



**Kemometria és FT-NIR spektroszkópia alkalmazása  
az élelmiszeranalitikában**

**Ph.D. értekezés**

**Rácz Anita**

**Budapest**

**2016**

**A doktori iskola**

**megnevezése:** Élelmiszertudományi Doktori Iskola

**tudományága:** Élelmiszertudományok

**vezetője:** **Dr. Vatai Gyula**  
egyetemi tanár, Dsc  
SZIE, Élelmiszertudományi Kar,  
Élelmiszeripari Műveletek és Gépek Tanszék

**Témavezetők:** **Dr. Fodor Marietta**  
habilitált egyetemi docens, PhD  
SZIE, Élelmiszertudományi Kar,  
Alkalmazott Kémia Tanszék

**Dr. Héberger Károly**  
tudományos tanácsadó, DSc  
MTA Természettudományi Kutatóközpont,  
Anyag- és Környezetkémiai Intézet

.....  
Az iskolavezető jóváhagyása

.....  
A témavezető jóváhagyása

.....  
A témavezető jóváhagyása

## Tartalomjegyzék

RÖVIDÍTÉSEK JEGYZÉKE .....	5
1. BEVEZETÉS .....	9
2. IRODALMI ÁTTEKINTÉS .....	11
2.1 A Q10 tartalmú étrendkiegészítők .....	11
2.2 Az energiailatok .....	13
2.3 Az antioxidáns kapacitás .....	16
2.4 Az FT-NIR spektroszkópia elméleti háttere .....	19
2.5 Az FT-NIR spektroszkópiai modellek .....	23
2.5.1 A spektrumok transzformációja .....	23
2.5.2 Modellek validálása .....	24
2.5.3 A modelleket leíró teljesítmény paraméterek .....	25
3. CÉLKITŰZÉS .....	28
4. ANYAG ÉS MÓDSZER .....	30
4.1 FT-NIR készülék felépítése .....	30
4.2 Nagyhatékonyságú folyadékkromatográfia .....	32
4.3 Mintaelőkészítés .....	33
4.4 Klasszikus mérési módszerek .....	33
4.4.1 Schoorl módszer .....	33
4.5 Kemometriai módszerek .....	34
4.5.1 Főkomponens-elemzés (PCA) .....	34
4.5.2 Hierarchikus fürtelemzés (klaszter analízis, HCA) .....	36
4.5.3 Lineáris diszkriminancia elemzés (LDA) .....	37
4.5.4 Parciális legkisebb-négyzetek módszere (regresszió és diszkriminancia elemzés) .....	38
4.5.5 Véletlen erdő módszere (Random forest, RF) .....	40
4.5.6 Fejlesztett fák módszere (Boosted tree, BT) .....	41
4.5.7 Változókiválasztási eljárások .....	41
4.5.8 Rangsorolás és összehasonlítás .....	43
5. EREDMÉNYEK .....	48
5.1 Kemometriai módszerfejlesztések .....	48
5.1.1 Új változóselektációs módszer létrehozása – az intervallum szelektivitási arány (iSR) ....	48
5.1.2 „n-class”, avagy több osztályos ROC görbék .....	49
5.1.3 Osztályozási módszerek összehasonlítása „n-class” ROC görbék segítségével .....	50
5.2 Q10 tartalmú étrendkiegészítők hatóanyagtartalmának vizsgálata .....	63
5.2.1 Az iPLS változókiválasztással kapott eredmények .....	65

5.2.2 iSR változó kiválasztással kapott eredmények .....	67
5.2.3 A genetikus algoritmus segítségével végzett változó kiválasztással kapott eredmények....	69
5.2.4 A kapott modellek összehasonlítása SRD módszerrel.....	72
5.2.5 A modellek külső validálása és az eredmények táblázatos összefoglalása .....	74
5.3 Az energiaitalok FT-NIR spektrumának kemometriai elemzése .....	75
5.3.1 Energiaitalok osztályozási lehetőségei .....	76
5.3.2 Energiaitalok koffein és cukortartalmának meghatározása .....	78
5.4 Antioxidáns kapacitás meghatározási módszerek csoportosítása, rangsorolása .....	84
5.4.1 HCA eredmények .....	85
5.4.2 PCA eredmények .....	85
5.4.3 SRD eredmények .....	87
5.4.4 GPCM eredmények .....	91
5.5 Regressziós modellek előrejelző képességét meghatározó paraméterek összehasonlítása .....	93
6. KÖVETKEZTETÉSEK ÉS JAVASLATOK .....	96
7. ÚJ TUDOMÁNYOS EREDMÉNYEK .....	98
8. ÖSSZEFOGLALÁS.....	100
9. SUMMARY .....	102
10. KÖSZÖNETNYÍLVÁNÍTÁS .....	104
11. IRODALOMJEGYZÉK.....	105
12. MELLÉKLETEK.....	116
M1 A teljesítményparaméterek összehasonlításához felhasznált paraméterek listája .....	116
M2 A HPLC mérésekhez használt eluensek táblázata .....	118
M3 A ROC görbék készítésének sémája .....	119
M4 A többszintű ROC görbék megalkotásának programkódja .....	120
M5 Az energiaitalok PCA adatkészlete alapján történő optimalás a fejlesztett fák módszerénél	121
M6 Publikációs lista .....	122

## RÖVIDÍTÉSEK JEGYZÉKE

ABTS	2,2'-azinobisz-(3-ethylbenzo-tiazolin)-6-szulfonsav
ACL	Zsírolható antioxidáns kapacitás
ACN	Acetonitril
ACW	Vízoldható antioxidáns kapacitás
APCI	Atmoszférikus nyomású kémiai ionizáció
AUC	Görbe alatti terület érték
<b>B</b>	A regressziós koeficiensek mátrixa (PLS)
<i>b</i>	Regressziós koeficiens
Bias	Torzítás
BT	Fejlesztett fák módszere (Boosted tree)
CCC	Egyezési (egybehangzási) tényező
CE	Kapilláris elektroforézis
CV	Kereszt-ellenőrzés
D	Folytonos bináris osztályozó változó
DAD	Diódasoros detektor
DLLME	Diszperzív folyadék-folyadék mikroextrakció
DPLS	PLS diszkriminancia elemzés
DPPH	2,2-difenil-1-pikrilhidrazil szabad gyök
<b>E</b>	hibamátrix (PCA, PLS)
ELSD	Elpárologtatásos fényszórás detektor
ESI	Elektronspray ionizáció
ET	Elektronátviteli reakció
<i>F</i>	Fisher érték
<b>F</b>	Hibamátrix (PLS)
FRAP	Vas redukálóképességén alapuló antioxidáns kapacitás
FT-NIR	Fourier-transzformált közeli infravörös spektroszkópia

GA	Genetikus algoritmus
GPCM	Általánosított pár-korrelációs módszer
HAT	Hidrogén atom átvitelével megvalósuló reakció
HCA	Hierarhikus fürtelemzés
HPLC	Nagyhatékonyságú folyadékkromatográfia
HPTLC	Nagyhatékonyságú vékonyrétegekromatográfia
iPLS	Intervallum PLS módszer
$K_x$	Változók közötti inter-korreláció, PCA alapján
LC	Folyadékkromatográfia
LDA	Lineáris diszkriminancia elemzés
LOD	Kimutatási határ (Limit of detection)
LOF	Fredman-féle „lack of fit” kritérium
LOO	Egy elem kihagyásos kereszt-ellenőrzés
LOQ	Meghatározási határ (Limit of quantification)
MAE	Átlagos abszolút hiba
MALDI	Mátrix segített lézer deszorpciós ionizáció
MLR	Többváltozós lineáris regresszió
MS	Tömegspektrometria
MS/MS	Tandem tömegspektrometria
MSC	Arányos (többszörös) szóródás korrekció
NIPALS	Nemlineáris iteratív parciális legkisebb-négyzetek algoritmus
NMR	Nukleáris mágneses rezonancia ( $^1\text{H}$ NMR - proton/hidrogén NMR)
ORAC	Oxigéngyök abszorpciós kapacitás
$P'$	a főkomponens-együttható mátrix transzponáltja
PCA	Főkomponens-elemzés
PEG	Polietilén glikol
PLS DA	Parciális legkisebb-négyzetek diszkriminancia elemzés
PLSR	Parciális legkisebb-négyzetek regresszió

PRESS	Előrebecslési hiba négyzetösszege
$Q'$	A PLS komponens-együttható mátrix transzponáltja
Q10	Q10 koenzim rövidítése
$Q^2$	A kereszt-ellenőrzésre vonatkozó determinációs koefficiens négyzete
Q2-F1, Q2-F2, Q2-F3	A kereszt-ellenőrzésre vonatkozó determinációs koefficiens négyzetének változatai a külső validáláskor (teszt készletre)
$Q^2_{\text{LOO}}$	Az egy elem kihagyásos kereszt-ellenőrzésre vonatkozó determinációs koefficiens négyzete
QCAR	Kvantitatív összetétel-hatás összefüggés
QSAR	Kvantitatív szerkezet-hatás összefüggés
R	(reakcióegyenletben) a szénváz, gyakran alkil-gyök
$R^2$	Determinációs koefficiens négyzete
RF	Véletlen erdő módszere (Random forest)
RMSE	Átlagos négyzetes hiba
RMSEC	Átlagos négyzetes hiba a kalibrációs adatkészletre nézve
RMSECV	Átlagos négyzetes hiba a validálási adatkészletre nézve
RMSEP	Átlagos négyzetes hiba a teszt adatkészletre nézve
ROC	Vevő-működtető jelleggörbe (Receiver operating characteristic curve)
RP	Fordított fázis
RPD	Korrigált tapasztalati szórás
RSS	Reziduális eltérés négyzetösszeg
SD	Korrigált empirikus szórás
SEP	A validálás torzítással korrigált becslési hibája
SIMCA	Osztályanalógiák közvetett modellezése
SR	Szelektivitási arány
SRD	Rangsámkülönbségek összegének módszere
$T$	a főkomponesek mátrixa (PCA), ill. a PLS komponensek mátrixa (PLS)
TAC	Teljes antioxidáns kapacitás
TBARS	Tiobarbitursav reaktív anyag mérési módszer

TEAC	Troloxra vonatkoztatott antioxidáns kapacitás
THF	Tetrahidrofurán
TPC	Összes polifenol tartalom meghatározás Folin-Ciocalteu reagenssel
TRSC	Összes antioxidáns meghatározása kemilumineszcencia elvén
TSS	Teljes eltérés négyzetösszeg
<i>U</i>	A PLS komponensek mátrixa (rejtett változók)
UPLC	Ultranagy hatékonyságú folyadékkromatográfia
VIP	Változó fontosság paramétere
<i>X</i>	A független változó(k) általában
<i>X</i>	A független változók mátrixa
<i>Y</i>	A függő változó általában
<i>Y</i>	A függő változó(k) mátrixa



*"Felfedezni valamit, annyit tesz, mint látni, amit mindenki lát, és közben arra gondolni, amire még senki." (Szent-Györgyi Albert)*

## **1. BEVEZETÉS**

Az élelmiszertudományi és élelmiszeripari vizsgálatokban előszeretettel használnak olyan analitikai eszközöket, amelyek gyors, költséghatékony és lehetőség szerint roncsolásmentes vizsgálatokat biztosítanak. A doktori munkám során használt közeli infravörös spektroszkópia is ezt az irányelvet követi, amely hatékonyan alkalmazható az élelmiszerek beltartalmi paramétereinek ellenőrzésére, mennyiségi és minőségi vizsgálatára. Az élelmiszerek minőségellenőrzése nem csupán a fogyasztók, de a gyártók érdeke is, hiszen hatalmas a verseny az élelmiszerek piacán, ahol a megfelelő minőség garanciája már egy elengedhetetlen feltételnek számít. Bár a magyar, élelmiszervizsgálatra vonatkozó szabványok sok esetben "elavultnak", hosszadalmasnak és drágának tekinthetők, a feltörekvőben lévő gyors analitikai eljárások – köztük a közeli infravörös spektroszkópia – alkalmazásával egy új, gazdaságilag sokkal kedvezőbb és hatékonyabb irányt mutathatunk meg.

Az ezredfordulót megelőző évtized óta a modern analitikai kémia és ezzel együtt az élelmiszeranalitika jelentősége is rohamosan nő. Az egyre modernebb és gyorsabb számítógépek megjelenése hatalmas szerepet játszott az analitikai módszerek fejlődésében, főként a spektroszkópia területén felhalmozódó nagymennyiségű adat feldolgozásában. Az utóbbi évtizedek "adatrobbanásával" járó hosszadalmas számolásokat a nagy teljesítményű számítógépek (akár szuperszámítógépek) használata nem csak lehetővé tette, de meg is könnyítette. A kemometria, mint a statisztika többváltozós kémiai adatokra történő alkalmazási területe már az előző évszázad közepétől egyre népszerűbbé vált, de az igazi áttörést a 80-as és 90-es évek számítógépes fejlődése hozta meg. Olyan cikkek születtek akkor, amelyek a mai napig a kemometriai tudomány alappillérei.

Az élelmiszertudományban is a kemometria számos alkalmazására lelhetünk. Segítségével felderíthetjük az adatainkban rejlő mintázatokat, hasonlóságokat és különbségeket, olyan rejtett összefüggéseket, amelyek mennyiségi és minőségi kiértékelések alapjául szolgálhatnak. A kemometriát előszeretettel társítják a spektroszkópiái kiértékelésekhez, hiszen a spektroszkópiái módszerek közül jó néhány az általuk generált hatalmas adatkészletek miatt elképzelhetetlen lenne az adatok kemometriai értékelése nélkül. A doktori munkám során alkalmazott legfontosabb analitikai eszköz, a közeli infravörös spektroszkópia elengedhetetlen

feltétele is a sokváltozós adatelemzés, mondhatni társtudománya a kemometria. A komplex élelmiszer-mátrixok spektrumainak kiértékelése csakis a kemometria segítségével történhet meg.

A doktori értekezésemben élelmiszer és élelmiszerekhez köthető minták mennyiségi és minőségi analitikai elemzését hajtottam végre, olyan kemometriai modelleket létrehozva, amelyek a későbbiekben az élelmiszerek minőségellenőrzésében, beltartalmi paramétereinek vizsgálatában elengedhetetlenek. Mindemellett olyan új kemometriai módszereket fejlesztettem, amelyek nem csak az élelmiszeranalitikában, de az analitikai kémia más területein is hasznosak lehetnek. A dolgozat célja tehát nem egy adott termék/komponens hatékonyságának vagy hasznosságának vizsgálata volt, hanem a változatos fizikai-kémiai paraméterekkel rendelkező élelmiszer minták komplex vizsgálata és a használt módszerek továbbfejlesztése, átfogó képet nyújtva a kemometria és FT-NIR spektroszkópia használati és továbbfejlesztési irányairól és lehetőségeiről.

A dolgozat a vizsgált mintamátrixok alapján három nagyobb részre tagolható: i) a Q10 tartalmú étrendkiegészítők FT-NIR és kemometriai vizsgálata ii) az energiatalok komplex mintázatfelismerési és mennyiségi meghatározása kimondottan a minták koffein és cukortartalmára összpontosítva, valamint iii) az antioxidáns kapacitás meghatározási technikák rangsorolása, osztályozása kemometriai módszerekkel. Míg az első két témakör az FT-NIR spektrumok alapján történő mintázatfelismerésre és regressziós modellépítésekre, azok fejlesztésére fekteti a hangsúlyt, addig a harmadik téma az élelmiszeranalitikai módszerek rangsorolási, összehasonlítási lehetőségeit mutatja be. A kemometriai módszerek fejlesztéseit a dolgozatomban az előbbieken felsoroltaktól külön tárgyalom.

## 2. IRODALMI ÁTTEKINTÉS

### 2.1 A Q10 tartalmú étrendkiegészítők

A Q10 koenzim számos egyéb néven is ismert a szakirodalomban, pl. ubikinon, ubidekarenon, koenzim-Q, 1,4-benzokinon. A Q10 egy zsírolható bezokinon származék, amelynek rövidítésében található Q a szerkezetében található kinon csoportra utal, a tizes szám pedig a kinon csoporthoz kötődő izoprenil alegységek számára. Az izoprenil alegységek a molekula mindkét féle oxidációs állapotában hat és tíz között változhatnak, de az emberben és a legtöbb állatban is a tíz izoprenil alegységet tartalmazó oldallánc forma található meg.

A Q10 koenzim nagy szerepet játszik, mint elektron átvivő (transzporter) a mitokondriális légzési láncban. Ezenkívül a redukált formája (ubikinol), antioxidánsként is hatásos (Ernster and Dallner, 1995). A Q10 koenzim redukált és oxidált formája is megtalálható a sejtmembránban, például a vérben vagy a szérum lipoproteinekben, de a mennyiségük az öregedéssel együtt csökken. A Q10 napjaink egyik legismertebb koenzime (talán „csodaszere” is), amely számos étrendkiegészítő aktív komponense.

A Q10 koenzim népszerűségének oka, hogy olyan betegségek megelőzésében és kezelésében segíthet, mint például a mitokondriális megbetegedések (Hargreaves, 2014) és az emlőrák (Lockwood et al., 1995). A koszorúerek megbetegedésében szenvedőknél képes csökkenteni az oxidatív stresszt és növelni az antioxidáns enzimek aktivitását (Lee et al., 2012). Más tanulmányok ajánlása szerint, a Q10 koenzim pangásos szívelégtelenség esetén növelheti az ejekciós frakciót (ez egy százalékos érték, amely megadja hogy a szív bal kamrájának összehúzódási periódusában a benne lévő vér térfogatának hány százaléka jut a főverőérbe) (Fotino et al., 2013). A szájon át alkalmazott Q10 koenzim krónikus szívbetegségek esetén megnövelheti a funkcionális kapacitást és javíthatja az endotéliális funkciókat mellékhatások fellépése nélkül (Belardinelli et al., 2006). López-Llunch és munkatársai (2010) kifejtették publikációjukban, hogy a Q10 koenzim az öregedési folyamatok egyik kulcsfaktora, így ezen tápanyag bevitele a jó egészség megőrzési stratégia része lehet különösen az életkor előrehaladtával. Q10 koenzimet tartalmaz például a sertés- és marhahús, a csirkehús, a hal, a repceolaj, illetve nagyobb mennyiség bevitelének szükségessége esetén számos Q10 tartalmú étrendkiegészítő is a fogyasztók rendelkezésére áll (Mattila és Kumpulainen, 2001).

Az elmúlt évek során a Q10 tartalmú étrendkiegészítők piaca rohamosan növekedésnek indult, rengeteg gyártó termékpalettáján található már meg, így a termékek minőségellenőrzésére nagy hangsúlyt kell fektetni. Az erre a célra legáltalánosabban használt analitikai vizsgálatokat

különböző típusú mintamátrixok esetén az **1. táblázatban** foglaltam össze. A szakirodalomban a módszerfejlesztések nem csak az étrendkiegészítők vizsgálatára térnek ki, hanem sok esetben az emberi szervezetben található Q10 koenzim mennyiséget is szeretnék detektálni, ezért a táblázatban felsorolt módszerek között néhány az ilyen típusú vizsgálatokra vezethető vissza.

**1. táblázat:** A Q10 koenzim tartalmú mintamátrixok analitikai vizsgálatainak összefoglalása

Módszer	Mintamátrix	Komponens	Név
RP-HPLC/DAD és LC-(APCI)MS az azonosításra	Étrend-kiegészítők	A és E vitamin, Q10 koenzim	D. E. Breithaupt és mts. (2006)
RP-HPLC-UV és elektrokémiai detektor	Vér, plazma, egyéb szövetek	E vitamin izomerek, Q10 és Q9 koenzimek (oxidált és redukált forma)	J. K. Lang és mts. (1987)
RP-HPLC-kulometriás detektor	Humán plazma	Q10 koenzim (redukált és oxidált forma)	P. H. Tang és mts. (2001)
Első derivált UV spektroszkópia és HPLC/DAD (PEG kötött kolonna)	Oldat	Ellagsav és Q10 koenzim	D. V. Ratnam és mts. (2006)
HPLC-UV	Nyers anyagok és étrend-kiegészítők	Q10 koenzim	S. Lunetta és mts. (2008)
UPLC-ESI-MS/MS	Természetes kozmetikumok	Q10 koenzim	J. H. Lee és mts. (2014)
HPLC-MS	Növényi olaj	Q10 és Q9 koenzimek	R. Rodríguez-Acuña és mts. (2008)
<sup>1</sup> H NMR	Étrend-kiegészítők	Q10 koenzim	Y. B. Monakhova és mts. (2013)
HPLC-UV	Emberi vizeletminta	Q10 koenzim	D. Yubero és mts. (2015)
HPLC-ESI-MS/MS	Étrend-kiegészítő	Q10 koenzim (redukált forma)	A. Vass és mts. (2014)

A felsorolt példák alapján szembejövő, hogy a leggyakrabban alkalmazott technikák között leginkább a folyadékkromatográfias módszerek különböző fajtái találhatóak, amelyek viszont drágák, időigényesek és hosszabb-rövidebb mintaelőkészítést is igényelnek. Így célom a Q10 tartalmú étrendkiegészítők hatóanyagtartalmának mennyiségi meghatározására FT-NIR spektroszkópiával; azaz egy új, roncsolásmentes és gyors analitikai eljárás kidolgozása volt, amely kiválthatja az eddig alkalmazott, sokkal időigényesebb módszereket.

## 2.2 Az energiatalok

Az energiatalok olyan funkcionális üdítő italok, amelyek bizonyos ideig képesek fokozni az emberi szervezet teljesítőképességét és anyagcseréjét. Az energiatalok az egyedi ízvilágukkal, színükkel és megjelenésükkel (és nagy koffein tartalmukkal) az utóbbi évtizedekben meghódították az egész világot. Magyarországon is hatalmas népszerűségnek örvendenek. Másrészt mértéktelen fogyasztásuk veszélyes mellékhatásokhoz vezet. Megtévesztő külsejük, az üdítőitalokhoz hasonló ízük és hatalmas akár 1,5 literes kiszerelésük miatt különösen veszélyesek a fiatalok, iskoláskorú gyermekekre, akik sokszor nagy mennyiségben fogyasztják. A legtöbb országban az energiatalok eladása kiskorúaknak is engedélyezett, így bárki kontrollálhatatlanul fogyaszthatja. Magyarországon az energiatal normál italként árusítható, forgalmazható a boltokban, de vannak olyan országok is, mint pl. Dánia és Norvégia, ahol ez nem megengedett. Bár hazánkban a márkák száma több mint százra tehető, a legpontosabb megfogalmazásuk a dobozon feltüntetett „koffein tartalmú (alkoholmentes) szénsavas üdítőital”. Ez a megfogalmazás azóta vált elterjedtté, hogy 2013-ban bevezették a „chipsadó”-ként elhíresült népegészségügyi adót, amely miatt az energiatalok jelentős összetevőbeli változásokon mentek keresztül. Az energiatal fogalma nem csak Magyarországon, de Európában sincs egyértelműen rögzítve, hazánkban a Magyar Élelmiszerkönyv se tartalmaz külön bejegyzést erre a kategóriára.

Az utóbbi évtizedben számos publikáció számolt be az energiatalok két legnagyobb veszélyt jelentő összetevőjének, a koffeinnek és a cukornak az egészségügyi hatásáról. Egy a témában íródott összefoglaló cikk szerint az 1980-tól egészen 2014-ig bezárólag végzett kutatásuk szerint összesen 2097 cikk foglalkozott (a PubMed és Google-Scholar eredményei alapján) az energiatalok káros egészségügyi hatásival. Arra a következtetésre jutottak, hogy az energiatal fogyasztás egy olyan problémakör, ami főként a fiatalokat, gyerekeket valamint a felnőtt férfiakat érinti. Emellett kapcsolódik a megnövekedett droghasználathoz és a kockázatvállalásos viselkedéshez is (Ali et al., 2015). A koffein túlzott bevitele a szervezetben magas vérnyomáshoz és szívritmuszavarhoz vezethet. Hosszú távú fogyasztás esetén az általános túladagolási tünetek: pl. hányinger, hányás, pánikroham mellett, akár máj és vese problémákat is okozhat (Reissig et al., 2009). Az energiatalok megnövelhetik a szisztolés vérnyomást, megváltoztathatják a szervezet elektrolit háztartását, amely változások kumulálódva repolarizációs abnormalitásokhoz vezethetnek. A szervezet ezen élettani válaszai szívritmuszavarokat és egyéb szív- és érrendszeri elváltozásokat okozhatnak (Kozik et al., 2016). A kontrollálatlan koffein bevétel a gyerekek és fiatal felnőttek körében nem csupán a szív működésének abnormalitásához, de hangulat és viselkedészavarhoz is vezethet (Seifert et al.,

2011). Heckman és munkatársai (2010) felhívták a figyelmet arra is, hogy a koffein bevitel, és így az energiatalok fogyasztása, a terhesség idején szintén veszélyekkel jár. Káros a magzat fejlődésére nézve, valamint csökkenti a termékenységet. A cukortartalom egészségügyi hatásai kapcsán is számos probléma merült fel, amelyek közül a két legkiemelkedőbb az elhízás és a kettes típusú cukorbetegség kialakulásának kockázata (Malik et al., 2006). Az energiatalok átlagos cukortartalma 10 g / 100 ml. Egy új „trend” is megjelent az utóbbi években, amely nagyon gyorsan terjed a fiatalok, egyetemisták körében: az energiatalok kombinálása alkohollal (Malinauskas et al., 2007). Ez viszont újabb súlyos problémákhoz vezethet, pl. a szervezet dehidratálásához, amely az alkohol és a koffein együttes hatásának köszönhető. Ferreira és munkatársai (2006) kimutatták, hogy az energiatal kombinációja alkohollal azt az érzetet kelti, hogy a motorikus képességek nem csökkennek. Egy másik kísérlet egyetemi diákok energiatal fogyasztási szokásait vizsgálta, amelyben összefüggést találtak az alkohol és energiatal kombinációját fogyasztó diákok és az alkoholos befolyásoltsághoz kapcsolódó eseményekben való részvétel között (O'Brien et al., 2008).

Az energiatalok legfontosabb összetevői általában a következők: víz, cukor, koffein, étkezési sav, savanyúságot szabályzó anyagok, színezék, aroma, inozitol, B vitaminok és taurin. Magyarországon a népegészségügyi adó bevezetése óta már csak elszórtan lehet találni taurinos italokat a rájuk kirótt nagy adóteher miatt. A taurint gyakran pótolják a nevesebb gyártók is argininnel, de a legtöbb esetben – főként az olcsóbb italoknál – meg se próbálják pótolni, hanem egyszerűen csak kihagyják az összetevők közül. Mivel az energiatalok fogyasztása egy mindennapos és növekvő problémának nevezhető (különös tekintettel a fiatalkorúakra nézve), így a cukor és koffein tartalom ellenőrzésének kérdése nagyon fontosnak számít mind a fogyasztók, mind a gyártók számára. Minden országnak eltérő szabályrendszere van, viszont a több száz energiatal márka között, mintha ez valahogy nem megfelelően történne. Számos módszerrel találkozhatunk a koffein koncentráció meghatározására a szakirodalomban és ugyanígy a cukortartalom meghatározásra is rengeteg publikáció áll a rendelkezésünkre. Általánosságban az alkalmazott módszereket két csoportra tudjuk bontani: spektroszkópiai és kromatográfiás technikákra. Az első csoportból kiemelhető Armenta és munkatársainak (2005) kutatása, amelyben Fourier-transzformált Raman spektroszkópiát használnak a kereskedelemben kapható energiatalok vizsgálatára, de ugyanígy találkozhatunk UV/Vis derivatív spektroszkópiával, amelyet szilárd fázisú extrakció kíséretében végeznek el (Pieszko et al., 2010). A másik csoportból példaként lehet említeni a HPTLC-UV denzitometriás elemzést (Abourashed és Mossa, 2004) vagy a diszperzív folyadék-folyadék mikroextrakciót (DLLME) gáz kromatográfiás elemzéssel összekötve (Seresheti et al., 2014). További példákat a koffein,

cukor és egyéb jellegzetes komponensek meghatározására a **2. táblázatban** foglaltam össze részletesen.

**2. táblázat:** A koffein, cukor és egyéb komponensek meghatározásának összefoglalása

Módszer	Mátrix	Komponens	Egyéb információ	Szerző
HPTLC-UV denzitometriás elemzés	Energiaital, gyógynövény termékek	koffein	visszanyerés= 98.90±3.46 pontosság= 99.84±2.87	E. Abourashed és mts. (2004)
Szilárd fázisú FT-Raman spektroszkópia	Energiaital	Koffein	LOD= 18 mg L <sup>-1</sup>	S. Armenta és mts. (2005)
UV/VIS derivative spektrofotometria + szilárd fázisú extrakció	Energiaital	Koffein, taurin	LOD=0.21 LOQ= 0.63 µg mL <sup>-1</sup>	C. Pieszko és mts. (2010)
Diszperzív folyadék-folyadék mikroextrakció (DLLME) + gáz kromatográfia	Tea, kávé, egyéb italok	Koffein	LOD=0.02 LOQ= 0.05 µg mL <sup>-1</sup>	H. Sereshti és mts. (2014)
felületaktív anyag segítségével végzett MALDI-TOF-MS	Energiaital	Koffein, riboflavin, nicotinamide, pyridoxine	RSD < 20 %	D. C. Grant és mts. (2008)
Rövid kapillárisal végzett elektroforézis csatlakozás nélküli vezetőképességi detektorral	Energiaital	Cukor tartalom: szacharóz, glükóz, fruktóz	LOD=15 LOQ=52 mg L <sup>-1</sup> szacharózra	B. Vochyánová és mts. (2012)
Folytonos szilárd fázisú extrakció + UV-Vis és ELSD detektorok	Üdítőitalok	Összes cukortartalom, IV. osztályú karamel, koffein	RSD=2.6 % szacharózra RSD=4 % koffeinre	R. Lucena és mts. (2005)
Planáris kromatográfia – többszörös detektálással	Energiaital	Riboflavin, piridoxin, nikotinamid, koffein, taurin	RSD=0.8-1.5 %	M. Aranda és mts. (2006)
Kapilláris elektroforézis (CE)	Energiaital	Koffein, C vitamin, PP és B6	Relatív hiba: 1.45 - 2.65 %	V. V. Khasanov és mts. (2013)

Habár a **2. táblázatban** felsorolt módszerek sikerrel alkalmazhatóak, legtöbbjük hosszadalmas és költséges eljárás. Ez főként az elhagyhatatlan mintaelőkészítési lépések

valamint az eljárások során alkalmazott oldószerek, vegyületek miatt van. Az FT-NIR spektroszkópia ezzel szemben kevésbé költséges és ugyanakkor gyors eljárást tud biztosítani. A módszer könnyen használható és a legtöbb esetben nem, vagy csak minimális mintaelőkészítést igényel. FT-NIR spektroszkópiával történő koffein tartalom meghatározásra ugyan található példa az irodalomban, de csak kávé mintákra (Huck et al., 2005; Zhang et al., 2013).

### 2.3 Az antioxidáns kapacitás

Az antioxidánsok és ezzel összhangban az élelmiszerek, zöldségek, gyümölcsök antioxidáns kapacitásának vizsgálata rendkívül népszerűvé vált az utóbbi években, köszönhetően az egészségtudatos életmód elterjedésének. Bár a meghatározási módszerek száma rohamosan növekszik, jelenleg nem képesek még az antioxidáns kapacitás pontos meghatározására *in vivo* környezetben (Cornelli, 2009).

A sejtekben különféle természetes biokémiai folyamatokban szabadgyökök keletkezhetnek. A szabad gyökök egy része ugyanakkor szükséges a szervezet megfelelő működéséhez, feleslegük viszont károsíthatja a fehérjéket, nukleinsavakat, zsírokat, és szénhidrátokat (Cadenas, 1989; Djuric et al., 1998). A szabadgyök képződés különböző belső és külső folyamatok hatására történik, amelyben szerepet játszhat a dohányzás, stressz, alkohol fogyasztás stb. Az antioxidánsok képesek a szabadgyökök koncentrációját csökkenteni, vagy gátolni a képződésüket. Másrészt pedig az antioxidánsok részei a szervezet integrált védő rendszerének (Benzie, 2000; Halliwell és Gutteridge, 1995). Néhány jól ismert antioxidáns például az E vitamin, a C-vitamin, a polifenolok és a karotinoidok.

Az antioxidánsokkal kapcsolatos kutatások, az antioxidánsok szabadgyökfogó képességének azaz antioxidáns kapacitásnak meghatározása rendkívül népszerű téma a szakirodalomban. Az antioxidáns kapacitást meghatározó módszerek száma száz fölött van (Cornelli, 2009). Az ok, amiért ennyi módszer létezik visszavezethető arra a tényre, hogy egyik módszer sem képes a természetes *in vivo* folyamatokat pontosan mérni és modellezni. Minden módszer néhány antioxidáns komponensre és reakcióra szelektív, de egyik sem képes mérni az összantioxidáns kapacitást kizárólagosan, ill. teljes pontossággal. Néhány példát az alkalmazott módszerek változatosságát és sokszínűségét különféle mintákon is bemutatva a **3. táblázatban** foglaltam össze.

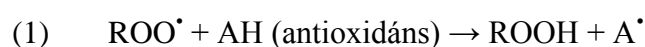


### 3. táblázat: Az antioxidánsok vizsgálatának szemléletes példái

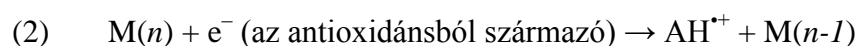
(a rövidítések a dolgozat elején találhatóak)

Minta/vegyület	Alkalmazott módszerek	Szerző
Trópusi gyümölcsök	TEAC (ABTS, DPPH)	Moo-Huchin és mts. (2014)
Vadon termő gombák	QCAR (mennyiségi összetétel-aktivitás összefüggés)	Froufe és mts. (2011)
Változatos biológiai minták	HPLC, HPTLC, UHPLC, stb.	Cervinkova és mts. (2016) (review)
Természetes vegyületek, új potenciális antioxidáns vegyületek	QSAR (mennyiségi szerkezet-hatás összefüggés)	Martinčič és mts. (2015)
Kumarin származékok	DPPH	Kotali és mts. (2016)
Lenmag kivonat	TPC	Slavova-Kazakova és mts. (2015)
Kínai rizsbor	TPC, TAC, FT-IR, Raman spektroszkópia	Wu és mts. (2015)
Áfonya minták	FRAP, ORAC, TPC, TAC, stb.	Kraujalyte és mts. (2015)
Citrusfélék	DPPH, TBARS, stb.	Cardeñosa és mts. (2015)
Mediterrán gyógynövény	DPPH	Boudkhili és mts. (2015)

A különböző antioxidáns kapacitás meghatározási módszerek többféle módon is csoportosíthatóak. Attól függően, hogy milyen típusú reakciót hangsúlyozunk, két főbb csoportra bonthatóak: hidrogén atom átvitelrel megvalósuló (HAT) vagy elektronátviteli (ET) reakciók. A HAT módszerek (pl. ACL, ACW, vagy ORAC) a szabadgyökök elfogásához köthetőek. A legtöbb ilyen módszer versengő (kompetitív) reakció sémát követ, ahol az antioxidáns és a szubsztrát vegyület versenyez a termikusan generált peroxil gyökökért. A reakcióegyenlet a következőképp írható fel (Huang et al., 2005):



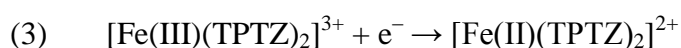
Ezzel szemben az elektronátviteli reakciók közé sorolt módszerek (pl. FRAP, TPC, TEAC) az oxidálószer redukálása közben mérik az antioxidáns kapacitást. Az oxidálószernek különböző színe van a redukált és oxidált formájában. A színváltozás milyenségéből és annak mértékéből következtetni lehet az antioxidáns kapacitásra (Apak et al., 2007). A reakció mechanizmusa a következő (Huang et al., 2005):



Minden antioxidáns kapacitás meghatározási módszer ugyanazt a tulajdonságot, a szabadgyökfogó képességet méri. Mégis, a különböző módszerek különböző eredményeket adnak. Ez annak köszönhető, hogy másféle modell vegyületeket alkalmaznak, másféle

mellékreakciók játszódnak le és egyéb zavaró tényezők is bekövetkezhetnek: a mátrixok is rendszerint különbözők. A lehetőségek száma szinte végtelen, vannak olyan kutatók is, akik megrögzötten egyfajta reakciót és módszert pártolnak a szisztematikus és véletlen hibák figyelembe vétele nélkül. Így egy átlagosnak mondható, kevésbé torzított eredményt produkáló módszer kiválasztása a sok közül egy fontos célként fogalmazható meg. A következőkben néhány, a doktori munkám során végzett statisztikai elemzés alapjául szolgáló módszert részletesebben is bemutatok:

– FRAP: A módszer alapja, hogy az antioxidánsok a minta  $\text{Fe}^{3+}$  ionjaikat  $\text{Fe}^{2+}$  ionokká képesek redukálni. A technika kifejlesztése Benzie és Strain (1996) nevéhez köthető. A reakció során a  $\text{Fe}^{3+}$ -TPTZ (tripridil-triazin) komplex redukálódik  $\text{Fe}^{2+}$ -TPTZ komplexszé az antioxidánsok hatására savas körülmények között (pH=3,6):



A redukált komplex kék színű, így a reakció spektrofotometriásan könnyen nyomon követhető 593 nm-en (Huang et al., 2005).

– TPC: Az összes polifenol tartalom meghatározása Folin-Ciocalteu reagenssel történik a Singleton és Rossi (1965; Singleton et al., 1998) által kidolgozott módszer alapján. Az alkalmazott reagensben volfrám és molibdén oxidok vannak, amelyek sárga színt adnak az oldatnak, míg a fénoxidok redukációjában az oldat színe kékre változik. Az elektronátviteli reakcióban a molibdén(VI) molibdén(V)-té redukálódik. A reakció 760 nm-en spektrofotometriásan detektálható. A nevével ellentétben a módszer nem szelektív a polifenol komponensekre nézve (Apak et al., 2007).

– TEAC: A troloxra vonatkoztatott antioxidáns kapacitás módszerének kidolgozása Miller és munkatársai (1993) nevéhez fűződik. A kulcs reakcióban az antioxidáns redukálja a 2,2-azino-bisz-(3-etilbenzotiazolin)-6-szulfonsavat (ABTS szabad gyök). Az  $\text{ABTS}^{+\cdot}$  gyök kation sötétzöld színű, de ha antioxidáns van a reakcióközegben, a gyök kation  $\text{ABTS}^{2-}$ -ra változik és elveszti a színét. A reakció spektrofotometriásan követhető 734 nm-en.

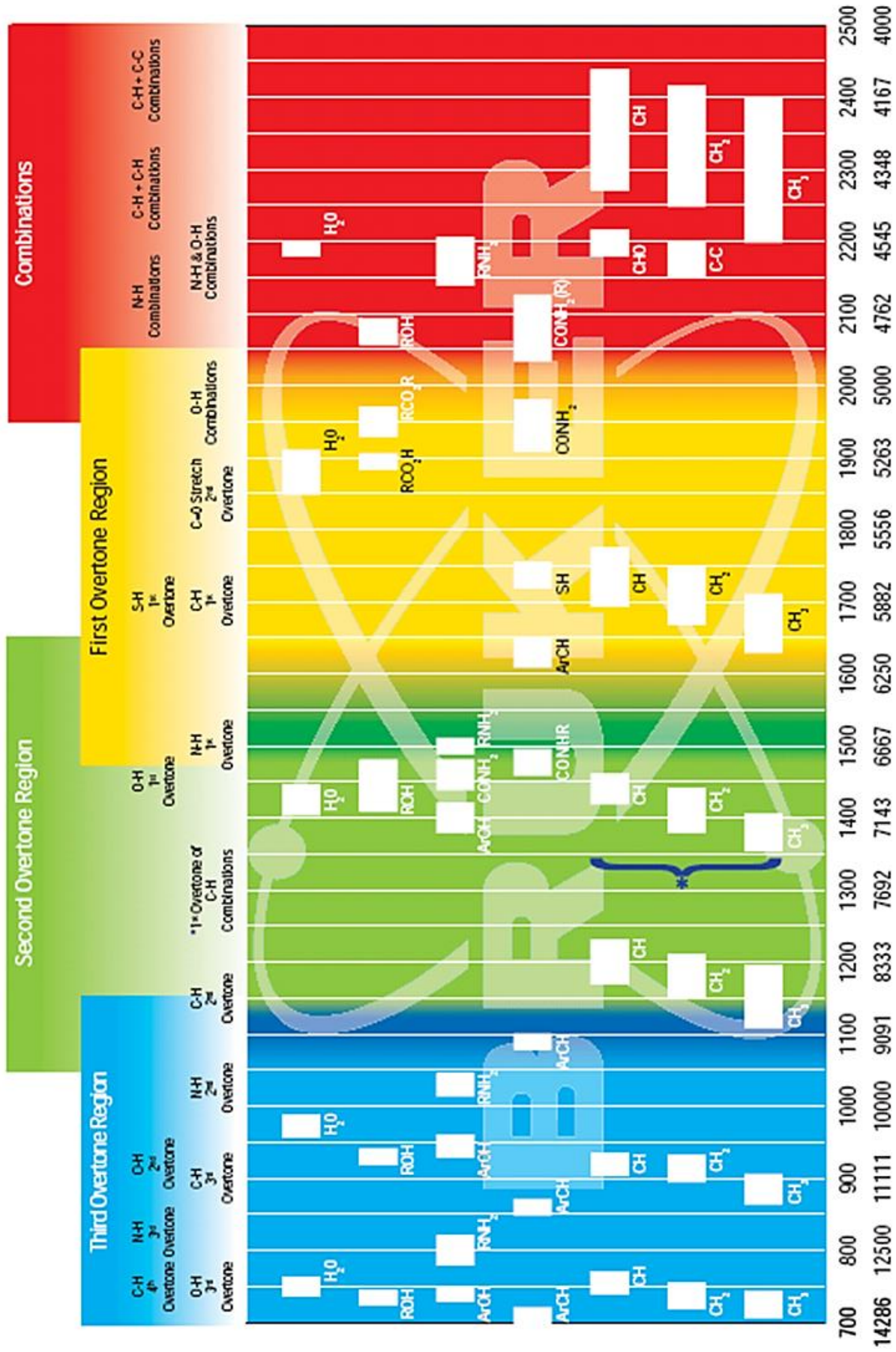
– DPPH és TRSC: A DPPH módszer (Blois, 1958) egyike a legrégebben fejlesztetteknek és 2,2-difenil-1-pikril-hidrazilt használ stabil gyökként (kereskedelemben is kapható). A redukció során az oldat színe eltűnik (a DPPH gyök eredetileg liláspiros színű). A reakció 734 nm-en detektálható spektrofotometriásan. A TRSC módszer (Blázovics et al., 1999) a  $\text{H}_2\text{O}_2/\text{OH}$  mikroperoxidáz-luminol rendszer gátlásán alapszik. Lúgos közegben a rendszer fényt emittál és

a vas komplex (mikroperoxidáz) hatására  $\text{OH}^\bullet$  gyök keletkezik a hidrogénperoxidból Fenton típusú reakcióban. A reakció spektrofotometriásan 420 nm-en detektálható.

