



TDK DOLGOZAT

Elsőéves VBK hallgatók teljesítményének
vizsgálata a Covid előtti és utáni időszakból

Köller Donát Ákos & Vlaszov Artúr

BME Matematikus MSc

Adattudomány szakirány

Témavezető: Szilágyi Brigitta

Geometria Tanszék



BME Matematika Intézet

Budapest

2022

Tartalomjegyzék

1. Bevezetés	1
2. A kognitív tesztről	1
3. Adatprepació	1
3.1. Az adatok beszerzése, adattáblák jellemzése	1
3.2. Adattisztítás	1
4. Felderítő adatelemzés a két évben	4
4.1. Általános ábrák	4
4.2. Folyamatábrák (Sankey-diagramok)	4
4.3. Klaszterezés	6
5. Prediktív analitika	6
5.1. Modellek és metodológia	6
5.2. Osztályozó algoritmusok	7
5.2.1. Lineáris regresszió	7
5.2.2. Naive Bayes	7
5.2.3. Gradient Tree Boosting	7
5.2.4. Logisztikus regresszió	7
5.2.5. SVM	8
5.3. Implementálás és optimalizálás	8
6. Modellek kiértékelése	9

1. Bevezetés

2. A kognitív tesztről

3. Adatprepació

3.1. Az adatok beszerzése, adattáblák jellemzése

A kutatáshoz használt adatokat Szilágyi Brigitta tanárnő bocsátotta rendelkezésünkre.

3.2. Adattisztítás

A felderítő és modellezési fázisban használt adattáblát a rendelkezésre álló táblák megfelelő mértékű tisztításából, majd ezt követő összeillesztéséből nyertük. Ezen műveleteket

Python-ban valamint R-ben végeztük el.

A kezdeti adattáblák közül a matematika és kognitív eredményeket tartalmazó adatsor tisztítása során először egységesítettük a szakmegnevezést ("Vegyésszmérnöki", "Biomérnöki", "Környezetmérnöki"), ahol pedig nem volt egyértelmű a szak, azt "UNKNOWN"-ra állítottuk. Ezt követően a teszten elért matematika és kognitív eredményt százalékos formátumról (tizedesjegy?) formátumra változtattuk, valamint létrehoztunk 3 új oszlopot, amely a matematikateszt 3 blokkjában helyesen megválaszolt kérdések számát mutatja. Végül a vizsgálataink szempontjából irreleváns oszlopokat eltávolítottuk, valamint azon oszlopokat is kiszűrtük, amelyek értéke a többitől egyértelműen kikövetkeztethető volt.

A kumulált átlagokat, felvételi pontszámokat illetve 0. ZH eredményeket tartalmazó tábláknál egy-egy sor eltávolításán kívül nem volt szükség adattisztításra, a matematika jegyeket tartalmazó tábla esetén pedig minden hallgató esetén meg kellett határozni azt az érdemjegyet, amellyel a tárgyat elvégezte.

A táblák összetűzése Neptun-kód alapján történt belső illesztés alkalmazásával. Azonban voltak olyan hallgatók, akikről nem minden táblában volt adat (vagy azért, mert nem írtak kognitív tesztet, vagy azért, mert az első félév vége előtt elhagyták a szakot), így az összeillesztés során 10-20 fős sorvesztéssel kellett számolnunk mindkét évben. Az így kapott adattáblák, amelyek a 2019-es és 2021-es évet reprezentálják, rendre 200 és 220 adatrekordot tartalmaztak, valamint mindkét tábla 14 attribútumaml rendelkezett, amelyek rendre a ['Neptun'], ['Hallgató neve'], ['Szak'], ['Matematika 1. blokk'], ['Matematika 2. blokk'], ['Matematika 3. blokk'], ['Matematika eredmény'], ['Kognitív eredmény'], ['0. ZH pontszám'], ['Tanulmányi pont'] ['Érettségi pont'], ['Többletpont'], ['Matematika A1a jegy'], ['Kumulált átlag'].

