

# SZAKDOLGOZAT



MISKOLCI EGYETEM

Szakdoga címe

**Készítette:**

Ferencsik Dorina Kinga

Évfolyam. szak-szak

**Témavezető:**

Egyik konzulens neve

Másik konzulens neve...

MISKOLC, 2019

**MISKOLCI EGYETEM**  
Gépészmérnöki és Informatikai Kar  
Alkalmazott Matematika Tanszék

**Szám:**

## **SZAKDOLGOZAT FELADAT**

Ferencsik Dorina Kinga (YXI8DJ) programtervező informatikus jelölt részére

**A szakdolgozat tárgyköre:** numerikus eljárások tesztelése

**A szakdolgozat címe:**

**A feladat részletezése:**

Nemlineáris egyenletrendszer megoldására szolgáló ABS-módszerek áttekintő jellegű ismertetése.

**Témavezető(k):** Dr. (vagy nem doktor) Valamilyen Valaki beosztás (pl. egyetemi adjunktus)

**Konzulens(ek):** (akkor kötelező, ha a témavezető nem valamelyik matematikai tanszékről való; de persze lehet egyébként is)

**A feladat kiadásának ideje:**

.....  
szakfelelős

## EREDETISÉGI NYILATKOZAT

Alulírott .....; Neptun-kód: .....  
a Miskolci Egyetem Gépészmérnöki és Informatikai Karának végzős .....  
szakos hallgatója ezennel büntetőjogi és fegyelmi felelősségem tudatában nyilatkozom  
és aláírással igazolom, hogy .....  
című szakdolgozatom/diplomatervem saját, önálló munkám; az abban hivatkozott szak-  
irodalom felhasználása a forráskezelés szabályai szerint történt

Tudomásul veszem, hogy szakdolgozat esetén plágiumnak számít:

- szó szerinti idézet közlése idézőjel és hivatkozás megjelölése nélkül;
- tartalmi idézet hivatkozás megjelölése nélkül;
- más publikált gondolatainak saját gondolatként való feltüntetése.

Alulírott kijelentem, hogy a plágium fogalmát megismertem, és tudomásul veszem,  
hogy plágium esetén szakdolgozatom visszautasításra kerül.

Miskolc, ..... év ..... hó ..... nap

.....

Hallgató

1.

szükséges (módosítás külön lapon)

A szakdolgozat feladat módosítása

nem szükséges

.....

dátum

.....

témavezető(k)

2. A feladat kidolgozását ellenőriztem:

témavezető (dátum, aláírás):

konzulens (dátum, aláírás):

.....

.....

.....

.....

.....

.....

3. A szakdolgozat beadható:

.....

dátum

.....

témavezető(k)

4. A szakdolgozat ..... szövegoldalt

..... program protokollt (listát, felhasználói leírást)

..... elektronikus adathordozót (részletezve)

.....

..... egyéb mellékletet (részletezve)

.....

tartalmaz.

.....

dátum

.....

témavezető(k)

5.

bocsátható

A szakdolgozat bírálatra

nem bocsátható

A bíráló neve: .....

.....

dátum

.....

szakfelelős

6. A szakdolgozat osztályzata

a témavezető javaslata: .....

a bíráló javaslata: .....

a szakdolgozat végleges eredménye: .....

Miskolc, .....

.....

a Záróvizsga Bizottság Elnöke

# Tartalomjegyzék

<b>1. Bevezetés</b>	<b>6</b>
<b>2. Téma elméleti kifejtése</b>	<b>8</b>
2.1. Mobil alkalmazások . . . . .	8
2.1.1. Bike calculator . . . . .	8
2.1.2. Strava . . . . .	8
2.1.3. Endomondo . . . . .	9
2.1.4. Best Bike Split . . . . .	9
2.1.5. Összefoglalás és célkitűzés . . . . .	9
2.2. Megválaszolandó kérdések . . . . .	11
2.2.1. Nehézségi osztályok meghatározása - eddigiek . . . . .	11
2.2.2. Nehézség meghatározása - új útvonalra . . . . .	12
2.2.3. Időtartam becslés . . . . .	12
2.2.4. Edzésterv . . . . .	12
2.3. Gépi tanulás . . . . .	13
2.3.1. Regresszió . . . . .	13
2.3.2. Osztályozás . . . . .	14
2.3.3. Klaszterezés . . . . .	14
2.4. A gépi tanulás alapköve - adatok . . . . .	15
2.4.1. Adatgyűjtés . . . . .	15
2.4.2. Adatstruktúra . . . . .	16
<b>3. Fejlesztői dokumentáció</b>	<b>18</b>
3.1. Adathalmaz előkészítés . . . . .	18
3.1.1. Adatstruktúra kialakítása . . . . .	18
3.1.2. Adattisztítás . . . . .	19
<b>4. Összefoglalás</b>	<b>20</b>
<b>Irodalomjegyzék</b>	<b>21</b>
<b>Adathordozó használati útmutató</b>	<b>22</b>

# 1. fejezet

## Bevezetés

Először 2018. nyarán kezdtem el az adatbányászattal és gépi tanulással foglalkozni és ahogy mind többet értettem meg belőle inkább érdekelt a terület, a kihívásai. A lelkesedés kitartott így ősszel a témaválasztás határidejének közeledtével az lebegett a szemem előtt hogy a területhez szorosan kapcsolódó szakdolgozatot tudjak elkezdeni.