– ACL és ACW: A víz- és zsírolható antioxidáns kapacitás meghatározása Popov és Lewin (1994, 1996) nevéhez köthető. A szuperoxid anion szabadgyök ( $\text{O}_2^{\bullet-}$  fotokémiaiag jön létre és detektálásra kemilumineszcenciát alkalmaznak. A szuperoxid anion szabadgyök reakcióba tud lépni a mintában lévő antioxidáns vegyületekkel, miközben a mennyisége csökken. A luminol a maradék szuperoxid gyök anion által aktiválódik és lumineszcenciát bocsát ki. A reakciók és módszerek megfelelőek a szabadgyökfogó képesség leírására. Az eredmények Troloxra vonatkoztatva adják meg a zsírolható, és C vitaminra a vízoldható antioxidáns kapacitás mérése esetén.

## 2.4 Az FT-NIR spektroszkópia elméleti háttere

A Fourier-transzformált közeli-infravörös spektroszkópia egyike az élelmiszer-analitikában leggyakrabban alkalmazott roncsolásmentes gyorsanalitikai és molekuláspektroszkópiai eljárásoknak. Míg az infravörös spektroszkópia az anyag és az elektromágneses sugárzás kölcsönhatását vizsgálja 780 nm és 300  $\mu\text{m}$  között, a közeli infravörös tartomány ennek az első harmada, 780 nm és 2500 nm között ( $12820\text{ cm}^{-1}$  és  $4000\text{ cm}^{-1}$  között) értékelhető. A spektrum folytonos, amely annak köszönhető, hogy a molekuláknak az infravörös tartománynak megfelelő gerjesztési energia esetén többféle rezgési és forgási átmenet gerjesztődik, miközben az átmenetekre jellemző frekvenciájú fotonok elnyelődnek. (Pokol et al., 2011). Aktív átmenetnek akkor tekintünk egy rezgést, ha a molekula dipólusmomentuma megváltozik a rezgés által. A közeli infravörös spektrumban jellemzően a vegyértékrezgések (a rezgések a kötések hossz tengelye mentén jönnek létre, és szimmetrikusak, ill. aszimmetrikusak) és deformációs rezgések (a vegyértékszög változik) felharmonikusai (felhangsáv) és kombinációs rezgései láthatóak (Subramanian és Rodriguez-Saona, 2009). A felharmonikusok, a normál rezgések egész számú többszöröse, míg a kombinációs rezgések során a normál rezgések csatolódnak és azok összege vagy különbsége jelenhet meg. Az jellegzetes rezgések összefoglalását az **1. ábrán** szemléltetem. Az FT-NIR spektroszkópiában jellemzően az infravörös fény transzmisszióját, reflexióját vagy transzreflexióját (az előbbi kettő kombinációját) detektálják.



### 1. ábra: A jellegzetes NIR rezgések összefoglalója

*A harmadik (kék színnel), második (zöld színnel) majd első (sárga színnel) felharmonikusok sávjai, végül pedig a kombinációs rezgések (piros színnel). Az alsó két tengely felülről lefelé haladva rendre a hullámhosszakat (nm-ben megadva) valamint a hullámszámokat ( $\text{cm}^{-1}$ -ben megadva) mutatja. Kék csillaggal jelölve a C-H kombinációs rezgések első felharmónikusa. (A kép forrása: (Bruker Optik GmbH))*

A NIR alkalmazásának előnyei, hogy nem csak a mintában lévő fő alkotóelemek, de a kis koncentrációjú komponensek mennyiségének illetve fizikai tulajdonságaik becslésére is használható. Emellett rendkívül gyors, minimális vegyszerigényű, a folyamatos elemzés is kivitelezhető és a méréshez felhasználandó mintamennyiség is nagyon kicsi. A módszer hátránya, hogy mindenképpen szükséges a mérésekhez egy kalibrációs modell megalkotása, ami azt jelenti, hogy szükség van referencia adatsorra a kalibráció végrehajtásához. Ettől eltekintve előszeretettel alkalmazzák a klasszikusnak mondható módszerek kiváltására az FT-NIR készülékeket számos területen, például mezőgazdaságban, élelmiszeriparban, gyógyszeriparban, stb. Bár az élelmiszertudományban nagy hagyománnyal rendelkező technika, az utóbbi évtizedekben már a gyógyszergyártás is „felfedezte” a módszerben rejlő lehetőségeket. Manapság egyre inkább elterjedté vált az on-line NIR technika (Huang et al., 2008a), használata során már a gyártósoron – vagy akár a búzamezőkön „on combine” módon – megvalósulhat a valós idejű ellenőrzés. NIR vizsgálatokat végezhetünk kiindulás, közti- és végtermékekre is. A vizsgálatok célja általában a minőségellenőrzés, a minták beltartalmi paramétereinek vizsgálata, esetleges szennyeződések kimutatása. De használható eredetvizsgálatra, a minták osztályozására is. A **4. táblázatban** összefoglaltam néhány szemléletes példát arra vonatkozóan, hogy milyen szerteágazó a szakirodalom az alkalmazások területén. Az on/in-line mérési módok mellett nagy hangsúlyt fektetnek manapság a hiperspektrális képfeldolgozásra is. Ebben az esetben a NIR spektroszkópiát digitális képkalkotó eljárásokkal ötvözik, amely során egy „spektrális hiperkockát” kapunk. A kocka első két dimenziói a képpontok geometriai információi, a harmadik dimenzió pedig maga a pontokhoz tartozó NIR spektrum. Ezekhez a mérésekhez minden esetben speciális szoftverekre van szükség (Balázs et al., 2011). Hasonlóképp elterjedt a hordozható NIR készülékek használata, hiszen így a labortól távol, a „terepen” is végre lehet hajtani a méréseket gyorsan és egyszerűen. Hátrányos tulajdonságuknak szokták tekinteni azonban azt, hogy ezeknek a készülékeknek a hordozhatósággal csökken az érzékenységük. Manapság viszont egyre több cikk olvasható arról, hogy valójában az asztali készülékekkel egyenértékű eredményt adnak (Lopo et al., 2016).

**4. táblázat:** FT-NIR spektroszkópia alkalmazási lehetőségeinek bemutatása néhány szemléletes (elsősorban élelmiszertudományi) példán keresztül.

Minta	Komponens	Cél	Eszköz	Szerző
Ehető olajok, zsírok	-	Mintázatfelismerés	FT-IR, FT-NIR, FT-Raman	Yang és mts. (2005)
Almalé	-	Mintázatfelismerés	NIR transzfelxiós üzemmódban	León és mts. (2005)
Számos élelmiszerfajta	Beltartalmi paraméterek	Mintázatfelismerés és regressziós modellépítés	On/in-line NIR	Huang és mts. (2008b) review
Sertéshús	Zsír	Regressziós modellépítés	NIR hiperspektrális képfeldolgozás	Huang és mts. (2016)
Kínai rizs bor	Antioxidáns kapacitás, $\gamma$ -aminovajsav	Regressziós modellépítés	FT-NIR	Wu és mts. (2015)
Gyógyszeripai alkalmazások	Hatóanyag tartalom, Nedvességtartalom stb.	Mintázatfelismerés és regressziós modellépítés	FT-NIR, on-line NIR	Roggo és mts. (2007) review
Méz	Antioxidáns tartalom	Regressziós modellépítés	NIR	Escuredo és mts. (2013)
Méz	Antioxidáns tulajdonságok	Mintázatfelismerés és regressziós modellépítés	NIR	Tahir és mts. (2016)
Szőlő	-	Mintázatfelismerés	Hordozható NIR készülék	Gutiérrez és mts. (2015)
Rizs	Fehérje, amilóz tartalom	Regressziós modellépítés	NIR	Bagchi és mts. (2016)
Szárastészta	Zsír, tojástartalom	Regressziós modellépítés	FT-NIR	Fodor és mts. (2011)

Az **5. táblázatban** olyan szakirodalomban megjelent példákat gyűjtöttem össze, amelyek az általam elvégzett kísérletek szempontjából is relevánsak, hiszen akár ugyanazon vagy hasonló kémiai komponensek FT-NIR spektroszkópiás vizsgálatára irányulnak különböző élelmiszer mintamátrixokban.

**5. táblázat:** FT-NIR spektroszkópiai modellek koffeinre, cukorra és vitaminokra különböző élelmiszer mintamatrixokban

Minta	Komponens	Szerző
Körte	Cukor tartalom	Ying és Liu (2008)
Pékáru	Cukor tartalom	Szigedi és mtsi. (2011)
Kávé	Koffein, teobromin, teofilin	Huck és mtsi. (2005)
Zöld tea por és granulátum	Koffein tartalom	Sinija és Mishra (2009)
Gyümölcslé	Cukor tartalom	Rodriguez-Saona és mtsi. (2001)
Lucerna, növényi olajok	E vitamin	Cozzolino (2009)
Cukorrépa	Cukor tartalom	Salgó és mtsi. (1998)
Bor	Cukor tartalom	Kaffka és Norris (1976)
Narancs	C vitamin	Magwaza és mtsi. (2013)

## 2.5 Az FT-NIR spektroszkópiai modellek

Az FT-NIR spektroszkópiára jellemző, hogy – ellentétben az IR spektroszkópiával – önmagában a spektrumok nem értékelhetők ki megfelelően, a kalibrációs és osztályozási modellek építéséhez is minden esetben valamilyen kemometriai módszert kell segítségül hívni. Erre számos lehetőség áll a rendelkezésre, így a megfelelő referencia adatok és spektrumok birtokában változatos kiértékelési lehetőségekkel kereshetjük meg a referencia és a spektrumok közötti összefüggéseket. A regressziós modellek mellett viszont elterjedtek a mintázatfelismerést szolgáló kemometriai módszerek is, amelyek nemcsak a kiugró értékek kiszűrésére, de különböző csoportosításokra, osztályozásokra is kiválóan alkalmazhatóak. Az említett módszereket részletesen tárgyalom az „Anyag és módszer” című fejezetben.

### 2.5.1 A spektrumok transzformációja

A NIR mérések során kapott spektrumokat, készülék típustól függően is, transzformálnunk kell a megfelelő kiértékeléshez. Gyakran adódik olyan eset, hogy egy osztályozási modell létrehozásához nem szükséges például a főkomponens-elemzés során használt standardizáláson kívül más komolyabb transzformációkat bevetnünk. Viszont a regressziós modellek esetén szinte kivétel nélkül szükségünk van a spektrum transzformációjára. Az erre alkalmas módszereknek szinte végtelen tárházát kínálja már a legismertebb

szoftvercsomagok is. Természetesen a legjobb módszer kiválasztása optimalás kérdése, ugyanakkor adatkészlet függő, hogy melyik technika lesz a legcélravezetőbb. Az általam is gyakran alkalmazott módszereket részletesebben is ismertetem:

- Standardizálás: Ez a legáltalánosabb adatelőkezelési módszer, nemcsak a spektrumok kiértékelésénél, de a kemometria egész területét tekintve is. Általában kombinációként alkalmazható például deriválás után. A spektroszkópia területén standard normál változóként (vagy vektor normalizációként) is ismert, viszont ebben az esetben a standardizálást soronként kell az adatkészletben végezni.

- Arányos (többszörös) szóródás korrekció (MSC): Eredetileg a fényszóródás hatásainak kompenzálására lett kifejlesztve (Geladi et al., 1985). A transzformáció az alapvonal eltolódás és az arányos (többszörös) szóródás hatását veszi figyelembe. Szintén gyakran kombinálják más transzformációs technikákkal is.

- Deriválás: A deriválás során a spektrumokat jelgazdagabbá tudjuk tenni, illetve kiemelhetünk vele spektrumrészleteket és absztrpciós csúcsokat. Hátrányként megemlíthető, hogy a csúcsok kiemelésével együtt a zajt is megnöveljük. Az első derivált a leggyakrabban alkalmazott fajtája, amelynek számolása során a lokális maximumok, ill. minimumok az eredeti spektrum csúcsinflexiós pontjainak felelnek meg. Az eredeti spektrum lokális maximum és minimum értékei a derivált spektrumban a nulla értékét veszik fel. Az alkalmazása előnyös lehet, ha a függőleges alapvonal eltolódást akarjuk kiküszöbölni, de veszélyes is lehet, ha pont az alapvonal-eltolódás hordozza a minta spektrumában a variancia egyik fő összetevőjét. A második derivált alkalmazásakor az eredeti spektrumban lévő csúcsok völgyként jelentkeznek. Leginkább a diszperziós NIR készülékek esetében érdemes használni, az FT-NIR készülékek esetén saját tapasztalataim alapján nem javított a modellek optimalásakor. Ahogy az első derivált, a második is, de még intenzívebben, felerősíti a spektrumban lévő zajt. A deriválás során felerősödött zaj miatt az eljárást célszerű kombinálni akár a standardizálás, akár az MSC technikákkal is.

## **2.5.2 Modellek validálása**

A NIR spektrumok regressziós modelljeinek megalkotásakor vizsgálandó a modellek jósága és előrebecslési képessége is. Hogyan dönthetjük el, hogy egy modell megfelelően érvényesített (validált) vagy egyáltalán alkalmazható-e a kívánt célra?



A regressziós modellek validálása többféleképpen is történhet. A legelterjedtebben alkalmazott módszernek a kereszt-ellenőrzés tekinthető. Ekkor a mintáinkat több alcsoportra osztjuk fel, amelyek közül egyet kiválasztunk és félre teszünk addig, amíg a maradék mintákkal megalkotjuk a kalibrációs modellt. Ezt követően a félretett mintacsoportra végezzük el a modell előrebecslését. Ezt a folyamatot meg kell ismételni az összes mintacsoporttal, majd az eredményeket összegezni. A kereszt-ellenőrzés két fontos paramétere az alcsoportok mennyisége és a kiválasztásának szabálya. Az előbbi paraméter meghatározására kiváló lehetőség a Hastie és munkatársai (2001) által írt könyvben szereplő ajánlás, miszerint a felosztásnak érdemes  $n=5$  és  $10$  közé esnie. Ezzel összhangban nagyon elterjedt a hét részre osztott kereszt-ellenőrzés alkalmazása. Az alcsoportok kiválasztására szintén számos lehetőség létezik. A minták besorolása történhet egyszerűen, csak sorrendben meghúzva a határokat, de ennél mindenképp kedvezőbbnek mondható, ha szisztematikusan vagy éppen véletlenszerűen (randomizáltan, többszöri ismétléssel) végezzük el az alcsoportba sorolásokat (Bro et al., 2008). A kereszt-ellenőrzés speciális esete, ha  $n=K$ , vagyis a csoportok száma ( $K$ ) megegyezik a mintaszámmal ( $n$ ). Ezt egy elem kihagyásos kereszt-ellenőrzésnek, vagy az angolból átvett „leave-one-out”, LOO kereszt-ellenőrzésnek nevezzük. Ez a verzió nagyon elterjedt, egyes tudományterületeken standard technika, pedig nagyobb véletlen hibával rendelkezik a több részre osztott kereszt-ellenőrzéshez képest és a legtöbb esetben optimistább becslést ad a modellünk jóságára. A belső validálás mellett fontos hangsúlyozni a külső, teszt mintákon történő validálást is, de önmagában, a kalibrációs modell belső kereszt-ellenőrzése alapján kapott eredmények nélkül nem szabad következtetéseket levonni belőle a modell jóságát illetően (Gütlein et al., 2013). A teszt validálásnak egy keverék változata, mikor nem külső, hanem belső, már az adatkészletben meglévő mintákat jelöljük ki az ellenőrző csoportba. Ez a belső validálás annyiban különbözik a kereszt-ellenőrzéstől, hogy nincs „keresztelés”, a mintákat csak egyszer használjuk fel, vagy kalibrációs, vagy ellenőrző mintaként. Megemlítendő a véletlenszerűségi (randomizációs) teszt is, amely nem csak a regressziós, de az osztályozási modellek esetén is rendkívül hasznos. A randomizációs teszt alapja, hogy vagy az  $\mathbf{X}$  vagy az  $\mathbf{Y}$  változóinkat összekeverjük az adatkészletben, és az így kapott modell teljesítmény paramétereit kiszámoljuk. Ha a modell jósága nem romlik el az eredetihez képest jelentős (szignifikáns) mértékben, akkor az eredeti modell hibásnak tekinthető (elterjedt kifejezéssel: „nincs validálva”).

### **2.5.3 A modelleket leíró teljesítmény paraméterek**

A kalibrációs modellek jóságát leíró teljesítmény paraméterek szakirodalmában is meglehetősen nagy. A legelterjedtebben alkalmazott paraméterek mellett időről-időre megjelennek az újabb fejlesztések is. Vitatott téma az ezzel foglalkozó kutatók között, hogy

mely paraméterek írják le legjobban a modellünk becslési képességét, avagy melyeknek lehet hinni a legjobban. A spektroszkópiai kiértékelésekben (és egyéb regressziós modellépítések során is) a determinációs koefficiens négyzete (magyarázott varianca), avagy a Pearson féle korrelációs koefficiens ( $R$ ) négyzete stabil helyet foglal el: a kemometriai kiértékelő szoftverekbe épített első számú teljesítmény jellemzőként használják. Az  $R^2$  érték mindig 0 és 1 között változik (de százalékban kifejezve is használják). Annál jobb a modell, minél jobban közelítik az 1-st. Az definíciója a következőképp írható fel:

$$(4) \quad R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\text{RSS}}{\text{TSS}},$$

ahol  $y_i$  a referencia érték,  $\hat{y}_i$  a becslt vagy számított érték,  $\bar{y}$  pedig a referencia értékek átlaga. Az RSS a reziduális eltérés négyzetösszeg a TSS pedig amely a teljes eltérés négyzetösszeg (Tóth et al., 2013). Az  $R^2$  szabadsági fokokkal korrigált verziója ( $R_{adj}^2$ ) a következő egyenlet szerint számolható:

$$(5) \quad R_{adj}^2 = R^2 - (1 - R^2) \times \frac{p}{n-p-1},$$

ahol a  $p$  a modellben lévő változók mennyiségét, az  $n$  pedig a mintaszámot jelenti. A determinációs koefficiens kiszámolható a kereszt-ellenőrzött modellre is, ilyenkor  $Q^2$ -tel jelöljük. Minden kereszt-ellenőrzés vagy egyéb validálási esetben kiszámolható, az egyenlet alap koncepciója ugyanaz marad. Több elem kihagyásos (több részre osztott) kereszt-ellenőrzésként ( $Q_{LMO}^2$ ) az egyenlet a következő:

$$(6) \quad Q_{LMO}^2 = 1 - \frac{\sum_{j=1}^m \sum_{i=1}^n (y_i - \hat{y}_{i/j})^2}{\sum_{i=1}^n (y_i - \bar{y})^2},$$

ahol  $y_i$  a referencia érték,  $\bar{y}$  a referencia értékek átlaga,  $\hat{y}_{i/j}$  pedig becslt érték az  $i$ -dik mintára, ha a  $j$ -dik alcsoport ki van hagyva a kalibrációs modellépítésből. A teljes adatkészlet  $m$  részre lett osztva. Hasonlóképp elterjedten alkalmazott paraméter az átlagos négyzetes hiba, amelyet RMSE rövidítéssel jelölünk. Ki lehet számolni a kalibrációs, a validálási és a teszt adatkészletre is, ilyenkor rendre RMSEC (calibration), RMSECV (cross-validation) és RMSEP (prediction) jelölést kapnak. Az RMSE érték általánosan a következő egyenlettel írható fel (Naes et al., 2002):

$$(7) \quad \text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}},$$

ahol  $y_i$  a referencia érték,  $\hat{y}_i$  a becslt vagy számított érték,  $n$  pedig a mintaszám. A különböző adatkészletekre történő számolás (RMSEP, RMSECV, RMSEC) leginkább csak a becslt értékben és a nevezőben tér el egymástól. Az RMSE értékek mértékegysége megegyezik a referencia adat mértékegységével, és természetesen minél kisebb annál jobb a modellünk. A spektroszkópiában szintén gyakran használt teljesítmény paraméter az RPD érték, amely a relatív korrigált tapasztalati szórás (Williams, 2001). A számolását tekintve a referencia értékek szórásának és a becslés vagy kereszt-ellenőrzés standard hibájának a hányadosa:

$$(8) \quad \text{RPD} = \text{SD}/\text{SEP},$$

$$(9) \quad \text{SD} = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}},$$

ahol  $y_i$  a referencia érték,  $\bar{y}$  a referencia értékek átlaga,  $n$  pedig a mintaszám. SD a korrigált empirikus szórás, a SEP a validálás torzítással korrigált becslési hibája:

$$(10) \quad \text{SEP} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i - \text{Bias})^2}{n-1}},$$

ahol  $\hat{y}_i$  a becslt vagy számított érték,  $y_i$  a referencia érték,  $n$  pedig a mintaszám. A Bias jelentése magyarul torzítás, amely számolása a következő egyenlet szerint történik (Naes et al., 2002):

$$(11) \quad \text{Bias} = \frac{\sum_{i=1}^n (\hat{y}_i - y_i)}{n},$$

ahol  $\hat{y}_i$  a becslt vagy számított érték,  $y_i$  a referencia érték,  $n$  pedig a mintaszám. Az eddigiekben ismertetett teljesítmény paraméterek csak a leggyakrabban alkalmazottak és ezeknek is van még számos egyéb változata, mint például az  $R^2$  és  $Q^2$  megvalósítása Y-randomizációs tesztekkel. További néhány, a regressziós modellekre jellemző teljesítmény-paraméter rövidítése csak felsorolásként szerepel itt: PRESS érték (előrebecslt hiba négyzetösszege),  $F$  (Fisher érték), CCC (egyezési együttható, coefficient of concordance) (Lin, 1989), MAE (átlagos abszolút hiba), LOF (Friedman-féle illeszkedés hiánya, lack of fit, kritérium) (Friedman, 1991a),  $K_x$  (változók közötti interkorreláció, a PCA alapján) (Todeschini et al., 1999). További ismert teljesítmény paraméterek és azok leírása az „Eredmények” fejezetben és a „Mellékletek” **M1** táblázatában találhatóak meg.

### 3. CÉLKITŰZÉS

A doktori munkám során a következő célokat fogalmaztam meg az egyes fejezetekhez tartozóan:

Q10 koenzim tartalmú étrendkiegészítők:

- A Q10 koenzim tartalmú étrendkiegészítők hatóanyagtartalmának vizsgálata FT-NIR spektroszkópiával, kalibrációs modell készítése és validálása kereszt-ellenőrzéssel és külső teszt mintákkal. Az eddig elterjedt kromatográfias módszer kiváltása FT-NIR spektroszkópia segítségével.
- A különböző változókiválasztási módszerrel létrehozott modellek összehasonlítása a rangszám különbségek összegének módszerével (SRD)

Energiaitalok vizsgálata:

- A koffein tartalom meghatározása HPLC-vel, mint referencia módszerrel. A kromatográfias módszer energiaital mintákra történő optimalizálása. Ezt követően FT-NIR módszerrel a regressziós modell megalkotása, valamint belső és külső validálása.
- Cukortartalom meghatározása Schoorl referencia módszerrel, majd FT-NIR regressziós modellépítés a Schoorl referencia módszerrel kapott, illetve a nominális (dobozon feltüntetett) értékek használatával. A létrehozott modellek belső és külső validálása új teszt mintákkal is.
- Az energiaitalok osztályozása az FT-NIR spektrumaik alapján. Elsőként cukortartalom szerint, majd a taurinos, arginines és normál (taurint és arginint nem tartalmazó) minták elkülönítése. Az osztályozás elvégzése többféle mintázatfelismerő módszerrel is.

Antioxidánsok vizsgálata:

- Már meglévő antioxidáns kapacitás adatok alapján, két adatkészlet (esettanulmány) felhasználásával, összehasonlító elemzés készítése a meghatározási módszerekről: Az antioxidáns kapacitás módszerek osztályozása és rangsorolása.

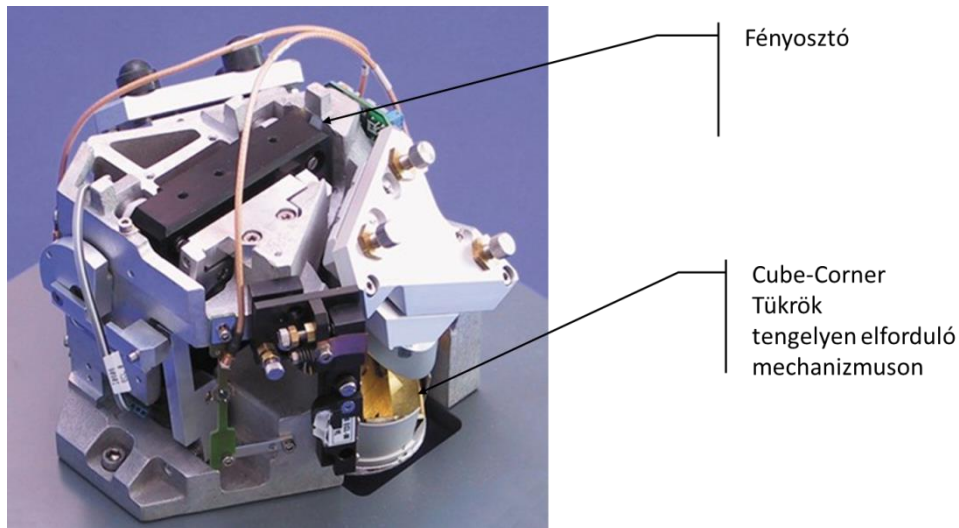
## Kemometriai módszerfejlesztések:

- Új változóselektációs módszer létrehozása és tesztelése a Q10 koenzim tartalmú étrendkiegészítők FT-NIR spektrumának felhasználásával. A modellek jóságának növelése a közel 2300 változószám csökkentésével.
- Olyan, „*n*-class” vevő-működtető görbék (ROC görbék) megalkotása, melyekkel kettőnél több osztályos elemzések is kiértékelhetővé válnak. A módszer megfelelő validálása és tesztelése az energiatalok osztályozása során. A mintázatfelismerési módszerek összehasonlítása az „*n*-class” ROC görbék segítségével.
- A véletlen fák (random forest) használata során történő optimalizálási lépés kidolgozása a minél jobb osztályozási modellek létrehozására.
- Regressziós modellek teljesítmény paramétereinek összehasonlítása és rangsorolása. Mennyiségi szerkezet-hatás összefüggések segítségével kapott regressziós modellek adatkészletei alapján, hogy szélesebb körű és általánosabban megfogalmazható eredményt kapjunk.

## 4. ANYAG ÉS MÓDSZER

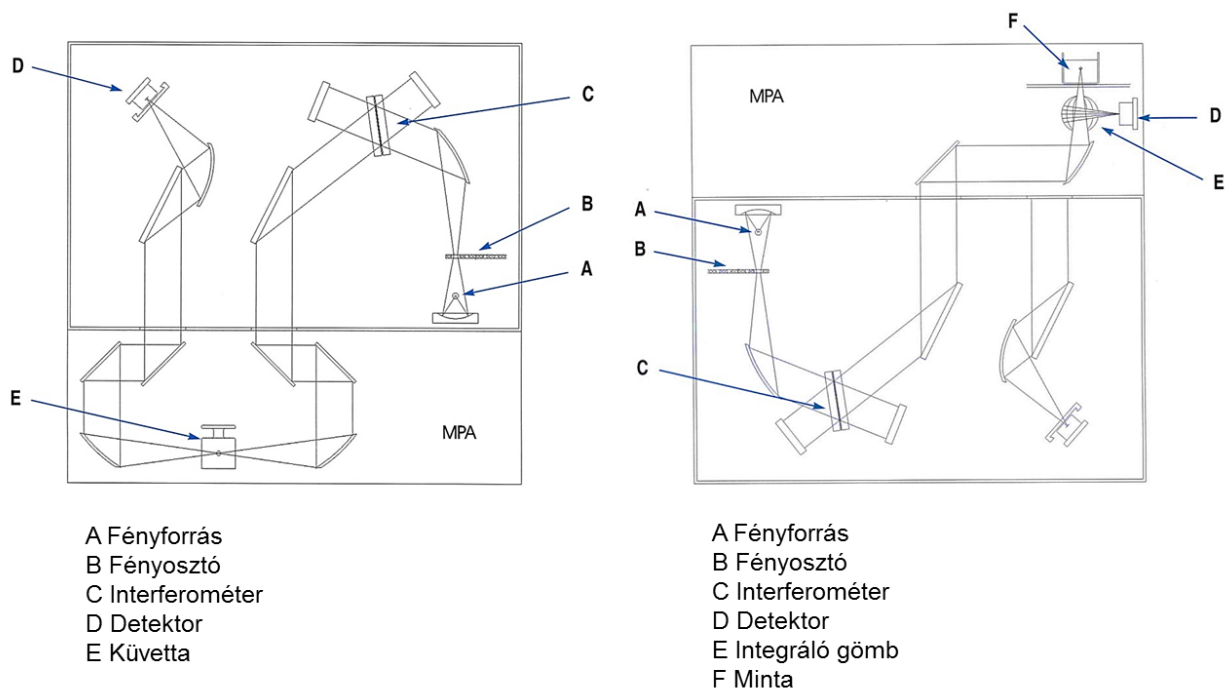
### 4.1 FT-NIR készülék felépítése

Az általam használt Bruker MPA NIR/NIT készülékbe (Bruker Optik GmbH, Ettlingen, Németország) kvarc fényosztót és Rocksolid™ interferométert építettek be, amely a Michelson típusú interferométerhez képest nem síktükröket, hanem ún. Cube Corner tükröket alkalmaz. A **2. ábra** a készülékbe épített Rocksolid interferométert mutatja be.



**2. ábra:** A Rocksolid interferométer. A kép forrása: (Internet, 1)

Ezeknek az az előnye az előbbivel szemben, hogy biztosítani képesek a fénysugarak párhuzamosságát. A készüléket a doktori munkám során diffúz-reflexiós és transzmissziós mérési módban is alkalmaztam. Míg a Q10 koenzim tartalmú étrendkiegészítők vizsgálatánál a szilárd mintákhoz megfelelően a forgó mintatartóval ellátott reflexiós egységet, az energiailokhoz az átfolyó küvettás (1 mm rétegvastagság) termosztálható transzmissziós egységet használtam. A készülék mindazonáltal szilárd, folyékony és kolloid minták vizsgálatára is képes. A Bruker MPA NIR/NIT spektrofotométer felépítése az átfolyó küvettás transzmissziós valamint a reflexiós rendszer szerinti beállításban a **3. ábrán** látható.



**3. ábra:** Az általam használt FT-NIR készülék átfolyó kuvettás (bal oldal) és reflexiós (jobb oldal) rendszere. A kép forrása: (Bruker Optik GmbH, 2003)

Diffúz-reflexiós mérési módban a mintára beérkező infravörös fény egy része a minta felületén (1-2 cm mélységig behatolva) elnyelődik, egy része pedig visszaverődik, és a detektorba jut. A detektorhoz egy arannyal bevont integráló gömbön keresztül jut a fénysugár, ami a készülékgyártó saját fejlesztése. Az integráló gömb segítségével tudjuk a referencia spektrumokat is felvenni. Reflexiós mérési módban PbS detektor áll a rendelkezésünkre. A spektrumokat  $12500\text{ cm}^{-1}$  és  $3800\text{ cm}^{-1}$  közötti tartományban rögzítettem.

Transzmissziós mérési módban a fénysugár áthatol a minta teljes rétegvastagságán és így jut el a detektorhoz. Egy része ilyenkor is elnyelődik a mintában, amely így intenzitásváltozást okoz. Alapvető különbségnek mondható a két mérési mód között, hogy míg a transzmissziós esetben a minta a bejövő fénysugár és detektor között helyezkedik el, a diffúz-reflexiós elrendezés szempontjából a beérkező fénysugár és a detektor is a minta ugyanazon oldalán helyezkedik el. A transzmissziós spektrumokat eredetileg csak  $12500$  és  $9000\text{ cm}^{-1}$  között tudnánk felvenni a mintákról, mert ennél alacsonyabb hullámszámnál nem tud már áthatolni a fénysugár az alacsony energiája miatt, de ez csak a szilárd és koloid mintákra igaz. A híg oldatokon sokkal könnyebben tud áthatolni a fény, így az energiatalok spektrumát  $12500\text{ cm}^{-1}$  és  $4000\text{ cm}^{-1}$  hullámszám között is fel tudtam venni. A transzmissziós mérési módhoz Te-InGaAs (Tellúr-Indium-Gallium-Arzén) detektor van beépítve a készülékbe.

A készülék mindkét mérési módjában a spektrális felbontás  $8\text{ cm}^{-1}$  volt, a szkennelési sebesség pedig  $10\text{ kHz}$ . A készülék ezekben az esetekben 32 „alspektrumot” rögzít, és ezek átlagát kapjuk egy-egy felvétel eredményeként.

## 4.2 Nagyhatékonyságú folyadékkromatográfia

A nagyhatékonyságú folyadékkromatográfias (HPLC) referenciamérésekre, két esetben, a Q10 koenzim tartalom meghatározásánál és az energiatalok koffein tartalmának meghatározásánál volt szükségem.

A Q10 koenzim HPLC-s referenciamérései az általam elvégzett FT-NIR kísérletek előtt megtörténtek. A mérések egy Agilent 1200 HPLC készülék használatával (Agilent Technologies, Santa Clara, CA, USA) történtek izokratikus módban Agilent Zorbax XDB C18 HPLC kolonna segítségével ( $2,1\text{ mm} \times 50\text{ mm} \times 3,5\text{ }\mu\text{m}$ ) UV detektorral  $275\text{ nm}$ -en nyomon követve. A kolonnatér hőmérséklete  $30\text{ }^\circ\text{C}$  volt. Eluensként acetonitril (ACN) : tetrahidrofurán (THF) : víz  $65:30:5$  térfogat százalékos (v/v %) elegye volt használva. Az áramlási sebesség  $0,35\text{ ml / perc}$  volt, az injektálási térfogat pedig  $10\text{ }\mu\text{l}$ . A fejlesztett módszer hibája  $10\%$ , detektálási limitje pedig  $0,05\text{ mg Q10 koenzim / g}$  volt.

Az energiatalok HPLC mérései esetén kiindulási módszernek egy nemzetközi standard eljárást választottunk (ISO 20481:2008), amely kávék koffein tartalmának meghatározására szolgál. Ezt a módszert fejlesztettem tovább. A készülék hasonlóan az előző esethez az Agilent 1200 HPLC volt UV detektorral felszerelve. A szabványmódszer módosításon átesett részletei a következők voltak: kollokálként egy Agilent Zorbax XDB C18 kolonnát ( $4,6\text{ mm} \times 150\text{ mm} \times 5,0\text{ }\mu\text{m}$ ) választottam. A kromatogramok felvétele izokratikus módban  $40\text{ }^\circ\text{C}$ -on történt. Áramlási sebességnek  $1\text{ ml / perc}$ et választottam, az injektálási térfogat pedig  $20\text{ }\mu\text{l}$  volt. A mérés futási ideje  $18\text{ perc}$  volt. Az UV detektálást  $273\text{ nm}$ -en történt, de párhuzamos mérésként  $260\text{ nm}$ -en csúcstisztaság vizsgálatot is végeztem. Erre azért volt szükség, hogy kimutatható legyen, ha az adott mintában bármilyen egyéb komponens esetleg zavaró tényezőként jelentkezne hasonló retenciós időnél. A vizsgálatokhoz használt eluensek és vegyületek listája a „Mellékletek” M2-es táblázatában található meg.



### 4.3 Mintaelőkészítés

A Q10 koenzim minták HPLC méréseinek mintaelőkészítése a Vass és munkatársai által közölt publikációban (2014) részletesen ismertetett.

Az energiatalok HPLC méréseihez ultrahangos fürdő (T2MODX; VWR, Radnor, PA, USA) segítségével szén-dioxid mentesítettem a mintákat 20 percen keresztül. Ezt követően 50 µl-t hígítottam 1600 µl-re desztillált vízzel. A külső kalibrációs mintasort használtam a koffein koncentráció meghatározására, amelyben a pontok rendre a következők voltak: 2,5; 5,0; 10,0 és 20,0 ppm.

Az FT-NIR mérések esetén a Q10 koenzim tablettákat dörzsmozsár segítségével megfelelően elporítottam és homogenizáltam. Az energiatalok mérése során az ultrahangos fürdő segítségével gázmentesített mintákból 10-10 ml-t használtam fel.

### 4.4 Klasszikus mérési módszerek

#### 4.4.1 Schoorl módszer

A Schoorl módszer az élelmiszerek cukortartalmának egyik legelterjedtebb meghatározási módszere, mely során a redukáló cukrokat, illetve hidrolízis segítségével redukáló cukrokká lebomló cukrokat lehet vizsgálni. A meghatározást az OÉTI módszergyűjteményének ÉLK 4.009 szabványa alapján végeztem el. A meghatározás megtalálható szabvány (standard) módszerként a nemzetközi szakirodalomban is (Horovitz, 1975).

Az élelmiszer mintákban található szénhidrátokat tömény sósavas hidrolízissel bontjuk le (68-70 °C vízfürdő), ezt követően 33 %-os nátrium-hidroxid oldattal semlegesítjük a mintákat. A nem szénhidrát komponenseket derítés segítségével tudjuk eltávolítani, amelyhez Carrez I. és Carrez II. oldatot használtunk. A meghatározás további lépései egy titrálási folyamat részét képezik. A mintát Schoorl A ( $\text{CuSO}_4$ ) és B oldattal ( $\text{NaOH}$ ,  $\text{K}_2\text{Na}_2\text{C}_2\text{O}_4$ ) forraljuk, amely során az oldat kékes vöröses színűvé válik. A redukáló cukrok a réz(II)-ionokat nátrium-hidroxid jelenlétében réz(I)-ionná redukálják. A lehűtött oldathoz 25 %-os kénsavat és kálium-jodidot hozzáadva a felszabadult jódot nátrium-tioszulfát mérőoldattal halványsárga színig titráljuk, majd keményítőoldatot hozzáadva a titrálást teljes fehér színig folytatjuk. A szénhidráttartalom kiszámításához szükségünk van a reagens oldatok pontos réz(II)-ion tartalmának ismeretére, ezért minden mérési sorozathoz egy-egy „vak” oldat titrálását is elvégezzük. A vak oldatra és a mintákra mért fogyás-különbségek alapján tapasztalati táblázat vagy empirikus képlet (Lásztity

és Törley, 1987) alapján történik a minta cukortartalmának kiszámítása. A mérés kimutatási határa 0,1 g / 100 g.

Az energiatalok cukortartalmának meghatározása során invert cukorra számoltam ki az eredményeket, tekintve, hogy savas hidrolízissel invertáltam a mintákat. Az energiatalokon csak nagyon kevés esetben van feltüntetve a konkrét vegyület, gyakran változnak a cukor, invertcukor, glükóz stb. megfogalmazások. Így az értékek invertcukorra történő megadása egy átlagos, jó közelítésnek mondható.

## 4.5 Kemometriai módszerek

A doktori munkám során egy- és többváltozós kemometriai eszközöket is alkalmaztam. Az egyváltozósak közül a különböző paraméteres és nem paraméteres tesztek, például  $t$ -próba vagy az előjel (Sign) és a páros Wilcoxon-próba, valamint doboz-bajusz ábrák használata volt gyakori. Ezeket többnyire a különböző összehasonlító elemzések során alkalmaztam. Ezen módszerek részleteire bővebben nem térek ki. A többváltozós elemzések világában azonban számos mintázatfelismerő valamint regressziós eljárást is alkalmaztam, amelyek viszont részletesebb kifejtést is igényelnek a gyakori használat és a sokkal speciálisabb alapelvek miatt.

### 4.5.1 Főkomponens-elemzés (PCA)

A főkomponens-elemzés (PCA) (Wold et al., 1987) manapság az egyik leggyakrabban alkalmazott mintázatfelismerő eljárás, a népszerűsége a 80-as 90-es évektől kezdve töretlen. Felügyeletlen tanítású módszerként (a minták csoportba sorolását a modellépítés során nem használjuk) nem használható tipikus osztályozási módszerként, de segítségével jellegzetes mintázatokat, csoportosulásokat, esetleg kiugró értékeket is észlelhetünk.