A kognitív eredményeket tartalmazó adattábla részletesen tartalmazott információt egyrészt minden hallgatóról (hova valósi, emelt érettségit tett-e matematikából, reál tagozatos volt-e, milyen szakra és tankörbe jár), másrészt a hallgató teszteredményéről is (mennyi idő alatt töltötte ki a tesztet, mely kérdéseket válaszolta meg jól, milyen lett a nyelvi és matekos teljesítménye), illetve tartalmazott néhány, a teszthez kapcsolódó egyéb információt is (például a teszthez használt edubase jelszó, felhasználónév). Természetesen nekünk ennyi adat nem kell, úgyhogy ebből az adattáblából jó pár irreleváns oszlopot ki kellett szűrni. Amelyeket meghagytunk, azok az alábbiak: a hallgató neve (ez már elég volt, hogy csak ez alapján fűzzük össze a táblákat) és Neptun kódja; emelt érettségit tett-e matematikából; reál/matematika tagozatos volt-e; szak és tankör; az elért pont és százalékos teljesítmény a nyelvi és matekos részben, valamint összességében. Problémát jelentett még, hogy a 'Szak' mezőben mindenki másképp írta be azt, hogy melyik szakon tanul, így ezt szabványosítani kellett, ha később szakok szerint akartunk vizsgálni. Erre a feladatra külön Python kódot írtunk, és ha volt olyan mező, ahol nem tudtuk eldönteni, hogy mi lenne az oda tartozó érték (például mert 'VBK' volt odaírva), arra bevezettünk egy globális 'UNKNOWN' változóértéket, azonban szerencsére ilyenből kevés volt. Kicsit

még tisztítani kellett a 'Tankör' értékeken is, de mivel ilyenből kevés volt, ezt manuálisan is meg tudtuk tenni. A 0. ZH eredményeket tartalmazó tábla szerencsére ennél jóval kisebb volt, csak a hallgató nevét, Neptun kódját, képzés nevét, illetve kódját, felvétel évét, valamint a ZH eredményt tartalmazta. Ebből értelemszerűen csak a névre és az eredményre volt szükségünk, a többi elhagyható.

A felvételi pontszámokat bontva (hozott pont, érettségi, többletpont) tartalmazó tábla hallgatók nevén és születési dátumán kívül tartalmazott még pár, a felvételi eljáráshoz és felvételi döntéshez kapcsolódó adatot, illetve a ponthatárt. Ebből a táblából csak a név és a pontokat tartalmazó oszlopok kellettek, a többit elvethettük.

A Matematika A1a jegyeket tartalmazó táblával már több dolgunk volt, mint az előző kettő esetben. Először is minden személyhez több rekord tartozhatott, legalább egy az A1a jegyhez, lehetett még egy az A2c jegyhez (nyilván csak azoknak, akik elvégezték az A1a-t, és ott maradtak az egyetemen), illetve, ha egy korábbi, nem a végleges jegyet eredményező vizsgán egy hallgató megbukott, akkor ahhoz is tartozott egy rekord. Az attribútumok között a hallgató nevén, Neptun kódján és az osztályzatán kívül szerepelt még a felvétel éve, képzés neve, kódja, státusz ID (Aktív, elbocsátott stb.), Pénzügyi státusz (állami/önköltséges), a tantárgy neve, kódja, kreditértéke, jegy típusa, bejegyzés dátuma, illetve, hogy elismert és hogy érvényes-e az adott jegy. Ezekből az adatokból nekünk csak a hallgató nevére és jegyértékeire volt szükségünk, ráadásul olyan formában, hogy minden sor egy hallgatóhoz tartozzon, és az oszlopok a tantárgyakból szerzett jegyeket tartalmazzák. Ehhez először Pythonban kiszűrtük az irreleváns oszlopokat, majd 'crosstab'-eléssel a kívánt formára hoztuk az adatokat, ahol még ügyelni kellett arra, hogy a korábbi vizsgajegyek ne kerüljenek bele, tehát minden hallgatóhoz tárgyként csak egy jegy tartozzon. Ezen kívül még, mivel az érdemjegyek szövegesen voltak megadva, azokat számszerűvé alakítottuk, hogy majd a későbbiekben könnyebb legyen velük dolgozni.

