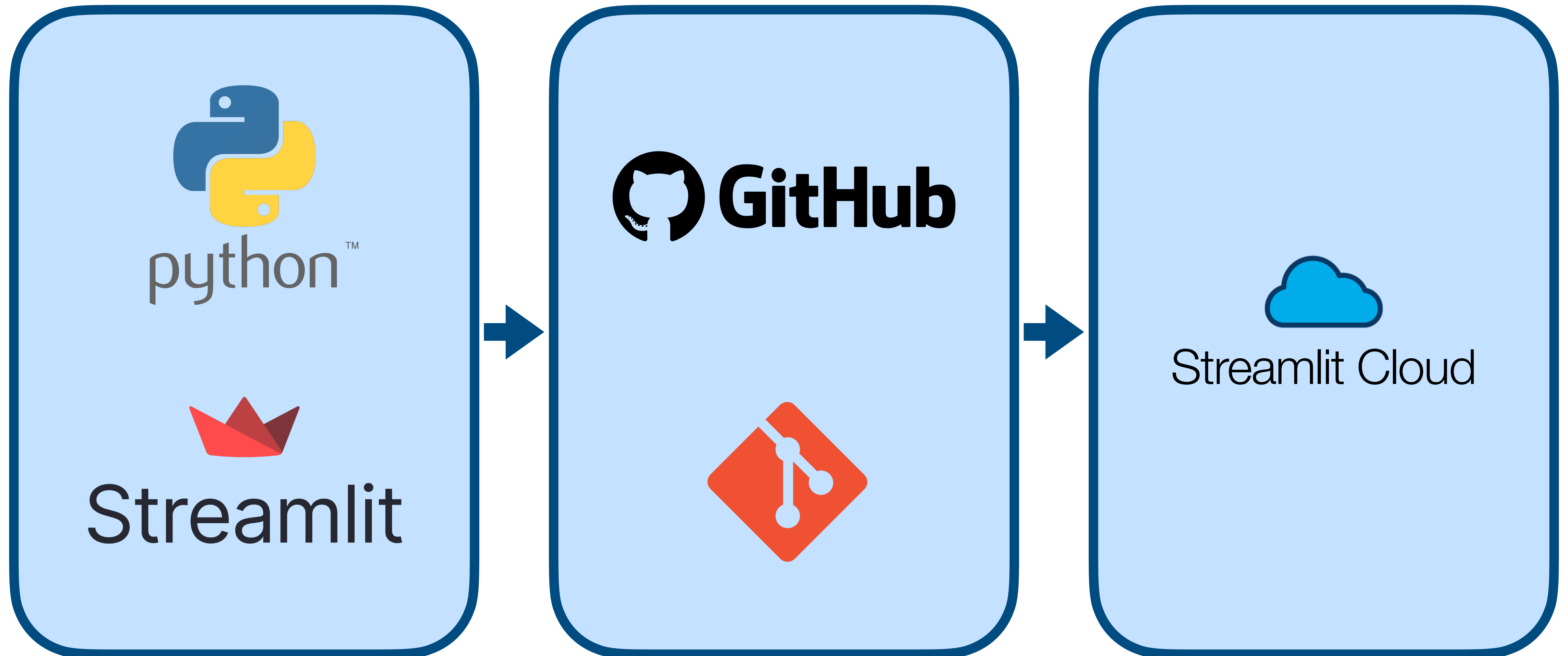
Alkalmazás / Vizualizáció

- Évi / összes eredmények kiértékelése
- Napi eredmények
- Időjárás adatok
- Osztályozó modellek és pontosságuk



- Előrejelzés
- Új évi tesztadat feltöltése
- Információk (BSZM, elérhetőség, stb)

Alkalmazás elérhetősége



Továbbfejlesztési lehetőségek

- Mindig lehet tovább pontosítani, hiperparaméter hangolás
- Új célok -> pl.: Eredmény közelítő becslés
- További együttműködés a szervezőkkel, új kimutatások
- Alkalmazás deployment



Köszönöm a figyelmet!

Csizmazia Máté - XI32IS - Programtervező informatikus MSc

- Téma ismertetése ✓
- Munkafolyamatok és problémák ismertetése ✓
- Megoldások és eredmények ismertetése ✓
- Kiértékelő alkalmazás bemutatása ✓
- Továbbfejlesztési lehetőségek ismertetése ✓

2022 Április 25.