Az irány megszabása után egy konkrét célt kellett találni, egy alkalmazási területet ahol a gépi tanulás jól hasznosítható, olyan problémák, feladatok vannak amihez érdemes lehet gépi tanulás algoritmusokat használni. Hosszas mérlegelés után a különböző sportok közül a kerékpározásra esett a választás. Ennek több oka is volt. Elsősorban a kerékpározás népszerű, elterjedt mozgásforma, így esélyesnek tűnt hogy viszonylag nagy mennyiségű statisztikai adatot lehet róla gyűjteni, változatos sportolói körből. Továbbá mikor időm engedi én is kerékpározok és ezeket az utakat rögzítem, így néhány kiinduló adat már elérhető volt ami alapján el tudtam kezdeni a tervezést. Az utak nyomon követésére a Strava mobil és webes alkalmazás szolgált, ami részletesen tárolja az egyes útvonalak statisztikai adatait. Az elérhető adatokat elemezve behatároltam több feladatot mint pl.:

- Egy adott útvonal jellemzőinek ismeretében a teljesítéshez szükséges várható időtartam meghatározása. Ez az időtartam két részre választható szét, megkülönböztetjük a mozgással töltött időt a teljes eltelt időtartamtól ami a mozgási időből és a járulékos idővesztésekből adódik mint a forgalmi lámpánál várakozás, pihenés, hosszabb utak esetén étkezés.
- Már megtett utak csoportosítása, nehézségi fokozatok meghatározása.
- A nehézségi szintek alapján edzéstervek készítése.

Azonban a sportteljesítmény becslése nehéz feladat, mivel egyénenként nagyon eltérő, valamint az út nyers adatain kívül az egyén pillanatnyi fizikai és szellemi állapota is nagyban befolyásolja a teljesítményét (álmoság, kimerültség, kedvtelenség stb.). Ezért el kell döntenie hogy egy általános modell megépítése a cél, ami esetleg életkor, nem, edzettségi szint szerint finomhangolható azonban kisebb pontossággal rendelkezik, vagy személyre szabott, pontos becsléseket kell elérni ahol viszont a gépi tanulás során nehézséget okozhat a kis méretű adathalmaz, hiszen kevesen rendelkeznek akár csak egy-két ezer rögzített úttal, míg a megfelelő pontossághoz az egyes algoritmusoknak ennél nagyságrendekkel több adatra van szüksége.

Az elsődleges cél általános modellek építése, ahol egy feladat megoldásához különböző gépi tanulási algoritmusokat lehet alkalmazni, ezeket egyenként javítani, ma-

---

ximalizálni a pontosságukat. Utána pedig az algoritmusokat össze kell hasonlítani a teljesítményük, felhasznált módszertanaik alapján.

## 2. fejezet

# Téma elméleti kifejtése

### 2.1. Mobil alkalmazások

Napjainkban rengeteg különböző kerékpáros alkalmazás érhető el amik más más területre fókuszálnak. A következő pontokban ezekre láthatunk néhány példát az egyszerű mérnöki alapú számítástól az útvonalak nyomon követésén keresztül a fejlett, összetett tervező rendszerekig. Látható hogy a felhasználók körében nagy népszerűségnek örvend a különböző alkalmazások közösségi funkciói ahol lehetőség van az ismerősök megtalálására, útvonalak megosztására.

#### 2.1.1. Bike calculator

A "Bike Calculator" [1] weboldal és mobil alkalmazás biciklis becsléseket kínál a felhasználói számára, teljesen mérnöki szempontból. A modell figyelembe veszi a sebesség, erő kifejtés, szélellenállás, gravitáció és egyéb tényezők összefüggését azonban egy általános eredményt ad, nincs lehetőség a felhasználó egyéni teljesítményének figyelembe vételére. Ezzel együtt is egy érdekes megoldás ami rengeteg paramétere révén jól skálázható.

#### 2.1.2. Strava

A Strava [2] az egyik legelterjedtebb alkalmazás a sporttevékenységüket nyomon követni kívánók körében, széles funkcionalitással és hatalmas felhasználói bázissal. Lehetőséget nyújt többek között biciklizés, futás, úszás, hegymászás és rengeteg egyéb sportág eredményeinek rögzítésére és korlátozott szintű tervezésére. A tevékenység során rögzített adatokból alapvető statisztikákat készít a felhasználók számára, mint a sebesség és a tengerszint feletti magasság változása valamint a sportoló pulzusának mozgása. Népszerűsíti a Strava-t hogy lehetőséget kínál az ismerősök megtalálására és az útvonalak, tevékenységek megosztására valamint kommentek és "kudo"-k egyfajta elismerés, tetszés nyilvánítás adására. Érdekes funkciói közé tartozik a szegmensek, rövid útvonal szakaszok kezelése ahol a felhasználóknak lehetőségük van rövid távon versenyezni, előrébb kerülni a szegmenst teljesítők listáján valamint számon tartani a saját rekordjaikat.

A Strava tervezője, jövővel kapcsolatban segítő funkciója egy útvonal tervezőben merül ki, ami segíti a felhasználókat a legnépszerűbb út megtalálásában A és B pont között, térképpel és részletes navigációs leírásokkal. Az útvonal kijelölésénél az alapvető



adatokon kívül egy becslést is kapunk a várható időre nézve. A várható idő számolása a bejelentkezett felhasználó legutóbbi 4 heti teljesítménye alapján történik. Előrelépés hogy egyáltalán figyelembe veszi a felhasználó eddigi útvonalait azonban ennek a módszernek is vannak hiányosságai így csak egy közepes kiindulási alapot ad. Alapvető gond ezzel a becsléssel hogy az útvonalak rögzítésénél a felhasználó hiába állítja be az út típusát (verseny, edzés, ingázás stb.) a Strava számítás során ezeket nem veszi figyelembe, így könnyen előfordulhat hogy a felhasználó kiugróan pontatlan, akár használhatatlan eredményt kap. Például az előző hetekben a saját határait feszegetve versenyre készült, most viszont szeretne egy nyugodt családi biciklizésre elindulni - ekkor nyilván nem fog rekordokat dönteni és a Strava által számolt várható idő semmilyen tényleges információval nem fog szolgálni neki.