Egy külön adattábla tartalmazta még a kognitív eredményeknél a matekos eredményt blokkokra lebontva, amely valójában az elsőként tekintett adattáblának volt egy egyszerűsített, kevesebb attribútummal bíró változata. Ebből az adattáblából csak a hallgatók nevére, illetve a blokkonkénti teljesítményre volt szükségünk, a többit elhagytuk.

Így már rendelkezésünkre állt az összes tábla, egyenként tisztítva, és már csak az összefűzés volt hátra, amit R-ben könnyen meg tudtunk tenni, valamint még a végén rendeztük az oszlopok sorrendjét, hogy az adathalmaz logikus szerkezetű legyen. Fontos megjegyezni azonban, hogy nem minden elsőéves írt abban az évben kognitív tesztet, így összefűzés során (ahol valójában 'inner-join'-oltunk) kevesebb sorunk lett, mint ahányan abban az évben a BME VBK karára felvételt nyertek. A legvégén az így kapott adathalmaz 231 rekorddal és 21 oszloppal rendelkezett, amelyek között még esetlegesen szűrtünk különböző algoritmusok használata során.

4. Felderítő adatelemzés a két évben

A felderítő elemzés célja az adatok ábrázolása egyszerűen és gyorsan, annak érdekében, hogy az adattípusok egyszerű viszonyai szemléletesen áttekinthetőek legyenek. Erre a célra alkalmasak a oszlopdiagramok, szórásdiagramok, folyamatábrák és más grafikonok.

4.1. Általános ábrák

4.2. Folyamatábrák (Sankey-diagramok)

A Sankey-diagramok olyan folyamatábrák, amelyekben az ábrázolt folyamatok szélessége arányos a megfelelő folyamértékekkel. Ezáltal könnyen ábrázolhatóak adott állapotok közötti vándorlások, illetve ezek egymással vett aránya. Egy hátránya van, hogy egy ábrán csak két osztálycsoport között olvashatjuk le az áramlásokat. Ha egy harmadik osztálycsoportot is belehelyezünk, és csak a másodikkal kötjük össze, akkor az első és a harmadik közötti vándorlások nem jelennek meg az ábrán.

Dolgozatunkban négy opciót vizsgáltunk a két évben az elérhető adatoknak megfelelően. Mindkét esetben az állapotok a hallgatók valamely eredményét jelentették, például felvételi pontszámot.

1. 2019. matematikai teszt 1.,2. és 3. részeiben helyesen megoldott feladatok száma.
2. 2019. kognitív teszt eredmény - *A1* jegy - *A2* jegy.
3. 2019. és 2021. években: felvételi pontszám - nulladik ZH - kognitív teszt eredmény.
4. 2019. és 2021. években: felvételi pontszám - kognitív teszt eredmény - első féléves átlag - kumulált átlag.

Az első két esetben csak a 2019-es évfolyamról volt adat. Az utolsó esetben 2021-ben megegyezett a kutatás idejében az első féléves és a kumulált átlag, így ennél az utóbbi értelemszerűen nem lett még egyszer beletéve.

A matematikai teszt részeinek eredményein és az *A1*, *A2* jegyeken kívül a többi változó folytonos volt, így szükség volt ezek diszkrétizálására. Öt-öt osztály lett létrehozva minden esetben, de nem feltétlenül került mindegyikbe hallgató. Az átlagok kerekítve lettek, a felvételi, kognitív teszt és nulladik zh pontszámok másképp lettek osztályozva.

A felvételi pontszámok terjedelmük alapján öt azonos hosszúságú intervallumra lettek osztva. A harmadik változatban a nulladik zh és a teszt eredményei nullától az elérhető pontszámig lettek felosztva 20%-os lépésközökkel. A negyedik változatban, mivel a kutatás későbbi szakaszában készült, már a teszt eredmény is a terjedelem szerint lett felosztva, nem az elérhető pontszám alapján.

Ezután for ciklus használatával elkészültek a tranzíciós mátrixok. Egy ilyen mátrix i -edik sorának j -edik eleme azt mutatta, hogy az egyik változó i -edik osztályában hány olyan hallgató volt, aki a másik változó j -edik osztályába került. Ez a folyamatok szélességének megadásához volt szükséges.