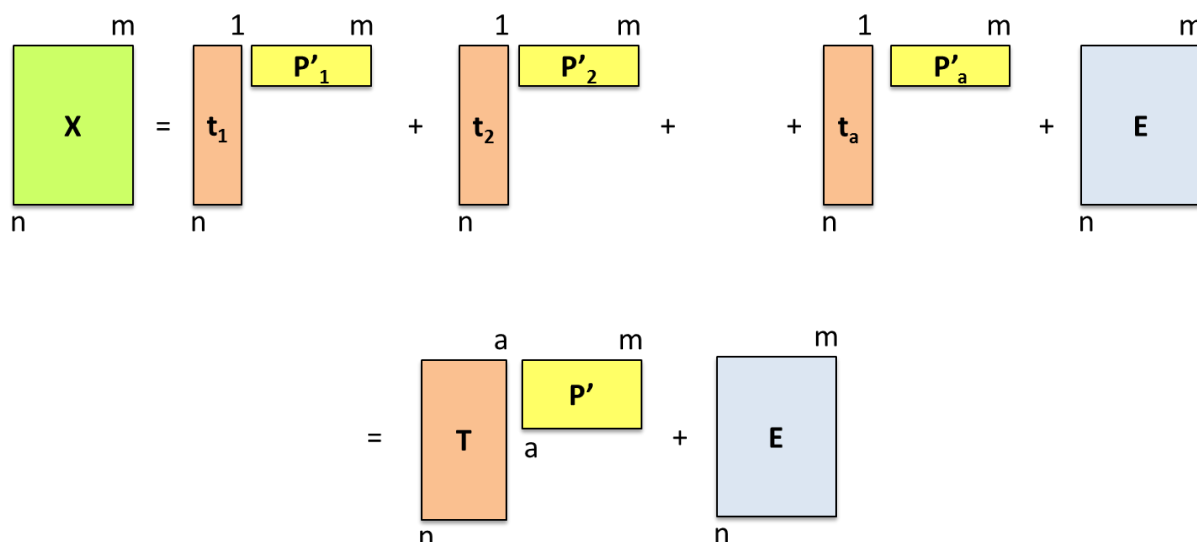
A PCA egy dimenzióredukciós módszer, alapelve szerint pedig az eredeti változók (tulajdonság vektorok, amelyek lehetnek pl. spektrum hullámszámok, koncentráció értékek stb.) lineáris kombinációjából új „látens” változókat hoz létre. Másképp megfogalmazva, az eredeti adatmátrixunk két mátrix szorzatára bontható fel: a főkomponens ( $T$ ) és főkomponens-együttható ( $P$ ) mátrixokra, amelyek ortonormáltak (vagyis egységnyi hosszúak és egymásra merőlegesek).

A két mátrix mellett egy hiba mátrix ( $E$ ) is felírható, melynek értéke nulla, ha az összes főkomponens és főkomponens-együttható vektort kiszámoljuk. A **4. ábra** az előbbi felbontást

szemlélteti. A főkomponensek úgy jönnek létre, hogy azok a lehető legjobban leírják az adatainkban rejlő varianciát. A PCA alapegyenlete a következőképp írható fel:

$$(12) \quad X = TP' + E,$$

ahol  $X$  az eredeti mintamátrix,  $T$  a főkomponens-,  $P'$  a főkomponens-együttható mátrix transzponáltja és  $E$  a hibamátrix.



**4. ábra:** A főkomponens-elemzés elvi sémája (Geladi and Kowalski, 1986)

A főkomponens-elemzés során első lépésben az adatmátrixot vagy standardizáljuk vagy centráljuk attól függően, hogy korrelációs vagy kovariancia mátrixból akarjuk kiszámolni a főkomponenseket. Ilyenkor az adatainkat beskálázzuk a változók által leírt koordináta rendszer nulla középpontjába (nulla közepűvé). A főkomponens értékek a mintákra (a változók hibáinak négyzetösszegének minimalizálásával) legjobban illeszkedő egyenesre (főkomponensre) vetített pontok nullától való távolsága. A főkomponens-együttható értékek pedig a főkomponensnek az eredeti változókkal bezárt szögének koszinusz értékei a korrelációs mátrix elemzése esetén. A főkomponensek kiszámítása a modern számítógépes eszközök birtokában néhány másodperc alatt megtörténik. A főkomponensek a főkomponens-együtthatókkal páronként, az iteratív NIPALS (nem-lineáris iteratív parciális legkisebb-négyzetek) algoritmussal számolhatók ki. A főkomponens-együtthatókhöz tartozik egy magyarázott variancia érték (sajátérték) is, ami azt mutatja meg, hogy az adott főkomponens hány százalékát magyarázza az adatkészlet teljes varianciájából. Minél nagyobb ez az érték, annál fontosabbnak, annál „hasznosabbnak” tekinthető az adott főkomponens. Természetesen ettől független adódhatnak olyan egyedi főkomponensek is, amelyek alapvetően csak egy adott változóval függnek össze, így ezekre külön figyelmet kell szentelni, mivel gyakran pont az ilyen főkomponensek használata (kombinációja) adja meg a választ az adott problémára.

A modellhez szükséges látens változó-, avagy a főkomponens szám meghatározásra több lehetőségünk is van. Lehet használni hegyomlás ábrát, ahol a főkomponenshez tartozó sajátértékek vannak ábrázolva a főkomponens szám függvényében. Ahol töréspont látszik a görbe lefutásában, annyi komponenst kell megtartani a modellben. De vannak egyéb meghatározási módszerek is, például a sajátérték nagyobb, mint egy kritérium, vagy a főkomponensszámot a magyarázott variancia bizonyos százalékában is rögzíthetjük.

A főkomponens-elemzés gyakran előzetes eszköz a változók számának csökkentésére, így alkalmazható olyan esetekben, amikor a túl sok változó miatt nem használhatunk egy adott módszert. Ennek megfelelően a doktori munkám során is szinte minden témakörben előkerült a használata. A PCA-val kapott eredményeket és speciális beállításokat az adott téma eredményeinek bemutatása során részletezem a szükséges ábrákkal szemléltetve.

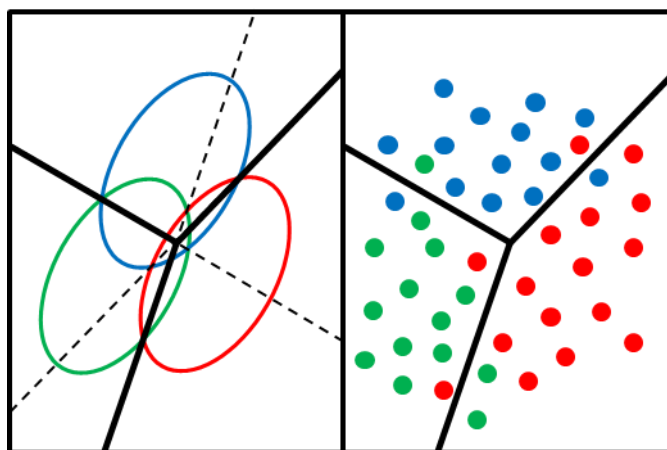
#### **4.5.2 Hierarchikus fürtelemzés (klaszter analízis, HCA)**

A hierarchikus fürtelemzés a PCA-hoz hasonlóan egy felügyeletlen tanítású mintázat felismerő eljárás (Hastie et al., 2001b; Otto, 1999). Egy rendkívül egyszerű és illusztratív módszer, amit nem csak a kémia, de a tudomány más területein is előszeretettel alkalmaznak. Többnyire a PCA-val együtt alkalmazzák, a kapott eredményeknek egymással történő megerősítésére (Melgarejo-Sánchez et al., 2015). Használata során a változók vagy vizsgált tulajdonságok közötti kapcsolatok különböző távolságmértékeken (minták vagy módszerek között vizsgálva) és kötődési szabályokon alapulnak. A mintákat vagy változókat ezen szabályok alapján csoportokba rendezzük. A kötődési (kapcsolási) szabályok adják meg a csoportok közötti távolságok számolásának módját, míg a távolság metrikák megadják a minták vagy módszerek közötti távolságokat (csoporton belül). Ez utóbbi lehet pl. Euklidészi-, Minkowski-, Mahalanobis- vagy Manhattan távolság. A kapcsolódási szabályra alkalmazhatjuk az egyszerű-, összetett-, teljes- vagy Ward módszert. A fürtelemzés hátránya, hogy a különböző kapcsolódási szabályok és távolság mértékek gyakran eltérő eredményekhez vezetnek. A legjobb gyakorlat, ha többféle kombinációt is kipróbálunk és a végső döntésünket így alapozzuk meg egy mintázattal kapcsolatban. A leggyakrabban használt távolságmérték és kapcsolódási szabály az Euklidészi távolság alkalmazása a Ward módszerrel.

A doktori munkám során a fürtelemzést az antioxidáns kapacitási módszerek összehasonlítására használtam, amikor az előbbieken említett kombinációt alkalmaztam. További részletek az „Eredmények” fejezet megfelelő részében találhatóak.

### 4.5.3 Lineáris diszkriminancia elemzés (LDA)

A lineáris diszkriminancia elemzés a nagyon hasznos és jól ismert osztályozási módszer, melyet a felügyelt tanítású mintázatfelismerési módszerek közé sorolunk. Ennek értelmében a minták csoportba sorolását előzetesen ismerni kell, és ezt az információt a modellépítés során fel is használjuk. Hasonlóképp dimenziócsökkentő technika, mint a főkomponens-elemzés, viszont ebben az esetben kanonikus változókat generálunk az eredeti változókból (hasonlóan látens változók). A minták csoportjai köré ellipszisek (vagy hiperellipszoidok) húzhatók, amelyek meghatározzák a csoportok által behatárolt teret, a diszkrimináló függvény pedig egy egyenes (vagy hipersík) által adható meg, amely az ellipszisek metszéspontjaiba húzható. Az LDA algoritmusának lényege, hogy a csoportok varianciái azonosak.  $N$  csoport esetén a kanonikus változók száma  $N-1$ . Az LDA alapelvét jól szemléltetik az **5. ábrán** látható ellipszisek (Hastie et al., 2001c).



**5. ábra:** A csoportok közötti határvonalak (diszkrimináló függvények) és ellipszisek

A modellépítés során sokféle lehetőségünk van a fontos és kevésbé lényeges változók kiszűrésére. Ilyen módzatok lehetnek a teljes változósám („all effects”), a lépésenkénti változó hozzáadás (forward stepwise), a lépésenkénti változó törlés („backward stepwise”) vagy a legjobb alrendszer („best subset”) kiválasztása. A teljes változósám választása esetén gyakran ütközhetünk túlillesztésbe. A lineáris diszkriminancia elemzésre és ezzel együtt minden felügyelt tanítású mintázatfelismerési eljárásra igaz, hogy a modellek validálására megfelelő hangsúlyt kell fektetni, hogy az eredmények a modell előrebecslő képességét mutassák, ne pedig egy mesterségesen létrehozott, túloptimált modellt. Célszerű a randomizációs teszt használata a modellek túlillesztésének vizsgálatára..

Az LDA módszert a doktori munkám energiatalokra vonatkozó részében használtam. A részleteket az „Eredmények” fejezetben fejtem ki bővebben.

#### 4.5.4 Parciális legkisebb-négyzetek módszere (regresszió és diszkriminancia elemzés)

A parciális legkisebb-négyzetek módszerét (PLS) alapvetően regressziós módszerként tartják számon, amely csak annyiban különbözik a diszkrimináló PLS-től, hogy az  $y$  függő változó folytonos vektorként referencia, mérési eredményeket tartalmaz, vagy csoportosítást. Ez utóbbi esetén beszélünk PLS diszkriminancia elemzésről (PLS DA vagy DPLS). A PLS módszer egyike a leggyakrabban alkalmazott többváltozós kemometriai módszereknek. Alapjában hasonló a többváltozós lineáris regresszióhoz (MLR), tekinthető akár az MLR általánosításának is. A PLS egyik nagy előnye az MLR technikával szemben, hogy olyan esetekben is alkalmazható, amikor több változónk van az adatkészletben, mint mintánk (ilyenkor az MLR megoldás azonosságra vezet).

A PLS alkalmazása során az  $X$  független és  $Y$  függő változó közötti kapcsolatot derítjük fel. Ez a kapcsolat lehet külső és belső típusú is. A külső kapcsolat felírásához PLS komponenseket hozunk létre hasonló módon, mint a főkomponens-elemzés során, csak külön-külön az  $X$  illetve  $Y$  mátrixunkra. Az eredeti mátrixok felbonthatók a  $T$  és  $U$  mátrixokkal jelzett PLS komponensekre, valamint a  $P'$  és  $Q'$  mátrixokkal jelzett PLS komponens-együttható mátrixok transzponáltjaira. A PLS komponensek szintén lineáris kombinációval jönnek létre az eredeti változókból. Ezen kívül a hiba mátrixok ( $E$  és  $F$ ) is hozzá adódnak, amelyek értéke nulla, ha az összes PLS komponens felírjuk. Ezt a folyamatot jól szemlélteti a **6. ábra**. A PLS módszer új, rejtett, „látens” változókat ( $T$  és  $U$ ) használja az  $Y$  értékek előrebecslésére (Brereton et al., 2014). A belső kapcsolat alapja a regressziós koefficiens ( $b$ ), amely segítségével a  $T$  és  $U$  mátrixok közötti összefüggés megadható. A regressziós koefficiens az  $u$  és  $t$  vektorok által ábrázolt pontokra húzható egyenes meredekségével ( $\tan\alpha$ ) egyezik meg.

$$(13) \quad \hat{u}_h = b_h t_h$$

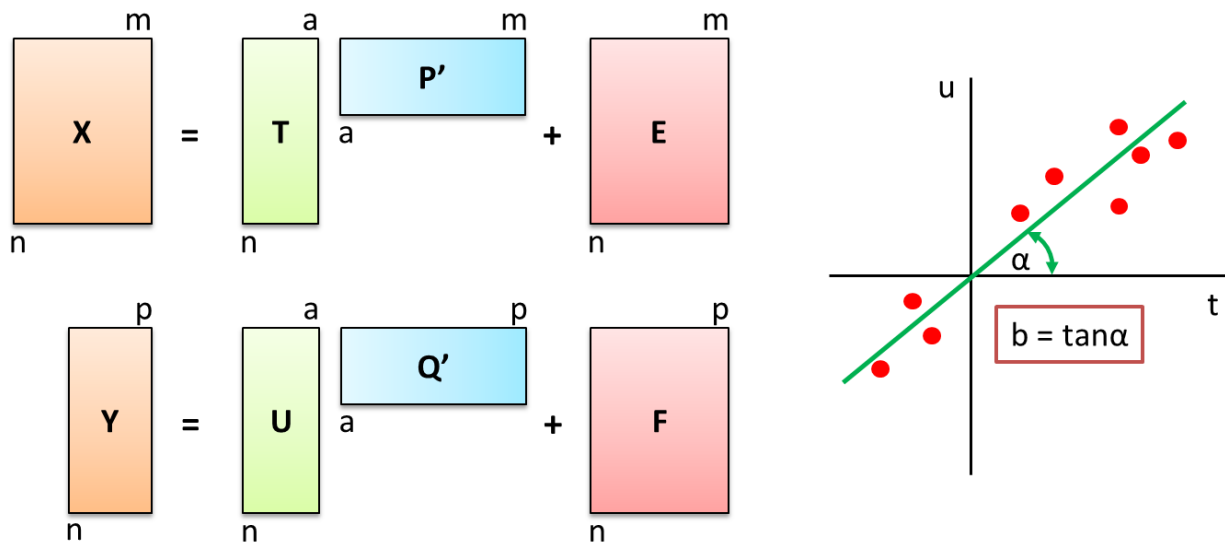
$$(14) \quad b_h = u'_h t_h / t'_h t_h ,$$

ahol  $\hat{u}$  a becsült (közelített) PLS komponens érték,  $t$  az  $X$  mátrixhoz tartozó PLS komponens vektor,  $u$  pedig az  $Y$  mátrixhoz tartozó PLS komponens vektor.

Végül a kombinált kapcsolat a két előbbi alapján a következőképp írható fel:

$$(15) \quad Y = TBQ' + F ,$$

ahol  $T$  a PLS komponens mátrix,  $B$  a regressziós koefficiensek mátrixa,  $Q'$  a PLS komponens együttható mátrix és  $F$  a hibamátrix.



**6. ábra:** A PLS alapelvének grafikus ábrázolása

*A mátrixok esetében  $n$  a mintaszámot,  $m$  a változós számot,  $a$  pedig a PLS komponensek számát jelenti. A regressziós koefficiensek jelölése  $b$ -vel történik (Geladi és Kowalski, 1986).*

A belső és külső kapcsolatok alapján a NIPALS algoritmus segítségével a két blokk PLS komponens vektorait kicserélhetjük, amellyel a belső kapcsolaton javíthatunk. Ehhez bevezetésre kerül egy  $w$  súlytag, amely a  $t$  és  $u$  vektorok közötti korrelációt maximalizálja. A PLS módszer, legyen szó regresszióról vagy diszkriminancia elemzésről, rendkívül könnyen megérthető a Geladi és Kowalski méltán híres publikációjából (1986).

A PLS modellépítések során a látens változók száma egy igen fontos paraméter. A PLS komponensek száma megadja a modell dimenzióját, és a legtöbb esetben nincs szükség az összes PLS komponens felhasználására a modellépítéshez. Ha túl keveset választunk alulillesztett lesz a modell, ha túl sokat, túlillesztetté válhat. A modell komplexitásának növekedésével együtt nő a becslési hiba, viszont csökken a modell hibája.

A megfelelő komponens szám kiválasztása többféle módon történhet. Gyakran alkalmazott az RMSECV vagy  $R^2$  alapján történő kiválasztás, amikor szoftvertől és módszertől függetlenül választhatjuk a PLS komponens szám függvényében kapott görbe lokális vagy globális minimumát is. Szintén használható a PRESS, azaz az előrebecslési hiba négyzetösszege, amelynek a PLS komponens szám függvényében való ábrázolásakor ismételt a görbe minimumát kell megkeresni. A látens változók kiválasztásának igen kiterjedt szakirodalma van, amelyben találhatunk kevésbé elterjedt de hasznos verziókat is, ilyen pl. a randomizációs teszt, az érzékenység korrekció vagy a sajátértékek alapján történő választás (Bro et al., 2008; Wiklund et al., 2007).

A PLS segítségével diszkriminancia elemzést is végezhetünk (PLS DA), amelyre szintén számos példát találhatunk a szakirodalomban, viszont a véletlen osztályozás növekvő valószínűségét mindig szem előtt kell tartani ilyen esetekben (Barker és Rayens, 2003; Kjeldahl és Bro, 2010).

A PLS módszer alkalmazása a doktori munkám majdnem minden témakörében előfordult vagy regressziós módszerként (PLSR) vagy diszkriminancia elemzéseként (PLS DA). A PLS elemzések során alkalmazott speciális paramétereket és beállításokat a hozzájuk tartozó ábrákkal együtt az eredmények részben szemléltetem.

#### **4.5.5 Véletlen erdő módszere (Random forest, RF)**

A véletlen erdő módszere a rekurzív felosztáson alapuló módszerek körébe tartozik, melyek regressziós és osztályozási feladatokra is kiválóan alkalmazhatóak. A fa alapú módszerek alapelve, hogy bináris osztályozásokat hozunk létre a csomópontjain keresztül, vagyis a kiinduláskor még egy csoportba tartozó teljes mintakészletből minden csomópont esetén egy eredeti változó beállított (limit) értéke alapján igen/nem kérdéseket teszünk fel, amely segítségével el tudjuk dönteni, hogy az adott minta melyik csoportba kerüljön.

A véletlen erdő módszere azzal fejlesztette tovább ezt a gondolatmentet, hogy nem csak egy, hanem rengeteg fát épít fel az osztályozás becsléséhez. A végső osztályozás a fa szekvenciák közötti szavazással dől el. A véletlen erdő módszerének betanító (training) algoritmusá egy általános használt agglomeráló „bagging” technikát alkalmaz. A megalkotott fák külön-külön gyenge előrebecslési képességűek lennének, de együtt „zenekarként” már képesek megfelelő becslést létrehozni (Breiman, 2001).

A véletlen erdő módszerének két fő paramétere, amit optimalizálni kell a modellépítés előtt, a tulajdonságváltozók száma és a fák mennyisége. A doktori munkám során ezt úgy választottam ki, hogy alapesetként kiválasztottam azt a legkisebb változós számot, ahol a helyes osztályozási százalék (CC %) nem változik láthatóan újabb változók bevonásakor (a CC % görbéje egy telítési (plató) állapotot ér el). A fák mennyiségének megválasztása ezt követően történt a már kiválasztott tulajdonságváltozó szám mellett az előbbiekhöz hasonló módon. Részletesebb leírás a módszer használatáról és annak eredményeiről az „Eredmények” fejezet energiatalokra vonatkozó részénél található.



#### **4.5.6 Fejlesztett fák módszere (Boosted tree, BT)**

A fejlesztett fák módszerét eredetileg osztályozási problémákra találták ki, majd később kiterjesztették regressziós feladatokra is. Alapelvét tekintve hasonló az RF módszerhez, itt is számos bináris fát hozunk létre: a fák csomópontjaiban kétfelé vágjuk az adott változót. Minden egyes lépésben létrejön a „fejlesztés”, amely eredményeképpen az adatok (legmegfelelőbb) részekre bontása megtörténik. A referencia értékektől (csoportosítástól) való eltérések, illetve hiba értékek kiszámolását is elvégezzük. Ez teljesen egyértelmű a regressziós mód esetén, az osztályozási feladatoknál a következő lépések történnek: a csoportok számának megfelelően felosztjuk az adatkészletet alcsoportokra, majd logisztikus transzformáció segítségével kiszámoljuk a súlyozott félreosztályozási hibákat az alcsoportok számára (fejlesztési lépés), majd végül a teljes félreosztályozási százalékot is (Hastie et al., 2009). Nagyon fontos tulajdonsága a fejlesztett fák módszerének a minták súlyozása aszerint, hogy melyik mintát milyen nehezen tudja osztályozni: a félreosztályozott minták büntető súlytagot kapnak az egyes lépésekben.

A sztohasztikus gradiens fejlesztés esetén minden egyes fa véletlenszerűen kiválasztott mintákra épül fel a teljes adatkészletből (betanító halmaz). A véletlenszerű kiválasztás modellépítésbe való bevonása erős védőeszközként működik a túlillesztés kockázata ellen (mivel minden fa különböző mintacsoportok segítségével épül fel és alkotja meg a végső modellt). A fejlesztett fák algoritmusát részletesen megtalálható Hastie és munkatársai által írt könyv ide vonatkozó fejezetében (Hastie et al., 2009).

A doktori munkám során az energiaitalok témakörében alkalmaztam ezt a módszert, az ekkor használt optimalizálási lépéseket a fejlesztett fák módszeréhez az „Eredmények” fejezetben tárgyalom.

#### **4.5.7 Változó kiválasztási eljárások**

A változó kiválasztási eljárásoknak jelentős szerepe van mind a regressziós mind pedig az osztályozási modellek megalkotásakor. A szakirodalomban a növekvő adatmennyiségek („big data”) kezelhetőségé tétele miatt egyre többféle módszert találni a modellépítéshez szükséges változó kiválasztására. FT-NIR spektrumok elemzésekor sok esetben ütközünk abba a problémába, hogy a teljes spektrum annyi fölösleges információt (zajt) tartalmaz, hogy a becslésünk eltorzul, nem kapunk megfelelően használható modelleket. Az ilyen esetekre nyújt kiváló megoldást a változó kiválasztási eszközök használata. Segítségükkel adott algoritmusok alapján kiválasztható az a néhány tíz-száz változó tartománya, amely a modellépítéshez szükséges információkat tartalmazza. Andersen és Bro (2010) tanító cikke részletesen és

közérthetően azt taglalja, hogy a spektrális adatok feldolgozása során milyen eszközök vethetők be, és hogyan. A cikkben megfogalmazott gondolatok szerint a változó kiválasztási módszerekkel nem csak a modellünk jóságát tudjuk növelni, de ezzel együtt könnyebb/jobban értelmezhetőséget kapunk, valamint csökkenteni tudjuk a kísérleti költségeket is. Természetesen, mint minden kemometria módszer, a változó kiválasztás sem képes csodát tenni, tehát ha a koncentráció és hullámszámok között nincs semmilyen összefüggés, a változó kiválasztásával sem fog működni a modellépítés. Nagyon fontos kitérni arra is, hogy a változó kiválasztási eljárásokat a spektrumok és egyéb adatkészletek esetén is az adatelőkezelési lépések után kell megtenni. Két fajtáját különböztethetjük meg: i) a változó csökkentést (állandó, variancia nélküli változókat vagy erősen korrelált változókat távolítunk el a kiindulási  $\mathbf{X}$  mátrixból) és ii) az adott  $y$ , referencia változóhoz ( $\mathbf{Y}$  változókhöz) válogatunk leíró és lényegi (releváns) változókat (ez utóbbi a felügyelt tanítású változó kiválasztás).

A változó kiválasztására egy könnyű és gyors út a regressziós koefficiensek ( $b$ ) használata. Azokat a változókat, amelyekre vonatkozó  $b$  érték nem számottevő (nem különbözik szignifikánsan a nullától), méltán elhagyhatjuk az adatkészletből.

Az előbbinél kissé komplexebb megoldás a változó-fontosság paramétere (VIP érték). Ez az érték megmutatja, mennyi a hozzájárulása a változónak az  $\mathbf{X}$  és  $\mathbf{Y}$  közötti összefüggéshez. Ha a VIP érték kisebb, mint egy, a változó elhanyagolható. Vigyázni kell viszont vele, mert sok esetben egy-egy változó meghagyása a korrelált spektrumban rosszabb eredményhez vezethet. Az előbbi ok miatt érdemes a változókat csoportosan, sávszerűen szelektálni, mintsem egyesével. Az alapelv szerint először a legkisebb VIP értékű tagokat távolítjuk el, majd addig folytatjuk, amíg a modellünk fejlődik.

A szelektivitási arány (SR) esetében, a változó magyarázott varianciája és a reziduális hiba értéke közötti arányt adjuk meg (Rajalahti et al., 2009). Minél nagyobb a szelektivitási arány értéke, annál fontosabb az adott változó.

Az intervallum PLS regresszió (iPLS) szintén kiváló változó kiválasztási eljárás. Ez az egyik legáltalánosabb választás, főként a spektrum adatkészletek esetén, mivel a spektrumok, ahogy már korábban említettem, erősen korrelált adatkészletek. Így a változó „ablakok” használata sokkal jobb lehetőség, a változó egyenkénti vizsgálatánál (Di Anibal et al., 2011; Kristensen et al., 2010; Nørgaard et al., 2000). Ekkor a változókat azonos nagyságú intervallumokra bontjuk fel (történhet manuálisan is), majd a változó szegmensekre külön-külön elvégezzük a PLS regressziót. A megfelelő intervallum(ok) kiválasztása történhet RMSECV,  $R^2$  vagy  $Q^2$  értékek alapján. Többféle felosztással is érdemes elvégezni az eljárást, jellemzően 10,

20 vagy 40 részre való felosztás tekinthető általánosnak (Andersen és Bro, 2010). A lépésenkénti változó kiválasztás (forward selection) az iPLS esetén is alkalmazható, csakúgy mint a lépésenkénti változó törlés (backward elimination). A további modellépítéshez a legkisebb hibával vagy legnagyobb  $R^2$  értékkel rendelkező intervallumok tekinthetők a legkedvezőbbeknek.

A genetikus algoritmus (GA) egy globális minimumkereső eljárás, amely viszont közelíteni tudja csak a globális minimumot. Az alapelve szorosan köthető a genetikához és annak működéséhez. Használható a minél jobb modellhez szükséges változók kiválasztására is. Ehhez több generáción keresztül kombinálódnak, mutálódnak a változó „alegységek” (egyedek). A kezdő lépésben kiválasztjuk a populáció méretét, amely általánosságban 20 és 500 egyedszám közötti érték (Andersen és Bro, 2010; Leardi, 2007), valamint az egyedekhez tartozó változók mennyiségét is (véletlenszerű kiválasztás). A következő lépések a reprodukció, kereszteződés és mutáció. A folyamat addig ismétlődik, amíg egy „leállási kritérium” meg nem állítja. Ilyen kritérium lehet, ha eléri a keresési fázis a célfüggvény minimumát, de lehet például a generáció szám felső korlátjának elérése vagy az előre meghatározott számolási idő is. A genetikus algoritmus szintén használható változó „ablakkal”. A PLS Toolbox szoftver (Eigenvector Research Inc., USA, 7.0.3 verzió) esetén ez az ablak maximum 50 változót tartalmazhat.

Összességében elmondható a változókiválasztási eljárásokról, hogy soha nem árt, ha többfélet is kipróbálunk a modellépítés során, valamint az is, hogy a „klasszikusnak” tekinthető módszerek gyakorlatilag szoftverfüggetlenek. A genetikus algoritmus az iPLS-hez viszonyítva mindenképpen számításigényesebb. Bár használatával javíthatunk a modelleken, önmagában nem építhetünk vele modellt.

A változó szelekciós eljárások a Q10 tartalmú étrendkiegészítők vizsgálatánál kaptak nagyobb szerepet, a modellépítések során továbbfejlesztettem a szelektivitási arány módszerét.

## **4.5.8 Rangsorolás és összehasonlítás**

### **4.5.8.1 A vevő-működtető jelleggörbék (ROC görbék)**

A ROC görbék (receiver operating characteristic curves) sajátos angol nevének eredete a II. világháború idejére nyúlik vissza, amikor katonai radaroknál egy ellenség-sajátgép felismerő rendszert (ellenséges objektumok felismerése a harctéren) dolgoztak ki. A ROC görbe használata ezt követően elsőként az orvostudomány és a pszichológia területén honosodott meg (Hanley és

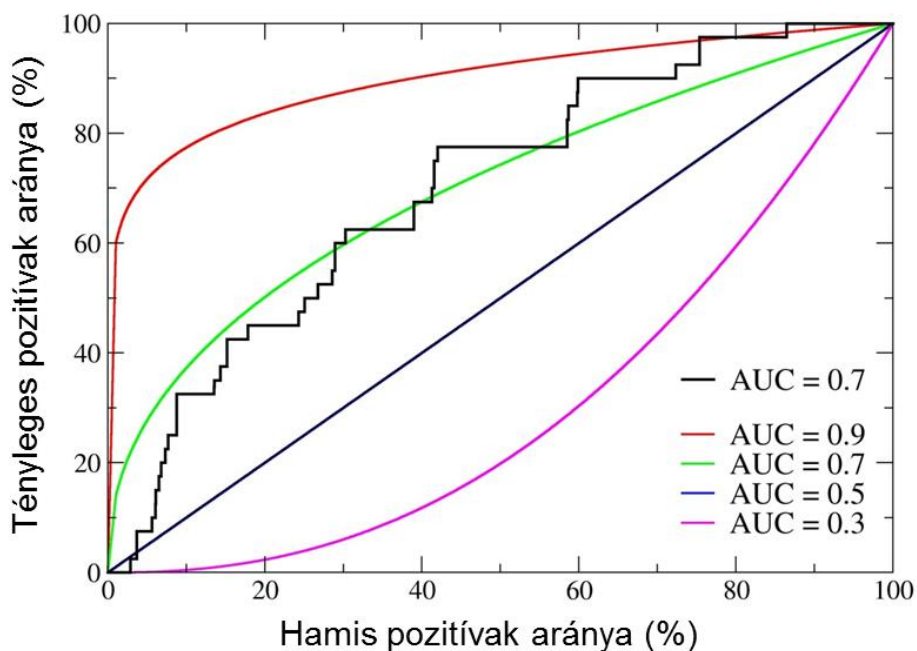
McNeil, 1982; Metz, 1978), majd később a gyógyszerkutatás és így a kémia területére is teret nyert.

Tételezzük fel, hogy  $D$  egy folytonos bináris osztályozó változó; az osztályok pedig  $+$  és  $-$  értékekkel vannak kódolva. Szeretnénk megbecsülni annak a valószínűségét, hogy a minta az 1. csoportba (más szóval a pozitívba) tartozik-e. Mi az a  $D$  korlát, ami fölött a mintát pozitívnak tekinthetjük? Minden  $D$  értékre kiszámoljuk a tényleges pozitívak arányát (true positive rate, TPR), amely a helyesen osztályozott pozitív minták aránya a teljes pozitív csoporthoz képest; hasonlóképp kiszámoljuk a hamis pozitívak arányát (false positive rate, FPR) is, amely a hamisan negatívnak osztályozott minták aránya az összes negatív mintához képest. A ROC, avagy a vevő-működtető jelleggörbe megalkotásakor a TPR értékeket ábrázoljuk az FPR értékek függvényében minden egyes korlát értékre nézve csökkenő sorrendben, amely végeredményeként egy monoton növekvő görbét kapunk ami  $(0;0)$ -tól  $(1;1)$  koordinátáig halad, ahogy a **7. ábrán** is látható.

A ROC görbék ábrázolásakor a  $(0;0)$ -tól  $(1;1)$ -ig terjedő diagonális a véletlenszerű osztályozást szimbolizálja: azok a módszerek vagy modellek, amelyek ROC görbéje az átló felett húzódik, jobbak a véletlen osztályozásnál, az átló alatt húzódók pedig rosszabbak. Így a ROC görbe egy gyors vizuális összehasonlítást tud adni módszerek, vagy modellek jóságáról. A ROC görbék teljesítmény paramétereként a görbe alatti területet (AUC érték) szokás használni. Az AUC érték 0 és 1 között változhat, de százalékban is szokás kifejezni. Hanley és McNeil (1982) publikációjukban bebizonyította, hogy az AUC érték megegyezik a Mann-Whitney-Wilcoxon statisztikával (U-teszt). Az AUC értékek varianciájának kiszámolása DeLong és munkatársai (2016) nevéhez köthető. A ROC görbék illesztésére számos módszer létezik a szakirodalomban, de kevésbé elterjedtek, mint az AUC értékek alapján történő ROC görbe illesztés, amely Hanley és McNeil (1982) munkásságához, illetve Nicholls (2014) nevéhez köthető:

$$(16) \quad Y = X^{\frac{1-AUC}{AUC}}$$

A **7. ábrán** látható illesztett ROC görbéket szintén a Hanley formula segítségével illesztettem meg. Az eredeti ROC görbével összehasonlítva, az illesztett görbék nem lépcsőzetesek, viszont ettől függetlenül tökéletes és gyors eszközt biztosítanak a vizualizációhoz.



**7. ábra:** A tényleges pozitívak aránya ábrázolva a hamis pozitívak arányának függvényében. Az eredeti ROC görbét fekete vonallal jelöltem, valamint különböző görbe alatti terület (AUC) értékeknél kiszámolt (Hanley-formula) görbék is fel vannak tüntetve példaként. A véletlenszerű osztályozás (amely az átló is egyben) fekete színnel van jelölve.

A ROC görbék megalkotásához segédletként található egy ábra a „Mellékletek” részeként (M3). A doktori munkám során megalkottam a ROC görbék tovább fejlesztett változatát a különböző mintázatfelismerési módszerek összehasonlítására, amelynek részleteit az Eredmények fejezetben tárgyalom.

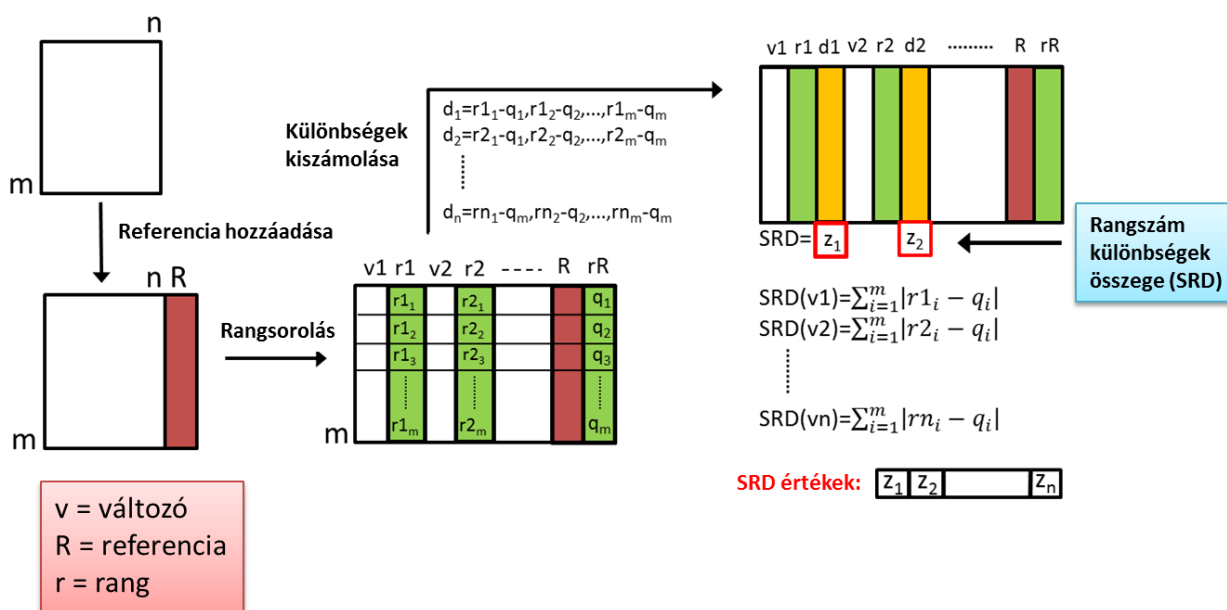
#### 4.5.8.2 Rangszám különbségek abszolút értéke összegének módszere (SRD)

A rangszámkülönbségek abszolút értékének összege rövidítve SRD (Héberger, 2010; Kollár-Hunek és Héberger, 2013) egy új, általánosan is használható módszer a modellek, módszerek stb. összehasonlítására (Moorthy et al., 2015; Nowik et al., 2013). Használata során az adatmátrixban a minták a sorokban a modellek (változók) pedig az oszlopokban helyezkednek el. Szükség van referencia oszlopra is, amely lehet konkrét referencia értékeket tartalmazó is, de ennek hiányában a sorátlag, minimum vagy maximum érték is használható az adatkészlettől és feladattól függően. Első lépésként a mintákat minden oszlopban (beleértve a referenciát is) rangsoroljuk növekvő sorrend szerint. Ezt követően a referencia oszlop és a modellek (módszerek) rangjai közötti különbséget adjuk meg. Végül a különbségek abszolút értékeit összegezzük minden változóra (modellre) nézve, így azok összehasonlíthatóakká válnak. Minél közelebb helyezkedik a nullához az adott módszer vagy modell SRD értéke, annál közelebb van

a referenciához, azaz annál jobbnak tekinthető. Az SRD módszer validálása randomizációs teszttel és kereszt-ellenőrzéssel (7-részre osztott vagy 14 minta alatt esetben egy elem kihagyásos módszerrel) történik. A **8. ábra** megfelelően szemlélteti az SRD értékek számolásának menetét.

Az SRD módszer végeredményeképpen kapott ábrán az SRD értékek mellett egy Gauss görbe (vagy ahhoz hasonló elméleti görbe) is található, amely a véletlenszerű rangsorolást mutatja. Azok a módszerek, amelyek a Gauss görbével átfednek, nem tekinthetők jobbnak a véletlenszámok által végzett rangsorolásnál; viszont jól jelezhetik a referenciától való különbözőséget.

Az SRD módszert a doktori munkám során az antioxidáns kapacitási módszerek, valamint a Q10 koenzim tartalmú étrendkiegészítők hatóanyagtartalmának becslésére alkotott modellek összehasonlítására használtam.



**8. ábra:** Az SRD értékek kiszámítása

A mátrixok esetében  $n$  a változó számot,  $m$  pedig a mintaszámot jelöli. A különbségek kiszámításánál a  $v$  megfelelő számozással ellátva jelöli az adott változókat, míg  $r$  a hozzá tartozó rangszámokat,  $R$  a referencia jelölése,  $d$  a különbségeket és  $z$  pedig a kiszámolt SRD értékeket jelöli.

#### 4.5.8.3 Általánosított pár-korrelációs módszer (GPCM)

A pár-korrelációs módszer hasonlóan az SRD-hez, szintén egy könnyű és gyors utat jelent módszerek vagy modellek (tulajdonságok) közötti rangsorolásra és kiválasztásra (Héberger és Rajkó, 2002; Rajkó és Héberger, 2001). A különbség, hogy a GPCM az SRD-vel szemben távolság független, viszont ugyanúgy nem paraméteres (robosztus) módszer. A GPCM technika csak az Y kiugró értékekre érzékeny. A kiindulási mátrix ugyanolyan elrendezésű, mint az SRD módszer esetében, tehát a minták a sorokban, a változók (tulajdonságok) az oszlopokban helyezkednek el. Szintén szükség van egy referencia oszlopra is az összehasonlítás során. A változók az összes kombinációs lehetőséget kihasználva páronként kerülnek összehasonlításra. Az összehasonlításoknak háromféle kimenetele lehet: győztes (ha az egyik változó az összehasonlított párból jobban „korrelál” a referenciával, mint a másik), vesztes (ha az adott változó az összehasonlított párból kevésbé korrelál a referenciával, mint a másik) és döntetlen (ha egyik változó se jobb vagy rosszabb szignifikánsan a másikonál). Az eredmények alapján a változókat, módszereket, modelleket háromféle módon lehet rangsorolni: egyszerű rangsorolás (a győzelmek száma számít csak), különbségi rangsorolás (kiszámoljuk a győzelmek és veszteségek különbségét), szignifikáns rangsorolás (a különbség rangsorolási valószínűséggel súlyozott változata). Különböző szelektációs kritériumok közül választhatunk az elemzés során, amely lehet a „feltételes egzakt Fisher próba” (Conditional Fisher’s exact test), a McNemar próba, a  $\chi^2$ -próba vagy a Williams  $t$ -próba. Az utóbbiak közül egyedül a Williams  $t$ -próba paraméteres.