### 2.1.3. Endomondo

Az Endomondo [3] webes és mobil alkalmazás sporttevékenység nyomon követését hivatott megkönnyíteni a felhasználói számára. Rengeteg különböző sportágat megkülönböztet, ezekről külön statisztikákat készít valamint egyéni célokat lehet ezekre vonatkozóan beállítani. Az egyes sportágakat GPS alapú jelkövetés segítségével lehet nyomon követni, illetve kézzel is létrehozhatóak új tevékenységek, továbbá lehetőség van különböző alkalmazásokból vagy fájlból importálni korábban teljesített útvonalakat. Az alkalmazás megvalósít közösségi funkciókat is, a felhasználóknak lehetőségük van megkeresni az ismerőseiket akikkel később megoszthatják legújabb sporttevékenységüket valamint versenghetnek egymással a közösen felállított célok eléréséért. Az Endomondo érdekes adaléka hogy mozgás közben audió visszajelzést biztosít a felhasználók számára az addigi teljesítményükről, így a felhasználók sportolás közben is pontos információval rendelkezhetnek anélkül hogy az alkalmazást futtató eszközüket kellene figyelniük.

### 2.1.4. Best Bike Split

A Best Bike Split [4] egy rendkívül összetett rendszer amely elsősorban versenyzők számára kínál széles funkcionalitást. Alapköve a versenytáv megtételének optimalizálása melynek keretei között személyre szabott, szakaszokra bontott tervet állít fel a felhasználó számára annak érdekében hogy a tervezett távot a lehető legjobb idő alatt teljesítse. Ez a terv alapul veszi a felhasználó biciklijének és felszerelésének a tulajdonságait (kerekek típusa, bicikli súlya, sisak kialakítása stb.) a preferált versenyzési módot (elhelyezkedés a verseny során, sík terep illetve emelkedő esetén) valamint már teljesített edzést, versenyeket és természetesen a teljesítendő verseny paramétereit. Prémium előfizetés használata során a kalkulációt befolyásolja a várható időjárás is. A rendszer használható tetszőleges útvonalakra azonban az átlag kerékpáros számára szükségtelesen bonyolult, összetett így elsősorban versenyzők számára hasznos.

### 2.1.5. Összefoglalás és célkitűzés

Több tucatnyi hasonló alkalmazás létezik, webes felülettel illetve Android és iOS lefedéssel hiszen a kerékpározás fellendülőben van napjainkban, így egyre több alkalmazás készül a felhasználói igények kielégítésére mind versenyzők mind pedig hétköznapi felhasználók számára akiknek a biciklizés inkább közlekedési mint sport eszköz. Ezen

alkalmazások felsorolása vagy részletezése nem célja ennek a szakdolgozatnak, azonban néhány példa alapján is könnyen meghatározhatóak általános csoportok amelyekbe a létező megoldások besorolhatóak.

### Útvonal rögzítés

A legalapvetőbb funkció amire különböző megoldások léteznek. A legegyszerűbb eszközök esetén például csak az aktuális sebesség és a távolság kerül megállapításra egy, a kerékre rögzíthető érzékelőnek a segítségével. Az érzékelő használatával számolni lehet hogy a kerék hányszor fordul körbe adott időablakon belül, ennek és a kerék méretének felhasználásával könnyen megállapítható a pillanatnyi sebesség illetve az útvonal kezdete óta megtett távolság. Azonban a növekvő igényekhez igazodva ma már lényegesen elterjedtebbek a GPS alapú mobil alkalmazások amik a pillanatnyi sebességen és távolságon kívül pontosan tudják rögzíteni az út egyéb paramétereit is mint pl.: kezdő és végpont koordinátái, mozgással töltött idő, megtett emelkedő, legalacsonyabb és legmagasabb pont, átlag sebesség, maximum sebesség. Ezek a jellemzők lényegesen pontosabb leírást adnak egy-egy útvonalról, letárolásukkal könnyen nyomon követhető a személyes fejlődés, kitűzhetőek heti, havi, éves célok. Az útvonal paramétereinek nyomon követése, struktúrája és tárolása megvalósításonként változó, azonban a cél mindig ugyanaz: minél több jellemző meghatározása a lehető legrészletesebb módon.

### Közösségi funkciók

A kerékpáros útvonalakat nyomon követő mobil alkalmazások elterjedésével egy új tendencia jelent meg; egyre többször figyelhető meg valamilyen jellegű közösségi funkció megvalósítása az alkalmazáson belül, amik nagy népszerűségnek örvendenek az alkalmazást használók körében. Ezen funkciók nagyban különbözhetnek, általánosságban lehetőség van az alkalmazást használó ismerősök megtalálására, megtett útvonalak megosztására illetve valamilyen formában tetszésnyilvánítás / elismerés adására. Ezeken kívül általában megjelenik valamilyen ösztönzés, lelkesítés ami segíthet a teljesítmény javításában, a kitűzött célok elérésében. A Strava esetében ez a szegmensek formájában jelenik meg, amik lényegében rövid, maximum néhány száz méteres útvonalak, általában valamilyen szempontból nehezebb szakaszok (pl.: emelkedő). Ezen szegmensek teljesítésével a felhasználó automatikusan felkerül több rangsorra is (globális, erre az évre vonatkozó, nem és korcsoport alapján lebontott) ahol lehetősége van ismerősökkel és ismeretlenekkel vetélkedni ezzel felébresztve a versenyszellemet.