A *plotly.graph_objects* könyvtár *go.Sankey* függvényével lettek elkészítve egyenként az ábrák. Minden esetben meg kellett adni számozva a kiinduló és beérkező állapotokat, így a középső csoportok mindkét listában szerepeltek. A program index alapján kapcsolta össze a két listából az állapotokat, így kellett egy harmadik lista, amelynek a megfelelő indexű eleme a két állapot közötti folyam mérete. Ezen kívül a folyamatokhoz és magukhoz az állapotokhoz is színeket kellett még rendelni, ismét egy-egy listában, figyelve az indexeket.

Végül a kirajzolódott ábrát kellett kézíleg igazítani, ugyanis az eredmények csoportjai nem feltétlenül álltak megfelelő sorrendben a függőleges tengelyen, illetve azokban az esetekben, amikor nem volt mind az öt osztályban hallgató, egy-egy csoport átcsúszott egy másik eredmény oszlopába. Ezek nem okoztak gondot, ugyanis a csoport mozgatásával a hozzátartozó sávok automatikusan mozogtak együtt.

A 2021-es gólyák a kutatás idejében még nem vehették fel az *A2* tárgyat, így értelem-szerűen csak a 2019-es évfolyamra készült ez az ábra.

Az ábrán rögtön látható, hogy a kognitív tesztben mindenki 20% felett teljesített, de nagyon kevés hallgatónak sikerült 80%-nál magasabb eredményt elérni. Legtöbbször 40% és 60% közötti pontszámot szereztek, a második legnagyobb osztály a 60-80%-os.

Ehhez képest az *A1* tárgyon elért jegyek eloszlása egyenletesebb. Néhány hallgatónak nem sikerült elvégezni a tárgyat, nagyobb részük kettest vagy hármaszt szerzett, és körülbelül a harmaduk 45/55 arányban négyest és ötöst szereztek.

A folyamatokat illetően, azok, akiknek nem sikerült a tárgy, a kettes és a hármas tesztosztályban voltak, azaz közülük valakinek nem sikerült 40%-ot elérni, de egyikőjük sem teljesített 60%-nál jobban. A tesztet legjobban megírok pedig többségben ötöst vagy négyest szereztek, azonban a négyes tesztosztályból több hallgató is csak hármaszt tudott szerezni. A hármas osztályból minden jegyhez megy folyam, a legtöbb ketteshez és hármas-hoz, illetve körülbelül a negyede egyenletesen oszlik el a négyes és ötös között. Egy érdekes eset van, a kettes tesztosztályból valakinek sikerült ötöst szerezni a tárgyból. Ennek ellenére ebből az osztályból a hallgatók kétharmada legfeljebb kettest ért el.

Ezekből sejthető, hogy egy hallgató teljesítménye a teszt során összefügg a további teljesítményével, de lehet jelentős romlás és javulás is.

Az *A2* tárgy jegyeinek eloszlása ismét más, több ötös és négyes, kevesebb hármas és körülbelül ugyanannyi kettes lett, és egy hallgatónak nem sikerült elvégezni. Érdekes, hogy az első körben kettest szerző hallgatók majdnem harmada javított négyesre vagy ötösre, a négyest szerző hallgatóknak pedig majdnem a fele kettesre rontotta. Hármasból is sokan javítottak, körülbelül a negyedik továbbra is hármaszt szerzett, viszont több, mint harmaduk rontott. Ötöst szerzőkből hatan kettesre rontottak (kb. 13%).

4.3. Klaszterezés

5. Prediktív analitika

5.1. Modellek és metodológia

Következő lépésként azt vizsgáltuk, hogy a két évben külön a bemeneti adatok alapján mennyire pontosan lehet előrejelezni egyes teljesítménymutatók értéket, illetve hogy a predikciókra nézve a különböző bemeneti változók milyen és mekkora szereppel bírnak. Ehhez többféle modellen többféle *gépi tanuláson alapuló osztályozó algoritmus* használatára volt szükség. A kutatás során kétféle eredmény alaposabb vizsgálatára összpontosítottunk mindkét évben: a "Matematika A1a - Analízis" tárgyból kapott érdemjegyre illetve az első félév végén megállapított kumulált tanulmányi átlagra. Ezen feladatok a felügyelt gépi tanulási problémák közé esnek, ahol a kitüntetett célváltozót szeretnénk a többi bemeneti változóval előrejelezni. Ekkor a gépi tanulás során a teljes adathalmazt tanító és teszhalmazra bontjuk, a tanító adathalmazon betanítjuk az algoritmusokat úgy, hogy a tanítóhalmazbeli adatok célváltozóját minél pontosabban prediktálják. Ezt követően valamilyen metrika szerint kiválasztjuk a betanított algoritmusok közül a legjobbat, majd az általánosítóképesség ellenőrzése végett a teszhalmazon is visszamérjük a teljesítményét.