Az SRD és a GPCM használatához is speciális programokat fejlesztettek, amelyet MS Excel segítségével lehet használni. A doktori munkám során a GPCM módszert kimondottan az antioxidáns kapacitás meghatározási technikák összehasonlítására használtam.

## 5. EREDMÉNYEK

### 5.1 Kemometriai módszerfejlesztések

A kemometriai módszerfejlesztéseimet az eredmények első szakaszában mutatom be, mert a továbbiakban a különböző vizsgálatoknak szerves részét képezik ezek a módszerek. Az alkalmazásukra az egyes adatkészletek értékelésekor külön ki fogok térni.

#### 5.1.1 Új változószelekciós módszer létrehozása – az intervallum szelektivitási arány (iSR)

A módszer alapötletét, a szelektivitási arány módszert, Rajalahti és munkatársai (Rajalahti et al., 2009) dolgozták ki. A cikk szerzői szerint ez a módszer megfelelő indikátor a változókiválasztásra. A szelektivitási arány minden egyes spektrum változóra külön-külön számolandó ki a következő egyenlet szerint:

$$(17) \quad SR_i = v_{exp,i}/v_{res,i}, \text{ ahol } i = 1, 2, 3 \dots m.$$

A  $v_{exp}$  megfeleltethető a magyarázott varianciának ( $R^2$ ) és a  $v_{res}$  pedig a reziduális varianciának (SEC). Az  $i$  a változók számát jelenti. A szelektivitási arány dimenziója így reciprok koncentráció négyzet lesz, de ez a használhatóságát nem befolyásolja. Az eredeti egyenlet némiképp meg kellett változtatni az „intervallumosítás” miatt: a nevezőbe egy négyzetgyök függvényt (RMSEC) még beillesztettem. Ez a hozzáadás (mint egyszerű algebrai transzformáció ugyanúgy, mintha osztást használnánk ugyanazon számmal végig) nem változtat a tendenciákon, amelyeket így meg tudunk figyelni. A módosított egyenlet a szelektivitási arány alapján a következő:

$$(18) \quad SR_i = R_i^2/RMSEC_i, i = 1, 2, 3 \dots m.$$

Ahol  $i$  az intervallumok számát jelöli. Az egyenletben az  $R^2$  esetén a szabadsági fokokat kihagytuk, amely viszont egyenlő mértékben befolyásolja az összes intervallumot (ekvidisztáns felosztás esetén, így a végső ábrán nem okoz aránytalan eltéréseket. A szelektivitási arány intervallumokra történő alkalmazására azért volt szükség, mert a szakirodalmi példák alapján tudni lehet, hogy az egyedi hullámszámok nagyon sok véletlen információt tartalmaznak és egyenként nem használhatók a modellek javítására. A szelektivitási arányhoz hasonlóan, az intervallum SR is minél nagyobb, annál fontosabb az adott intervallum a modellépítés szempontjából.



### 5.1.2 „n-class”, avagy több osztályos ROC görbék

A ROC görbék alapvetően kétosztályos adatkészletek esetén használható módszerek. A többosztályos kiterjesztésének alapját az „egy vs. összes” módszer képezte, amely Provost és Domingos (2001) munkájához köthető. Másrészt az alapelv nagyon hasonlít az osztályanalógiák közvetett modellezésének (SIMCA) módszeréhez (Wold, 1976). Ennek megfelelően az ROC görbe alatti terület (AUC) úgy számolható ki, hogy az egyik csoportot pozitívként kezeljük az összes többit pedig negatívként és ezt ismétljük annyiszor, ahány csoportunk van. A doktori munkámhoz kötődően, ezt a módszert három csoport esetén mutatom be.

Az AUC értékek súlyozott átlaga adja meg a teljes osztályozásra adható AUC értéket:

$$(19) \quad \overline{\text{AUC}} = \frac{\sum_{i=1}^n N_i \text{AUC}_i}{\sum_{i=1}^n N_i},$$

Ahol  $N_i$  a csoportokban (osztályokban) lévő minták száma. A Hanley formula segítségével vizualizálni tudunk egy „ROC-szerű” görbét az AUC értékek súlyozott átlagának felhasználásával (összesített AUC érték). A hibaterjedés törvényének alkalmazásával az előbbieken említett átlagos AUC érték varianciája is kiszámolható (Ku, 1966):

$$(20) \quad \text{Var}(\overline{\text{AUC}}) = \frac{\sum_{i=1}^n N_i^2 \text{Var}(\text{AUC}_i)}{\sum_{i=1}^n N_i^2}$$

Több módszer ROC görbékkel történő összehasonlításakor, ahhoz hogy el tudjuk dönteni, hogy az adott módszerek szignifikánsan különböznek-e egymástól, az átlag  $\pm$  az összesített AUC érték egyszeres standard deviációját alkalmaztam.

Bár az AUC értékek súlyozott átlagát korábban már alkalmazták (Provost és Domingos, 2000), ezt egy további lépéssel megtoldottam, felhasználva a Nicholls által közölt (2014) formulát (az eredeti cikk 22. és 23. egyenlete alapján), amely az AUC értékek szórásának kiszámítását vagy akár a 95 %-os konfidencia intervallumának megadását is lehetővé tette. Továbbá a hibaterjedés törvényének alkalmazásával, az egyes osztályokhoz tartozó külön-külön számított AUC értékek hibáit figyelembe véve, megadtam az átlagos vagy összesített AUC értékhez tartozó hibát is. Ezt követően a már meglévő átlagos AUC értékeket ábrázoltam. Ehhez segítségemre volt a Hanley formula (lásd. Rangsorolás és összehasonlítás fejezet). A vizualizációhoz létrehozott programkód a „Mellékletek” fejezet **M4** pontjában található

A rendelkezésre álló egyenletek és számolási lehetőségek alapján egy olyan új, többosztályos ROC görbe számolási és megjelenítési metodikát hoztam létre, amely már a

hibaterjedést és szórásokat is figyelembe tudja venni és kiválóan alkalmas módszerek értékelésére és az osztályozó képességük megjelenítésére kettőnél több csoport esetén is.

A továbbiakban a több osztályos ROC görbék egy olyan alkalmazását mutatom be, ahol mintázatfelismerő eljárásokat hasonlítottam össze. A ROC görbe generálásához az LDA, a véletlen erdő (RF) és a fejlesztett fák módszere (BT) során kapott osztályba sorolási valószínűségeket használtam fel. A PLS DA esetén kiszámoltam az abszolút értékben vett különbségeket a becsült csoport érték (folytonos) és az adott „pozitív” csoport (diszkrét) érték között: ekkor természetesen az alacsonyabb értékek tekinthetők jobbnak. Például ha kettes csoportot tartjuk a ROC görbe készítés során „pozitívnak”, egy 1,8 becsült értékkel rendelkező minta jobbnak tekinthető egy 2,4 becsült értékkel rendelkezőnél, hiszen az abszolút különbségük  $|1,8 - 2| = 0,2$  az első esetre és  $|2,4 - 2| = 0,4$  a második esetre nézve. A többi csoportra is hasonlóan történik a ROC görbe megalkotása, majd az átlagos (összesített) ROC és AUC értékek a csoportokra vonatkoztatott ROC görbékéből és AUC értékekből számolhatóak ki a korábban már részletezett módon.

### 5.1.3 Osztályozási módszerek összehasonlítása „*n*-class” ROC görbék segítségével

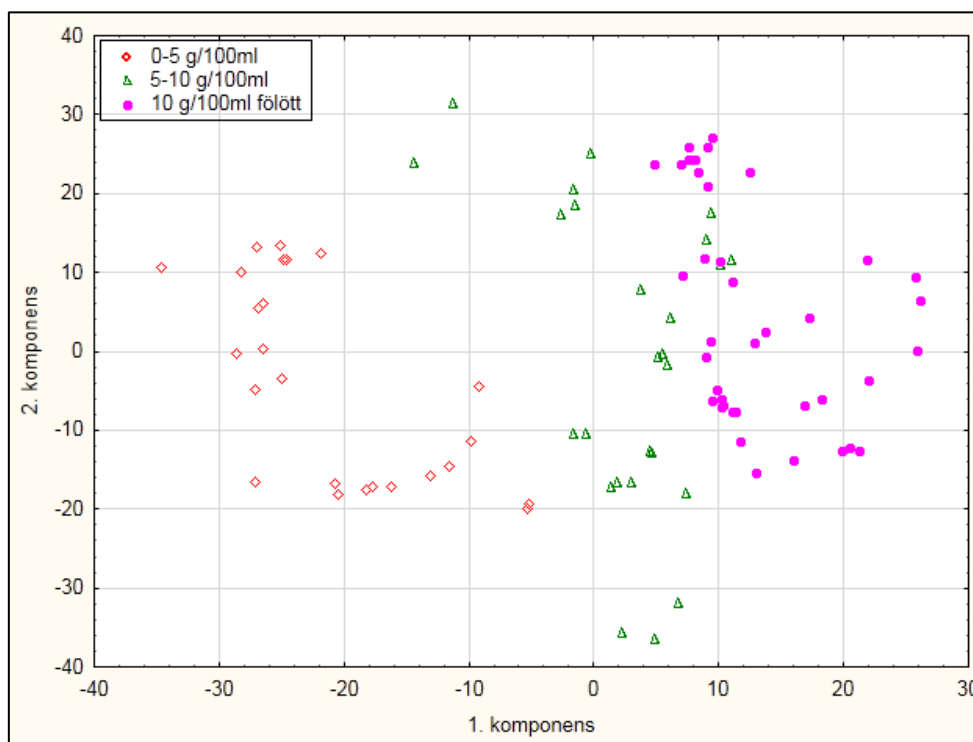
Az előző fejezetben említett véletlen erdők (RF), fejlesztett fák (BT), LDA és PLS DA módszerek teljesítményét hasonlítottam össze a több osztályos ROC görbék segítségével. Ezek a módszerek manapság nagyon elterjedt felügyelt tanítási mintázatfelismerő eljárások közé sorolhatók. Az LDA és PLS DA klasszikusoknak számítanak, míg az RF és BT módszerek – bár már nem napjaink fejlesztései – jóval újabbak, és használatuk egyre gyakoribb. A vizsgálatokhoz energiatalok FT-NIR spektrumait és azok cukortartalom szerinti csoportosíthatóságát osztályozási változóként használtam fel. Általánosságban az energiatalok cukortartalma 0 és 15 g / 100 ml közötti érték. Ezen belül a gyártók előszeretettel alkalmaznak 7,9 vagy 10,9 g / 100 ml-es értéket, valamint megfigyelhető volt, hogy a 4-5 g / 100 ml és 9 g / 100 ml körüli cukortartalmak hiányoznak az előzetesen folytonosnak gondolt skáláról. A csoportosítás tehát azon a tendencián alapult, miszerint az energiatalok cukortartalma három jól körülhatárolható tartományban változik: 1. csoport = 5 g / 100 ml alatti; 2. csoport = 5-10 g / 100 ml közötti; 3. csoport = 10 g / 100 ml fölötti értékűek. Ez az osztályozás bár önkényesnek tűnhet, valójában a mintakészletben rejlő „természetes” csoportosulást követte.

Két adatkészletet használtam fel az osztályozási módszerek összehasonlítására. Az első esetben 90 energiatal minta FT-NIR spektrumát, majd ugyanezen minták spektrumainak főkomponens-elemzés révén kapott főkomponenseit. Az FT-NIR spektrum 649 spektrális változót tartalmazott, amelyek 9000 és 4000  $\text{cm}^{-1}$  hullámszám értékek között változtak. A

spektrumokat transzmissziós üzemmódban vettem fel. Az osztályozási módszerek használatakor validálásként véletlenszerűségi tesztet (Y változó keverésével) és szisztematikus, három részre osztott kereszt-ellenőrzést használtam. Adatelőkezelésként standardizálást alkalmaztam (lásd. 2.5.1. fejezet).

### 5.1.3.1 PLS DA eredmények az eredeti spektrumok felhasználásával

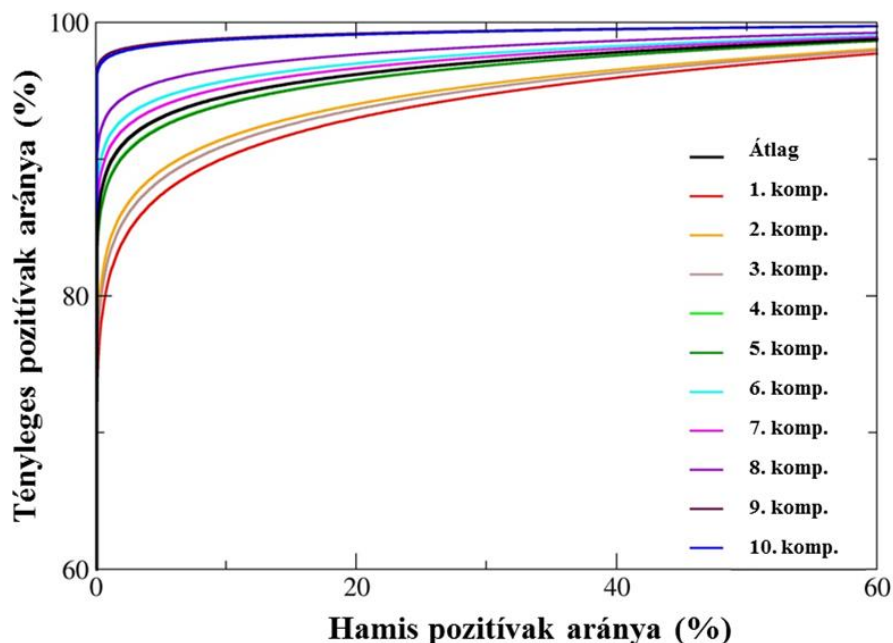
Elsőként a PLS DA módszert alkalmaztam, amely esetében a ROC görbék elkészítéséhez a becsült értékeket használtam. A kiértékeléshez szükséges PLS komponensek optimális száma öt volt a PRESS értékek alapján (globális minimum). Az első két PLS komponens egymás függvényében ábrázolva a **9. ábrán** látható, hogy bár a csoportok közel helyezkednek el egymáshoz, mégis jól elkülöníthetőek, csak néhány minta tekinthető félreosztályozottnak.



**9. ábra:** Az energiai talok csoportosítása PLS DA segítségével a cukortartalom szerint. A második PLS komponens ábrázoltam az első függvényében

Az összehasonlítás szempontjából szintén érdekesnek tartottam a különböző PLS komponens számok figyelembe vételével kapott eredményeket is. Ezért a PLS komponens számot egy és tíz között változtatva megfigyeltem, hogyan változik a modellek becslési képessége. A ROC görbék minden PLS komponens szám esetén a becsült értékek alapján határoztam meg, az előző fejezetben már említett módon. Szintén megalkottam az átlagos ROC görbét is a három csoport különálló ROC görbéihez tartozó AUC értékek felhasználásával. A **10. ábrán** a PLS komponens számokhoz tartozó átlag ROC görbék ábrázoltam (1 és 10 között), valamint a tíz ROC görbe átlagát is. Ez utóbbi alapján akár a PLS

komponensek számának meghatározását is elvégezhetjük (a.m. pszeudorang meghatározás, a modern kemometria egyik kulcskérdése).



**10. ábra:** A tényleges pozitívak aránya ábrázolva a hamis pozitívak arányának függvényében A PLS komponens számok függvényében alkotott átlagos ROC görbékét különböző színekkel jelöltem, míg a tíz ROC görbe átlagát fekete színnel. Az ábrát a jobb áttekinthetőség miatt a megfelelő részekre fókuszáltam.

A **10. ábrán** látható, hogy a ROC görbék átlagából 4-5 PLS komponens szám következik, ami azt jelenti, hogy öt PLS komponens szám kiválasztásával tudom a legjobban megközelíteni a becslési képességet úgy, hogy a modellem nem lesz sem túlillesztve, sem pedig alulillesztve. Ez a szám pedig megegyezik az eredetileg kiválasztott öt PLS komponenssel.

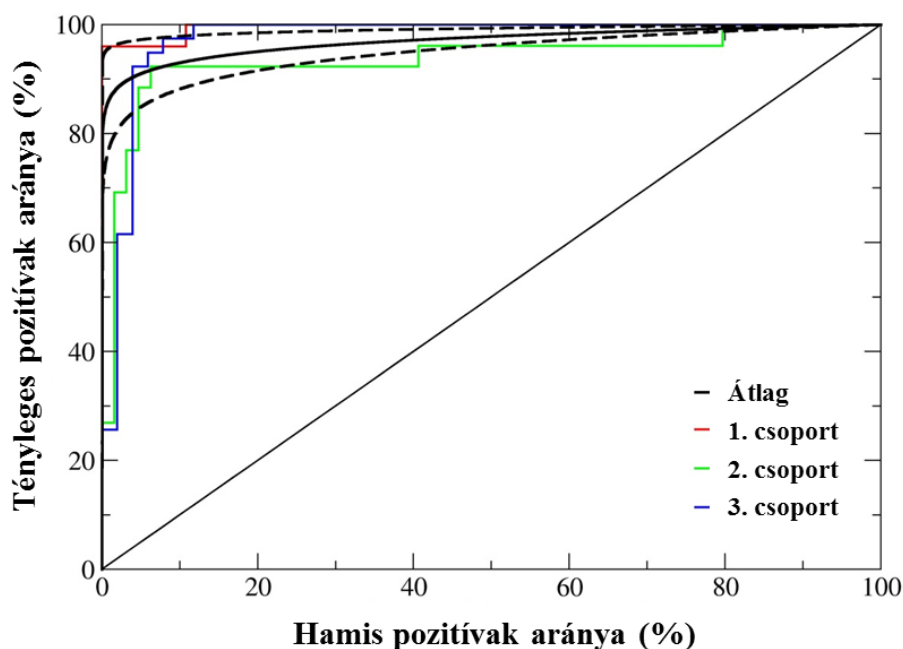
A csoportokra külön-külön megalkotott ROC görbék átlag AUC értékét felhasználva a Hanley-formulával ábrázoltam a végső ROC görbét, ami már a többi osztályozási módszer eredményeivel is összevethető. A későbbi összehasonlításhoz a PLS DA esetén az öt PLS komponenssel kapott modell ROC görbáját használtam fel.

### 5.1.3.2 RF eredmények az eredeti spektrum felhasználásával

A véletlen erdő (RF) módszere esetén első lépésben a modellépítés során szükséges paramétereket, a fák mennyiségét és a tulajdonságváltozók számát állítottam be az Anyag és módszer fejezetben már korábban részletezett módon (lásd. 4.5.5 fejezet). Ennek megfelelően a legrobosztusabb tulajdonságváltozó szám (a helyes osztályozási százalékok stabilitásán alapulva) negyven lett, amely esetben a helyes osztályozási százalékok (CC %) az első, második és harmadik csoportra nézve rendre 0,9600; 0,8077; és 0,8718 lettek. A folyamat során azt a

legalacsonyabb tulajdonságváltozó számot választottam ki, ahol a CC %-ok már nem változnak számottevően az újabb tulajdonságváltozók hozzáadásával. A kiválasztásra empirikus ajánlások is léteznek, pl.  $\sqrt{M}$  vagy  $\log_2(M+1)$ , ahol  $M$  a teljes változós szám, de ezek a lehetőségek csak kisebb változós számoknál működnek megfelelően. A fák mennyiségének 30-at választottam, ahol már a CC % értékek nem változtak láthatóan (lásd. 4.5.5 fejezet). Harminc fa használatával a harmadik csoport CC % értéke javult: 0,8718-ról 0,9487-re. Véletlenszerűségi tesztet valamint kereszt-ellenőrzést is használtam az osztályozási modell validálására.

A ROC görbék megalkotásához a helyes osztályozási százalékok helyett a mintákhoz és csoportokhoz tartozó valószínűség értékeket használtam fel. A valószínűség értékeket növekvő sorrendben rendeztem. A módszerek összehasonlításához a csoportokhoz tartozó ROC görbék átlagát használtam fel. A csoportok görbéi illetve az átlag ROC görbe is látható a **11. ábrán**.



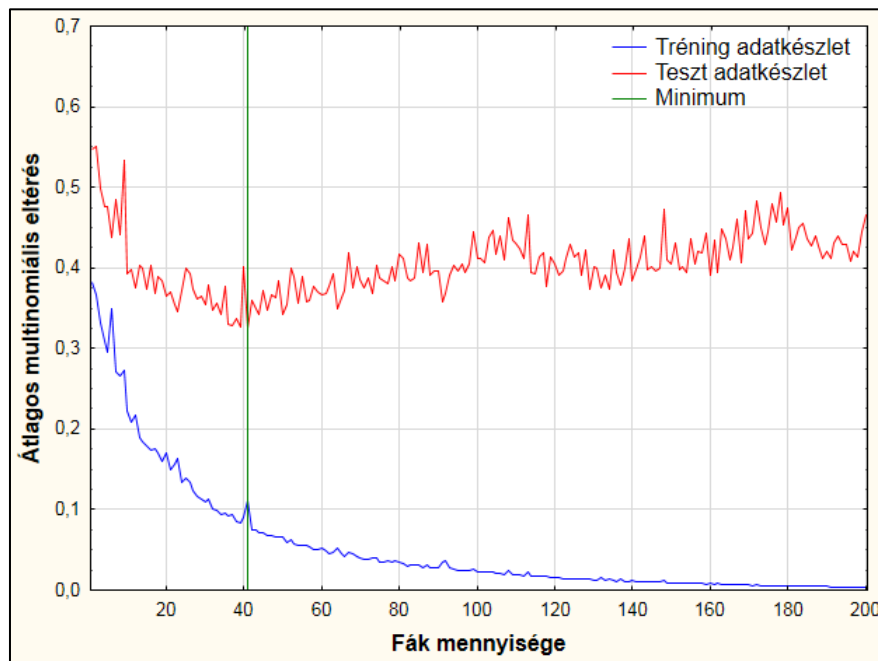
**11. ábra:** A tényleges pozitívak ábrázolása a hamis pozitívak arányának függvényében az véletlen erdő módszer esetén

*A csoportokhoz tartozó görbéket különböző színekkel, míg az átlag görbét fekete folytonos vonallal jelöltem. A szaggatott vonal jelenti az átlagtól való  $\pm 1$  SD eltérést.*

### 5.1.3.3 Fejlesztett fák alapján kapott eredmények az eredeti spektrumok felhasználásával

Ahogy azt már az „Anyag és módszer” részben ismertettem (lásd. 4.5.6 fejezet), a fejlesztett fák módszere nagyon hasonló az RF módszerhez, de ebben az esetben a modellépítés során a fák mennyiségét és a fa méretét kell optimalni. Ehhez az optimalizáláshoz, az átlagos multinomiális eltérések globális minimumát a fák mennyiségének függvényében kerestem meg.

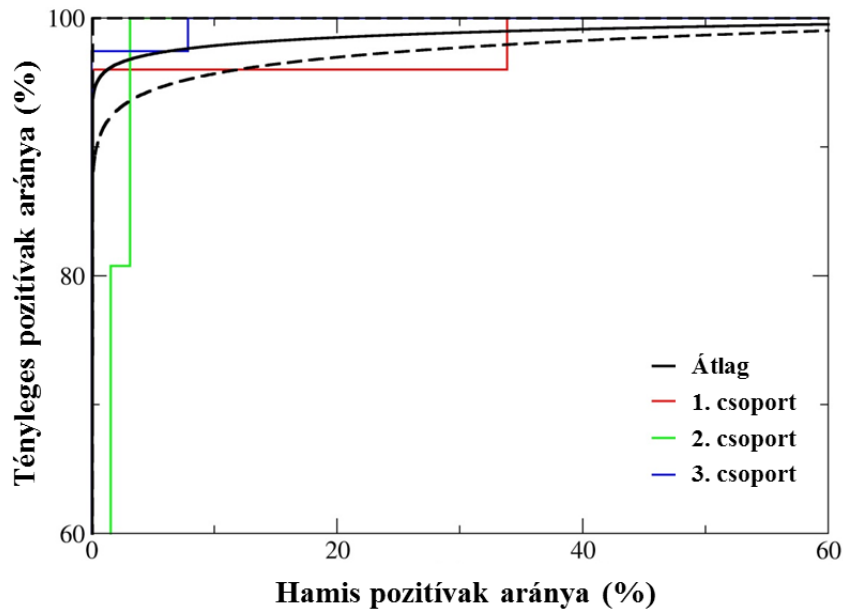
A **12. ábrán** ez a folyamat jól követhető. Esetemben az átlagos multinomiális eltérés alapján, a fák mennyisége 41 volt, amelyet a teszt adatkészlet figyelembe vételével határoztam meg.



**12. ábra:** A hibaértékek ábrázolása a fák mennyiségének függvényében  
*A hibaértékek globális minimumát zöld függőleges egyenes szemlélteti.*

Az előrebecslési (teszt) adatkészletet véletlenszerűen generáltam minden egyes lépésben, az adatkészlet 50 %-át felhasználva. A maximális fa méret három lett. Ebben az esetben is kereszt-ellenőrzést és véletlenszerűségi tesztet alkalmaztam validálásként.

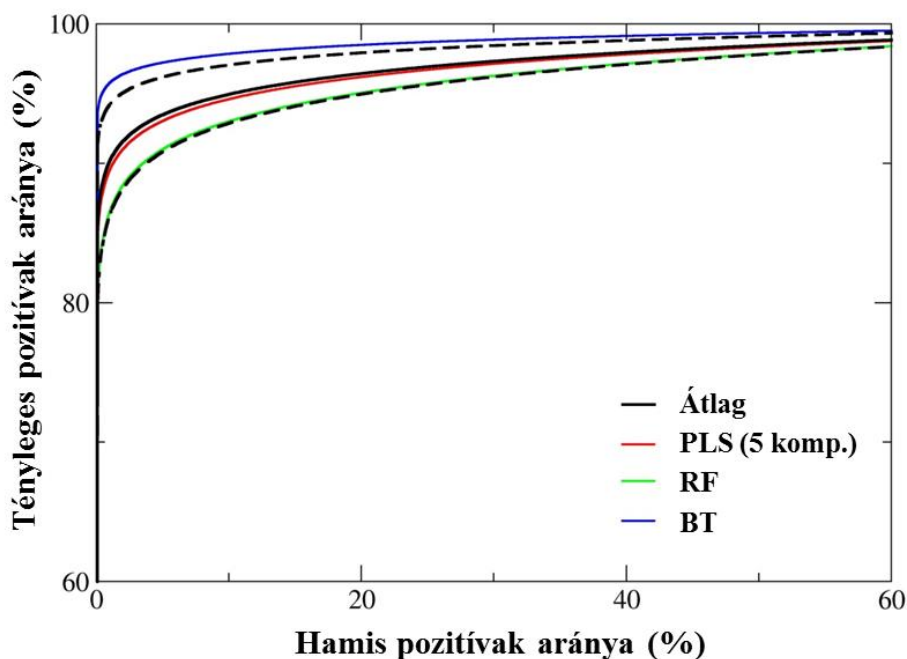
Hasonlóképpen a véletlen erdő módszeréhez, a ROC görbék megalkotásához itt is a valószínűségi értékeket használtam fel. A **13. ábra** alapján látható a csoportonkénti és az átlag ROC görbe is. Szembetűnő, hogy a módszer majdnem tökéletes eredményt adott, ezért is volt szükség, ebben az esetben is, a kapott ábra nagyítására.



**13. ábra:** A tényleges pozitívak arányának ábrázolása a hamis pozitívak függvényében a fejlesztett fák módszere esetén

*Folytonos fekete vonallal az átlagos ROC görbét, szaggatott vonallal az átlagos görbe  $\pm 1$  SD eltérését szemléltetem. A kép a megfelelő láthatóság érdekében fókuszált.*

Végül az összehasonlításhoz elkészítettem a különböző osztályozó módszerek (PLS DA az optimális 5 komponenssel, RF, BT) felhasználásával kapott átlagos ROC görbék ábráját is, valamint a végső átlagos ROC görbét is. A **14. ábra** alapján látható, hogy az energiatalok osztályozása során a fejlesztett fák módszere a legjobb. Bár a másik két módszer rosszabb végeredményt adott, összefoglalásképpen elmondható, hogy mind a három módszer jó választás lehet az energiatalok osztályozásakor, hiszen a ROC görbéik messze jobbak voltak, mint a véletlen osztályozással kapott egyenes (átló). Szintén fontos kiemelni, hogy bár a fejlesztett fák módszere korántsem annyira ismert, mint az LDA vagy PLS DA módszerek, megbízható és jobb eredményt adott az előbbieknél. A fejlesztett fák módszere még az átlagnál és annak  $\pm 1$  SD hibájánál is jobban teljesített, vagyis szignifikánsan ( $H_0$ = A ROC görbék átlaga és az adott görbe között nincs szignifikáns különbség,  $\alpha=0,05$ ) jobbnak tekinthető a másik két módszernél.



**14. ábra:** Az osztályozási módszerek ROC görbék alapján történő összehasonlítása  
*A fekete folytonos vonallal az átlagot, szaggatott vonallal pedig az átlaghoz tartozó  $\pm 1$  SD szórást jelöltem. Az eredeti ábra fókuszált verziója látható a jobb láthatóság érdekében. RF és BT jelölések rendre a véletlen erdő és fejlesztett fák módszerét jelentik.*

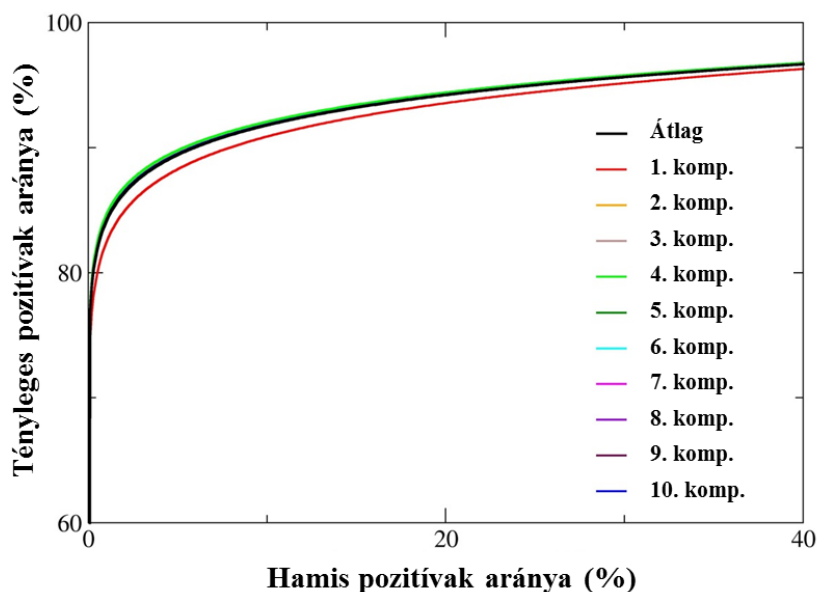
#### 5.1.3.4 PLS DA eredmények a PCA adatkészlet alapján

A második adatkészlet ugyanazon energiainformációkat tartalmazta, viszont az eredeti spektrum helyett a spektrumok PCA-ja során kapott első tíz főkomponenst. Ezzel a megközelítéssel már a lineáris diszkriminancia elemzés (LDA) is bekerülhetett az összehasonlítható módszerek közé, hiszen így a változók mennyisége már kevesebb volt, mint a mintaszám. (A korábbi esetben a konzisztencia miatt nem akartunk változókat kihagyni, ezért kellett az LDA nélkül elvégezni az elemzést.) Az első tíz komponens használata elégnak bizonyult, mivel összességében a magyarázott variancia százalékának 99,7 %-át lefedték.

A PLS DA módszer esetén a PLS komponens szám meghatározása még egyszerűbben történt, ugyanis az  $R^2$  és PRESS értékek is két komponens után már alig látható módon változtak. Ez a PLS komponens szám csökkenés már az eredeti adatmátrix redukciójából is várható volt. Ez azt jelentette, hogy az első két PLS komponens használata elég a modellépítéshez. Ettől függetlenül látni szerettem volna, hogy további PLS komponensek hozzáadásával növelhető-e a modell jósága, ezért a ROC görbéket egészen tíz PLS komponens számig ábrázoltam a korábban már ismertetett módon. A **15. ábrán** látható eredmény alapján két komponens használatán túl a



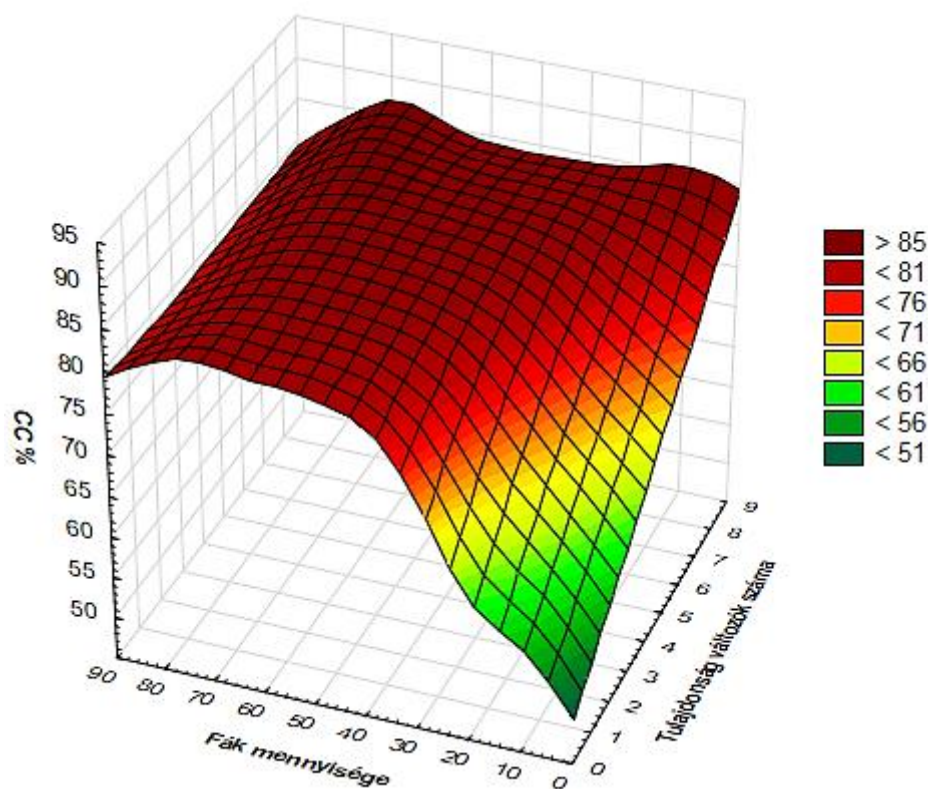
görbék már nem változnak, ami ebben az esetben is igazolta a PRESS értékek alapján történő két PLS komponenses kiválasztást.



**15. ábra:** Tényleges pozitívak arányának ábrázolása a hamis pozitívak függvényében. A különböző PLS komponens számmal kapott ROC görbéket különböző színekkel jelöltem, míg az átlagos ROC görbét fekete folytonos vonallal. Az ábra az eredeti fókuszált verziója.

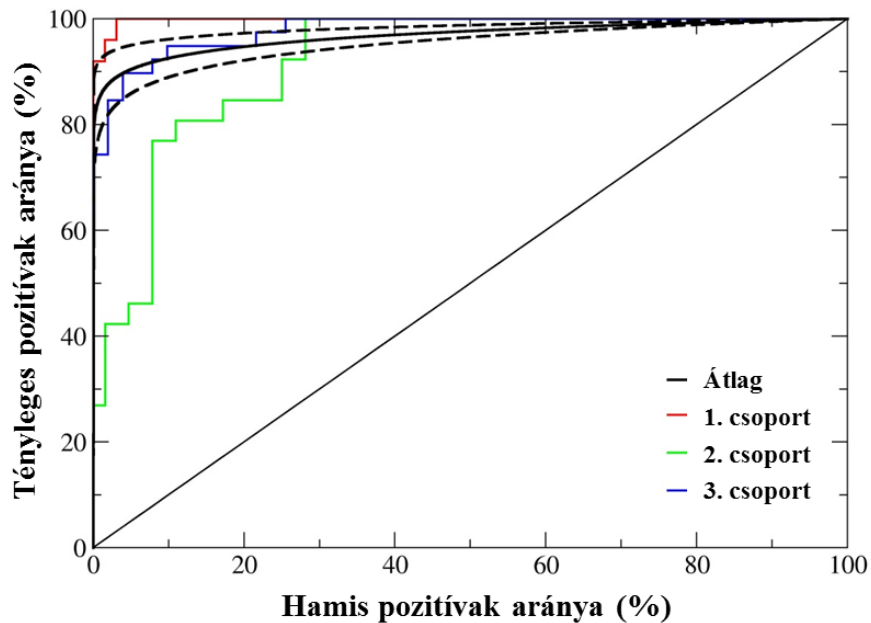
### 5.1.3.5 RF eredmények a PCA adatkészlet alapján

Első lépésként ebben az esetben is a fák mennyiségének és tulajdonság változók számának optimalítása volt a feladat. Viszont most egy háromdimenziós felületet illeszttem távolság-súlyozású legkisebb négyzetek módszere segítségével: a helyes osztályozási százalékokat (CC %) ábrázoltam a tulajdonságvektorok száma (1 és 8 között) valamint a fa mennyiségek (10 és 80 között) függvényében. A **16. ábra** ezt az illesztett felületet szemlélteti, amelyen látható, hogy a CC % a fák és a tulajdonság változók mennyiségének növelésével együtt emelkedik. Ez csak egy ideig tart, ugyanis a felület közepétől állandóvá kezd válni az érték, sőt hetven fa érték fölött egy kisebb csökkenés következik be. Az illesztés alapján harminc fát és öt tulajdonság változót választottam ki, amelyek az „elfogadhatónak” ítélt régió elején helyezkedtek el. Ez azt jelenti, hogy ezek az értékek a legkisebbek, amelyek használatával jó osztályozást hajthatok végre a további kiértékelés során és a túlillesztést is el tudom kerülni. A CC % ebben a kombinációban az egyik legnagyobb értéket vette fel.



**16. ábra:** A helyes osztályozási százalékok (CC %) ábrázolása a fák mennyisége és a tulajdonság változók mennyiségének függvényében távolság-súlyozású legkisebb-négyzetes illesztés segítségével

A modellépítést, valamint a valószínűségi értékek felhasználásával minden egyes csoport esetén a ROC görbék megalkotását, a megfelelő paraméterek segítségével végeztem el. A **17. ábrán** látható a csoportonkénti és az átlag ROC görbe is, melyek mindegyike jobb eredményt mutat a véletlen osztályozásnál. Összehasonlítva az eredményt a spektrumokat tartalmazó adatkészlet alapján kapottal, lényeges változás nem tapasztalható a két ábra között, ami várható is, ha nem lépett föl lényeges információvesztés a PCA során.



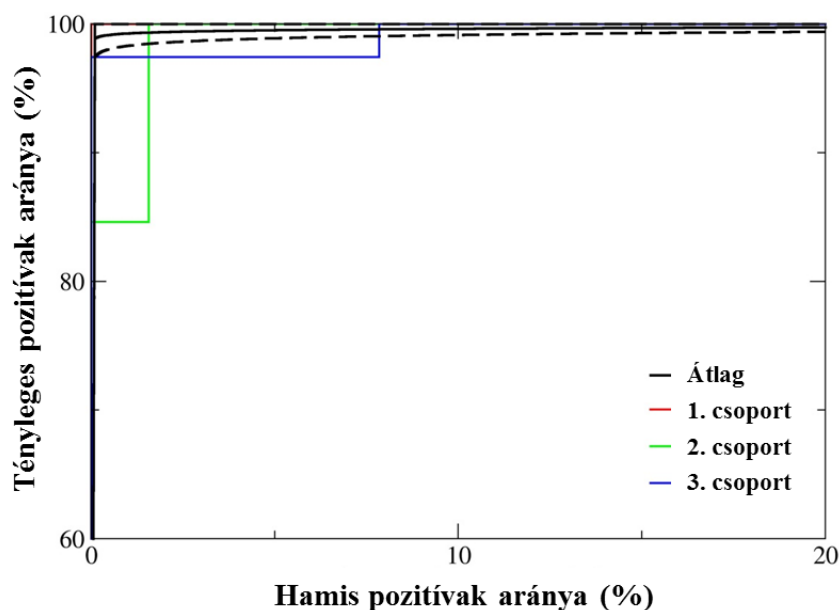
**17. ábra:** Az RF módszerrel kapott ROC görbék

*A tényleges pozitívak arányát a hamis pozitívak arányában ábrázoltam. Fekete folytonos vonal szemlélteti az átlagos ROC görbét a három csoportra nézve, szaggatott vonal pedig a hozzá tartozó  $\pm 1$  SD értéket.*

#### 5.1.3.6 BT eredmények a PCA adatkészlet alapján

A fejlesztett fák módszerénél az előző adatkészlethez hasonlóan itt is a fák mennyiségét és a fa méretét optimaltam. Az átlagos multinomiális eltérés alapján a globális minimum 190 fa értéknél volt (Melléklet **M5**). Az előző, és ezen eredmények tükrében felfedezhetünk egy inverz kapcsolatot a fák mennyisége és az eredeti változók száma között, mivel a spektrális (nagyobb) adatkészlet esetén egy jóval kisebb mennyiséget: 41 fát választottam ki optimumként. A modellek validálására ezúttal is véletlenszerűségi tesztet és kereszt-ellenőrzést használtam.