### Tervezés, becslés

A kerékpáros alkalmazások fejlődésének következő lépcsőfoka az útvonalak tervezése és a teljesítmény becslése. Útvonalak tervezésére kevésbé kerékpár specifikus alkalmazások is lehetőség adnak mint például a Google Térkép de sok, kifejezetten kerékpárosoknak készült alkalmazásban is léteznek megvalósítások. A tervezés különböző szempontok alapján történhet: egyik fontos szempont a távolság (ez a legfontosabb amennyiben a cél az út minél gyorsabb megtétele A és B pont között), egy másik pedig a népszerűség szóval foglalható össze. Az utóbbi főleg akkor fontos ha a kerékpározásra nem közlekedés szerűen tekintünk hanem inkább mint sport, kikapcsolódás hiszen ilyenkor előrébb kerül a jogos felhasználói igény az olyan útvonalakra ahol szép környezetben,

lehetőség szerint autós és gyalogos forgalomtól zavartalanul lehet biciklizni. Az útvonalak tervezéséhez szorosan kapcsolódik a teljesítmény, a várható idő meghatározása ami koránt sem egyszerű feladat. Ennek legalapvetőbb megoldása a Bike calculator által is megvalósított mérnöki, fizikai modell, ahol a teljesítmény tisztán a felszerelés és a megtenni kívánt út paraméterei alapján kerül kiszámításra

### Tervezés, becslés - egyéni teljesítmény alapján

Az útvonal tervezés és a várható teljesítmény becslésének továbbfejlesztésének, pontosításának alapjául az eddig megtett útvonalak és az elért eredmények szolgálhatnak, ez azonban egy olyan terület ahol még nem sok megoldás született. Ide sorolható a Strava által fejlesztett útvonal tervező, ahol a várható idő becsléséhez alapul veszi a bejelentkezett felhasználó elmúlt 4 heti teljesítményét, és míg ez egy jó kezdeményezés a valóságban nem mindig szolgál pontosabb információval mint egy egyszerű, fizikai összefüggéseken alapuló általános modell. Megnehezíti a pontos becslések készítését hogy egy adott felhasználó sokszor kevés adattal rendelkezik amelyek nem elég változatosak egy új útvonal esetén várható teljesítmény megbecsléséhez.

Ezen 4 fő csoport közül az első három által meghatározott problémákra, igényekre mára már rengeteg, változatos módon megvalósított megoldás létezik, így ezeknek a tárgyalásától a továbbiakban eltekintünk. Azonban az utolsó probléma, a felhasználó egyéni teljesítményének figyelembe vétele egy nyitott terület, ahova nem sokan merészkedtek még. **TODO** ez a szakdoga célja

minden sport esetén az adott ember, sportoló az alapvető tényező, az ő képességei fogják meghatározni legnagyobb részben az eredményeit. Természetesen a felszerelés és a terepviszonyok is nagy hatással vannak a végeredményre de az emberi tényező nem hagyható ki a számítások során amennyiben növelni akarjuk a pontosságot.

## 2.2. Megválaszolendő kérdések

A sportteljesítmény becslése nehéz feladat, mivel egyénenként nagyon eltérő lehet, valamint az út nyers adatain kívül (távolság, emelkedő, időjárás) a egyén pillanatnyi fizikai és szellemi állapota is nagyban befolyásolja (álmosság, kimerültség, kedvtelenség stb.). Ezért el kell döntenie hogy egy általános modell megépítése a cél, ami esetleg életkor, nem, edzettségi szint szerint finomhangolható azonban kisebb pontossággal rendelkezik, vagy személyre szabott, pontos becsléseket kell elérni ahol viszont a gépi tanulás során nehézséget okozhat a kis méretű adathalmaz, hiszen kevesen rendelkeznek akár csak néhány száz rögzített úttal, míg a megfelelő pontossághoz ehhez nagyságrendekkel több adatra van szüksége az egyes algoritmusoknak

### 2.2.1. Nehézségi osztályok meghatározása - eddigiek

A Strava-ról gyűjtött utak sok adatot tartalmaznak azonban nincsenek kategorizálva, így cél az eddig megtett utakat nehézség alapján több osztályba sorolni, azonban ennek megvalósítása kérdéses. Mivel a helyes label-ek hiányoznak klaszterezési módszereket kell használni. Várhatóan az egyén adatai alapján és a összes összegyűjtött adat alapján készült osztályok különbözőek lesznek, hiszen más útvonalakat tettek meg, eltérő eredményekkel.

Egy általános modell elkészítése az elsődleges feladat, mivel az adatok sokszínűsége és mennyisége miatt jobb kiindulópontot jelent. Amennyiben egy adott felhasználónak megfelelően sok és sokféle útadata van akkor érdemes lehet az egyéni adatokra is kipróbálni az osztályozást.

### 2.2.2. Nehézség meghatározása - új útvonalra

A már megtett útvonalak osztályozása után a következő lépés az új, megtenni kívánt útvonalak nehézségének meghatározása. Ez megvalósítható egyrészt a már betanított klaszterezési algoritmusok felhasználásával, másrészt a klaszterezés által meghatározott osztályok label-jeit felhasználva valamilyen osztályozással.