A két vizsgált probléma közül az előbbi alapvetően egy osztályozási probléma öt osztállyal, míg az utóbbi egy regressziós feladat. Mivel azonban viszonylag kevés adat állt a rendelkezésünkre, ezért az érdemjegyek prediktálásánál az adatrekordok esetén a pontos érdemjegy helyett érdemjegy csoportok előrejelzésére koncentráltunk. A matematika érdemjegycsoportok prediktálására így kétféle modell került felvázolásra: egy *3 csoport modell*, illetve egy *2 csoport modell*.

A 2 csoport modell esetén a két csoportot a $\{5,4,3\}$ illetve $\{2,1\}$ osztályok adják, míg a 3 csoport modellnél az osztályok $\{5,4\}$, $\{3,2\}$ illetve $\{1\}$ módon alakultak. Az előbbi esetben az osztályok intuitívan a lemorzsolódási veszélyeztetettség szerint formálódtak, míg az utóbbiban egy általánosabb "jól teljesítő", "rosszul teljesítő", "lemorzsolódott" csoporthármast kívántunk elérni. Mindkét modell esetén külön vizsgáltuk a teljesítményt összes hallgatóra vonatkozóan illetve szétbontva vegyészmérnök és biomérnök hallgatókra egyaránt (a környezetmérnökökről nem állt rendelkezésre elég adat, így őket a szakonkénti bontásban kihagytuk).

Az egyes modellek és almodellek esetén a tanítás előtt főkomponens analízist (Principal Component Analysis) is alkalmaztunk. Az eljárás lényege, hogy az attribútumok súlyozott kombinációjával új attribútumokat hozunk létre, kevesebb számút mint az eredeti attribútumok száma, oly módon, hogy az információvesztés minimális legyen. Ennek következtében az adatok egy kisebb dimenziós térbe transzformálódnak át, s az ehhez használt megfelelő súlyozás pedig az attribútumok tapasztalati kovariancia mátrixának

segítségével határozható meg. PCA használata során az algoritmusok gyorsabban tanulnak a kisebb dimenziószám miatt, és olykor jobb eredményt is érnek el. Jelen kutatásban másrészt azért is döntöttünk a PCA alkalmazása mellett, mert a korábbi, ugyanezen problémakört vizsgáló, általunk átnézett tanulmányok nem használtak PCA-t még nagyobb attribútum szám esetén sem, ugyanakkor mi a korai tanító-tesztelési fázisnál azt tapasztaltuk, hogy az áttanszformált adaton jobban teljesítenek az osztályozó algoritmusok, így potenciálisan megéri alkalmazni.

5.2. Osztályozó algoritmusok

5.2.1. Lineáris regresszió

A lineáris regresszió jellegéből adódóan alapvetően folytonos célváltozó prediktálására alkalmas, ugyanakkor kategorikus, ordinális címkéjű adatok osztályozására is fel lehet használni. Az algoritmus lényege, hogy a magyarázóváltozók mindegyikéhez egy-egy súlyt rendelünk, majd az adatpont attribútumértékeinek vesszük a súlyokkal való súlyozott összegét, és esetleg egy bias tagot hozzáadva az így kapott szumma lesz a prediktált címkeérték. A cél a tanítás során a súlyok optimális megtalálása, miközben a tanítóhalmazon minimalizáljuk a predikciós hibát.

5.2.2. Naive Bayes

A Naive Bayes algoritmus működése mögött álló alapelv az, hogy feltesszük az attribútumok feltételes függetlenségét, amennyiben a célváltozó értéke ismert. Osztályozás során azt vizsgáljuk, hogy mely címkeérték mellett a legnagyobb a valószínűsége annak, hogy az adott adatrekord attribútumai éppen a felvett értékeket kapták. A megfelelő valószínűségeket a tanítóhalmazbeli adatok attribútumértékeinek különböző címkék melletti relatív gyakoriságai adják.