Az osztályozási modell építése után a ROC görbéket az energiatartalom csoportokra és az átlagos AUC értékekre is ábrázoltam. A ROC görbék elkészítéséhez szintén a valószínűség értékeket használtam fel. A **18. ábrán** nagytáblázatosan látható a ROC görbék egy része a jobb láthatóság érdekében. Az eredmény nagyon érdekes, hiszen mind az átlagos görbe, mind a csoportonkénti ROC görbék rendkívül közel helyezkednek az AUC érték maximumhoz (majdnem 1), így ebben az esetben az osztályozás nagyon jónak tekinthető.

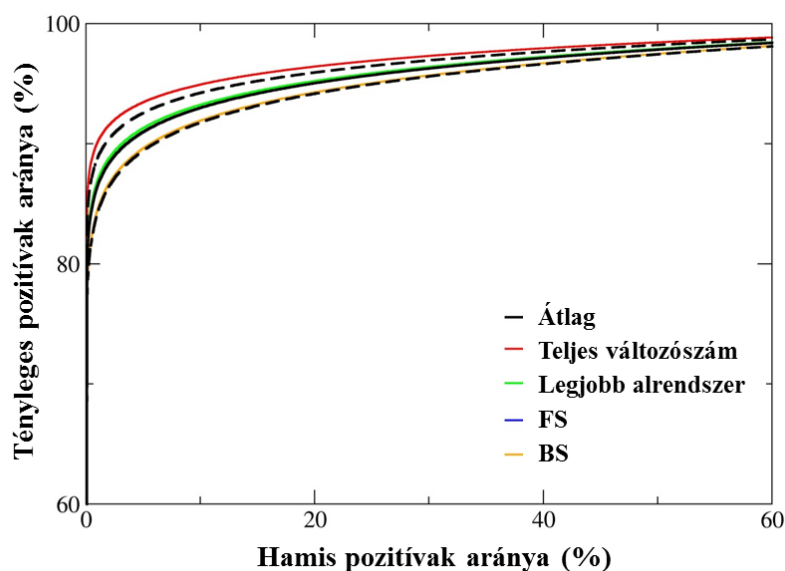


**18. ábra:** A tényleges pozitívák arányának ábrázolása a hamis pozitívák függvényében a fejlesztett fák módszere és a PCA adatkészlet esetén

*Fekete vonal szemlélteti az átlagos ROC görbét, a szaggatott vonal pedig a hozzá tartozó  $\pm 1$  SD értéket. Az ábra az eredeti fókuszált verziója.*

#### 5.1.3.7 LDA eredmények a PCA adatkészlet alapján

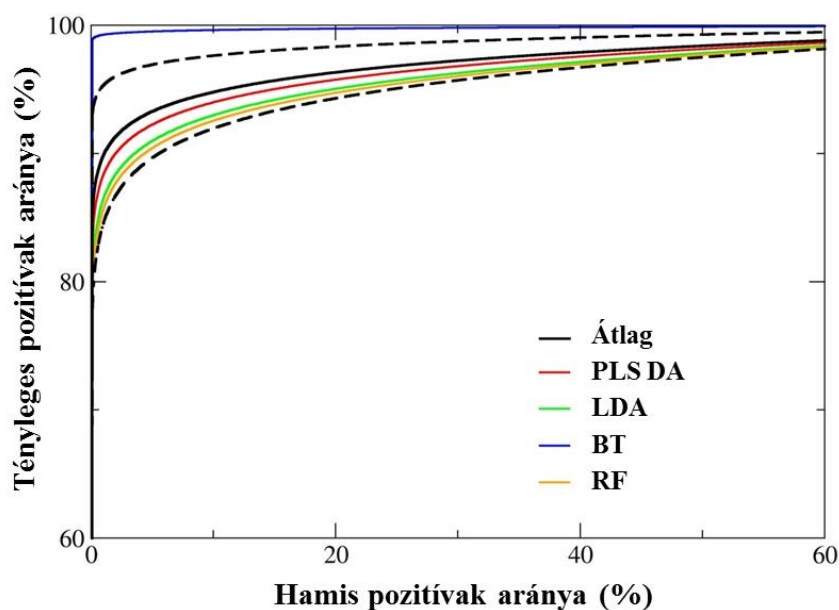
A lineáris diszkriminancia elemzés során szükséges hasznos változókat a lépésenkénti változó hozzáadással, lépésenkénti változó törléssel és a legjobb alrendszer módszerek segítségével választottam ki. Ezen felül építettem egy modellt a teljes (tíz) változós szám felhasználásával is. A legjobb alrendszer kiválasztásánál a maximális változós szám öt volt, valamint Wilk's lambda statisztikát alkalmaztam a szignifikáns komponensek meghatározására. A lépésenkénti változó hozzáadás és lépésenkénti változó törlés esetén egy változó akkor épült be, vagy törlődött ki, ha a hozzáadása vagy törlése 5 %-os szinten szignifikáns változást eredményezett a modellben. Összességében elmondható, hogy bármelyik kiválasztási módszerrel illetve a teljes változó mennyiséggel is megfelelő osztályozáshoz jutottam. A különböző változó kiválasztási módszerekkel megalkotott modellek és ROC görbéik alapján (**19. ábra**) elmondható, hogy a legjobb eredményt a teljes változós szám használatával értem el, de ezt követte a legjobb alrendszer módszere, majd végül a lépésenkénti változó hozzáadás és lépésenkénti változó törlés következett. Mind a négy modell AUC értéke a használhatóság szempontjából megfelelt az elvárásoknak.



**19. ábra:** Az LDA elemzés során kapott ROC görbék különböző változó szelektálási lehetőségek használatával

*Az átlag AUC értékhez tartozó ROC görbe feketével van jelölve, a hozzá tartozó  $\pm 1$  SD pedig szaggatott vonallal. FS a lépésenkénti változó hozzáadást, BS pedig a lépésenkénti változó törlést jelenti. Az ábrán az előbbi két modell ROC görbéi egymást fedik, ezért csak a piros és a zöldes szín látható. Az ábra az eredeti fókuszált verziója.*

A végső eredmény, amelyen mind a négy módszerhez tartozó ROC görbét feltüntettem, a **20. ábrán** látható. Ekkor minden használt módszer esetén az átlag ROC görbékét vettem alapul, amelynek szerepe a PLS DA és LDA módszereknél volt leginkább látható, hiszen az RF és BT módszerek esetén csak egy átlag ROC görbe volt a három csoportra nézve. Az eredmények tükrében a fejlesztett fák módszere lényegesen jobb eredményt produkált az összes többi módszerhez képest. Hozzá kell tenni viszont, hogy az energiatartalom eltérő cukortartalmú csoportjaira nézve minden módszer megfelelően jó AUC értékekkel és így megfelelő osztályozással rendelkezett.



**20. ábra:** A vizsgált négy osztályozási módszer átlagos ROC görbéi

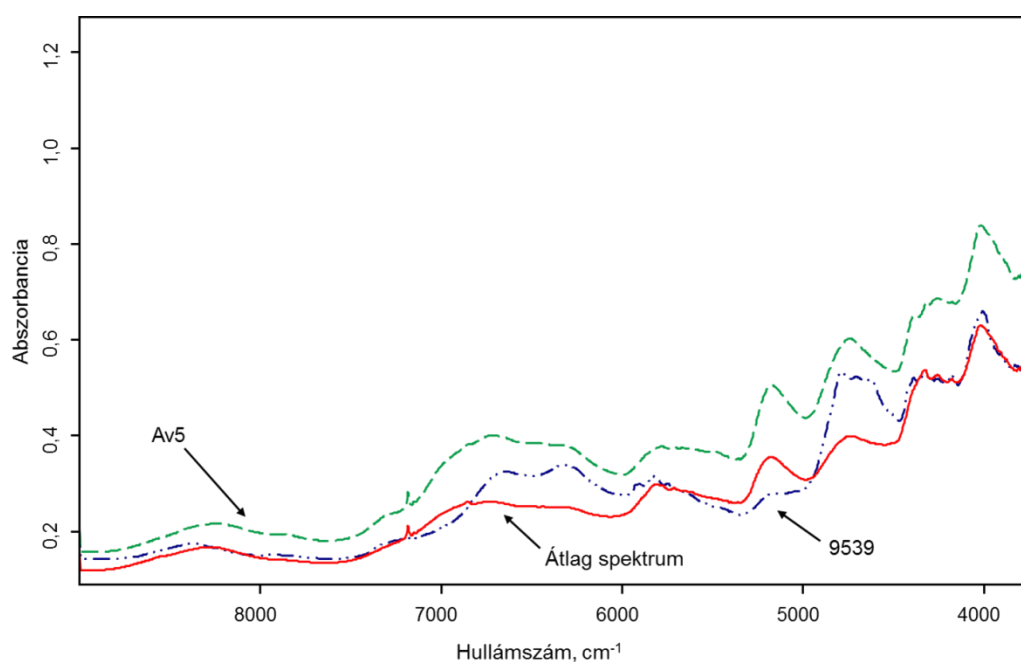
*A tényleges pozitívák arányát ábrázoltam a hamis pozitívák függvényében. Az átlag AUC értékhez tartozó ROC görbe feketével van jelölve, a hozzá tartozó  $\pm 1$  SD pedig szaggatott vonallal. Az ábra az eredeti fókuszált verziója a jobb láthatóság érdekében.*

Az általam továbbfejlesztett több osztályos ROC görbe módszere egy, a korábbi módszereket felülmúló lehetőség nemcsak osztályozási, de egyéb módszerek és modellek összehasonlítására is. Ezzel a lehetőséggel megoldást nyújthatunk olyan problémákra, mikor a CC % kiszámolása nehézkes vagy megoldhatatlan feladat lenne. Az AUC értékek ezzel szemben a valószínűségi értékekre vagy akár becsült értékekre is jól használhatóak. A nevében is szereplő „több csoportnak” megfelelően, nem csak három, hanem számos csoport esetén is kiválóan alkalmazható. Az átlagos ROC görbék vizualizációja pedig megfelelő lehetőség a gyors összehasonlításra.

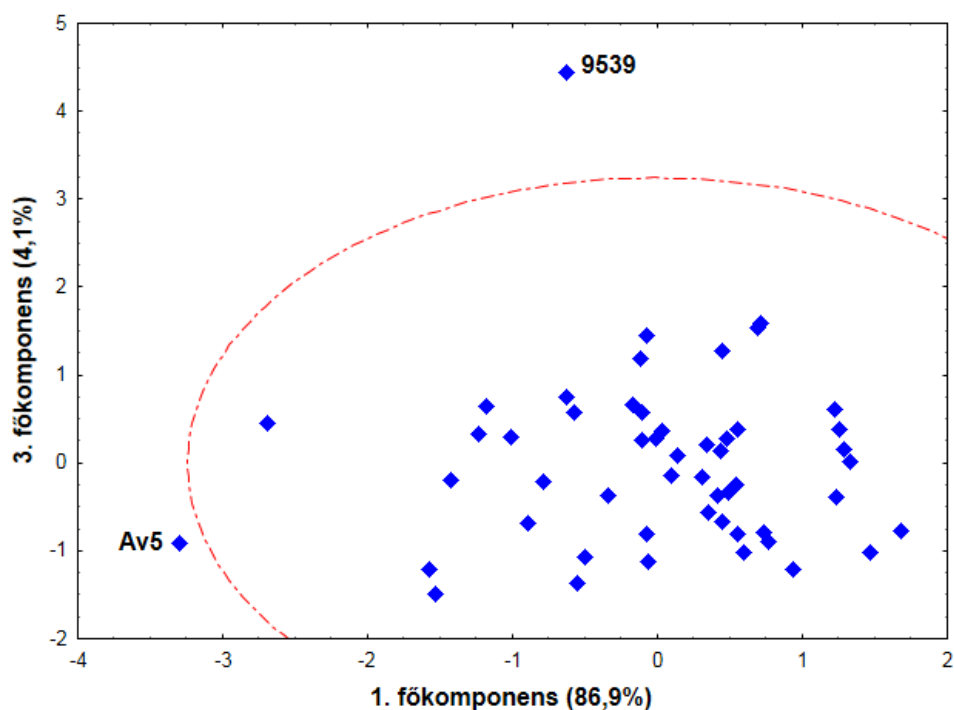
Az energiatalok esetében egyértelműen megállapíthattam, hogy a fejlesztett fák osztályozási módszere adta a legjobb eredményt, de összességében mindegyik módszer jól tudta az energiatalokat cukortartalom szerint osztályozni. A fejlesztett fák módszerének (előzetesen nem várt) sikere és az LDA módszerénél jobb eredménye arra ösztönözhet mindenkit, hogy bátran merjünk kevésbé ismert, vagy kevésbé bevált módszerekkel is próbálkozni, hiszen néha a megoldás kulcsa éppen bennük van az adott problémára.

## 5.2 Q10 tartalmú étrendkiegészítők hatóanyagtartalmának vizsgálata

Összesen 52 különböző Q10 koenzim tartalmú étrendkiegészítőt vizsgáltam meg, amelyek Q10 koenzim tartalma az előzetes HPLC mérések alapján a 6,9 és 118,1 mg / g közötti tartományban volt. A porított készítmények FT-NIR spektrumait diffúz reflexiós mérési módban  $12500 - 4000 \text{ cm}^{-1}$  (800–2500 nm) hullámszám tartományban vettem fel. Minden mintából két párhuzamos almintát készítettem, s ezek spektrumát rendre kétszer vettem fel. A további kemometriai elemzésekhez a minták átlagspektrumai használtam fel. A spektrum  $12500 \text{ cm}^{-1}$  és  $9000 \text{ cm}^{-1}$  (800–1111 nm) közötti részét kihagytam a kiértékelésből, mivel a spektrum ezen szegmense nem hordozott szisztematikus információt számomra. Az első lépésként a lehetséges spektrális kieső mintákat szűrtem ki az adatkészletből. Ehhez főkomponens-elemzést (PCA) használtam (STATISTICA 12, Statsoft Inc., Tulsa, OK, USA). A módszer segítségével két spektrális kieső mintát találtam, amelyek az átlag spektrummal összehasonlítva is jelentős eltérést mutatnak, így a **21. ábrán** látható Av5 és 9539 kódszámú minták spektrumainak a többi spektrumtól való szemmel látható eltérését kemometriai módszerrel is igazolni tudtam. A két említett minta a 99 %-os konfidencia intervallumon kívülre esett. A főkomponens-elemzés során kapott eredményeket a **22. ábrán** mutatom be.



**21. ábra:** A két spektrális kieső minta és az átlag FT-NIR spektrumának összehasonlítása  
*Az abszorbancia értékeket a hullámszámok függvényében ábrázoltam. A zöld szaggatott vonal az Av5 minta, a kék (duplán) pontozott-szaggatott vonal a 9539 minta, a piros folytonos vonal pedig az átlag spektrumot jelenti.*



**22. ábra:** Spektrális kieső minták (kiugró értékek) keresése főkomponens-elemzéssel  
 A 3. főkomponenst az 1. főkomponens függvényében ábrázoltam. A piros pontozott-szaggatott vonal a 99 %-os konfidencia intervallumot jelzi. A főkomponensek által magyarázott variancia értékei zárójelben találhatóak a nevük mögött.

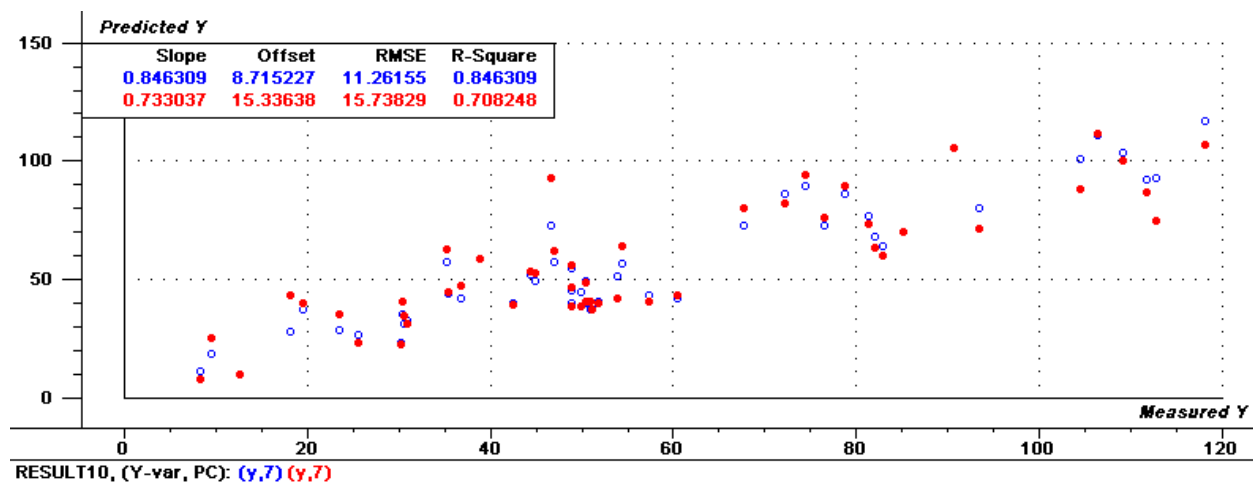
A két spektrális kieső minta elemzésekor kiderült, hogy ezek a minták nagy mennyiségben tartalmaztak egyéb vitaminokat is, amelyek jelentős változást okoztak a spektrumokban.

A kalibrációs modellek elkészítését az Unscrambler 9.7 szoftver (CAMO Software, Oslo, Norway) segítségével végeztem. Az első modellt az eredeti adatkészletre hoztam létre, amely a spektrális kieső minták kivételével összesen ötven mintát tartalmazott. Ebben az esetben a teljes spektrumot figyelembe vettem ( $9000\text{ cm}^{-1}$  és  $3800\text{ cm}^{-1}$ ). Ekkor még adatelőkezelési módszereket sem alkalmaztam a modellépítéshez, mert létre akartam hozni egy olyan modellt, amely összehasonlítás alapként szolgálhat a későbbiekben a különböző adatelőkezelési és változókiválasztási módszerek használatával kapott modellekhez. Ahogy azt a későbbiekben bemutatom, az adatelőkezelési módszerek használata nem javította számottevően a modellépítés jóságát. Ezen modellek végső paramétereit, a további modellek paramétereivel együtt a **6. táblázatban** foglaltam össze.

A modellépítések során az előzetesen HPLC segítségével megkapott koncentráció eredményeket használtam fel  $\mathbf{Y}$ , vagyis referencia változóként. Az első, még kezdeti kalibrációs



modell megépítéséhez 7 PLS komponensre volt szükségem. A **23. ábrán** az elkészített modell alapján becsült  $\hat{Y}$  értékeket a referencia  $Y$  értékek függvényében ábrázoltam.



**23. ábra:** Az eredeti, kiindulási Q10 koenzim koncentráció kalibrációs és validált modellje. Kékkel a kalibrációs, pirossal pedig a validált modellt jelöltem. A becsült  $\hat{Y}$  értékeket a referencia  $Y$  értékek függvényében ábrázoltam.

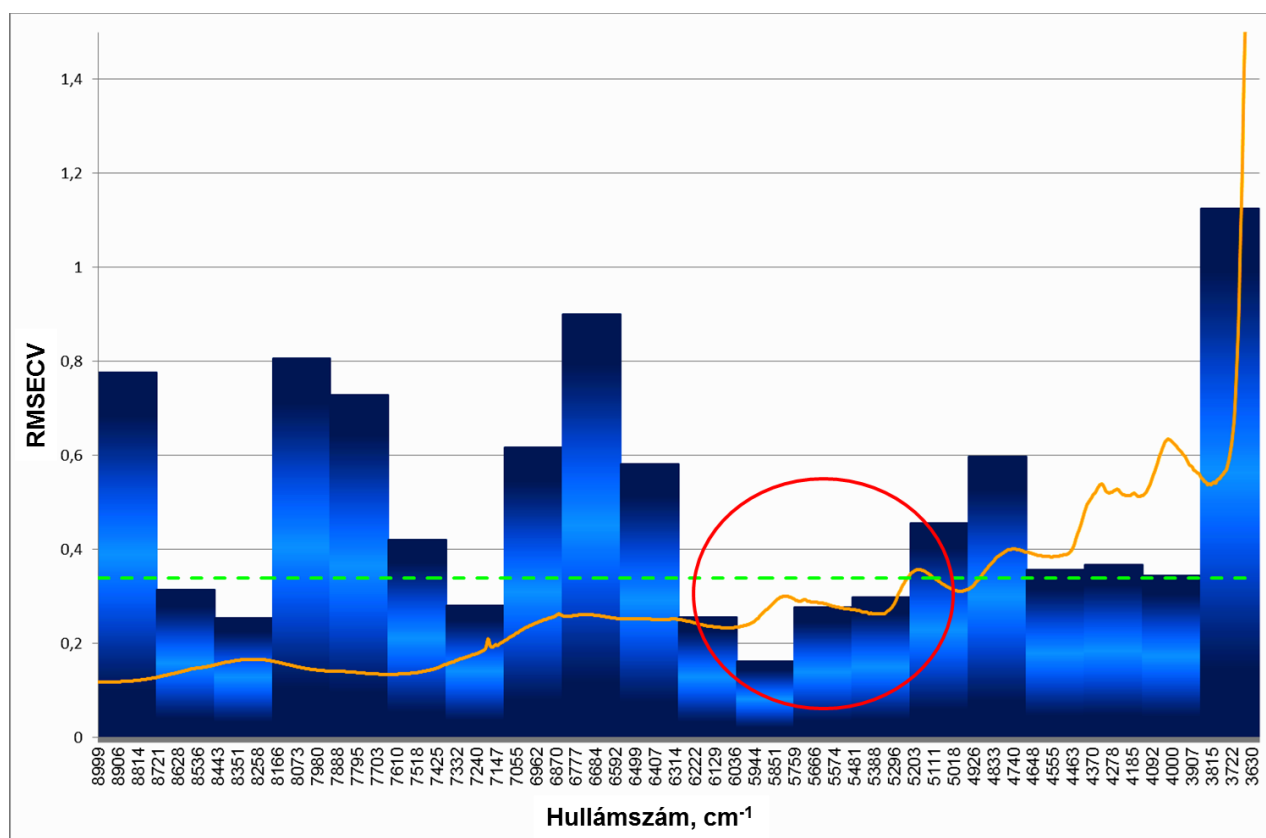
Az  $R^2$  érték 0,85 lett a kalibrációs modellre nézve a  $Q^2$  érték pedig 0,71 a validált modell esetén. Az átlagos négyzetes hiba a kalibrációra (RMSEC) 11,26 mg / g, míg a kereszt-ellenőrzéssel (RMSECV) kapott modellre nézve 15,74 mg / g lett. Öt részre osztott kereszt-ellenőrzést alkalmaztam validálási lépésként véletlen minta kiválasztással. Hangsúlyozandó, hogy ez a metodika a további modellek esetén is érvényes maradt. Ezek az értékek még messze nem érték el azt, ami elvárható, így mindenképp szükséges volt modellfejlesztési lehetőségeket bevetni.

Az ezt követő modellépítések során már többféle adatelőkezelési és változókiválasztási módszerrel dolgoztam a minél előnyösebb modell megalkotása érdekében.

### 5.2.1 Az iPLS változókiválasztással kapott eredmények

Mivel a kezdeti modellem semmiképp sem volt tökéletesnek mondható, a fejlesztési lépést a változók számának csökkentésével kezdtem meg. A változókiválasztási módszerek közül az iPLS egy rendkívül egyszerű és elterjedten alkalmazott technika. Segítségével kiválaszthatjuk azokat a változó intervallumokat, amelyek a modellépítés szempontjából fontosabbak. Elsőként az adatkészletet tíz részre osztottam fel, de ez nem vezetett megfelelő eredményre, mert a változók túl nagy csoportokat alkottak. A második esetben húsz részre osztottam a változókat és ez már megfelelőnek bizonyult a további elemzésekhez. Így az iPLS módszert az adatelőkezelés

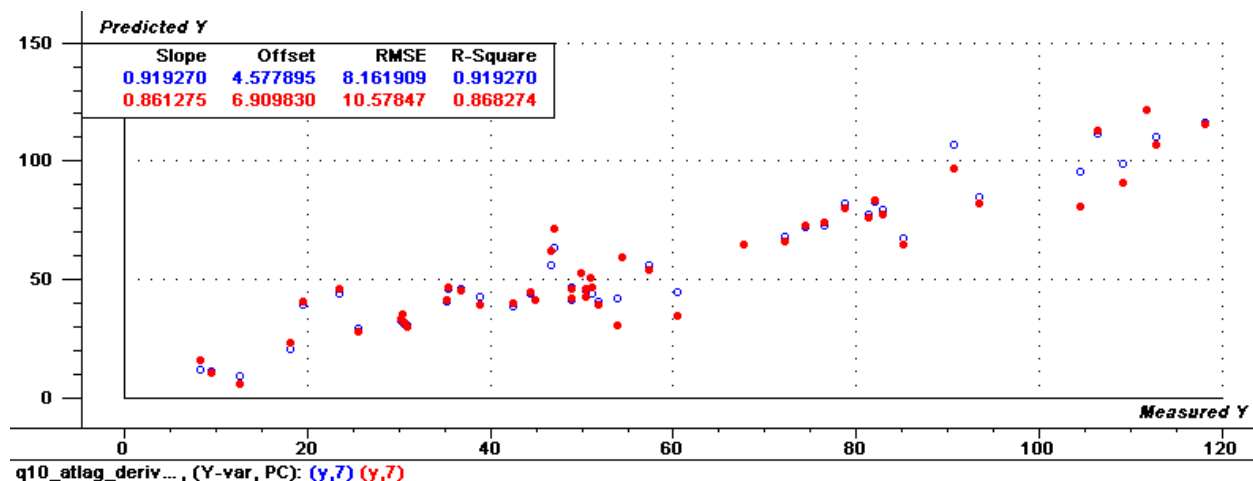
nélküli adatkészletre alkalmaztam húsz részre felbontva a mátrixot. Az iPLS alapelvénél fogva így 20 PLS elemzést hajtott végre egyenként az összes változó intervallumra. Az RMSEC és RMSECV értékeket minden részre külön-külön kiszámoltam. Végül az intervallumok függvényében ábrázoltam az RMSECV értékeket, amely segített a megfelelő intervallumok kiválasztásában, hiszen minél kisebb az adott intervallum RMSECV értéke, annál fontosabbnak tekinthető. A **24. ábra** megfelelően szemlélteti az iPLS végeredményét. A már említett RMSECV értékeket normalizált alakban az intervallumok függvényében ábrázoltam, valamint szaggatott vonallal jeleztem a teljes spektrum használatával kapott normalizált RMSECV értéket is. Az ábrán az eredeti átlagspektrumot is feltüntettem. Az RMSECV értékek skálázására a spektrummal együtt történő ábrázolás miatt volt csak szükség.



**24. ábra:** Az iPLS változókiválasztás RMSECV értékei a hullámszám függvényében  
Zöld szaggatott vonal jelzi a teljes spektrumhoz tartozó RMSECV értéket, illetve okkersárgával  
jelöltem az eredeti átlagspektrumot. A piros kör a kiválasztott öt intervallumokat szemlélteti.

Öt szomszédos intervallumot (összesen 350 változó) választottam ki az RMSECV értékek alapján a további modellépítéshez ( $6303 - 4953 \text{ cm}^{-1} = 1587 - 2019 \text{ nm}$ ). Az utolsó szomszédos intervallum ( $5003 - 4953 \text{ cm}^{-1}$ ) bár nagyobb RMSECV értékkel rendelkezik a többi intervallumhoz képest, de a szükségessége az  $R^2$  értékek és a spektrum fontosnak ítélt hullámszámjai alapján indokolt volt. Az eredeti adatkészletet deriváltam a modellépítéshez. A

kiválasztott intervallumokat alkalmazva a modellhez hét PLS komponensre volt szükség. A **25. ábrán** látható a modellezés végeredménye, amely alapján a kalibrációs modell  $R^2$  értéke 0,92 lett (kékkel jelölve), a  $Q^2$  pedig 0,87 lett (pirossal jelölve). A hiba értékek a következőképp alakultak: az RMSEC érték 8,16 mg / g lett, míg az RMSECV 10,58 mg / g. Elmondható tehát, hogy az iPLS módszerrel alkotott modell sokkal kisebb hibával jelez előre, mint a kiindulási modell főként ha figyelembe vesszük, hogy  $Q^2$  értéke 0,71-ről 0,87-re emelkedett.



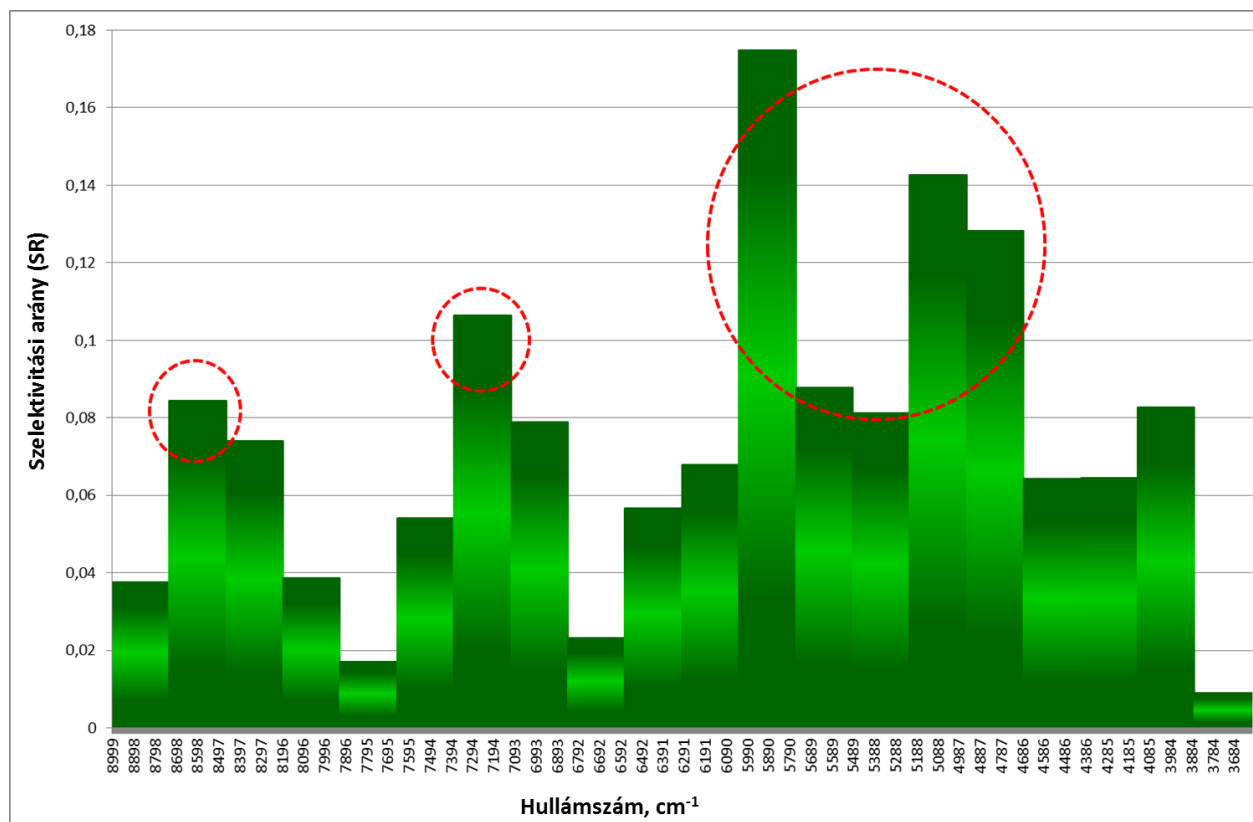
**25. ábra:** Az iPLS módszerrel kapott végső kalibrációs és validált modell

*A becsült Y értékeket a referencia Y értékek függvényében ábrázoltam. A kalibrációs modellt kékkel, a validált modellt piros színnel jelöltem.*

A kiválasztott spektrumtartományok különböző funkcionális csoportokhoz és kötésekhez rendelhetők. Ilyenek a metil csoport aszimmetrikus és szimmetrikus C–H nyújtó rezgéseinek, a metilén aszimmetrikus és szimmetrikus C–H nyújtó rezgéseinek, a triszubsztituált alkénhez tartozó C–H nyújtó vibrációs rezgésnek és az 1,4-kinon csoport rezgésének a felharmónikusai (Workman et al., 2007; Workman, 2000).

### 5.2.2 iSR változókiválasztással kapott eredmények

A második változókiválasztási módszerként a saját fejlesztésű intervallum szelektivitási arányt választottam. Az eljáráshoz ugyanazt a hús intervallumot használtam fel, mint amelyet az iPLS módszernél is alkalmaztam. Az RMSECV és  $R^2$  értékek segítségével meghatároztam a szelektivitási arányt minden intervallum esetében a 18. egyenlet alapján. A szelektivitási arányokat a **26. ábrának** megfelelően, az intervallumok számának és a spektrum hullámszámainak függvényében ábrázoltam.



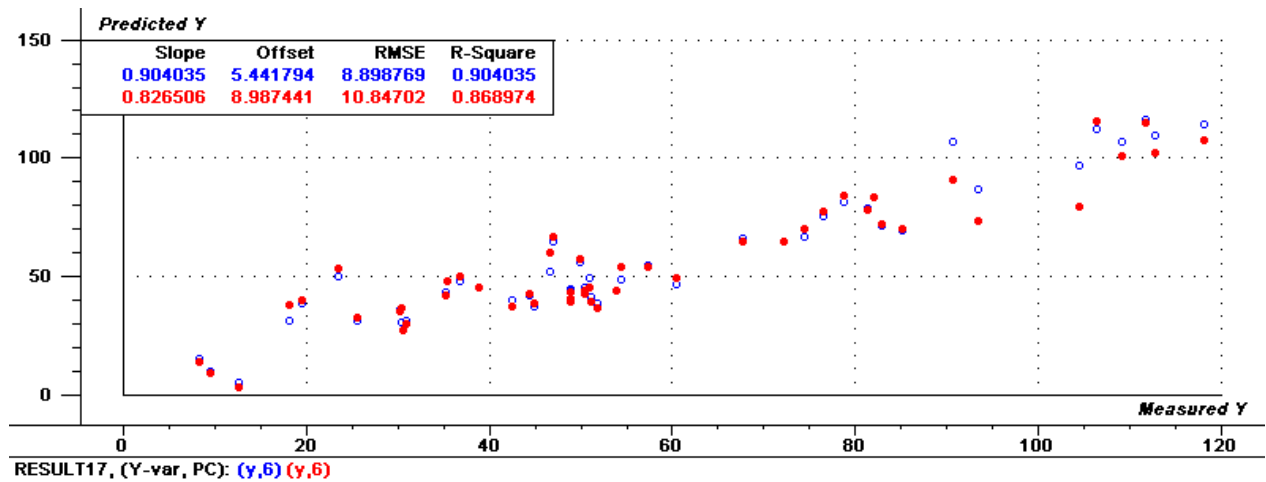
**26. ábra:** iSR értékek ábrázolása az intervallumok függvényében

Az *x* tengelyen a hullámszámokat tüntettem fel. A kiválasztott intervallumokat piros szaggatott körök jelzik.

Minél nagyobb a szelektivitási aránya egy intervallumnak, annál fontosabbnak tekinthető a kiértékelés szempontjából. A **26. ábra** alapján a következő spektrum intervallumokat választottam kik: 8728–8462, 7378–7112 és 6032–4682 cm<sup>-1</sup> (1146–1182, 1355–1406 és 1658–2136 nm). Az így kiválasztott 490 változóval elvégeztem a további modellépítést a derivált adatkészletre. Ehhez hat PLS komponensre volt szükség. Végül a becsült *Y* értékeket ábrázoltam a referencia *Y* értékek függvényében, amelynek eredményét a **27. ábra** szemlélteti. A kalibrációra vonatkozó *R*<sup>2</sup> érték 0,90 lett, míg a validálásra vonatkozó *Q*<sup>2</sup> érték 0,87. Az RMSEC 8,90 mg / g értéknek, míg az RMSECV pedig 10,85 mg / g értéknek adódott. A kapott eredmény hasonlóan jónak bizonyult, mint az iPLS változó szelektálással kapott eredmény és lényegesen jobbnak mutatkozott az eredeti adatkészlet alapján kapott kiindulási modellhez képest.

A kiválasztott spektrum szegmensek ezúttal is hozzárendelhetők a megfelelő funkciócsoporthoz. Az iPLS esetében is, a feltüntetetteken kívül, a következő felharmonikusok is fellelhetők voltak: metil aszimmetrikus és szimmetrikus C–H nyújtó rezgése (második felhangsáv), metilén aszimmetrikus és szimmetrikus C–H nyújtó rezgése

(második felhangsáv), valamint metil és metilén kombinációs nyújtó és hajlító rezgések felharmonikusai is fellelhetőek voltak (Workman et al., 2007; Workman, 2000).



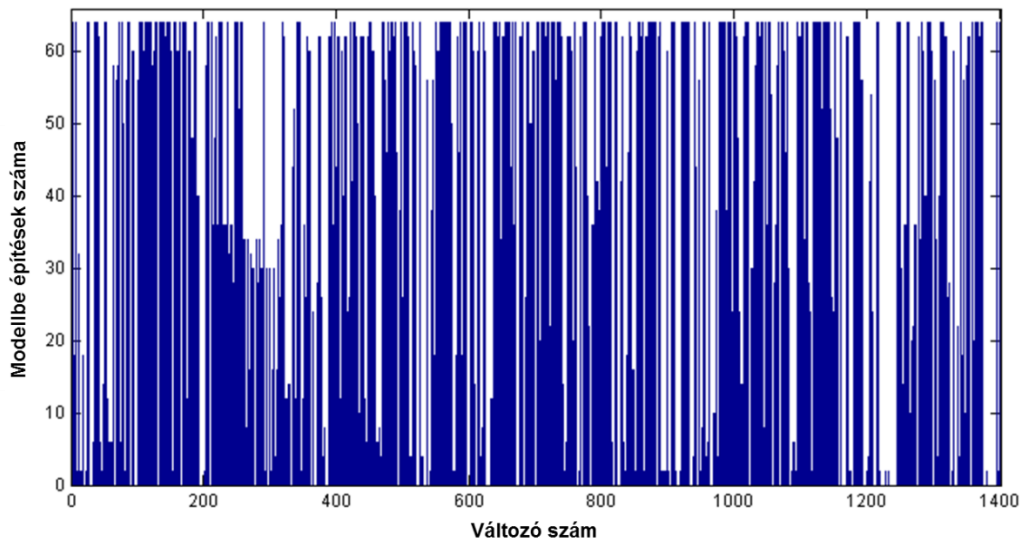
**27. ábra:** Az iSR változó szelektálással kapott modell

*Kék színnel a kalibrációs, pirossal pedig az öt részre osztott kereszt-ellenőrzéssel kapott validált modellt jelöltem. A becsült Y értékeket ábrázoltam a referencia Y értékek függvényében.*

### 5.2.3 A genetikus algoritmus segítségével végzett változókiválasztással kapott eredmények

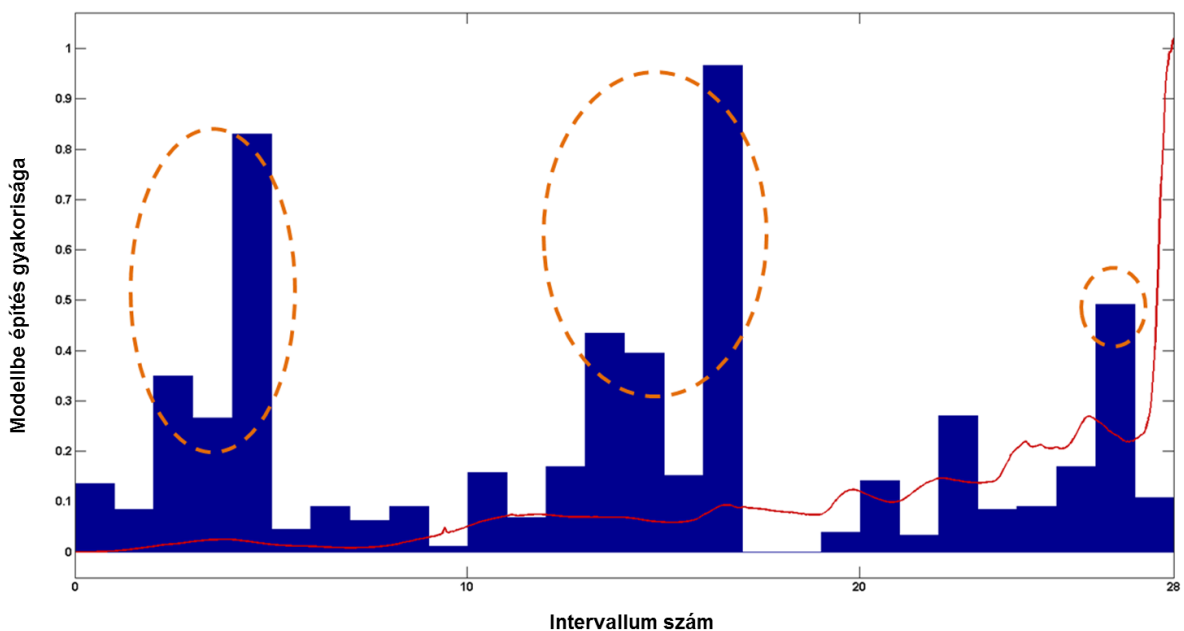
Utolsó változókiválasztási módszerként a genetikus algoritmust (GA) alkalmaztam. A kiértékelés ezen részéhez MATLAB R2012a-t (MathWorks Inc., Natick, MA, USA) illetve a PLS Toolbox 7.9 (Eigenvector Research Inc., Wenatchee, WA, USA) szoftvert használtam. Elsőként intervallumok nélküli esetben próbáltam ki a genetikus algoritmust. A populáció mérete 160 volt, a benne lévő egyedek kezdeti értéként a spektrum változói 30 %-át tartalmazták véletlenszerűen kiválasztva. A maximális generáció szám 100 volt, a mutációs arány pedig 0,005 és dupla kereszteződést használtam. A dupla kereszteződés használatával nagyobb az átfedés a „szülő” és a „utód” gének között. A folyamat során kereszt-ellenőrzést és a változók adatelőkezelését is elvégeztem.