### 2.2.3. Időtartam becslés

Az egyik legalapvetőbb feladat a várható idő meghatározása, hiszen ez az egyik legmeghatározóbb pontja a kerékpározásnak - mennyi időbe fog telni mire megtesz valaki adott kilométert? Ennek a kérdésnek a megválaszolására alapvetően a különböző regressziók alkalmasak, amik az út paramétereire alapján folytonos értéket adnak vissza. Azonban ez a pont kicsit bonyolultabb mint amilyennek tűnik, mivel alapvetően kétféle időtartamot különböztetünk meg: a mozgási és a tényleges időtartamot. Az első jelöli az összesen mozgással töltött időt, míg a második magába foglalja a különböző megállókat mint a pihenés, várakozás forgalmi lámpáknál, hosszabb utak esetén étkezéssel töltött időket stb. Ennek megoldására több lehetőség is kínálkozik.

A legegyszerűbb megoldás olyan regresszió használata amely egy inputhoz több eredmény, output-ot is tud társítani. Erre kínál lehetőséget a 2.3.1. pontban részletezett MLPRegressor. Ennek használatával rögtön mindkét idő előáll, nem szükséges további lépéseket tenni.

A második lehetőség több részből áll. Első lépésben tetszőleges regresszióval csak a mozgási időre készül becslés, mivel ez szorosabban összefügg az olyan alapvető adatokkal mint a távolság. Második lépésként pedig egy újabb regressziós modellt kell használni, ami a mozgási idő és egyéb paraméterek segítségével képes a teljes várható időre becslést készíteni.

### 2.2.4. Edzésterv

A felhasználók számára javasolt edzésterveknek több lehetséges megvalósítási módja is van. A legteljesebb kivitelezés a térképen jelölt konkrét útvonal, célul kitűzött teljesítési és részhidővel, nehézségi szinttel. Azonban a szakdolgozat során egy valamivel egyszerűbb feladat elérése a cél, amit két alrészre lehet bontani.

- Az első típus esetén a felhasználó ad meg egy útvonalat, aminek adataiból kell a teljes várható időre, részhidőkre, sebességre célkitűzéseket adni, alapul véve hogy a felhasználó milyen nehézségi szintet tűz ki. A részhidőket érdemes külön vizsgálni, hiszen egy meredeken emelkedő szakasz teljesítése sokkal nagyobb kihívás a kerékpározók számára mint egy viszonylag sík terepen való haladás.
- A második típus esetén az út adatait is az algoritmusnak kell kiszámolnia, ideértve az út hosszát és a szintemelkedéseket. Ezen adatok valamint a kitűzött teljesítési idő egyensúlyba hozásával kell egy-egy nehézségi szinthez racionális nehézségű útvonal tervet létrehozni.

## 2.3. Gépi tanulás

Rengeteg gépi tanulási módszer létezik amik jól használhatóak különböző modellek, becslések, osztályok készítésére. Ezek a módszerek ma már legtöbbször egy programozási nyelv adott könyvtárának, csomagjának segítségével egyszerűen használhatóak. Ennek megfelelően a Python nyelv scikit-learn [6] csomagja adta a gépi tanulási módszerek alapját a szakdolgozat elkészítése során.

Az egyes algoritmusok használatán túl fontos a megfelelő adatforrás megtalálása, valamint az adatkinyerési és feldolgozási folyamat leegyszerűsítése, átláthatóvá tétele.

Ebben a fejezetben találhatóak meg a szakdolgozat során használt gépi tanulási módszertanok részletezése, valamint az adatgyűjtéshez létrehozott rendszer eszközeinek kifejtése.

### 2.3.1. Regresszió

A regresszió egy felügyelt gépi tanulási módszer amely folytonos kimenetet generál. A tanuláshoz szükség van a bemeneti adathalmazra valamint a helyes kimeneti értékekre, amik elemzésével az algoritmus képes új bemeneti adatpárookra kimenetet generálni.

#### Ridge Regression

A Ridge regresszió sajátossága hogy a túlilleszkedés kiküszöbölésének érdekében módosítja a költség funkció számítását. Ehhez egy "büntető" tényezőt, egy plusz súlyt használ, ezzel ellensúlyozva a hibákat.

$$\min_w \|Xw - y\|_2^2 + \alpha \|w\|_2^2$$

Itt  $\alpha \geq 0$  egy komplexitási paraméter ami a csökkenés mértékét befolyásolja: minél nagyobb az  $\alpha$  értéke annál gyorsabb a csökkenés.

Az algoritmus bonyolultsága egy  $n \times p$  méretű  $X$  mátrix esetén  $O(np^2)$  amennyiben  $n \geq p$ .

#### Lasso

A *Least Absolute Shrinkage and Selection Operator* angol kifejezés rövidítése.

$$\min_w \frac{1}{2n_{samples}} \|Xw - y\|_2^2 + \alpha \rho \|w\|_1$$

#### Elastic Net

**TODO** szöveges kifejtés

$$\min_w \frac{1}{2n_{samples}} \|Xw - y\|_2^2 + \alpha \rho \|w\|_1 + \frac{\alpha(1 - \rho)}{2} \|w\|_2^2$$

#### MLPRegressor

Multi-layer Perceptron Regression, azaz mesterséges neuronokból felépülő többretegű struktúrán, más néven egy neurális hálón alapuló regressziót valósít meg. Ez egy felügyelt tanulási algoritmus, tehát inputként a feature-ok egy mátrixát ( $X$ ) valamint a

hozzájuk tartozó helyes label-eket ( $y$ ) várja. Ezen értékek alapján backpropagation-t használva a módszer folytonos értékek készletével tér vissza. Az MLPRegressor továbbá több kimenetű regressziót is támogat, így egy mintához több output érték is rendelkezhető.