5.2.3. Gradient Tree Boosting

A Gradient Boosting algoritmus egy Ensemble típusú osztályozó, amelyek lényege, hogy sok gyenge teljesítményű prediktor (*weak learner*) eredményét felhasználva hoz egy erős predikciót.

5.2.4. Logisztikus regresszió

A logisztikus regresszió alapvetően bináris osztályozási problémák megoldására alkalmas, de kiterjeszthető másfajta feladatok megoldására is. Lényege, hogy a lineáris regresszióhoz hasonlóan az adatrekord attribútumértékeinek súlyozott összegét használjuk egy szigmoid¹ függvény bemeneteként, amely azt utána leképzi az $(0, 1)$ intervallumra, és amennyiben

¹A szigmoid függvény: $\sigma(z) = \frac{1}{e^{-z} + 1}$, ahol $z = x_1w_1 + x_2w_2 + \dots + x_nw_n + b$ a súlyozott összeg.

az output 0.5-nél nagyobb, úgy a pozitív osztályba soroljuk az adott rekordot.

5.2.5. SVM

Az SVM (Support-Vector-Machine) egy lineáris szeparálást használó bináris osztályozó algoritmus, amely egyéb problémákra is kiterjeszthető. Lényege, hogy különböző magfüggvénye segítségével az adatrekordokat egy magasabb dimenziós térbe képzi le, ahol olyan szeparáló hipersíkot keres, amely maximalizálja a vele párhuzamos, adatpontot nem tartalmazó térrészt. A cél a megfelelő magfüggvény és a szeparáló hipersík megtalálása, amellyel a különböző címkéjű adatpontok lineárisan szeparálhatóak.

5.3. Implementálás és optimalizálás

A 2 és 3 csoport modellnél a célváltozó-értékeket rendre 1, 0-ra illetve 2, 1, 0-ra módosítottuk. Az osztályozó algoritmusok közül a Naive Bayes, kNN és Gradient Tree Boosting alaptól jól kezeli a kategorikus változókat, az SVM-nél valamint a regressziós eljárásoknál viszont minimális változtatásra volt szükség. Az SVM és logisztikus regresszió mivel csak bináris osztályozásra alkalmas, így ott *One-Vs-Rest* elvű osztályozást használtunk, lineáris regressziónál pedig az címkeértékek ordinális jellege miatt egyszerűen a prediktált értéket kerekítettük a legközelebbi csoportcímkére. Kumulált átlag előrejelzésénél csupán lineáris regressziót alkalmaztunk.

Valamennyi modellnél a teljes adathalmazt felbontottuk tanító- és teszhalmazra 70-30 arányban, továbbá a két adathalmazban a folytonos változók kvantilis alapú 0-1 skálázásnak lettek alávetve, hogy a távolság alapú osztályozók jobban teljesítsenek. Az algoritmusok hiperparamétereinek optimális megválasztására 5-szörös kersztvalidációt alkalmaztunk. Fontosabb optimalizációs lépések csak a 2-3 csoport modellek esetén történtek. A Naive Bayes és lineáris regresszió algoritmusoknál jellegük és/vagy implementálásuk végett nem történt optimalizálás, a többi algoritmus esetén az alábbi hiperparaméterek kerültek változtatásra:

- **Gradient Tree Boosting:**
 - Tanulórata (0.01 és 5 között változtatva 0.1-es lépésközzel)
 - Boosting fázisok száma (5 és 100 között változtatva 5-ös lépésközzel)
 - Vágási feltétel (négyzetes hiba, Friedman MSE)
 - Maximális famélység (3,4 és 5)
- **Logisztikus regresszió:**
 - Regularizációs paraméter (0.05 és 5 között változtatva 0.05-ös lépésközzel)
 - Optimalizálási módszer ("SAG", "SAGA")
- **SVM:**

- Regularizációs paraméter (0.1 és 5 között változtatva 0.1-es lépésközzel)
- Magfüggvény (lineáris, legfeljebb 3-ad fokú polinomiális, RBF)