Az egyedi változók használata intervallumok helyett nem hozott nagyobb sikereket a genetikus algoritmus alkalmazása során, ugyanis a kapott eredményekből nehezen lehetett volna levonni bármilyen fizikai következtetést, jelentést vagy kiválasztási lehetőséget. A **28. ábra** ezt megfelelően szemlélteti. Természetesen ez párhuzamba állítható Andersen és Bro azon következtetésével, hogy a genetikus algoritmus nem használható az egyedi spektrum hullámszámok kiválasztására (2010). A legnagyobb probléma, hogy a változók mennyisége túl nagy, a véletlen kiválasztás valószínűsége megnő, és így a fontosabb hullámszámok nehezen találhatók meg. A hullámszámok közötti kapcsolat nem vehető figyelembe megbízhatóan.



**28. ábra:** A genetikus algoritmus eredménye egyedi változókra nézve, intervallumok nélkül  
*A modellbe építés gyakoriságát ábrázoltam a változók számának függvényében (hullámszámok).*

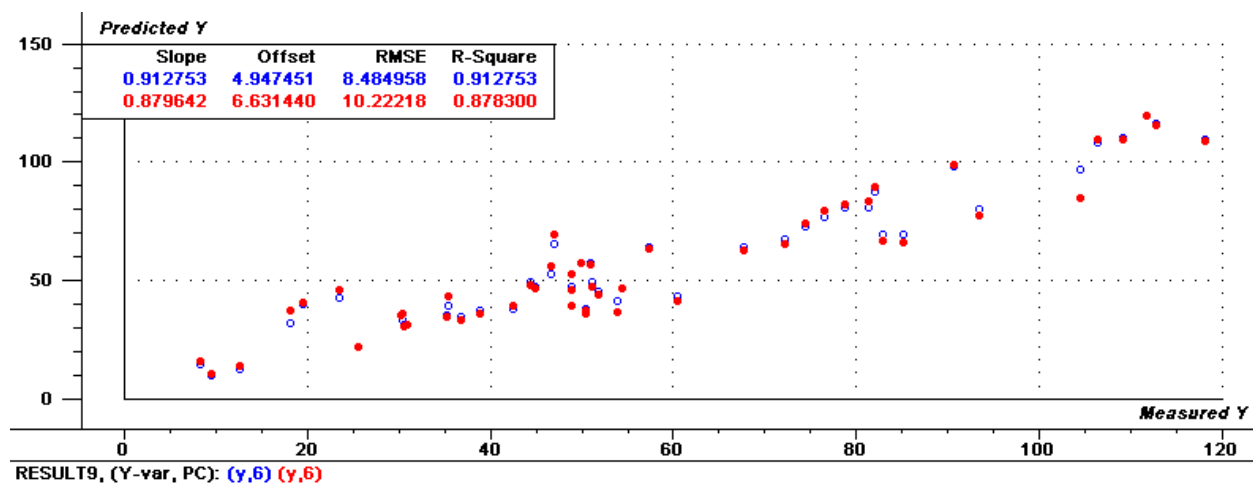
A második esetben már változtatott szélességű ablakokkal (intervallumokkal) végeztem el a genetikus algoritmus számításait, ahogy azt (Andersen és Bro, 2010) ajánlja. Az intervallum szélességét 50 változóra állítottam be, így összesen 28 intervallumot kaptam. A további paraméterek beállításai hasonlóképp történtek, mint az előbbi esetben, így azokat nem részletezem ismét. A **29. ábrán** látható, hogy a GA eredménye itt már jól értelmezhetővé vált.



**29. ábra:** A genetikus algoritmussal kapott változókiválasztás  
*A modellbe építés gyakoriságát az intervallumok függvényében ábrázoltam. Az átlag spektrum piros színnel látható az ábrán. Ezen kívül narancssárga szaggatott ellipszisekkel jeleztem a kiválasztott intervallumokat.*

Az ábra alapján minél nagyobb a modellbe építés gyakorisága, annál jobban használható a modellépítésre az adott változó intervallum. A modellbe építések gyakoriságát ábrázoltam az intervallum szám függvényében, valamint az átlagspektrum is látható az ábrán.

Összesen 350 változót választottam ki a modellépítéshez, amely spektrumtartományok a következők voltak: 8616–8038  $\text{cm}^{-1}$ , 6495–6109  $\text{cm}^{-1}$ , 5916–5724  $\text{cm}^{-1}$ , 3988–3800  $\text{cm}^{-1}$ . Adatelőkezelésként ezúttal is deriválást alkalmaztam. Hat PLS komponensre volt szükség a modell megalkotásához. Az  $R^2$  érték a kalibrációra 0,91 lett, és a  $Q^2$  érték pedig 0,88-nak adódott a validálás során. Az RMSEC érték 8,48 mg/g-nak, míg az RMSECV 10,22 mg/g értéknek adódott. A **30. ábra** szemlélteti a genetikus algoritmussal kapott végső modellt. A becsült  $Y$  értékeket ábrázoltam a referencia  $Y$  értékek függvényében. A kapott modell hasonlóképpen megfelelő, mint a korábbi változókiválasztási módszerekkel kapottak, és a kiindulás modellnél lényegesen jobb.



**30. ábra:** A genetikus algoritmus segítségével kapott modell

*A becsült  $Y$  értékeket ábrázoltam a referencia  $Y$  értékek függvényében. A kalibrációs modellt kékkel, a validált modellt pedig pirossal jelöltem.*

A kiválasztott spektrum részletek a már említetteknek megfelelően a fontosabb funkciós csoportok rezgéseinek feleltethetőek meg. Az előző esetekben leírtakon kívül a kiválasztott spektrum intervallumok kapcsolatban állnak a molekula redukált formájának fenol csoport és metoxi csoport rezgéseivel, valamint C–C/C–H kombinációs rezgésekkel (Workman és Weyer, 2007; Workman, 2000).

#### 5.2.4 A kapott modellek összehasonlítása SRD módszerrel

A végső modellek összehasonlítására a rangszám különbségek összegének módszerét használtam fel. Az SRD kiszámítható egy MS Excel VBA makro segítségével, amely a következő honlapon keresztül érhető el:

<http://aki.ttk.mta.hu/srd>

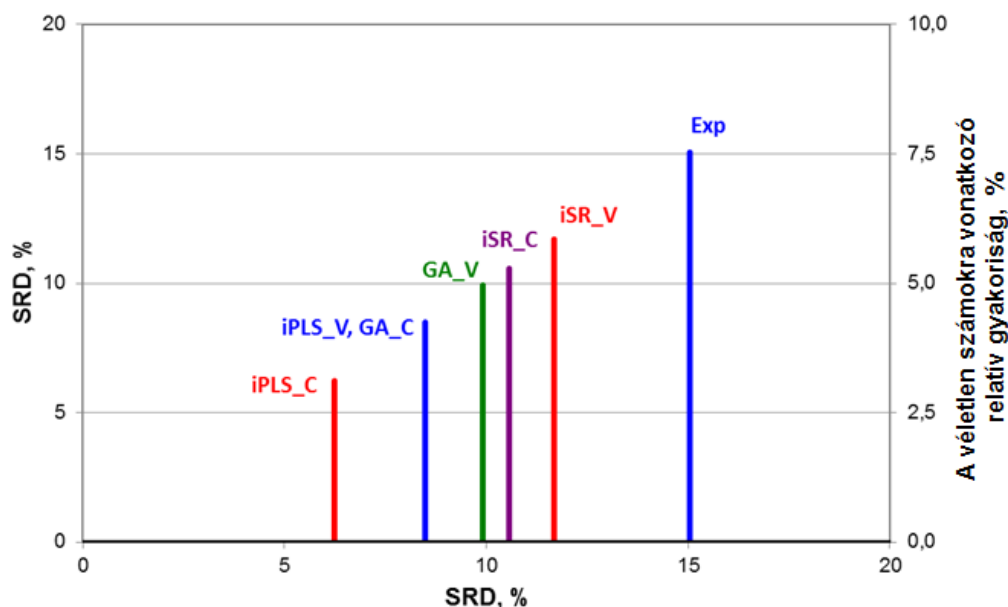
Az SRD-hez szükséges adatkészlet ötven sort (mintát) tartalmazott és hét oszlopot, amelyek a három végső modell becsült kalibrációs és validált  $Y$  értékeit foglalták magukba. Az utolsó oszlopban a HPLC-s referencia  $Y$  értékek szerepeltek. Azért helyeztem az összehasonlítandó modellek közé a referencia mérés eredményeit is, mert tudni akartam, hogy a modellek jobbak, vagy rosszabbak a kísérleti HPLC módszernél. Az SRD módszer referencia oszlopába így nem a HPLC-vel kapott eredmények kerültek, hanem az átlagértékek. A sorátlag értékek lényegében a „konszenzus”-t jelentik, a „legnagyobb valószínűség elve” alapján, amely szerint a legnagyobb valószínűséggel bekövetkező eseményt kell választanunk referenciaként (ez pedig az átlag (Hastie et al., 2001d)). Ekkor nem csak a véletlen hibák, de a módszerek torzítása is kiejti egymást (legalábbis részlegesen), mert a különböző laboratóriumok és mérési módszerek rendszeres hibáinak normális eloszlása tényként kezelendő megállapítás (Youden, 1975).

Az SRD értékek két skálán adhatók meg. Az első az eredeti, a második egy skálázott változat, amelyek  $SRD_{nor}$ -ként rövidítünk. A **31. ábrán** a skálázott SRD értékek vannak feltüntetve, amely segítségével így a modellek összehasonlíthatóvá válnak. A skálázott SRD értékek 0 és 100 közé esnek. A skálázás egyenlete a következő:

$$(21) \quad SRD_{nor} = 100SRD/SRD_{max}$$

Ahol  $SRD_{max}$ , az SRD érték maximumát jelenti az adott számú sor esetén.





**31. ábra:** Az SRD módszerrel kapott modell-összehasonlítás eredménye

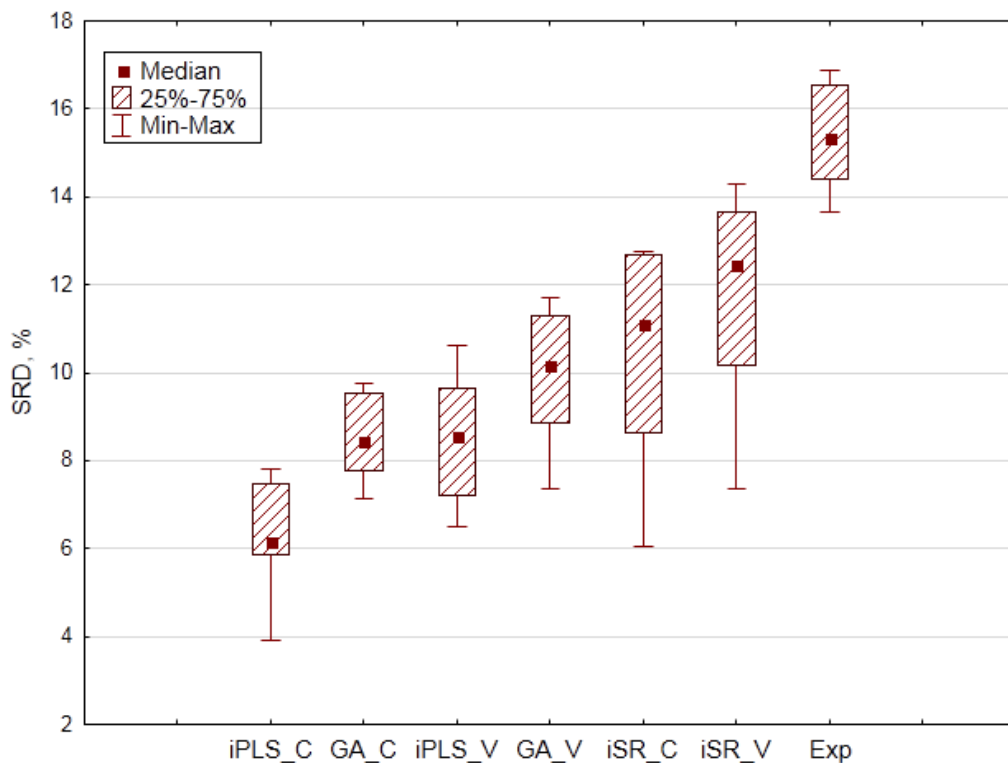
Az jobboldali Y tengelyen található véletlen számokra vonatkozó relatív gyakoriság a „Gauss-szerű” görbéhez kapcsolódik, amely az ábra szükséges nagyítása miatt nem látható a képen.

Rövidítések: *iPLS\_C* és *iPLS\_V* a kalibrált és validált *iPLS* változó szelektált modelleket jelentik, *iSR\_C* és *iSR\_V* a kalibrált és validált *iSR* változó szelektált modellekhez tartoznak és a *GA\_C* valamint *GA\_V* a genetikus algoritmussal változó szelektált kalibrációs és validált modellt jelentik. *Exp* jelentése a HPLC segítségével meghatározott referencia eredmények.

A **31. ábra** alapján elmondható, hogy az összes általam alkotott modell konzisztensebbnek bizonyul a kísérletileg meghatározott eredményeknél az átlag értékekhez viszonyítva, és az átlag értékeket legjobban közelítőnek az *iPLS* kalibrációs modellje mondható. Összességében az *iPLS* segítségével kapott modellek tekinthetőek a legjobbnak. Bár a véletlenszámokhoz tartozó Gauss-görbeszerű eloszlás az ábra nagyítása miatt nem látszik, viszont ez alapján az is elmondható, hogy az összes modell „átmegy” a véletlenszerűségi teszten, mindegyikük lényegesen jobb a véletlen számok használatánál.

Az SRD eredményének validálása szintén megtörtént, mégpedig hét részre osztott kereszt-ellenőrzéssel. Az SRD eljárást így hétszer egymás után a sorok mintegy 1/7-ét kihagyva végeztem el, hogy a kiértékelés bizonytalanságát is megkapjam. Az eredmény ábrázolásához doboz-bajusz ábrát használtam (Box and Whisker plot) mediánokkal, a **32. ábrán** látható módon. Ez alapján elmondható, hogy az *iPLS* kalibrációs modell igen messze helyezkedik el a többi modelltől és a validált *iPLS* valamint a kalibrált *GA* modell nagyon közel helyezkedik egymáshoz (ahogy ez az SRD ábrán is látszott). Ezt a Sign és Wilcoxon nem paraméteres próbák

is megerősítették, tehát az utóbbi két modell között nem található szignifikáns különbség, míg a többi között igen.



**32. ábra:** Doboz-bajusz ábra a normalizált és validált SRD értékekre (hétszeres kereszt-ellenőrzés)

*Az ábrán a modellek medián, 25-75 %-os percentilisei és minimum-maximum értékei vannak jelölve. A modellek rövidítései megegyeznek az SRD ábránál használtakkal.*

### 5.2.5 A modellek külső validálása és az eredmények táblázatos összefoglalása

A külső validálás elengedhetetlen lépés a modellek megbízhatóságának ellenőrzésére (Chirico és Gramatica, 2011; Esbensen és Geladi, 2010; Gramatica, 2007). Ha a lehetőség adott rá, új, még nem vizsgált minták használatát javasolják a modellek jóságának bizonyítására.

A három végső kalibrációs modellt hat külső forrásból származó, új mintával ellenőriztem, melyek a modellek megalkotása után érkeztek a laboratóriumba. Ugyanazon PLS komponens szám és beállítások mellett az iPLS során kapott modell esetén a  $Q^2$  értéke 0,93 lett és a validálás RMSEP (a külső validálásra vonatkozó közepes négyzetes hiba) értéke pedig 8,82 mg / g volt. Az iSR modell esetében a  $Q^2$  érték 0,83, az RMSEP pedig 13,74 mg / g lett. Végül a GA szelektált modell során a  $Q^2$  érték 0,89 lett az RMSEP érték pedig 11,12 mg / g-nak adódott.

A teljes 11 modellből álló sorozat részleteit a **6. táblázat** szemlélteti. A 11 modellből három javasolható későbbi használatra. Mindhárom változó kiválasztási módszer rendkívül hasznos volt a modellek javítása során. A három legjobb modell mindegyike 0,90 fölötti  $R^2$  értékkel és 0,86 fölötti  $Q^2$  értékkel rendelkezik. A külső teszt ellenőrzés szintén sikeres volt mindegyik végső modell esetén. Összességében a legjobb adatelőkezelésnek a deriválást tekinthető. Az SRD módszer pedig bebizonyította, hogy a három modell bármelyike kiválthatja az eddigiekben használt HPLC-s eljárást a Q10 koenzim hatóanyag tartalom mérése során.

**6. táblázat:** A Q10 koenzim koncentráció meghatározására szolgáló modellek. Vastagon kiemelve találhatóak a későbbi alkalmazásra javasolható (legjobb) modellek.

	$R^2$	$Q^2$	RMSEC	RMSECV	Skálázás	Komp. szám	Változó kiválasztás
1.	0,85	0,71	11,26	15,74	-	9	-
2.	0,90	0,53	8,91	19,93	deriválás	9	-
3.	0,88	0,80	9,94	13,14	standardizálás	8	-
4.	0,89	0,74	9,69	15,68	MSC	8	-
5.	0,91	0,72	8,63	15,45	MSC+deriválás	10	-
6.	<b>0,92</b>	<b>0,87</b>	<b>8,16</b>	<b>10,58</b>	<b>deriválás</b>	<b>7</b>	<b>iPLS</b>
7.	0,88	0,78	9,97	14,03	MSC	7	iPLS
8.	0,89	0,82	9,68	12,12	standardizálás	8	iPLS
9.	<b>0,90</b>	<b>0,87</b>	<b>8,90</b>	<b>10,85</b>	<b>deriválás</b>	<b>6</b>	<b>iSR</b>
10.	<b>0,91</b>	<b>0,88</b>	<b>8,48</b>	<b>10,22</b>	<b>deriválás</b>	<b>6</b>	<b>GA</b>

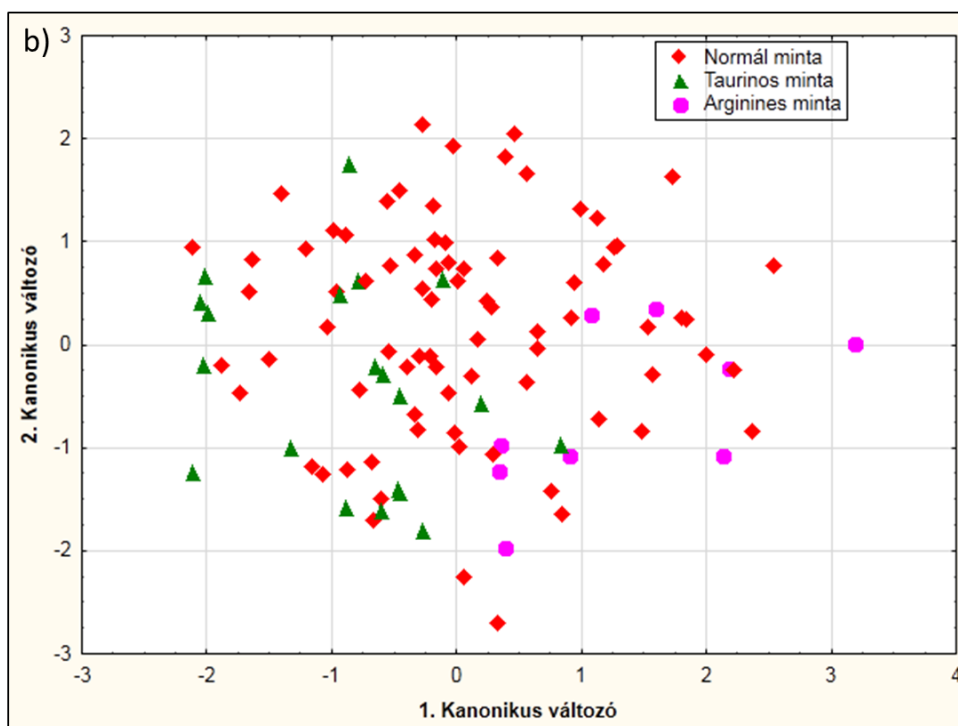
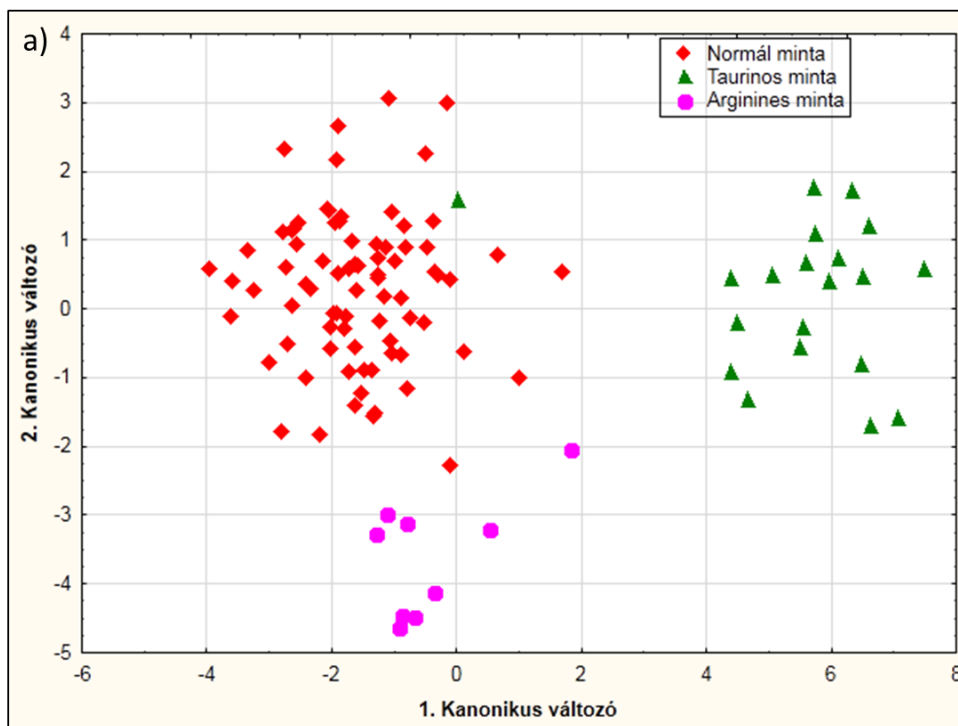
### 5.3 Az energiatalok FT-NIR spektrumának kemometriai elemzése

Az energiatalokról már korábban, a többosztályos ROC görbék fejezetben esett szó, viszont a minták között akadnak eltérők, hiszen a mérési folyamat gyakran sokkal több időt vett igénybe, mint amit a minták eltarthatósága engedett számomra. A kemometriai vizsgálatokhoz hasonlóképpen itt is az energiatalok FT-NIR spektrumát, valamint a különböző referenciamérésből származó adatokat (lásd. 4.2 és 4.4 fejezet) használtam fel. Az eredmények bemutatását két csoportra tagolva teszem meg: i) mintázatfelismerési/osztályozási modellek létrehozása valamint ii) regressziós modellépítések az energiatalok koffein és cukortartalmának becslésére, mivel ez a két szegmens kemometriai szempontból teljesen más módszert és szemléletet igényel.

### 5.3.1 Energiaitalok osztályozási lehetőségei

A többsztályos ROC görbék tárgyalása során az energiaitalok osztályozási lehetőségét már bemutattam, amely alapján cukortartalom szerinti elkülönítést végezhetünk rendkívül pontos osztályozási modellekkel. A cukortartalom mellett viszont osztályozni akartam a különböző magyar, szlovák és görög mintákat arginint, taurint és az előbbieket nem tartalmazó csoport szerint. Ehhez összesen 108 energiaital FT-NIR spektrumát vettem fel majd PCA és LDA módszerekkel végeztem el a kiértékelést. A PCA szerepe a kiértékelésben az adatkészlet „redukciója” volt, vagyis a változósám csökkentése, hiszen az LDA változósámra vonatkozó erős követelménye nem engedte volna meg, hogy több (vagy akárcsak közel annyi) változót használjak fel, mint mintát. Ezt a metodikát követtem a cukortartalom szerinti osztályozás egyik fázisában is (lásd. 5.1.3 fejezet). Ahogy az már az irodalmi áttekintés során is elhangzott, néhány gyártó Magyarországon kicseréli a taurint argininre, míg néhány gyártó egyik komponenst se használja a készítés során. Más országokban a taurinnal kapcsolatban nem olyan szigorú a szabályozás, így a mintákat akár e vegyületek jelenléte vagy hiánya alapján is lehet osztályozni. Hangsúlyozandó, hogy az egyes csoportokban lévő minták száma erősen különböző volt, ami a hagyományos osztályozási eljárásokkal nehezen kezelhető.

Első lépésben a minták FT-NIR spektrumait vettem fel 12500 és 4000  $\text{cm}^{-1}$  között, majd létrehoztam három párhuzamos mérés alapján az átlag spektrumokat. Ezt követően az adatmátrixot standardizáltam. Az így kapott adatkészletet PCA módszerrel értékeltem ki és a kapott 20 főkomponenst használtam a további LDA elemzésekhez. A kemometriai módszerek elvégzéséhez STATISTICA 12 (Tulsa, OK, USA) szoftvert használtam. A változók kiválasztásának módozatai közül a lépésenkénti változó hozzáadást választottam; a validálást három részre osztott kereszt-ellenőrzéssel végeztem el. A megfelelő validálásra itt különösen nagy szükség van, ugyanis könnyen kaphatunk a PCA-LDA kombináció során műtermékeket, túlllesztett modelleket. Ezért nemcsak kereszt-ellenőrzést, de az X változók (független tulajdonság változók) véletlenszerűségi tesztjét is használtam, mégpedig háromszor egymás után. A **33. ábra a)** és **b)** része alapján a végeredmény és annak X véletlenszerűségi (randomizációs) tesztje közötti különbségek jól láthatóak. Az **ábra a)** része alapján a három előzetesen meghatározott csoport jól elkülönül egymástól, amelyet az X véletlenszerűségi teszt megfelelő mértékben összekever a **b)** ábrán látható módon. A modell helyes osztályozási százaléka 95,68 % volt. Megjegyzendő, hogy az egyes csoportokban található minták számossága erősen különböző volt, ami erős követelményt támaszt minden osztályozási eljárással szemben. A sikeres osztályozás itt azt jelenti, hogy csak egy-egy mintát osztályoztunk félre.



**33. ábra:** Az energitalok osztályozási modellje: arginines, taurinos és normál (taurint és arginint nem tartalmazó) minták szerint

*Az a) rész a végső eredményt, míg a b) rész a véletlenszerűségi teszt eredményét szemlélteti.*

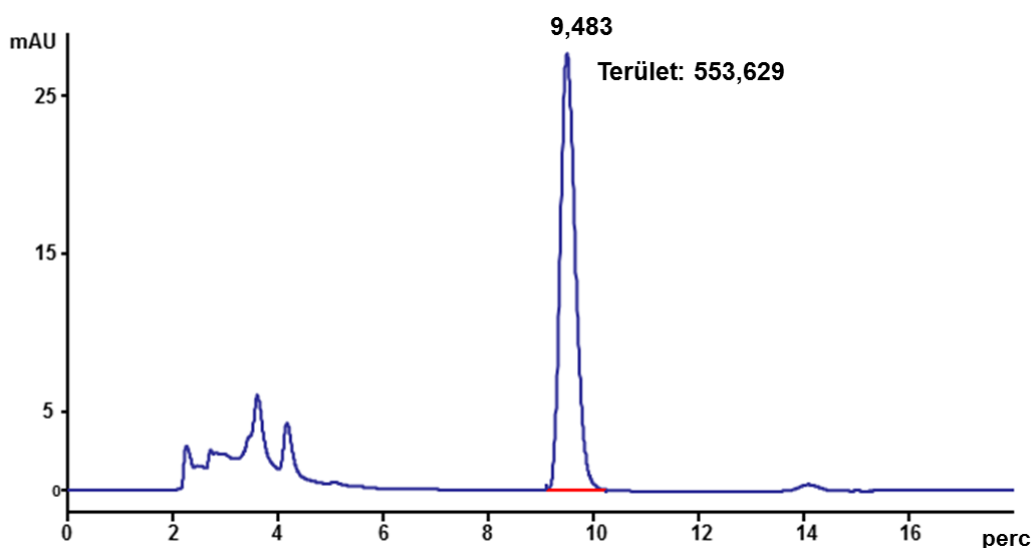
*Mindkét részen a második kanonikus változót ábrázoltam az első kanonikus változó függvényében.*

A kapott osztályozási százalék alapján a megalkotott modellt a további vizsgálatokhoz megfelelőnek tartom. A használt csoportosítási rendszer segítségével eredetazonosítást hajthatunk végre, amellyel így kiszűrhetővé válnak a hamisított energiaital minták is.

### 5.3.2 Energiaitalok koffein és cukortartalmának meghatározása

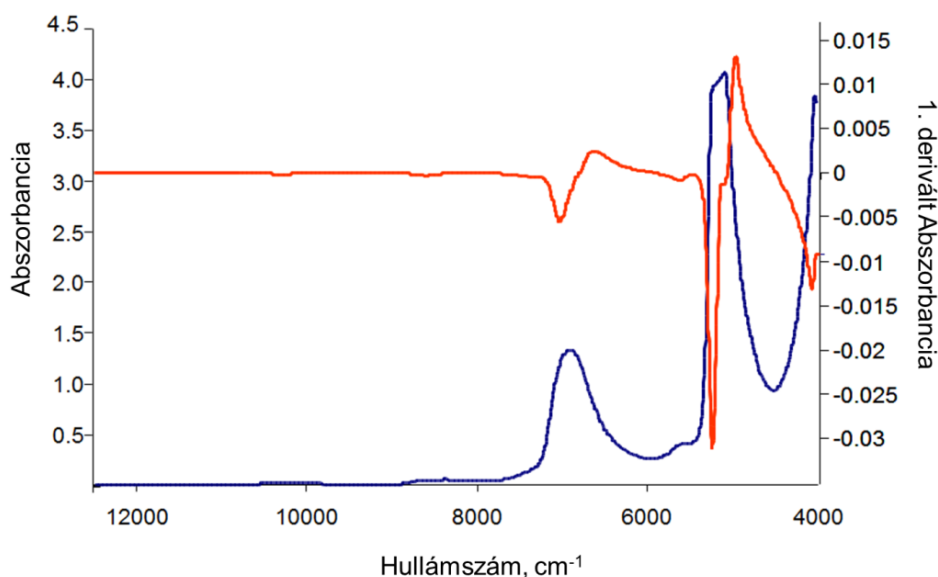
#### 5.3.2.1 A koffein tartalom meghatározása

Első lépésben HPLC-UV technika segítségével 42 eredeti energiaital minta koffein tartalmát határoztam meg. Ezen kívül 33 keverék mintát is létrehoztam az eredeti minták felhasználásával. Mivel az eredeti 42 minta koncentrációját pontosan ismert volt, a belőlük készített keverékek közül csak néhányat ellenőriztem ismételt HPLC-s vizsgálattal: a százalékos eltérés (RSD %) 5 % alatt volt. Minden mintából három párhuzamos mérést készítettem és az átlagos koffein koncentráció értéküket vettem figyelembe a későbbi modellépítés során referencia értéként. A csúcstisztaság vizsgálathoz a minták koncentrációjának meghatározását 260 nm-en is elvégeztem, majd összehasonlítottam a 270 nm-en kapott értékekkel. A kétmintás  $t$  próba használatával (5 % szinten) nem tapasztaltam szignifikáns eltérést a két különböző hullámszámon történt mérés eredménye között. A jellegzetes koffein csúcs a kromatogramon megközelítőleg 9,5 percnél látható. A vizsgált minták közül egy példa kromatogram látható a **34. ábrán**.



**34. ábra:** Egy tipikus példa-kromatogram a 270 nm-en történt meghatározások közül *A koffeinhez tartozó retenciós időt és a csúcs alatti területet a koffein csúcs fölött és mellett tüntettem fel. Az abszorbanca értékek (ezredrészét) az eltelt idő függvényében ábrázoltam.*

Az FT-NIR mérésekhez 10 ml mintát használtam fel minden mintából, és három párhuzamos mérést készítettem ekkor is. A példa-spektrum és annak derivált verziója a **35. ábrán** látható. A koffein koncentrációja a HPLC-UV mérések alapján 118 és 338 mg/100 ml között változott. A megfelelő kalibrációs modellhez a keverék minták koncentrációját úgy állítottam be, hogy azok a vizsgált koncentráció tartományt minél jobban lefedjék, és így minél folytonosabbá váljon a koncentrációeloszlás a két határérték között.

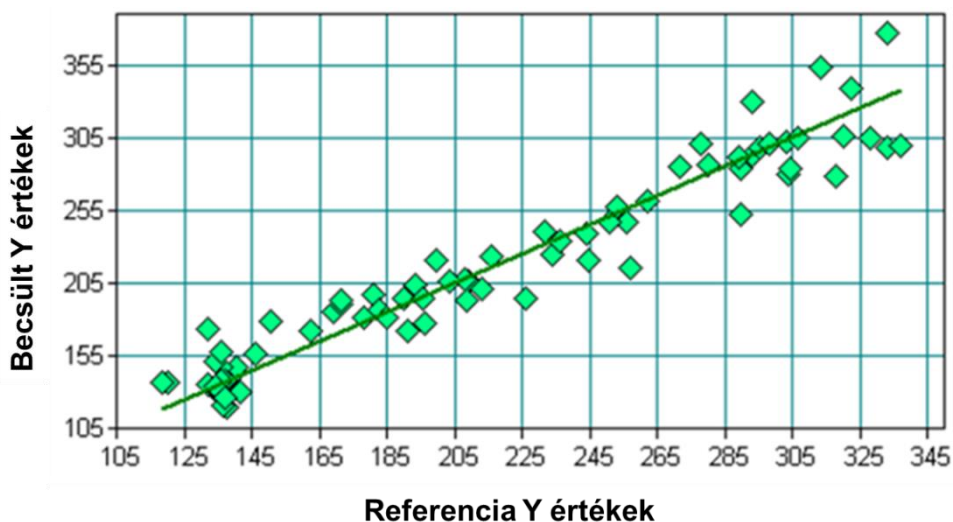


**35. ábra:** A vizsgált minták egy szemléletes példa spektruma, ill. annak derivált formája. Az abszorbancia értékeket ábrázoltam a bal Y tengelyen a hullámszámok függvényében, valamint az első derivált abszorbancia értékeket a jobb Y tengelyen. Az eredeti spektrumot kézzel, míg a derivált verzióját piros színnel jelöltem.

Az FT-NIR spektrumok felvételét követően PCA segítségével elvégeztem a spektrális kieső minták keresését. Ám nem kellett egy mintát se spektrális kiesőnek nyilvánítani, így a végső mintamennyiség 75 maradt. A modell optimalizálást a különböző hullámhossz választással és adatelőkezelés kombinációkkal az OPUS 7.2 (Bruker Corp., Ettlingen, Németország) szoftver segítségével végeztem el. Deriválást és standardizálást használtam adatelőkezelési módszerként. A kiválasztott spektrum tartományok a következők voltak: 12490-7498, 6102-5446 és 4605-4243  $\text{cm}^{-1}$ . A PLS regresszió során 8 PLS komponens elegendőnek bizonyult a modellépítéshez, amelyet az RMSECV értékek globális minimuma alapján határoztam meg.

A **36. ábrán** látható a végső, hétszeres kereszt-ellenőrzéssel kapott modell. Az  $R^2$  értéke a kalibrációra nézve 0,966 lett és a hozzá tartozó közepes négyzetes hiba (RMSEC) értéke pedig 13,4 mg / 100 ml volt. A kereszt-ellenőrzés során kapott modell  $Q^2$  értéke 0,928 lett, az RMSECV értéke pedig 18,3 mg / 100 ml. A külső validálást is elvégezve 13 kereskedelmileg is

kapható új mintára az  $Q^2$  értéke 0,898 lett és a becslés közepes négyzetes hibája (RMSEP) pedig 36,3 mg / 100 ml. Itt a kisebb szabadsági fok miatt előzetesen is magasabb előrebecslési hiba volt várható.



**36. ábra:** A végső, kereszt-ellenőrzéssel kapott koffein koncentráció modell  
*A becscült Y értékeket a mért Y értékek függvényében ábrázoltam.*

A kiválasztott spektrum tartományok a koffein molekulában található különböző funkciós csoportoknak és kötéseknek feleltethetők meg: a metil csoport aszimmetrikus és szimmetrikus első és második felharmonikusainak (Workman, 2000), valamint a C=O, N–H és a C=C vegyértékrezgések első felharmonikusainak és a CONH amid kombinációs rezgéseknek (Magalhães et al., 2016; Ribeiro et al., 2011).

### 5.3.2.1 A cukortartalom meghatározása

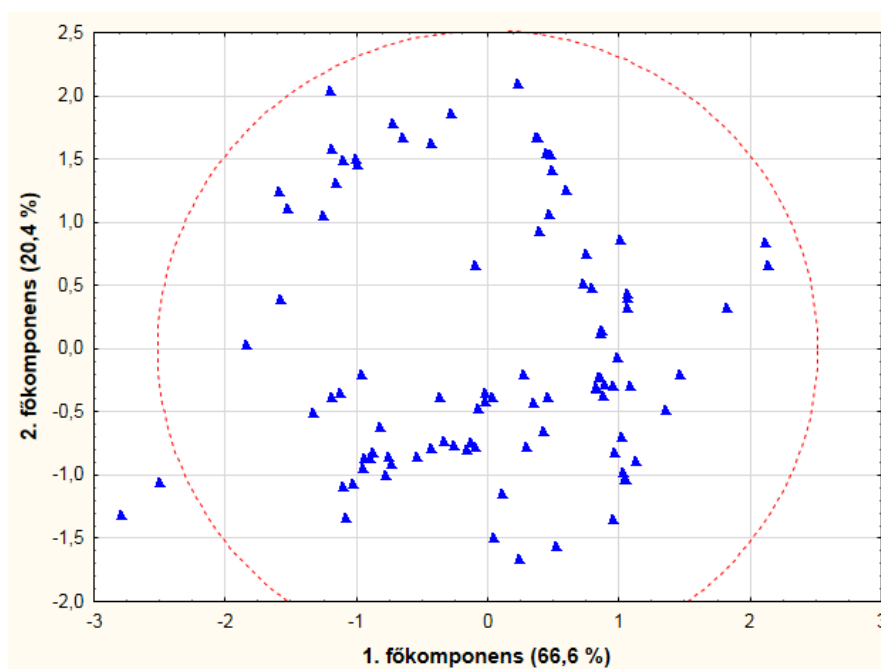
A cukortartalom meghatározásához 71 eredeti és 20 keverék mintát (összesen 91 mintát) használtam fel. Hasonlóan a koffein koncentráció meghatározáshoz, itt is az eredeti mintákból hoztam létre a keverékeket előre meghatározott keverési arányban. A keverékek itt is a cukortartalom jellegzetes értékeinek (gyártók által kedvelt értékek) bővítésére szolgáltak, azért, hogy a kalibrációs tartományt minél jobban le tudjuk fedni.

A referencia Y értékeként a Schoorl módszerrel meghatározott cukortartalom értékeket használtam fel. Ez a módszer rendkívül népszerű az élelmiszeranalitika területén. A teljes mintasorozatból ezzel a technikával 75 mintát tudtam megmérni. Sajnos a Schoorl módszer bár rendkívül népszerű, elég nagy torzítású és szórású (12,4 %), különösen az kisebb koncentrációk tartományában (1-2 g / 100 ml). Ezért úgy döntöttem, hogy az ezzel a módszerrel kapott értékek mellett a nominális (dobozon feltüntetett) értékekkel is végrehajtom a modellépítést, ugyanis a



nominális koncentráció értékek feltehetőleg sokkal kisebb hibájúak az egyszerű tömegmérésnek köszönhetően.

Ez utóbbi esetben minden minta FT-NIR spektrumát háromszor vettem fel, majd átlagspektrumot számítottam belőlük. A mérésekhez 10 ml mintára volt szükségem ekkor is. A spektrális kieső minták vizsgálatára PCA módszert használtam, amely eredményeként két mintát kihagytam a további elemzésekből. Az első két főkomponens egymás függvényében való ábrázolása a **37. ábrán** látható. Mint látható, az eredetileg 91 mintából kettő kívül helyezkedett a 95 %-os konfidencia intervallumon (Hotelling- $T^2$  ellipszis).