### 2.3.2. Osztályozás

Az osztályozás egy felügyelt gépi tanulási algoritmus, melynek lényege hogy adott input adatokból és helyes label-jeikből (amik megadják hogy az egyes adatok melyik osztályhoz tartoznak) tanulva az új adatokhoz megfelelő osztályt, csoportot rendeljen.

#### Nearest Neighbors Classification

A szomszéd-alapú osztályozás példányalapú tanulás: nem egy általános belső modellt próbál konstruálni, csak egyszerűen a training data példányait tárolja.

A scikit-learn csomag két különböző legközelebbi szomszéd módszert is kínál: a KNeighborsClassifier-t és a RadiusNeighborsClassifier-t

- **KNeighborsClassifier:** ez a módszer az egyes pontok  $k$  legközelebbi szomszédja alapján osztályoz, ahol  $k$  egy a fejlesztő által megadott egész szám. A  $k$  optimális értéke nagyban függ az adattól: általában egy nagyobb  $k$  elnyomja az anomáliákat, zajok hatását, azonban az osztály határok kevésbé lesznek egyértelműek. Ez a módszer elterjedtebb.
- **RadiusNeighborsClassifier:** az előző módszerrel ellentétben nem a  $k$  legközelebbi szomszédot veszi figyelembe, hanem a fejlesztő által megadott  $r$  sugárba eső pontokat csoportosítja egy osztályba. Amennyiben az adathalmaz nem egyenletesen mintavételezett a RadiusNeighborsClassifier jobb választás lehet.

#### Nearest Centroid Classifier

Egyszerű algoritmus, ahol minden osztályt a tagjai középpontja reprezentál

#### MLPClassifier

Multi-layer Perceptron Classification, azaz mesterséges neuronokból felépülő többretegű struktúrán, más néven egy neurális hálón alapuló osztályozást valósít meg. Bemenetként a feature-ok  $X$  mátrixát és a label-ek  $y$  vektorát várja. A modell backpropagation-t használ a tanulás során, pontosabban gradiens leereszkedést, ami backpropagation használatával kerül kiszámításra. Az osztályozáshoz a kereszt entrópia funkciót minimalizálja, az egyes  $x$  mintákhoz a  $P(y|x)$  valószínűséget becsülve. Támogatja a multi-class osztályozást a Softmax függvény output funkcióként való alkalmazásával, valamint a multi-label osztályozást, melynek segítségével egy minta több osztályba is sorolható.

### 2.3.3. Klaszterezés

A klaszterezés alapfeladata hasonló az osztályozáséhoz, az adatokat különböző csoportokba kell rendezni, azonban ez felügyelet nélküli módszer, a helyes label-ek nélkül próbál adott számú osztályt kialakítani.

### KMeans

**TODO** kifejtés

### MeanShift

**TODO** kifejtés

## 2.4. A gépi tanulás alapköve - adatok

A gépi tanulás módszerei önmagukban mit sem érnek, megfelelő méretű és minőségű adathalmazra van szükség amiből az algoritmusok tanulni tudnak. Azonban adatgyűjtés mindig érzékeny terület, sok időt, energiát felemészt és a kívánt eredmény ezek befektetésével sem mindig garantált. A szakdolgozathoz elsősorban biciklivel megtett utak statisztikai adataira van szükség, amihez a Strava mobil és webes alkalmazás ad kiindulási alapot. A fejlesztői munkát segíti a Strava saját API-ja, aminek segítségével egyszerűen megvalósítható a felhasználók autentikációja valamint a szükséges adatok elérése.

### 2.4.1. Adatgyűjtés

A Strava által gyűjtött és tárolt adatok elérésének alapját egy, a szakdolgozat keretein belül készített weboldal adja. Ezen oldal felkeresésével a felhasználók tájékozódhatnak a szakdolgozatról, az adatgyűjtés fontosságáról valamint alapvető funkciókat használhatnak. Az oldalon lehetőség van új felhasználó létrehozására, már létező fiókokba való bejelentkezésre. Bejelentkezés után a felhasználó hozzákapcsolhatja a Strava fiókját a weblapon létrehozott fiókjához, ezzel lehetővé téve az útvonalai statisztikai adatainak olvasását. A weboldalon továbbá lehetőséget fog kínálni az elkészült, betanított program letöltésére, amivel a felhasználók is megtekinthetik az egyes algoritmusok hatékonyságát. A weboldal az alábbi eszközök felhasználásával valósult meg:

### JQuery

A JQuery egy népszerű JavaScript könyvtár amely lehetőséget ad a HTML elemek könnyebb manipulálására, az események egyszerűbb kezelésére, animációk megvalósítására.

### Firebase

A Firebase egy Google által fejlesztett sokoldalú rendszer, amely egyszerű eszközöket kínál egy weboldal logikai részének felépítéséhez. Nem alkalmas robosztus rendszerek létrehozásához, azonban kisebb projektek menedzseléséhez megfelelő alapot kínál. A projekteket létrehozás után online lehet kezelni és a Firebase funkcióit weboldalakban és mobil alkalmazásokban is (mind Android és iOS) egyszerűen lehet használni. A rendszer alapját egy Node.js szerver és egy JSON alapú adatbázis adja valamint elérhető egy általános célú tárhely tetszőleges file-ok tárolásához.