A PCA esetében a főkomponensek számát 2 és 8 között változtattuk 2-es lépésközzel, az algoritmusokat minden főkomponenyszám mellett 5-szörös keresztvalidációval optimalizáltuk. Az optimalizálásnál használt metrikák:

- A 2 csoport modell esetén az **F1-érték**, amely a Precision és a Recall harmonikus közepe, melyek közül a Precision jelöli, hogy az 1-es címkéjűnek osztályozott
- A 3 csoport modell esetén a **Kiegyensúlyozott pontosság** (Balanced Accuracy), amely címkeosztályok kiegyensúlyozatlanságát figyelembe véve az osztályokhoz tartozó Recall értékek számtani közepét adja vissza.

A megfelelő metrikák mellett kiválasztottuk a legjobb algoritmusokat, amelyeket a teszt-halmazon visszamértünk.

6. Modellek kiértékelése

		Osztályozó algoritmusok					
		Grad. Boost.	Naive B.	Log. reg.	SVM	Lin. reg.	
3 csoport	Összesített	2 PC	66.67	62.67	52.00	58.67	61.33
		4 PC	70.67	60.00	57.33	53.33	65.33
		6 PC	65.33	62.67	54.67	68.00	64.00
		8 PC	64.00	68.00	53.33	62.67	64.00
	Szakonként	2 PC	78.71	80.36	73.85	67.24	80.36
		4 PC	75.41	80.36	47.80	65.59	78.71
		6 PC	82.02	80.36	75.47	67.21	82.02
		8 PC	80.36	78.71	37.42	75.44	78.71
2 csoport	Összesített	2 PC	59.42	69.57	69.57	72.46	69.57
		4 PC	59.42	71.01	73.91	69.57	73.91
		6 PC	63.77	69.57	73.91	71.01	73.91
		8 PC	68.12	69.57	71.01	73.91	72.46
	Szakonként	2 PC	67.31	67.31	65.69	62.41	64.03
		4 PC	64.00	64.00	67.31	67.31	63.97
		6 PC	64.03	65.69	64.07	65.65	67.27
		8 PC	64.07	62.51	57.49	67.34	68.96

1. táblázat. A 2019-es adatsor eredményei

Az 1. táblázatban a 2019-es adatokra optimalizált modellek teljesítménye látható, ahol az első 3 oszlop a modellt és a használt főkomponensek (PC) számát mutatja, míg a jobb oldali 5 oszlop az optimalizált algoritmusok teljesítményét szemlélteti. A szakonkénti bontásban szereplő teljesítménymutatók értékei a vegyész- és biomérnök hallgatók adatain elért pontszámok, az egyes szakokon tanuló hallgatók számával vett súlyozott átlagai. A 3 csoport modellek esetén a Gradient Boosting algoritmus, míg a 2 csoport modellek esetén a regressziós algoritmusok teljesítettek a legjobban, ugyanakkor a 2 csoport modellek teljesítménye többségében alulmúlja a 3 csoport modellekét. Mivel a 2 csoport modellben az osztályok eloszlása kiegyensúlyozottabb, ez arra enged minket következtetni, hogy a 2-es és 3-as érdemjegyet szerzők között a bemeneti adatok tekintetében nincsenek nagy különbségek.

		Osztályozó algoritmusok					
		Grad.	Boos.	Naive B.	Log. reg.	SVM	Lin. reg.
3 csoport	Összesített	2 PC	53.74	72.97	67.75	73.30	56.80
		4 PC	61.86	70.49	65.82	66.63	51.91
		6 PC	68.27	68.53	67.73	70.80	64.00
		8 PC	66.31	70.20	67.21	67.17	64.00
	Szakonként	2 PC					
		4 PC					
		6 PC					
		8 PC					
2 csoport	Összesített	2 PC	69.45	67.89	69.77	69.77	67.89
		4 PC	72.74	72.90	78.23	76.51	77.91
		6 PC	67.89	72.90	79.80	79.80	76.19
		8 PC	68.05	74.62	79.96	83.24	79.63
	Szakonként	2 PC					
		4 PC					
		6 PC					
		8 PC					

2. táblázat. A 2021-es adatsor eredményei

Hivatkozások