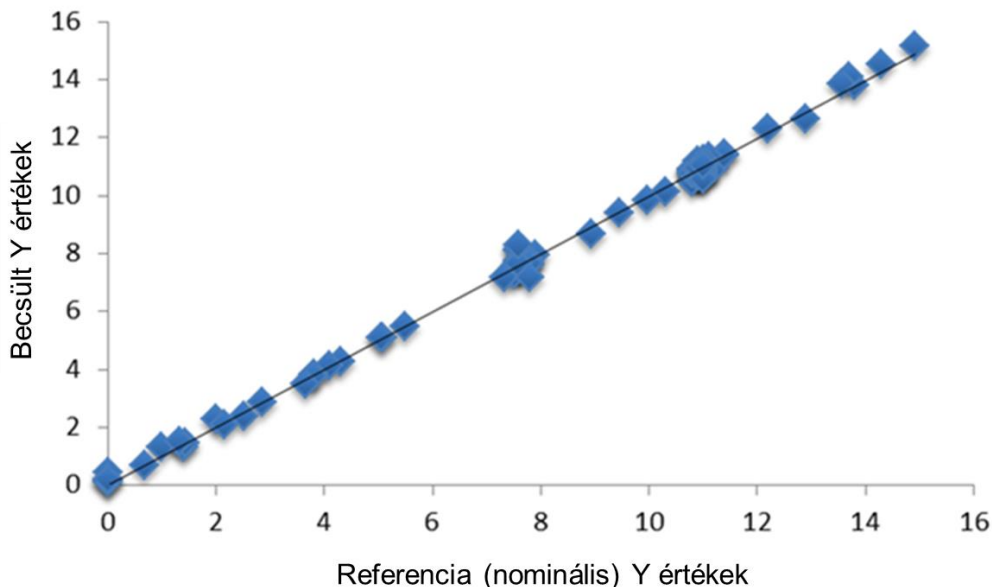
**37. ábra:** Spektrális kiesők vizsgálata

*Az első két főkomponenst ábrázoltam egymás függvényében. Piros szaggatott vonal jelzi a 95 %-os konfidencia sávot.*

A továbbiakban a maradék 89 mintával végeztem el a PLS regressziót (OPUS 7.2). Adatelőkezelésként első deriváltat és standardizálást használtam. A koncentráció tartomány 0,0 és 14,9 g / 100 ml között volt (a mintakészlet tartalmazott „light” energitalokat is). Hat PLS komponens elég volt a modell megalkotásához, amelyet az RMSECV értékek minimuma alapján határoztam meg. Két spektrális tartományt választottam ki a regresszióhoz: 7506–6796 és 4605–4243  $\text{cm}^{-1}$  (összesen 141 változó). Az  $R^2$  érték a kalibrációs modell esetén 0,997 lett, az RMSEC pedig 0,219 g / 100 ml.

A hét részre osztott kereszt-ellenőrzéssel kapott validált modellt a **38. ábra** szemlélteti. A kereszt-ellenőrzés mellett külső teszt ellenőrzést is végrehajtottam 12 új minta felhasználásával.

A kereszt-ellenőrzés során kapott modell  $Q^2$  értéke 0,995 lett az RMSECV értéke pedig 0,29 g / 100 ml. A külső validálás során, ugyanezen teljesítmény paraméterek rendre 0,996 valamint 0,26 g / 100 ml lettek. Összességében a modell közepes négyzetes hibája mindkét validálás esetén 0,30 g / 100 ml érték alatt volt.

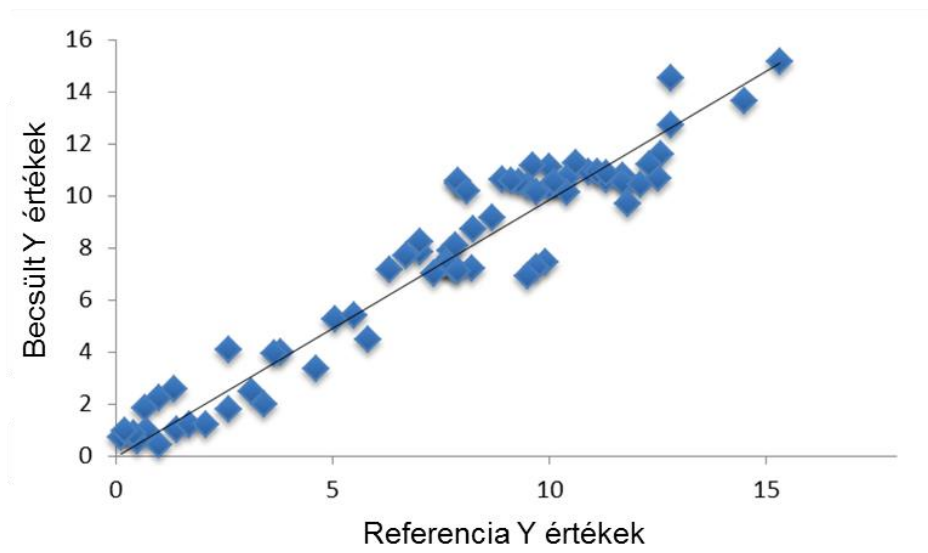


**38. ábra:** A végső kereszt-ellenőrzéssel validált modell a cukortartalom meghatározására  
*A becsült Y értékek lettek ábrázolva a referencia Y értékek függvényében.*

A kiválasztott spektrum tartományok, ebben az esetben is, megfeleltethetők a cukor molekula funkciós csoportjainak és kötéseinek, úgymint az OH csoport vegyértékrezgésének valamint a hidrogénhid kötésben lévő OH első felharmonikusának, illetve a CH vegyértékrezgések és a CH<sub>2</sub> deformációs rezgések kombinációjának (Workman, 2000).

A modellépítést megismételtem a Schoorl módszerrel kapott referencia Y értékek segítségével is. A két spektrális kieső ugyanúgy szerepelt a 75 mintából álló adatkészletben is, így végül 73 mintát használtam fel a kiértékelés e részében. A cukortartalom koncentráció tartománya 0,1 és 15,3 g / 100 ml között volt. A spektrumok adatelőkezelése az előző esethez hasonlóan a deriválás és a standardizálás volt. A kiválasztott spektrum tartományok pedig a következők voltak: 7506–5446 és 4506–4243 cm<sup>-1</sup>. A tartományok kissé eltértek a nominális Y értékek használatakor kapottakhoz képest, amely annak köszönhető, hogy az OPUS szoftver változókiválasztási módszere – illetve a PLS regresszió is – használja az Y (függő) változóban lévő információt. Az így kiválasztott spektrum részek szintén tartalmazzák az elméletileg várt és jellegzetes rezgési sávokat.

A PLS regresszió során hat PLS komponenszt használtam fel a modellépítéshez, amelyet a megszokott módon, az RMSECV értékek alapján határoztam meg. A **39. ábrán** látható a hét részre osztott kereszt-ellenőrzéssel kapott végső modell. A kalibráció során kapott  $R^2$  érték 0,943 volt az RMSEC érték pedig 1,00 g / 100 ml. A validált modellnél a  $Q^2$  érték 0,919 volt, az RMSECV értéke pedig 1,13 g / 100 ml-nek adódott. Végül 11 új mintát felhasználva a külső teszt ellenőrzés során a  $Q^2$  érték 0,935 lett az RMSEP pedig 1,23 g / 100 ml.



**39. ábra:** A Schoorl módszerrel kapott cukortartalmak felhasználásával kapott PLS regressziós modell 7 részre osztott kereszt-ellenőrzéssel

*A becsült Y értékeket a referencia értékek függvényében ábrázoltam.*

A kapott modell teljes mértékben elfogadható, de összehasonlítva a nominális értékek esetén kapottal láthatóan nagyobb hibával terhelt. Ez természetesen nem meglepő, hiszen a Schoorl módszer mérései során sokkal nagyobb torzítást és hibát vittünk be a rendszerbe, mint amekkora egy tömeg bemérése a nominális koncentrációk esetén. Ezt az bizonyítja, hogy a nominális értékek figyelembe vételével a modell sokkal kisebb hibájú volt.

Az energiaitalokra megalkotott három regressziós modell teljesítmény paramétereit a **7. táblázat** foglalja össze.

**7. táblázat:** Az energiai italok koffein és cukortartalmára vonatkozó regressziós modellek teljesítmény paramétereit

*N a mintaszámot jelöli, C a PLS komponensek számát, valamint az RMSECV és RMSEP értékek mértékegysége a cukor modellek esetén g / 100 ml volt, a koffein modell esetében pedig ppm. Az ext jelölés a külső teszt validálás  $Q^2$  értékére vonatkozik.*

	N	C	$R^2$	$Q^2_{ext}$	$Q^2$	RMSECV	RMSEP
<i>Koffein modell</i>	75	8	0,966	0,898	0,928	16,8	36,3
<i>Cukor modell (Schoorl)</i>	73	6	0,943	0,935	0,919	1,13	1,23
<i>Cukor modell (nominális)</i>	89	6	0,998	0,996	0,995	0,29	0,26

A három modell teljesítmény paramétereit figyelembe véve elmondható, hogy mind a koffein mind a cukortartalom meghatározása sikeresen megtörtént az FT-NIR spektrumok alapján. A  $Q^2$  értékek minden validált modellnél 0,90 fölöttiek. Ezért a megalkotott modellek képesek az energiai italok koffein és cukortartalom meghatározására szolgáló HPLC és egyéb elterjedten alkalmazott eljárások kiváltására. A megközelítőleg 100 energiai ital minta alapján készült modellek gyakorlatilag a Magyarországon kereskedelemben kapható energiai italok koncentráció tartományát teljes mértékben lefedik.

#### 5.4 Antioxidáns kapacitás meghatározási módszerek csoportosítása, rangsorolása

Az antioxidáns kapacitás meghatározási technikák összehasonlításához két adatkészletet használtam fel: az első 13 bogyós gyümölcs mintát (szamóca, málna, piros és fekete ribiszke) tartalmazott, a második pedig 12 féle meggy mintát. Összességében hét féle meghatározási módszerrel vizsgálták meg a mintákat, adatkészletekre lebontva a bogyós gyümölcsök antioxidáns kapacitását a FRAP, TPC, TRSC, DPPH, ACL és ACW módszerekkel, míg a meggy mintákét FRAP, TPC, TEAC, ACL és ACW módszerekkel mérték meg. Minden mintából kettőt vizsgáltak és mindegyiket háromszor mértek meg. A kemometria kiértékeléshez a párhuzamos mérések átlagát vettem figyelembe.

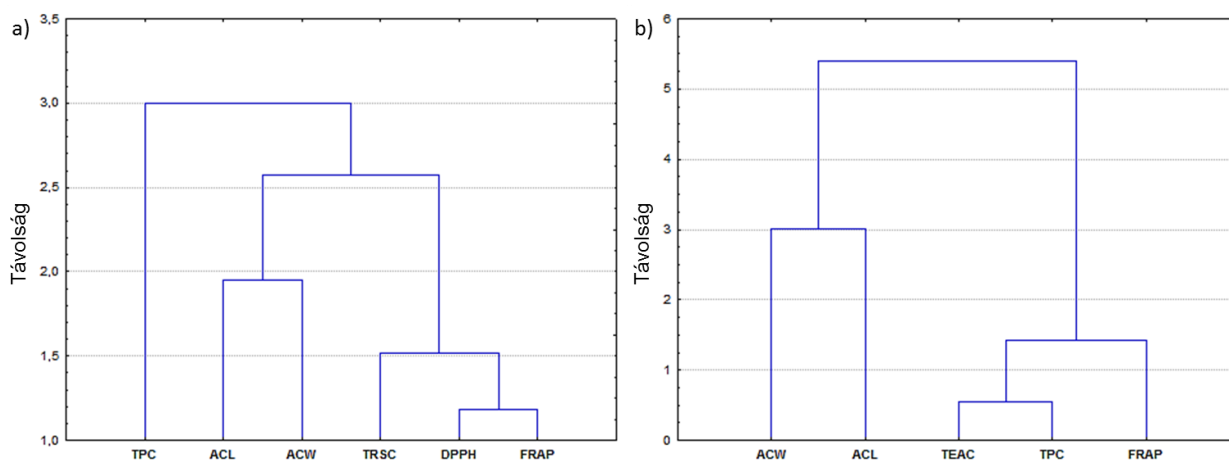
Mindkét adatkészletet standardizáltam a további kiértékelések előtt. Az adatelemzéshez STATISTICA 12 (Statsoft Inc., Tulsa, OK, USA) szoftvert, valamint MS Excel VBA makrókat használtam (SRD és GPCM makrók). Az SRD makró mellett a GPCM is ingyenes elérhető makró a következő honlapon:

<http://aki.ttk.mta.hu/gpcm>

A csoportosíthatóságokat és kapcsolatokat a különböző módszerek között a PCA és hierarchikus fürtelemzés (HCA) segítségével tártam fel. Az antioxidáns kapacitás meghatározási technikák rangsorolásához az SRD és GPCM technikákat használtam.

#### 5.4.1 HCA eredmények

A fürtelemzés során kapott különböző csoportokat és kapcsolatokat a **40. ábrán** mutatom be. Euklideszi távolságot használtam távolságmértékként és Ward módszert kapcsolási szabályként mindkét adatkészlet esetén. Az a) ábra bogyós gyümölcsökre vonatkozó részén látható, hogy két csoport teljesen elkülöníthető egymástól: az egyik az ACL és ACW technikákat, a másik pedig a TRSC, DPPH és FRAP módszereket foglalja magába. Az eredmény tükrében az említett módszerek a csoportokon belül sokkal jobban hasonlítanak egymásra (közelebb állnak egymáshoz). A b) rész a meggy fajtákhoz tartozó két csoportot mutatja. Az ACW és ACL módszerek itt is, az alkalmazott három másik módszerhez képest külön csoportba kerültek.



**40. ábra:** A HCA-val kapott nem felügyelt tanítású mintázat

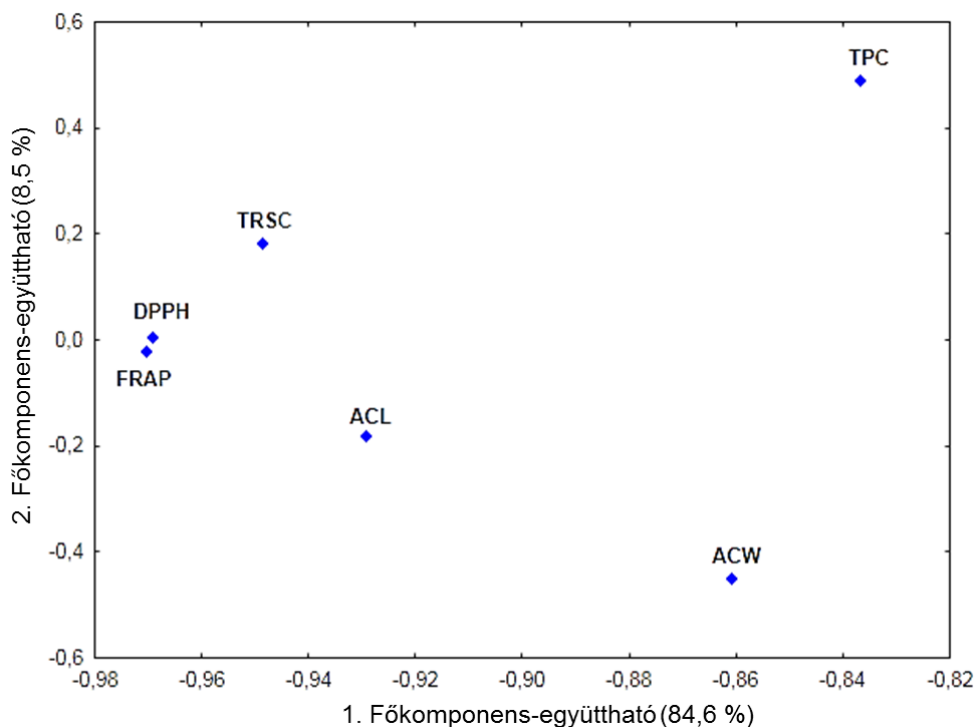
*Az a) és b) rész rendre a bogyós gyümölcsökhöz és meggy fajtákhoz tartozó adatkészletek eredményeit szemléltetik. Az Y tengelyen az Euklideszi távolság van ábrázolva az antioxidáns kapacitási technikák függvényében.*

Bár a fürtelemzés segítségével a hasonlóságokat felderíthettem, de módszerek rangsorolása így nem oldható meg.

#### 5.4.2 PCA eredmények

Mintázatfelismerési módszerként főkomponens-elemzést is végre kellett hajtani, hiszen ellenőrizni akartam a fürtelemzéssel kapott eredmények helyességét és megbízhatóságát. A bogyós gyümölcsök adatkészlete esetén két főkomponens elegendőnek bizonyult a magyarázott

variancia nagy részének leírására (több mint 90 %). Az első két főkomponens-együttható vektort egymás függvényében ábrázolva a **41. ábrán** látható eredményhez jutottam.

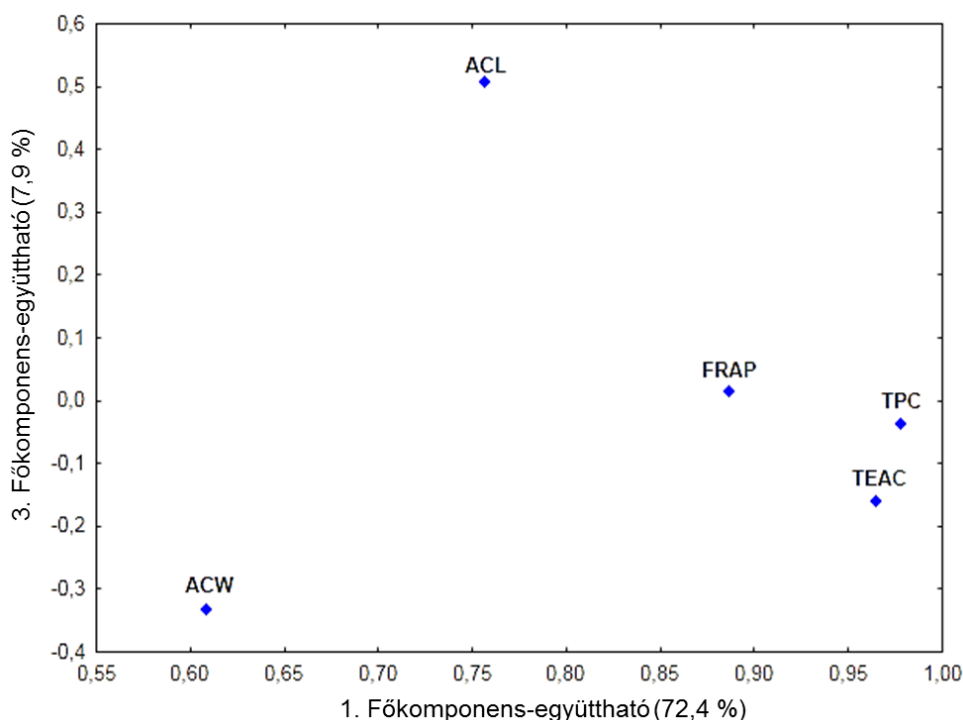


**41. ábra:** A főkomponenselemzéssel kapott mintázat a bogyós gyümölcsöket tartalmazó adatkészlet esetén

*A második főkomponens-együtthatót az első függvényében ábrázoltam. A magyarázott variancia százaléka zárójelben található a tengelyfeliratokon.*

A legtöbb módszer elszórtan található az ábrán, de a DPPH és FRAP módszerek rendkívül közel helyezkednek el egymáshoz. Az ACW, ACL és TPC módszerek az első két főkomponens-együtthatót egymás függvényében ábrázolva nem alkottak csoportokat.

A meggy fajták adatkészlete esetében a kapott főkomponensek közül az első hármat tartottam meg a későbbi elemzésekhez. Ezen három főkomponens összesen 98 %-át fedi le a magyarázott varianciának. A második és harmadik főkomponens egyedi főkomponens volt, vagyis egyetlen eredeti változó hordozta a variancia nagy részét. A legjobb elkülönülést az első és harmadik főkomponens-együttható egymás függvényében történő ábrázolása adta, amely a **42. ábrán** látható.



**42. ábra:** A harmadik főkomponens-együttható ábrázolása az első függvényében  
(meggy fajták adatkészlete)

*A főkomponens-együtthatókhöz tartozó magyarázott variancia százaléka zárójelben található a tengelyfeliratokon.*

A kapott mintázat igazolni tudja a fúrtelemzés során kapott eredmények helyességét, hiszen az ACW és ACL módszerek kiugró pontokként detektálhatóak és így külön csoportot képeznek a másik három egymáshoz viszonylag közel helyezkedő módszerhez képest.

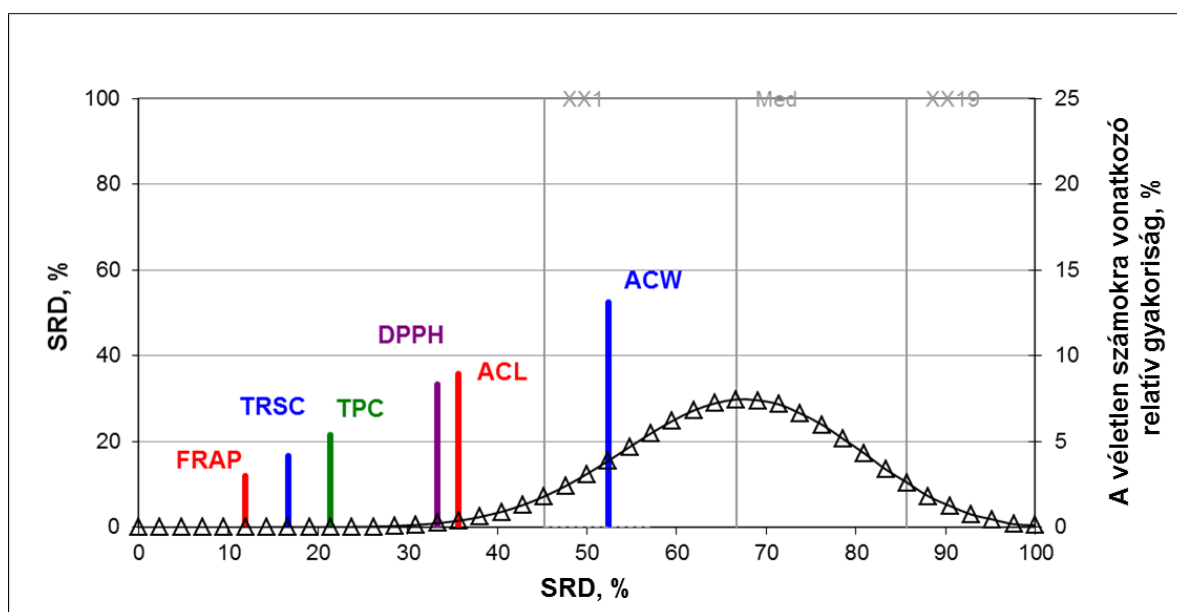
Bár az első adatkészlet még nem tűnt elegendőnek a következtetések levonására, a PCA és HCA eredményekből együttesen azt láthattam, hogy az ACW és ACL módszerek nagyon eltérőek a többi alkalmazott módszerhez képest. A TPC, FRAP és TEAC módszerek közötti kapcsolat erősebbnek bizonyult, amely a bogyós gyümölcs adatkészlet alapján ugyanígy elmondható a FRAP és DPPH technikákról is. Ezek a következtetések kapcsolódnak a módszerek alapelveiben rejlő különbségekhez is. Hiszen a TPC és egyéb ET módszerek sokkal közelebb helyezkedtek egymáshoz, míg a HAT elven működő ACW és ACL sokkal távolabb helyezkedett el az előbbiektől. A bogyós gyümölcs adatkészlet esetén a FRAP és DPPH módszerek erősen korrelálnak, a közöttük lévő korrelációs koefficiens 0,941-nek adódott.

### 5.4.3 SRD eredmények

Az SRD elemzés már kimondottan a módszerek közötti rangsorolásra irányult. Ekkor az adatmátrixok mindkét esetben a módszereket, mint változó oszlopokat, és a mintákat, mint

sorokat tartalmazták. Az adatkészletek kiegészültek egy referencia, sorátlag oszloppal is. Ez a megoldás az olyan esetekben kedvező, ha nincs egzakt referencia értékünk az adatomátrixban. Ilyenkor a konszenzusos megoldást célszerű választani a „maximum likelihood elv” értelmében, amely kimondja, hogy azon paraméter választása a legmegfelelőbb, amely a legnagyobb valószínűséggel következik be (ez esetben az átlag) (Hastie et al., 2001d).

Az SRD értékeket egy MS Excel VBA program segítségével számítottam ki, kétféleképpen: „nyers” adatokként majd azok eredeti 0 és 100 közé skálázottan (azaz normalizált változatban). A **43. ábrán** a bogyós gyümölcs adatkészletből számolt skálázott SRD értékek láthatóak. Azért a skálázott verziót használtam, mert így a két eltérő mennyiségű mintával rendelkező adatkészletet össze tudtam hasonlítani.



**43. ábra:** A skálázott SRD értékek (százalékos) ábrázolása a módszerek függvényében a bogyós gyümölcs adatkészlet esetében

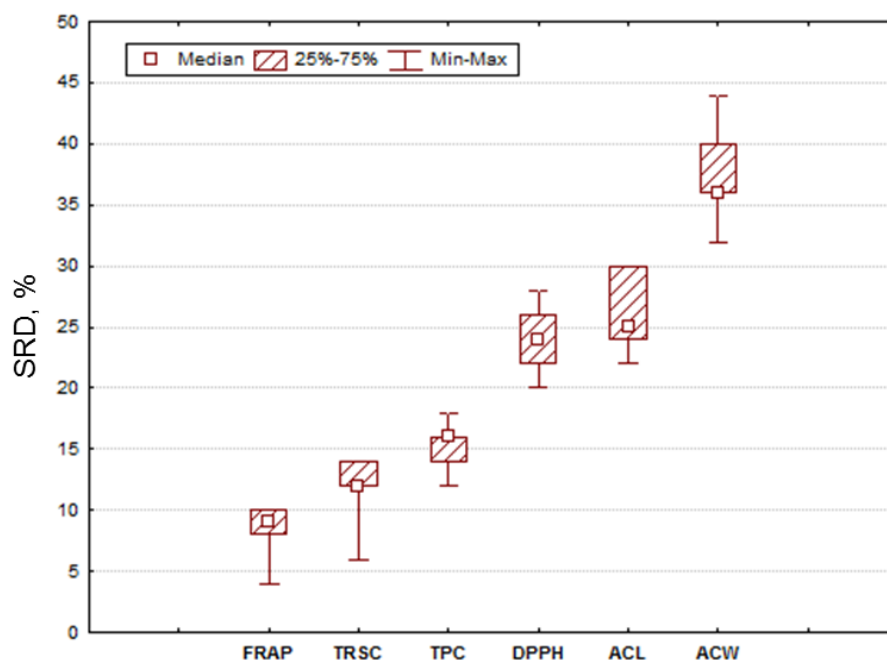
*A baloldali Y és X tengelyeken az SRD értékek, míg a jobboldali Y tengelyen a véletlen számokra vonatkozó relatív gyakoriság van feltüntetve. Az ábrán látható Gauss-szerű eloszlás az előbbi relatív gyakoriságokhoz tartozó görbe.*

A kiértékeléshez azzal a feltételezéssel éltem, hogy minden módszer véletlen és szisztematikus hibákkal terhelt, valamint feltételezhető az is, hogy a véletlen hibák és a torzítások is részben kiejtik egymást. Ahogy az az SRD fejezetben már szerepelt, minél kisebb a kapott SRD érték, annál közelebb van a referencia értékhez (jelen esetben az átlaghoz). Így az eredmény alapján elmondható, hogy a FRAP módszer esett a legközelebb a referenciához, vagyis ez a módszer képes a legkisebb hibával kiváltani a többi módszert. A véletlen számokra vonatkozó relatív gyakoriságok Gauss-szerű görbéjéből pedig láthatjuk, hogy a vizsgált módszer



jobb, hasonló, vagy rosszabb eredményt ér el a véletlen számok használatával végzett rangsorolásnál. Ebben az esetben majdnem mindegyik antioxidáns kapacitás meghatározási technika jobbnak mutatkozott a véletlen számoknál, kivéve az ACW technikát.

Az eredmény validálásához véletlenszerűségi tesztet és hétszeres kereszt-ellenőrzést használtam. Az utóbbi esetén így minden módszer esetén hét validált SRD értéket számoltam plusz az eredeti teljes adatmátrix felhasználásával kapott értékeket is. Erre azért volt szükség, mert az előbbi ábrán vonalakkal szemléltetett SRD értékekről nem tudjuk eldönteni, hogy mekkora a hibájuk, átfednek-e, vagy épp közelségüket csak a véletlen okozza-e. Ezért a kereszt-ellenőrzéssel együtt kapott SRD értékeket doboz-bajusz ábrán szemléltettem, amelyen a minimumot, maximumot a 25-75 %-os percentiliseket és medián értékét is feltüntettem. Az eredmény a **44. ábrán** látható.

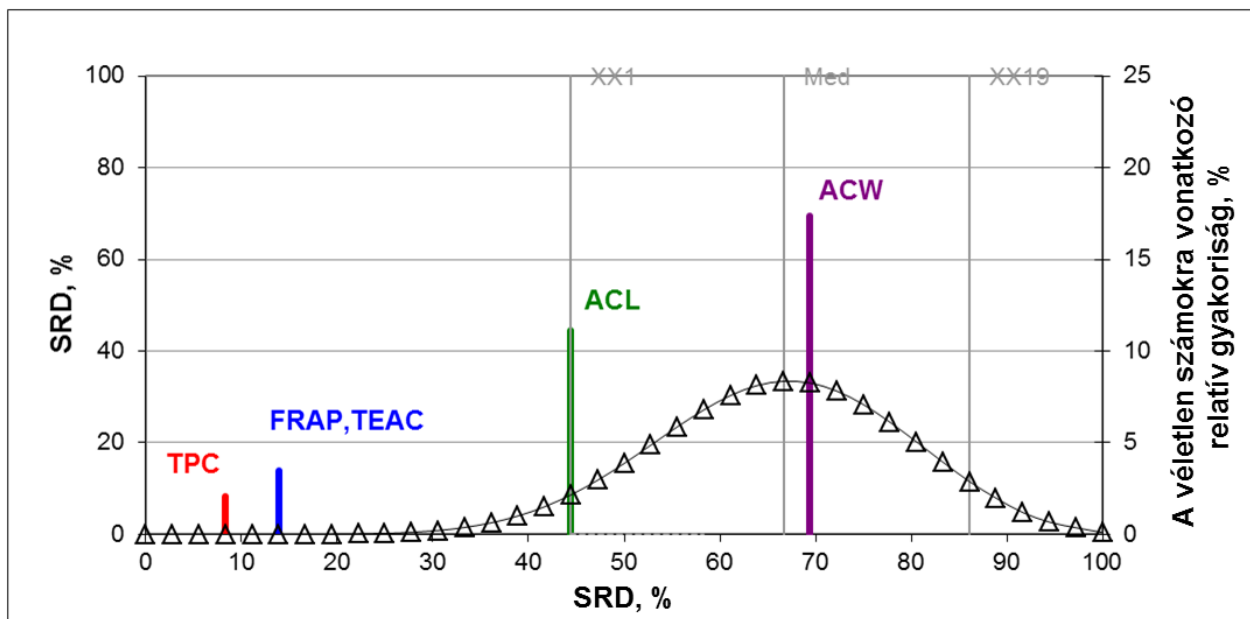


**44. ábra:** A bogyós gyümölcs adatkészlet validált SRD értékeinek doboz-bajusz ábrája

*Az SRD értékeket az Y tengelyen ábrázoltam. A doboz az adatok 50 %-át jelképezi, míg a belső négyzet a mediánt.*

A **44. ábra** alapján elmondható, hogy a DPPH és ACL módszerek igen közel helyezkedtek el egymáshoz képest. A null hipotézisem szerint a módszerek átlag értékei nem különböznek szignifikánsan egymástól. A két mintás *t*-próba segítségével összehasonlítottam az SRD átlagértékeket páronként minden módszerre. A hasonlóképpen nemparaméteres Sign és Wilcoxon próbát is alkalmaztam a mediánok közötti eltérések megállapítására (ezen tesztek esetén a normalitás nem feltétel). A *t*-teszt és a nemparaméteres tesztek is ugyanazt az eredményt adták, miszerint a DPPH és ACL módszerek nem különböznek szignifikánsan. A többi vizsgált módszer eredménye viszont szignifikánsan különböző volt.

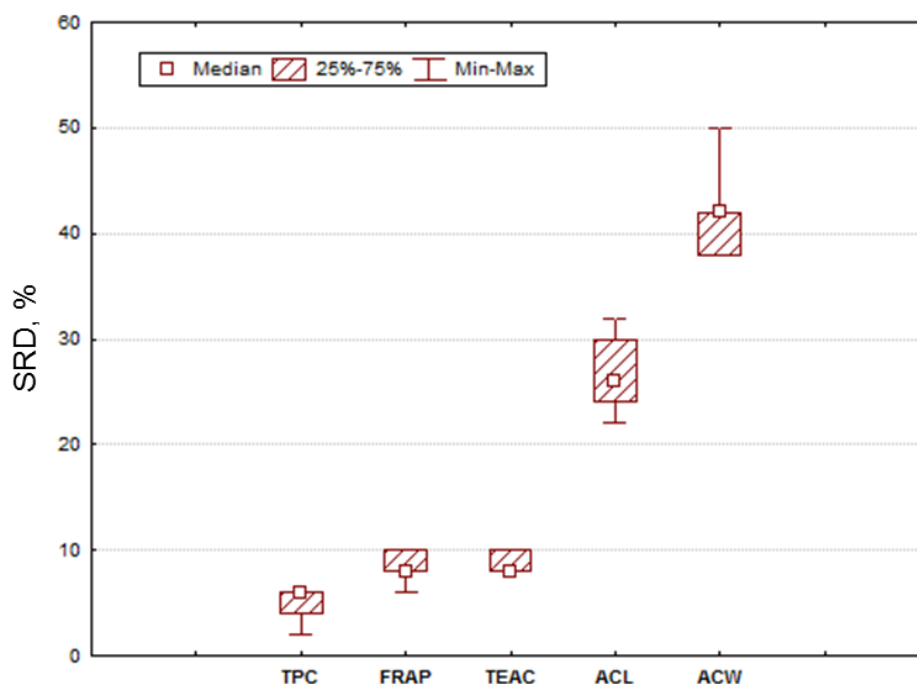
A meggy adatkészlet vizsgálatát hasonlóképpen elvégeztem az SRD segítségével hasonló beállítások mellett. Az öt módszerre kapott skálázott SRD értékeket a **45. ábrán** mutatom be.



**45. ábra:** A skálázott SRD értékek (százalékos) ábrázolása a módszerek függvényében a meggy adatkészlet esetében

A baloldali Y és X tengelyeken az SRD értékek, míg a jobboldali Y tengelyen a véletlen számokra vonatkozó relatív gyakoriság van feltüntetve. Az ábrán látható Gauss-szerű eloszlás az előbbi relatív gyakoriságokhoz tartozó görbe.

Az ábra alapján megállapítható, hogy a TPC eredményezte a legkisebb SRD értéket, így e módszer képes a többi vizsgált módszer kiváltására. Az ACW módszer viszont kívül helyezkedtek az elfogadható tartományon, míg az ACL határeset. A FRAP és TEAC módszerek ugyanazon SRD értékeket adták, és a validálás során kapott értékek mediánjai Sign és Wilcoxon próbával vizsgálva sem térnek el szignifikánsan a két módszer esetén. A kapott mintázat rendkívül hasonlít a két adatkészletre nézve a doboz-bajusz ábrák és az alkalmazott próbák alapján is. A **46. ábra** alapján látható, hogy az említett FRAP és TEAC módszerek mediánjai és a 25-75 %-os percentilisei tökéletesen megegyeznek.



**46. ábra:** A megyy adatkészlet validált SRD értékeinek doboz-bajusz ábrája

*Az SRD értékeket az Y tengelyen ábrázoltam. A doboz az adatok 50 %-át jelképezi, míg a belső négyzet a mediánt mutatja.*

Összegezve elmondható az SRD módszer képes volt rangsorolni az átlag értékekhez képest a különböző antioxidáns kapacitás meghatározási technikákat és kimuattani, hogy azok szignifikánsan különböznek-e egymástól. A vizsgált adatkészletek alapján a FRAP valamint TPC módszerek bizonyultak a legjobbaknak (legalacsonyabb SRD értékek), míg az ACW és ACL technikák rendelkeztek a legnagyobb SRD értékekkel.

#### 5.4.4 GPCM eredmények

Az általánosított pár-korrelációs módszer segítségével, ellenőrzési célból, a rangsorolás eredményeit kibővítettem. A használt adatmátrix ugyanaz volt, mint az SRD módszer esetén, tehát az antioxidáns kapacitási módszerek az oszlopokban, míg a minták a sorokban helyezkedtek el. A GPCM módszer segítségével sikerült is validálni az SRD során kapott eredményeket. A GPCM számolásokat egy MS Excel VBA makró segítségével végeztem. A referencia ebben az esetben is a sorátlagokat jelentette. A páronkénti összehasonlításokat minden lehetséges kombinációban elvégeztem. A kiértékelés során „feltételes egzakt Fisher próbát” valamint a különbségek alapján történő rangsorolás valószínűséggel súlyozott változatát (szignifikáns rangsorolás) használtam. A **8. táblázat** tartalmazza a bogyós gyümölcs adatkészletre vonatkozó eredményeket. A szignifikáns rangsorolás szerint a FRAP illetve TRSC módszerek kapcsolódtak leginkább az Y referencia változóhoz, amely a sorok átlaga volt. A

valószínűséggel korrigált győzelmek száma, ebben a két esetben, nagyon közel esett egymáshoz. A TPC harmadik lett a sorban, és az ACW volt az utolsó a vizsgált módszerek között.

**8. táblázat:** A GPCM módszer összefoglaló táblázata a bogyós gyümölcsök esetén.

*A pWinner a szignifikáns rangsorolással kapott győzelmek számát, a pLoser a hasonló módszerrel kapott vereségek számát adja meg. A rangsorolás a kettő különbsége alapján számolható ki. Az előre meghatározott hibalimit érték jelölése  $\alpha$  (user) volt, míg  $\alpha$  (emp.) az elméleti limitet jelentette. A „Crit. Sum.”-al jelölt kritikus összeg (Critical Sum) a szignifikáns rangsorolással kapott győzelmek összegét jelenti a vereségek száma nélkül (2. szám). Az első szám a konfidencia korlát, amely a Crit. Sum. szorzata  $(1 - \alpha)$ -val.*

	FRAP	TRSC	TPC	DPPH	ACL	ACW
<b>pWinner</b>	2,9968	2,9951	0,9999	0,9956	0,9927	0
<b>pLoser</b>	0	0	0	1,9965	1,9955	4,9882
<b>No Decision</b>	2	2	4	2	2	0
<b>Rank by:</b>	<b>1</b>	<b>2</b>	<b>3</b>	4	5	6
<b>pWin-pLos</b>	$\alpha$ (user)	0,05	$\alpha$ (emp.)	0		
		No Diff. in Y		<b>Crit. Sum</b>	6,65	7

A **9. táblázat** hasonlóképp a GPCM eredmények bemutatására szolgál, csak a meggy fajták adatkészletét vizsgálva. Ekkor a TPC és FRAP módszerek szinte megkülönböztethetetlenek voltak egymáshoz és a többi módszerhez képest. A TEAC harmadik lett a rangsorolásban, és az ACW az előző esethez hasonlóan itt is az utolsó volt.

**9. táblázat:** A GPCM módszere összefoglaló táblázata a meggy fajták adatkészletére kapott eredményekről

*A részletes ábramagyarázat a 8. táblázatnál található.*

	TPC	FRAP	TEAC	ACL	ACW
<b>pWinner</b>	1,9999	1,9984	1,9895	0,9997	0
<b>pLoser</b>	0	0	0	2,9878	3,9997
<b>No Decision</b>	2	2	2	0	0
<b>Rank by</b>	<b>1</b>	<b>2</b>	<b>3</b>	4	5
<b>pWin-pLos</b>	$\alpha$ (user)	0,05	$\alpha$ (emp.)	0	
		No Diff. in Y	<b>Crit. Sum</b>	5,7	6

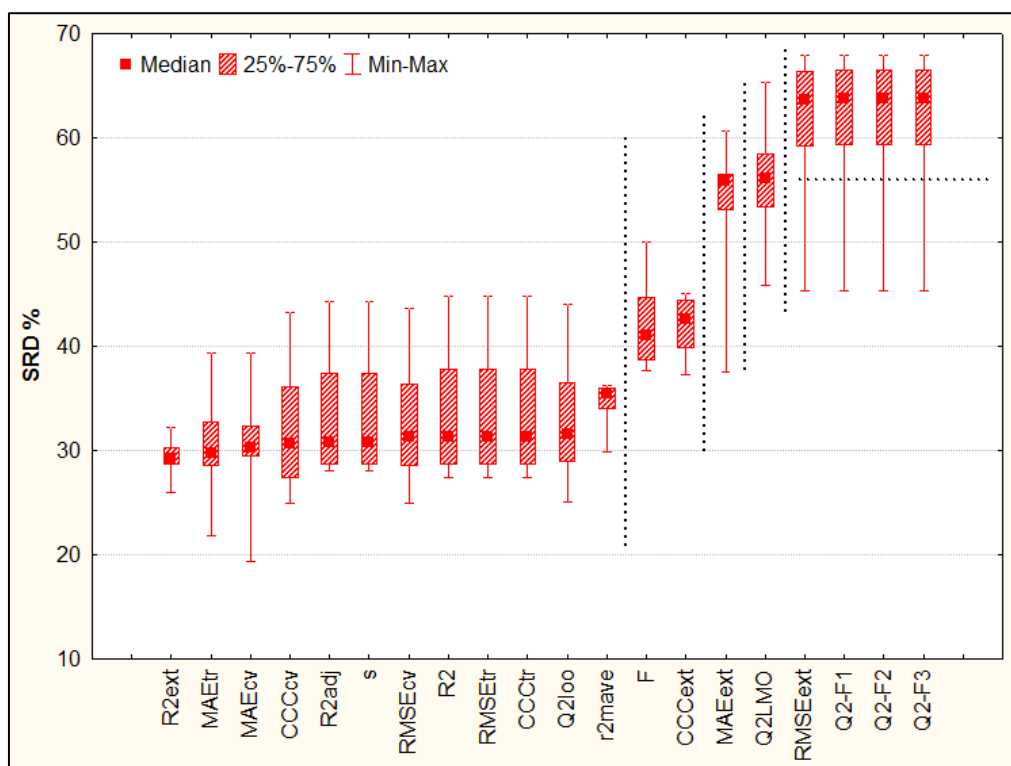
Az eredmények rendkívül hasonlóak lettek az SRD során kapottakhoz, viszont így a GPCM technikával a DPPH és ACL módszerek között az első adatkészlet esetén, míg a FRAP valamint TEAC módszerek között a második adatkészlet esetén különbséget tudtam tenni.