### Firestore - Autentikáció

A szakdolgozathoz készített weboldal szempontjából a Firestore egyik legfontosabb eleme az autentikációs rendszer. A Firestore projektben való engedélyezés után néhány sor JavaScript kóddal gyorsan és biztonságosan megvalósítható az új felhasználók regisztrálása, bejelentkeztetése. Az email alapú regisztráláson túl lehetőség van külső rendszerekkel való autentikációra is mint a Google fiók, Facebook, Github. Továbbá elérhető a már regisztrált felhasználók emailen keresztül történő értesítése is.

### Firestore - Adatbázis

A Firestore alapvetően ingyenes, emiatt az adatbázisnak vannak méret és sebességbeli korlátjai, azonban megfelelő kiindulási alapot ad. A JSON alapú struktúra megkönnyíti az adatok kliensen való kezelését valamint egy ilyen méretű alkalmazás esetén egyszerűen karban tartható. Fontos kiemelni az adatbázis hozzáférhetőségének konfigurálási lehetőségét. Lehetőség van az adatbázis egészére vagy a részfákra olvasási és írási beállításokat megadni, amik korlátozhatóak a bejelentkezett felhasználó azonosítója alapján is. Ennek felhasználásával könnyen megvalósítható hogy a rendszer összes felhasználója csak a saját adataihoz férjen hozzá ami fontos adatvédelmi szempont.

#### 2.4.2. Adatstruktúra

A szakdolgozathoz készített weboldal fő célja hogy lekérje és tárolja azon felhasználók útvonalait akik regisztrálás után kapcsolódtak a Strava fiókjukhoz. A Strava API-ja az utak lekérésekor SummaryActivity objektumok tömbjével tér vissza, ahol minden objektum 39 változóval jellemez egy utat, aktivitást. Ez a 39 jellemző magába foglalja többek között a felhasználó azonosítóját, a megtett távot, a teljesített kihívások számát, az ismerősök megjegyzéseinek a számát és rengeteg egyéb információt. A gépi tanulás során ezeknek a jellemzőknek csak egy részét érdemes használni, amik szorosabban kapcsolódnak az út valós leírásához. Ezek a jellemzők a következők:

- **average\_speed:** átlagsebesség (méter/másodperc)
- **average\_watts:** átlagos erő kifejtés az útvonal során. Amennyiben nem biztosított az értéke -1
- **commute:** egy tevékenység rögzítésekor lehetőség van azt 'ingázásként' felcímkézni, ezzel megjelölve hogy nem egy verseny vagy túra adat, hanem mindennapi, például munkába vagy boltba járási útvonal. Feltételezhető hogy ilyenkor az illető nem a minél jobb eredmények elérésére törekszik, inkább rutin szerűen, átlagos tempóban halad (Logikai érték, -1 amennyiben nem biztosított)
- **device\_watts:** azt jelöli hogy az erő kifejtési adatokat tényleges mérő készülék szolgáltatja vagy csupán becsült érték (logikai, Hamis amennyiben becsült érték)
- **distance:** a megtett távolság (méter)
- **elapsed\_time:** összesen eltelt idő (másodperc)
- **elev\_high:** az út legmagasabb pontja (méter)
- **elev\_low:** az út legalacsonyabb pontja (méter)



- **end\_latlng:** az útvonal végpontjának hosszúsági és szélességi helyzete
- **flagged:** az útvonal meg van e jelölve (logikai)
- **has\_heartrate:** elérhető e pulzus adat (logikai)
- **heartrate\_opt\_out:**
- **kilojoules:** összes kifejtett erő mennyisége (kilojoule)
- **max\_speed:** maximum sebesség az út során (méter/másodperc)
- **max\_watts:** maximális erő kifejtés az út során
- **moving\_time:** mozgással töltött idő (másodperc)
- **pr\_count:**
- **resource\_state**
- **start\_date\_local:** az út kezdési ideje ( dátum)
- **start\_latlng:** az útvonal kezdőpontjának hosszúsági és szélességi helyzete
- **total\_elevation\_gain:** összesen megtett emelkedő (méter)
- **trainer:** az útvonal gépen lett e megtéve
- **weighted\_average\_watts:**
- **workout\_type:** tevékenység típusa, lehet edzés vagy verseny (egész szám, ha nincs megadva akkor -1 ?? )

## 3. fejezet

# Fejlesztői dokumentáció

Az Elméleti kifejtés során meghatározott célokat különböző gépi tanulási módszerekkel lehet megoldani amik eltérő eredménnyel, pontossággal fognak működni. A nehézségi osztályok meghatározására elsősorban a 2.3.3. pontban részletezett klaszterezési módszerek használhatóak

### 3.1. Adathalmaz előkészítés

A fejlesztés legelső lépése az összegyűjtött adatok megfelelő struktúrára való átalakítása és megtisztítása, ezzel megkönnyítve a későbbi feldolgozást valamint növelve a gépi tanulási algoritmusok pontosságát

#### 3.1.1. Adatstruktúra kialakítása

Az adatgyűjtés elsődleges eszközeként a szakdolgozat keretein belül készített weboldal szolgált, amely minden információt egy JSON alapú adatbázisban tárolt le. Az adathalmaz előkészítésének a legelső lépése a JSON struktúráról való áttérés egy, a Python programozási nyelv által kezelt, könnyen használható adatstruktúrára. Erre a szakdolgozat során a Pandas [17] nevű Python csomag által megvalósított DataFrame nevű struktúra fog szolgálni, amely segítségével az adatokat táblázathoz hasonló formában lehet tárolni. Egy DataFrame oszlopai egyedi, a fejlesztő által definiált oszlopnevekkel érhetőek el, soraira index használatával lehet hivatkozni. Az oszlopok egyenként különböző típusúak lehetnek és tartalmazhatnak NULL értékeket. A DataFrame egyik legnagyobb előnye az oszlopok nevesítése, amivel könnyen nyomon követhető hogy az aktuális értékek melyik feature-nek felelnek meg, mit reprezentálnak a valóságban

**3.1. Program részlet.** A parancs segítségével a JSON struktúra (jsonData) könnyen átalakítható DataFrame-é (dataset). A dataset oszlopai megfelelnek a 2.4.2. pontban részletezett adatoknak.

```
import pandas

jsonData = pandas.read_json("2019_02_08_10h.json", type='series')
cleanData = []
for u in userData:
    if userData[u]['connectedToStrava'] == True:
        for i in range(0, len(userData[user]['activities'])):
```

```
userData[u]['activities'][i]['ageGroup'] = userData[u]['ageGroup']
userData[u]['activities'][i]['sex'] = userData[u]['sex']
userData[u]['activities'][i].pop('external_id', None)
userData[u]['activities'][i].pop('map_id', None)
userData[u]['activities'][i].pop('map_resource', None)
userData[u]['activities'][i].pop('map_summary', None)
cleanData.append(userData[u]['activities'][i])
else:
    print('User does not have any activities')
dataset = pandas.DataFrame(cleanData)
```

A fenti parancs segítségével a JSON struktúra (jsonData) könnyen átalakítható DataFrame-é (dataset). A dataset oszlopai megfelelnek a 2.4.2. pontban részletezett adatoknak.

### 3.1.2. Adattisztítás

A megfelelő adatstruktúra kialakítása után a következő lépés az adatok megtisztítása. Ez a lépés azért szükséges mert a legfigyelmesebben gyűjtött adathalmaz is tartalmazhat rossz értékeket illetve előfordulhat hogy több helyen hiányzik a valódi érték. Ezeknek a hibáknak a megtalálása és kijavítása több fázisból áll.

#### NULL értékek

Egy általános adatbázis esetén gyakran előfordul hogy egyes oszlopok NULL értékeket tartalmaznak - ebben az esetben például NULL jelöli ha egy útvonal nincs elérhető adat egy adott jellemzőről. Azonban a gépi tanulási algoritmusok számokat képesek feldolgozni így elengedhetetlen a NULL értékek kiküszöbölése valamilyen formában.

A NULL értékek kiküszöbölésére két elterjedt megoldás létezik:

- **NULL értékek törlése:** az egyszerűbb megoldás a NULL értékeket tartalmazó sorok törlése, azonban ennek a módszernek nagy hátránya hogy sok NULL-t tartalmazó adathalmaz esetén jelentős mértékben megcsappan az adathalmaz mérete
- **NULL értékek feltöltése:** összetettebb, azonban sok esetben célravezetőbb megoldást jelent a NULL értékek helyettesítése valamilyen számított értékkel. Az új értékek számítása különböző módokon történhet, ez nagyban függ az adott jellemző jellegétől

Amennyiben egy oszlop nagy mértékben tartalmaz NULL értékeket érdemes megfontolni az elvetését vagy külön esetként kezelni mikor van tényleges érték

#### Kiugró értékek

Gyakori jelenség hogy egy jellemző tartalmaz néhány magasan kiugró értéket, amelyek sokszor érvényesek azonban fakadhatnak mérési / rögzítési hibából is. Akár érvényes értékek, akár valamilyen hibából fakadnak érdemes kiküszöbölni őket mivel könnyen eltorzíthatják az eredményeket.

## 4. fejezet

# Összefoglalás

Ebben a fejezetben kell összefoglalni a szakdolgozat eredményeit, sajátosságait és a témában való elhelyezkedését. A fejezet címe nem módosítható! Lehet benne több alfejezet is, de nem ajánlott. Min 1 max 4 oldal terjedelemben

# Irodalomjegyzék

- [1] <http://bikecalculator.com/>
- [2] Strava sporttevékenység követő alkalmazás. Strava.com
- [3] Endomondo
- [4] <https://www.bestbikesplit.com/>
- [5] Firebase
- [6] @articlescikit-learn, title=Scikit-learn: Machine Learning in Python, author=Pedregosa, F. and Varoquaux, G. and Gramfort, A. and Michel, V. and Thirion, B. and Grisel, O. and Blondel, M. and Prettenhofer, P. and Weiss, R. and Dubourg, V. and Vanderplas, J. and Passos, A. and Cournapeau, D. and Brucher, M. and Perrot, M. and Duchesnay, E., journal=Journal of Machine Learning Research, volume=12, pages=2825–2830, year=2011
- [7] Scikit-learn Ridge regression
- [8] Scikit-learn Lasso
- [9] Scikit-learn Elastic net
- [10] Scikit-learn MLPRegressor
- [11] Scikit-learn KNeighbors Classifier
- [12] Scikit-learn Radius Neighbors Classifier
- [13] Scikit-learn Nearest Centroid Classifier
- [14] Scikit-learn MLPClassifier
- [15] Scikit-learn KMeans
- [16] Scikit-learn MeanShift
- [17] Pandas: Python csomag

# Adathordozó használati útmutató

Ebben a fejezetben kell megadni, hogy a szakdolgozathoz mellékelte adathordozót hogyan lehet elérni, milyen struktúrát követ. Minimum 1 max 4 oldal. Lehet több alszakkasz. Fejezet címe nem módosítható