## **5.5 Regressziós modellek előrejelző képességét meghatározó paraméterek összehasonlítása**

A regressziós modellek nemcsak az élelmiszertudomány, de a más tudományterületeken is, pl. a gyógyszerkémia területén is legalább olyan fontosnak bizonyultak, mint pl. az utóbbi évtizedekben, különböző szerkezet-hatás összefüggések (QSAR) megalkotásában. A munkám során egyfajta kitekintésképpen, két olyan adatkészletet használtam fel a becslés jóságára vonatkozó paraméterek összehasonlítására, amelyek bár a gyógyszerkémia területéről származnak, de sokkal szélesebb körű, nagyobb mennyiségű paraméter számolását tették lehetővé. A munka további alkalmazhatóságát ezzel is növelni tudtam. Természetesen az élelmiszertudomány és azon belül is NIR spektroszkópia modellekre jellemző tipikus becslési paraméterek is szerepelnek a vizsgálatban, pusztán a szélesebb alkalmazhatóság vezetett az adatkészletek (és szoftverek) bővítésére.

Az összehasonlított paraméterek teljes listáját a Melléklet **M1** része tartalmazza, belső és külső validáláshoz tartozó teljesítményparaméterek egyaránt szerepeltek közöttük. A két vizsgált adatkészlet egyike különböző benzol származékok toxicitás értékeit tartalmazta független (Y) változóként, míg a másik N-szubsztituált maleimidek inhibitor aktivitás értékeit ( $IC_{50}$ ).

A kiértékelésekhez többváltozós lineáris regressziót (MLR) használtam. A kapott modellek teljesítményparamétereit SRD módszerrel rangsoroltam, és hasonlítottam össze. Az első adatkészlet esetén 20 paraméter szerepelt az oszlopokban és összesen hatvan modellt készítettem el ekkor. A sorátlagokat használtam referenciaként, a kapott értékeket pedig később hét részre osztott kereszt-ellenőrzéssel validáltam. A **47. ábra** a validálási eredményeket szemlélteti egy doboz-bajusz ábra segítségével.

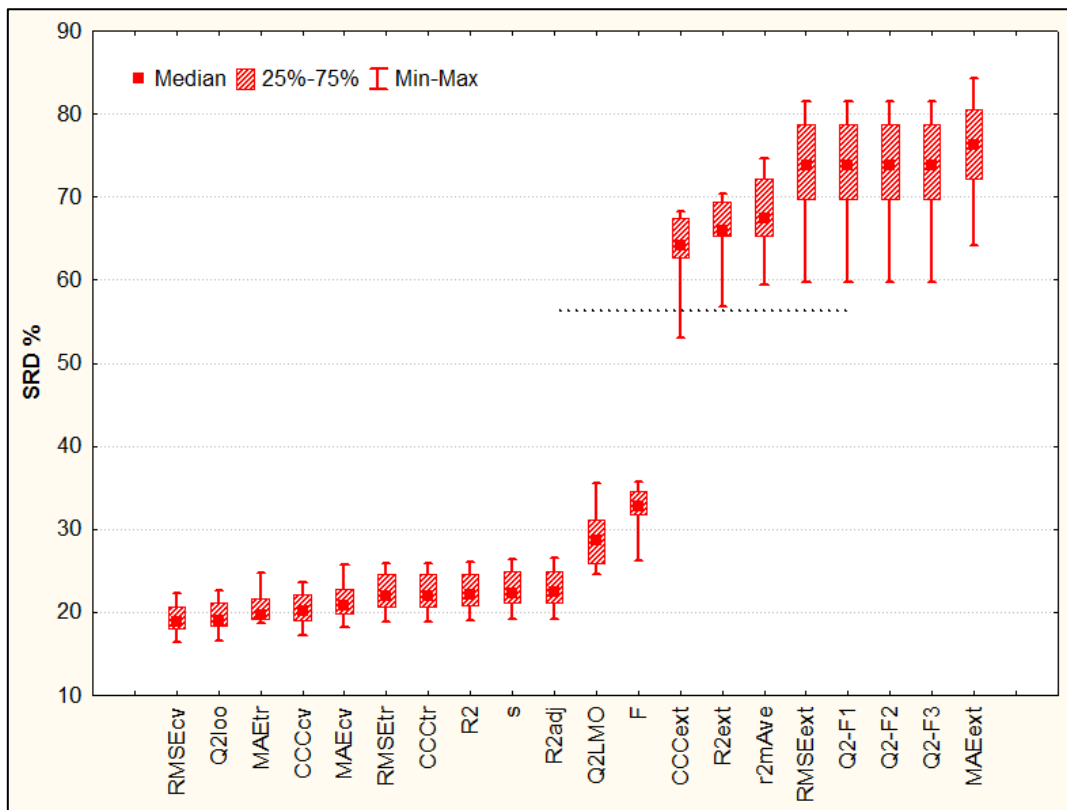


**47. ábra:** A toxicitási modellek SRD elemzése során kapott eredmények doboz-bajusz ábrája. Az SRD értéket az Y tengelyen ábrázoltam a teljesítményparaméterek függvényében. A vízszintes szaggatott vonal a véletlen számok SRD értékeinek eloszlásához tartozó 5 %-os hibahatárt szemlélteti. A függőleges vonalak a páros Wilcoxon próba alapján történő szignifikáns elkülönüléseket jelöli.

A **47. ábra** alapján néhány külső validálási paramétert leszámítva ( $RMSE_{ext}$ ,  $Q^2-F1$ ,  $Q^2-F2$ ,  $Q^2-F3$ ) a többi paraméter, a véletlen számok SRD eloszlásához tartozó 5 %-os hiba határ alatt van ( $H_0$  hipotézis: az adott paraméter SRD értéke nem különbözik szignifikánsan a véletlen számok SRD értékétől). Az első 12 teljesítmény paraméter a páros Wilcoxon próba alapján nem különbözik szignifikánsan egymástól ( $\alpha = 0,05$ ,  $H_0$  hipotézis: az adott paraméterek medián értékei nem különböznek szignifikánsan egymástól).

Ugyanezt a kiértékelést elvégeztem a második adatkészletre is, amely eredményeképp nagyon hasonló rangsorhoz jutottam. A vizsgált teljesítmény paraméterek száma ekkor 20 volt, és 70 modellt használtam fel.

A **48. ábra** alapján bár kissé eltérő sorban, de nem számottevő különbséggel ugyanazon paraméterek kapták a legkisebb SRD értékeket. Az  $RMSE_{cv}$  volt a leginkább konzisztens (az átlaghoz legjobban közelítő), majd ezt követték a  $Q^2_{LOO}$ ,  $MAE_{tr}$ ,  $CCC_{cv}$  és  $MAE_{cv}$  paraméterek.



**48. ábra:** Az aktivitási modellek SRD elemzése során kapott doboz-bajusz ábra. Az SRD értékeket az Y tengelyen ábrázoltam a teljesítményparaméterek függvényében. A vízszintes szaggatott vonal a véletlen számok SRD értékeinek eloszlásához tartozó 5 %-os hibahatárt szemlélteti.

A külső validálási paraméterek ebben az esetben is sokkal távolabb helyezkedtek el, és 5 %-os hibahatár fölött voltak ( $H_0$  hipotézis: az adott paraméter SRD értéke nem különbözik szignifikánsan a véletlen számok SRD értékéhez képest) az első, nagyobb, belső validálási paramétereket tartalmazó csoporthoz képest. Ez természetesen nem jelenti azt, hogy a külső validáláshoz tartozó teljesítmény paraméterek, amelyek a hibahatár fölé kerültek, tehát nagy SRD értékekkel jellemezhetők, használhatatlanok lennének, hiszen pont a nagy különbségek miatt hordozhatnak értékes információt (de a véletlen rangsorolástól nem különböztethetők meg). Összességében elmondható, hogy a betanító készletre vonatkozó, vagy a kereszt-ellenőrzéssel (vagy a kereszt-ellenőrzéssel egy elem kihagyásos változatával) kapott paraméterek alkották a leginkább reprezentatív és leginkább konzisztens csoportot. A két adatkészletet figyelembe véve a leginkább az RMSECV, CCC<sub>cv</sub> és  $Q^2_{LOO}$  és MAE paraméterek ajánlhatók.

## 6. KÖVETKEZTETÉSEK ÉS JAVASLATOK

A doktori munkám során számos példán keresztül bemutattam az FT-NIR kutatásokban és a kemometriában rejlő, még máig is gyakran kiaknázatlan lehetőségeket. Az elvégzett kísérleteimmel és elemzéseimmel törekedtem a kemometriai tudományág és az élelmiszertudomány területén minél több fejlesztést és újítást létrehozni.

A Q10 koenzim tartalmú étrendkiegészítők piaca még a mai napig is növekvőben van, hiszen sokáig egyfajta „csodaszerként” kezelték ezt a koenzimet. A manapság kiélezett piaci versenyhelyzetben a különböző formában forgalomba hozott termékek minőségellenőrzése kihangsúlyozandóan, nem csak a fogyasztók, de a gyártók érdeke is. Ebben a szegmensben az általam megalkotott kalibrációs modellek tökéletesen ki tudják váltani a korábban alkalmazott, idő és vegyszerigényes HPLC-s technikát. A modellépítések során kifejlesztett változókiválasztási technikával pedig a kemometria világában új lehetőséget nyújtok a korábban használt változókiválasztási lépésekhez képest. A modellépítések e része gyakran kulcsfontosságúnak bizonyul, ahogy ez a Q10 koenzim tartalmú étrendkiegészítők vizsgálatokor is bebizonyosodott számomra. Mivel az iSR változószelektálás megfelelő hatékonyságot mutatott – akárcsak a másik két alkalmazott módszer – így bátran javasolható a későbbi regressziós modellek során történő alkalmazása is.

Az energiatalok területén végzett kutatásaim során szembesültem azzal, hogy az energiatal kereskedelemben betöltött szerepe és fogyasztók nagy száma miatt mekkora jelentőségre tett szert Magyarországon. Ahogy egy évtizeddel ezelőtt hazánkban még éppen csak néhány márka volt forgalomban, úgy az utóbbi években már ez a szám több, mint százra volt tehető. Sajnálatos módon a legnagyobb fogyasztói köre az energiataloknak mindmáig a fiatalok, beleértve a 18 és sokszor 14 év alatti korosztályt is. A rohamos fejlődésnek indult energiatal piacon így szintén fontos – már csak a fogyasztói kört figyelembe véve is – a termékek minőségellenőrzése, különös tekintettel a két legnagyobb egészségügyi kockázattal járó összetevőre: a koffeinre és a cukorra. A kutatásaim során mindkét komponensre nézve sikerült olyan regressziós modelleket építenem, amelyek gyors és környezetbarát módon percek alatt képesek kellően pontos mennyiségi meghatározásra. Ezek a modellek így ismételten kiválthatják a területen gyakorta alkalmazott nagyműszeres analitikai eljárásokat. A vizsgálatok során közel száz energiatalt használtam fel, így a magyarországi piacon található energiatalokra mindkét modell kellőképpen optimáltnak tekinthető. Az osztályozás területén pedig sikerült megalkotni egy olyan modellt, amely megfelelő elkülönítést biztosít a taurinos, arginines és normál (taurint és arginint nem tartalmazó) minták között. Ez a modell az eredetazonosítás és hamisítványok



kiszűrése során juthat nagy szerephez. Az energiainformációk spektrumainak felhasználásával pedig kifejlesztettem és teszteltem a többosztályos ROC görbék egy új formációját. Az energiainformációk adatait felhasználva rendkívül jó megjelenítést biztosítva tudtam összehasonlítani többféle osztályozási módszert is. Az ebből származó eredmények következtetéseként érdemes sokszor a még nem túlzottan elterjedt, akár új kemometriai módszereket is tesztelni, hiszen gyakran a jól bevált és több évtizedes múlttal rendelkező módszerek is csődöt mondhatnak. Esetemben a kifejlesztett fák módszere (BT) tűnt ki a többi alkalmazott módszer közül, így ennek a technikának az alkalmazását feltétlenül javaslom osztályozási feladatokra.

Az antioxidáns kapacitás meghatározási technikák világában is számos lehetőséggel találkozhatunk, amelyek közül nehéz eldönteni melyiket érdemes választani a vizsgálatainkhoz. Gyakran ajánlott lehetőség a minél többféle technika használata. A kutatásom során én azt a kérdést tettem fel, hogy melyik módszer tudná visszaadni leginkább a legkisebb hibájú (konszenzusos) értékeket, melyik az, amelyik kiválthatja a többit olyan esetekben, ha nincs elég időnk, vagy több mérés elvégzése túl költségesnek bizonyul. A vizsgálatok során erre a célra legmegfelelőbbnek a FRAP és TPC módszerek bizonyultak. Így e két módszer használata az előbbieken felvázolt esetekben erősen ajánlott, hiszen ezek a módszerek tudták a két vizsgált adatkészlet esetén az összes mérési módszer átlagát figyelembe véve a konszenzusos eredményt legjobban visszaadni.

A regressziós modellépítések során állandó problémát okoz, a helyes teljesítményparaméterek megfelelő figyelembe vétele, legyen szó FT-NIR spektrumokról vagy gyógyszerkémiai adatokról. A modellnek a különböző teljesítmény paraméterek alapján történő eltérő megítélése miatt nagyon fontos kérdés, hogy melyiket érdemes használni. Az általam elvégzett számítások során bebizonyosodott, hogy bár a különböző külső validálási paraméterek fontos plusz információt adhatnak a modellekről (és bizonyos esetekben elengedhetetlen a külső validálás megléte), a belső validáláshoz tartozó RMSECV vagy éppen  $CCC_{CV}$  értékek sokkal inkább egymással konzisztens modelleket válogatnak ki. Az általam elvégzett vizsgálatok így azt tükrözik, hogy a belső validálásra, kereszt-ellenőrzésre és annak eredményeire érdemes nagyobb hangsúlyt fektetni.

## 7. ÚJ TUDOMÁNYOS EREDMÉNYEK

1. A doktori munkám során a vevő-működtető jelleggörbék (ROC görbék) alapelvét továbbgondolva, többosztályos („n-class”) ROC görbéket hoztam létre, melyek a különböző módszerek összehasonlítására, értékelésére – beleértve a hibabecslést is – valamint az osztályozó képességük ábrázolására, megjelenítésére is kiválóan alkalmasak. Ezen görbék segítségével energiatalok FT-NIR spektrumai alapján többféle mintázatfelismerési (osztályozási) eljárást hasonlítottam össze. A véletlen erdő módszerével történő osztályozás modellépítésére kétféle paraméteroptyimálási eljárást is kidolgoztam és sikeresen alkalmaztam. Megállapítottam, hogy a vizsgált adatkészleten a fejlesztett fák módszere képes a legjobb modellépítésre. A módszerhez saját fejlesztésű programot használtam, amelynek programkódja a dolgozat mellékletében megtalálható.

*(A. RÁCZ, D. BAJUSZ, M. FODOR and K. HÉBERGER (2016): Comparison of classification methods with "n-class" receiver operating characteristic curves: A case study of energy drinks. Chemometrics and Intelligent Laboratory Systems, Vol. 151. pp. 34-43. (IF=2.321))*

2. A Q10 koenzim tartalmú étrendkiegészítők vizsgálata során három regressziós modellt is építettem, melyek segítségével az étrendkiegészítő minták FT-NIR spektrumaiból gyorsan, és kellő pontossággal, környezetkímélő módon meghatározható a Q10 koenzim koncentrációja. A modelleket belső és külső validálásnak is alávettem. A három PLS regressziós modell segítségével kiválthatók az idő- és költségigényes, kevésbé környezetbarát HPLC eljárások is.

*(A. RÁCZ, A. VASS, K. HÉBERGER and M. FODOR (2015): Quantitative determination of coenzyme Q10 from dietary supplements by FT-NIR spectroscopy and statistical analysis. Analytical and Bioanalytical Chemistry, Vol. 407:(10). pp. 2887-2898. (IF=3.436))*

3. A Q10 koenzim étrendkiegészítők vizsgálata során, a modellépítésekhez továbbfejlesztettem a szelektivitási arány változó kiválasztási módszert, amely az intervallum szelektivitás arány kiszámításán alapul (iSR). Az így kifejlesztett változó szelekciós eljárást sikeresen alkalmaztam a modellépítés során.

*(A. RÁCZ, A. VASS, K. HÉBERGER and M. FODOR (2015): Quantitative determination of coenzyme Q10 from dietary supplements by FT-NIR spectroscopy and statistical analysis. Analytical and Bioanalytical Chemistry, Vol. 407:(10). pp. 2887-2898. (IF=3.436))*

4. Az energiatalok vizsgálata során a cukor és koffein tartalom meghatározására FT-NIR spektrumok felhasználásával PLS regressziós modelleket dolgoztam ki. A fejlesztett modellek mindkét komponensre nézve megfelelő pontossággal és robusztussággal rendelkeznek. A modellek mindkét esetben belső és külső validáláson estek át, így a továbbiakban alkalmazhatók az eddigi időigényes, és kevésbé gazdaságos módszerek helyett az energiatalok mennyiségi meghatározására. Ezen kívül megfelelő osztályozási képességgel rendelkező osztályozási modellt hoztam létre a taurinos, arginines és normál (taurint és arginint nem tartalmazó) minták elkülönítésére. Ez utóbbi modell a minőségellenőrzés, eredetazonosítás vagy hamisítványok kiszűrése során használható.

(A. RÁCZ, K. HÉBERGER, M. FODOR (2016): *Quantitative determination and classification of energy drinks using near-infrared spectroscopy. Analytical and Bioanalytical Chemistry, Doi: 10.1007/s00216-016-9757-8 (IF=3.125)*)

5. Az antioxidáns kapacitási módszerek vizsgálata és összehasonlítása során megállapítottam, hogy az azonos kémiai háttérrel rendelkező módszerek a PCA és HCA eredmények alapján sokkal hasonlóbb (egységesebb) eredményt mutatnak. Az ACW és ACL módszerek, mindkét adatkészlet vizsgálata során, jelentős eltérést mutatnak a többi módszerhez képest. Statisztikai szempontból leginkább konzisztens módszernek a FRAP és TPC technikákat tartom, amelyek kiválthatják a többi módszer alkalmazását, amellyel így szűkös időkeret esetén is a legkisebb hibával tudjuk megadni az antioxidáns kapacitás értékeket.

(A. RÁCZ, N. PAPP, E. BALOGH, M. FODOR and K. HÉBERGER (2015): *Comparison of antioxidant capacity assays with chemometric methods. Analytical Methods: Advancing Methods and Applications, Vol. 7. pp. 4216-4224. (IF=1.821)*)

6. A regressziós modellek teljesítményparamétereinek széleskörű összehasonlítása során, jelentős eltérést tapasztaltam a legtöbb külső és belső validálásból származó paraméter között. A belső validáláshoz tartozó paraméterek általában sokkal jobbnak (konzisztensnek) bizonyulnak. Ezek közül is kiemelten jónak tekinthetők az RMSECV,  $CCC_{cv}$  és  $Q^2_{LOO}$  paraméterek. A külső validálási paramétereknek a modellépítések során tapasztalt nagy eltérései viszont új információt hordozhatnak.

(A. RÁCZ, D. BAJUSZ, K. HÉBERGER (2015): *Consistency of QSAR models: Correct split of training and test sets, ranking of models and performance parameters. SAR and QSAR in Environmental Research, Vol. 26:(7-9). pp. 683-700. (IF=1.596)*)

## 8. ÖSSZEFOGLALÁS

Kutatómunkám során többféle kemometriai módszerfejlesztést, összehasonlítást és FT-NIR adatokra történő modellépítést végeztem. A módszerfejlesztések területén kiemelhető a többcsoportos vevő-működtető jelleggörbék (ROC görbék) megalkotása, valamint azok osztályozási modellek értékelésére történő használata, beleértve a hibabecslést, és az osztályozó képességük megjelenítését is. Az elemzések során parciális legkisebb négyzetek módszerével végzett diszkriminancia elemzést (PLS DA), véletlen erdő (RF), fejlesztett fa (BT) és lineáris diszkriminancia elemzés (LDA) módszereit hasonlítottam össze a megalkotott ROC görbék alapján. A felhasznált adatkészlet egyike 90 energiatartalmú minta FT-NIR spektrumát, a másik pedig az eredeti adatkészlet főkomponens-elemzésének főkomponenseit tartalmazta. Az osztályozási modelleket az energiatartalom szerinti elkülönítésére alkottam meg, mely modelleket felhasználtam az osztályozási módszerek összehasonlítására. A kapott eredmények alapján a legjobb módszernek egyértelműen a fejlesztett fák (BT) módszere bizonyult, de mind a négy eljárás meglehetősen jó előrebecslő képességű modelleket eredményezett. Egy másik módszerfejlesztés során továbbfejlesztettem a szelektivitási arány változó kiválasztási módszert, melynek lényegi paraméterét intervallum szelektivitási aránynak neveztem el.

A Q10 koenzim tartalmú étrendkiegészítők vizsgálata során a minták hatóanyag-tartalmát becsültem: PLS regressziós modelleket építettem. Háromféle változó kiválasztási módszert alkalmaztam, amelyek az intervallum PLS-t, az intervallum szelektivitási arányt, illetve a genetikus algoritmus módszerét használták. Mindhárom változó kiválasztási módszer segítségével sikerült megfelelő kalibrációs modelleket létrehozni. A modell által magyarázott variancia, az  $R^2$  értékek, a kalibrációs modellekre mindhárom esetben 0,90 fölöttiek voltak az analóg, de validációra jellemző  $Q^2$ -el jelölt mértékek pedig a mindig 0,87 fölöttiek voltak. A kapott modelleket a rangszámkülönbségek (abszolút) értékeinek összegén alapuló módszer (SRD) segítségével összehasonlítva, a becsült értékek alapján, a genetikus algoritmus alkalmazásával kapott modell tekinthető a legjobbnak.

Az energiatartalom FT-NIR spektrumait felhasználva PLS regressziós modelleket hoztam létre a koffein- és cukortartalom meghatározására. A koffein tartalom meghatározásához referencia módszerként egy HPLC-UV szabvány módszer a saját kísérleteimre optimalizált változatát használtam. Az elkészített koffein regressziós modell  $R^2$  értéke 0,9667 volt, a keresztellenőrzést elvégezve pedig a validált modell  $Q^2$  értéke 0,9279. A cukortartalom meghatározásához két modellt hoztam létre a nominális (89 mintával) és a Schoorl módszerrel (73 mintával) meghatározott referencia értékek alapján. Mindkét modell rendkívül jónak tekinthető, hiszen bár a nominális értékek felhasználásával jobb modellhez jutottam, mindkét

esetben a modellek 0,90 fölötti  $R^2$  és  $Q^2$  értékekkel jellemezhetőek. Az energiatalok adatkészletének felhasználásával megalkottam egy osztályozási modellt is PCA majd LDA módszerek használatával 108 mintára, amely segítségével a taurinos, arginines és normál (taurint és arginint nem tartalmazó) mintákat szinte 100 %-os pontossággal el lehet különíteni. A vizsgálat során a magyar minták mellett, szlovák és görög mintákat is felhasználtam.

Az antioxidáns kapacitás meghatározási technikák elemzése során hétféle különböző technikát hasonlítottam össze. Mindkét vizsgált adatkészlet esetén a főkomponens-elemzés és hierarchikus fürtelemzés (HCA) nevű módszerek azt mutatták, hogy a különböző technikák a meghatározási mód kémiai háttérétől függően csoportosulnak. Mindkét esetben megfigyelhető volt az ACL és ACW módszerek külön csoportba helyeződése vagy elkülönülése a többi módszertől. (az antioxidáns kapacitás mérés rövidítéseinek feloldását a rövidítések jegyzéke tartalmazza). Az SRD eljárás és az általánosított párkorrelációs módszer (GPCM) használatával kimutattam, hogy az átlagértékekhez legjobban közelítő eredményt a FRAP és TPC jelű technikák adták, így ezek a módszerek tekinthetőek a leginkább konzisztensnek. Az általam használt adatkészletek alapján ez utóbbi két módszer kiválthatja a többi vizsgált meghatározási technikát. Az ACL és ACW jelű módszerek teljes mértékben eltérő eredményt adtak, ha a többi módszerhez hasonlítjuk őket, így ezek használata inkább csak a plusz információhordozás (pl mellékreakciók megléte) szempontjából lehet érdekes.

Szintén a kemometriai módszerfejlesztések közé sorolható az az összehasonlítás, amelyet a regressziós modellek során használt teljesítmény paraméterek rangsorolásával végeztem el. A 20 különböző teljesítmény paraméter és két különböző adatkészlet alapján elvégzett kiértékelés segítségével megállapítottam, hogy a belső és külső validálásból származó teljesítmény paraméterek között nagy eltérések tapasztalhatóak, valamint azt, hogy a belső validáláshoz tartozó paraméterek bizonyultak leginkább konzisztensnek és ezáltal jobban ajánlhatók. Ezen teljesítmény paraméterek közül is kiemelkedően jónak tekinthetők: az átlagos négyzetes hiba kereszt-ellenőrzésre vonatkozó értéke ( $RMSE_{CV}$ ), az egyezési tényező kereszt-ellenőrzésre vonatkozó értéke ( $CCC_{CV}$ ), az egy elem kihagyásos validálásra vonatkozó determinációs koefficiens négyzete ( $Q^2_{LOO}$ ) és az átlagos abszolút hiba (MAE) paraméterek.

## 9. SUMMARY

In my PhD thesis several achievements regarding chemometric method development, method comparison and model building for FT-NIR spectra are reported. One of the method developments was a novel implementation of multi class ( $n$ -class) receiver operating characteristic curves (ROC curves) and the application of this process to the evaluation and visualization (with error prediction) of various models based on their classification abilities. Based on the created ROC curves I compared the partial least-square regression (PLS R), random forest (RF), boosted tree (BT) methods and linear discriminant analysis (LDA). I used a dataset with 90 energy drink samples and their FT-NIR spectra. The other dataset contained the PCA components resulting from the principal component analysis of the original matrix. The grouping system was based on the sugar content of the energy drinks and the classification models were used for the comparison of the different classification methods. The boosted tree technique was the best one among them, however all of the used methods produced excellent classification models.

Another method development was the introduction of a new variable selection technique, the interval selectivity ratio (iSR) – based on an existing variable selection technique, the selectivity ratio.

In the study of coenzyme Q10 dietary supplements I predicted the Q10 coenzyme content of the samples with PLS regression models. I applied three variable selection methods: interval PLS, interval selectivity ratio (iSR) and genetic algorithm (GA). All of the three variable selection methods provided good calibration models. The explained variance ( $R^2$ ) was above 0.90 for all of the three models, whereas the same parameter for the validation ( $Q^2$ ) was over 0.87 in every case. The models were compared using sum of (absolute) ranking differences (SRD), based on the predicted values. The result showed that the model produced by the genetic algorithm was the best in my case.

In the case of energy drinks, PLS regression models were built for the determination of sugar and caffeine content of energy drinks based on their FT-NIR spectra. The reference determination of caffeine was based on a standard HPLC-UV procedure with some optimization. The  $R^2$  of the final regression model was 0.967 and the  $Q^2$  value of the validated model was 0.928. I built two models for the determination of sugar concentration: the first was based on the nominal concentrations (with 89 samples), while the second was based on the Schoorl method as a reference measurement (with 73 samples). Both of the models were very good, however the model with nominal concentration values was slightly better than the other one. The  $R^2$  and  $Q^2$

values were above 0.90 in both cases. Three types of energy drinks that contain (i) taurine, (ii) arginine, and (iii) none of these two components were also classified correctly using principal component analysis and linear discriminant analysis. I used 108 Hungarian, Slovakian and Greek samples for the model building process. The correct classification was almost 100 %.

Seven antioxidant capacity methods were compared and evaluated using several statistical methods. The results of principal component analysis (PCA) and hierarchical cluster analysis (HCA) prove that the antioxidant capacity assays based on similar principles are connected to each other closely. In both cases the ACL and ACW techniques were different from the other methods. SRD and the generalized pair correlation method (GPCM) showed that the FRAP and TPC methods were recommended to substitute all the other antioxidant capacity methods for both datasets. These techniques were the most consistent ones. ACL and ACW gave really different results compared to other methods, thus the use of these techniques is recommended in those cases, where they can carry plus information.

Another part of my work that can be considered as chemometric method development is the comparison and ranking of performance parameters of regression models. Based on a comparison of twenty performance parameters on two datasets, I have detected significant differences between performance merits based on cross- and external validation, and concluded that those based on cross-validation are more consistent with the consensus ranking and thus more preferable. In particular,  $RMSE_{CV}$ ,  $CCC_{cv}$ ,  $Q^2_{LOO}$  and MAE are highlighted as particularly good ones.

## 10. KÖSZÖNETNYÍLVÁNÍTÁS

Hálás köszönettel tartozom témavezetőimnek, Dr. Fodor Mariettának és Dr. Héberger Károlynak mindenekelőtt azért, hogy a diákjuk lehettem és azért a rengeteg szakmai segítségért, ami nélkül ez a doktori értekezés nem jöhetett volna létre. „Tudományos szüleimként” egyengették a pályafutásomat, amiért mindig hálás leszek.

Köszönettel tartozok az Alkalmazott Kémia tanszék összes dolgozójának a doktori munkám során nyújtott sok segítségért, illetve külön köszönet illeti Firisz Zsuzsannát is, az örök hangulatfelelőst, aki nélkül unalmasak lettek volna a labor előkészítőben töltött idők.

Szeretném megköszönni Dr. Dernovics Mihálynak, Vass Andreának és Nagy Renátának a HPLC készülék használatában nyújtott segítségüket.

Szintén köszönet illeti az MTA TTK Plazmakémiai Csoportjának összes tagját a támogatásukért és Károly Zoltán csoportvezetőt, hogy lehetővé tette a feltételeket az MTA TTK-ban való munkavégzésemhez. Külön köszönet illeti Dr. Demeter Attilát is a doktori értekezésem gondos javításáért.

Szeretném megköszönni Geréné Radványi Dalmának és Dr. Gere Attilának a doktorandusz éveim felejthetlenné tételét, valamint a sok tanácsot és útmutatást.

Hálás köszönettel tartozok páromnak, Bajusz Dávidnak azért a hihetetlen kitartásért és rengeteg biztatásért, ami nélkül nem tudtam volna végig csinálni az elmúlt három évet. Köszönöm a sok segítséget és tanácsot, amit nem csak a disszertáció elkészítésében, de az összes közös és nem közös munkánk során is nyújtott.

Szeretném megköszönni a Családom összes tagjának, hogy kiálltak mellettem és mindig támogattak, bátorítottak a döntéseimben a doktorandusz éveim alatt is.

Végül köszönet illeti minden barátomat is a biztató és együttérző szavakért, ami miatt soha nem adtam fel.



## 11. IRODALOMJEGYZÉK

- Abourashed, E.A. and Mossa, J.S. (2004), "HPTLC determination of caffeine in stimulant herbal products and power drinks.", *Journal of Pharmaceutical and Biomedical Analysis*, Vol. 36 No. 3, pp. 617–620.
- Ali, F., Rehman, H., Babayan, Z., Stapleton, D. and Joshi, D.D. (2015), "Energy drinks and their adverse health effects: A systematic review of the current evidence", Taylor and Francis Inc., Vol. 127 No. 3, pp. 308–322.
- Andersen, C.M. and Bro, R. (2010), "Variable selection in regression-a tutorial", *Journal of Chemometrics*, Vol. 24, pp. 728–737.
- Di Anibal, C. V, Callao, M.P. and Ruisánchez, I. (2011), "1H NMR variable selection approaches for classification. A case study: the determination of adulterated foodstuffs.", *Talanta*, Elsevier B.V., Vol. 86, pp. 316–323.
- Apak, R., Güçlü, K., Demirata, B., Özyürek, M., Çelik, S.E., Bektaşoğlu, B., Berker, K.I., et al. (2007), "Comparative evaluation of various total antioxidant capacity assays applied to phenolic compounds with the CUPRAC assay", *Molecules*, Vol. 12 No. 7, pp. 1496–1547.
- Aranda, M. and Morlock, G. (2006), "Simultaneous determination of riboflavin, pyridoxine, nicotinamide, caffeine and taurine in energy drinks by planar chromatography-multiple detection with confirmation by electrospray ionization mass spectrometry.", *Journal of Chromatography. A*, Vol. 1131 No. 1-2, pp. 253–260.
- Armenta, S., Garrigues, S. and de la Guardia, M. (2005), "Solid-phase FT-Raman determination of caffeine in energy drinks", *Analytica Chimica Acta*, Vol. 547 No. 2, pp. 197–203.
- Bagchi, T.B., Sharma, S. and Chattopadhyay, K. (2016), "Development of NIRS models to predict protein and amylose content of brown rice and proximate compositions of rice bran", *Food Chemistry*, Vol. 191, pp. 21–27.
- Balázs, G., Bugyi, Z., Gergely, S., Hegyi, A., Hevér, A., Salgó, A. and Tömösközi, S. (2011), "Közeli infravörös spektroszkópiai eljárások", *Élelmiszeranalitika Gyors és Automatizált Módszerei*, Nemzeti Tankönyvkiadó, available at: [http://www.tankonyvtar.hu/hu/tartalom/tamop425/0011\\_2A\\_5\\_modul/1384/index.html](http://www.tankonyvtar.hu/hu/tartalom/tamop425/0011_2A_5_modul/1384/index.html).
- Barker, M. and Rayens, W. (2003), "Partial least squares for discrimination", *Journal of Chemometrics*, Vol. 17 No. 3, pp. 166–173.
- Belardinelli, R., Muçaj, A., Lacalaprice, F., Solenghi, M., Seddaiu, G., Principi, F., Tiano, L., et al. (2006), "Coenzyme Q10 and exercise training in chronic heart failure", *European Heart Journal*, Vol. 27, pp. 2675–2681.
- Benzie, I.F.F. (2000), "Evolution of antioxidant defence mechanisms", *European Journal of Nutrition*, Vol. 39 No. 2, pp. 53–61.
- Benzie, I.F.F. and Strain, J.J. (1996), "The ferric reducing ability of plasma (FRAP) as a measure of 'antioxidant power': The FRAP assay", *Analytical Biochemistry*, Vol. 239 No. 1, pp. 70–76.
- Blázovics, A., Kovács, A., Lugasi, A., Hagymási, K., Bíró, L. and Fehér, J. (1999), "Antioxidant defense in erythrocytes and plasma of patients with active and quiescent Crohn disease and ulcerative colitis: A chemiluminescent study", *Clinical Chemistry*, Vol. 45, pp. 895–896.

- Blois, M.S. (1958), “Antioxidant determination by the use of stable free radicals”, *Nature*, Vol. 4617, pp. 1198–2000.
- Boudkhili, M., Greche, H., Misbahi, H., Giovanelli, S., Noccioli, C., Pistelli, L. and Aarab, L. (2015), “Isolation and antioxidant activity of flavonoids from *Coriaria myrtifolia* methanolic extract”, *Chemistry of Natural Compounds*, Vol. 51, pp. 141–142.
- Breiman, L. (2001), “Random Forests”, *Machine Learning*, Vol. 45 No. 1, pp. 5–32.
- Breithaupt, D.E. and Kraut, S. (2006), “Simultaneous determination of the vitamins A, E, their esters and coenzyme Q10 in multivitamin dietary supplements using an RP-C30 phase”, *European Food Research and Technology*, Vol. 222, pp. 643–649.
- Brereton, R.G. and Lloyd, G.R. (2014), “Partial least squares discriminant analysis: taking the magic away”, *Journal of Chemometrics*, Vol. 28 No. 4, pp. 213–225.
- Bro, R., Kjeldahl, K., Smilde, A.K. and Kiers, H.A.L. (2008), “Cross-validation of component models: A critical look at current methods”, *Analytical and Bioanalytical Chemistry*, Vol. 390 No. 5, pp. 1241–1251.
- Bruker Optik GmbH. (2003), *Bruker MPA User Manual*, Ettlingen, Németország.
- Cadenas, E. (1989), “Biochemistry of oxygen toxicity”, *Annual Review of Biochemistry*, Vol. 58, pp. 79–110.
- Cardeñosa, V., Barreira, J., Barros, L., Arenas-Arenas, F., Moreno-Rojas, J. and Ferreira, I. (2015), “Variety and Harvesting Season Effects on Antioxidant Activity and Vitamins Content of *Citrus sinensis* Macfad.”, *Molecules*, Vol. 20, pp. 8287–8302.
- Cervinkova, B., Krcmova, L.K., Solichova, D., Melichar, B. and Solich, P. (2016), “Recent advances in the determination of tocopherols in biological fluids: from sample pretreatment and liquid chromatography to clinical studies”, *Analytical and Bioanalytical Chemistry*, Vol. 408, pp. 2407–2424.
- Chirico, N. and Gramatica, P. (2011), “Real external predictivity of QSAR models: How to evaluate It? Comparison of different validation criteria and proposal of using the concordance correlation coefficient”, *Journal of Chemical Information and Modeling*, Vol. 51 No. 9, pp. 2320–2335.
- Consonni, V., Ballabio, D. and Todeschini, R. (2009), “Comments on the definition of the Q2 parameter for QSAR validation.”, *Journal of Chemical Information and Modeling*, American Chemical Society, Vol. 49 No. 7, pp. 1669–78.
- Consonni, V., Ballabio, D. and Todeschini, R. (2010), “Evaluation of model predictive ability by external validation techniques”, *Journal of Chemometrics*, Vol. 24 No. 3-4, pp. 194–201.
- Cornelli, U. (2009), “Antioxidant use in nutraceuticals”, *Clinics in Dermatology*, Vol. 27 No. 2, pp. 175–194.
- Cozzolino, D. (2009), “Near infrared spectroscopy in natural products analysis”, *Planta Medica*, Vol. 75 No. 7, available at:<http://doi.org/10.1055/s-0028-1112220>.
- DeLong, E.R., DeLong, D.M. and Clarke-Pearson, D.L. (2016), “Comparing the Areas under Two or More Correlated Receiver Operating Characteristic Curves : A Nonparametric Approach”, *Biometrics*, Vol. 44 No. 3, pp. 837–845.
- Djuric, Z., Depper, J.B., Uhley, V., Smith, D., Lababidi, S., Martino, S. and Heilbrun, L.K.

- (1998), “Oxidative DNA damage levels in blood from women at high risk for breast cancer are associated with dietary intakes of meats, vegetables, and fruits”, *Journal of the American Dietetic Association*, Vol. 98 No. 5, pp. 524–528.
- Ernster, L. and Dallner, G. (1995), “Biochemical, physiological and medical aspects of ubiquinone function”, *Biochimica et Biophysica Acta*, Vol. 1271, pp. 195–204.
- Esbensen, K.H. and Geladi, P. (2010), “Principles of proper validation: Use and abuse of re-sampling for validation”, *Journal of Chemometrics*, Vol. 24 No. 3-4, pp. 168–187.
- Escuredo, O., Carmen Seijo, M., Salvador, J. and Inmaculada González-Martín, M. (2013), “Near infrared spectroscopy for prediction of antioxidant compounds in the honey”, *Food Chemistry*, Vol. 141 No. 4, pp. 3409–3414.
- Ferreira, S.E., de Mello, M.T., Pompéia, S. and de Souza-Formigoni, M.L.O. (2006), “Effects of energy drink ingestion on alcohol intoxication.”, *Alcoholism, Clinical and Experimental Research*, Vol. 30 No. 4, pp. 598–605.
- Fodor, M., Woller, A., Turza, S. and Szigedi, T. (2011), “Development of a rapid, non-destructive method for egg content determination in dry pasta using FT-NIR technique”, *Journal of Food Engineering*, Vol. 107, pp. 195–199.
- Fotino, A.D., Thompson-Paul, A.M. and Bazzano, L.A. (2013), “Effect of coenzyme Q10 supplementation on heart failure: a meta-analysis”, *The American Journal of Clinical Nutrition*, Vol. 97, pp. 268–275.
- Friedman, J.H. (1991a), “Multivariate Adaptive Regression Splines”, *The Annals of Statistics*, Vol. 19, pp. 1–67.
- Friedman, J.H. (1991b), “Multivariate Adaptive Regression Splines”, *The Annals of Statistics*, Vol. 19 No. 1, pp. 1–67.
- Froufe, H.J.C., Abreu, R.M. V and Ferreira, I.C.F.R. (2011), “QCAR models to predict wild mushrooms radical scavenging activity, reducing power and lipid peroxidation inhibition”, *Chemometrics and Intelligent Laboratory Systems*, Vol. 109 No. 2, pp. 192–196.
- Geladi, P. and Kowalski, B.R. (1986), “Partial least-squares regression: a tutorial”, *Analytica Chimica Acta*, Vol. 185 No. C, pp. 1–17.
- Geladi, P., MacDougall, D. and Martens, H. (1985), “Linearization and Scatter-Correction for Near-Infrared Reflectance Spectra of Meat”, *Applied Spectroscopy*, Vol. 39 No. 3, pp. 491–500.
- Gramatica, P. (2007), “Principles of QSAR models validation: Internal and external”, *QSAR and Combinatorial Science*, Vol. 26 No. 5, pp. 694–701.
- Grant, D.C. and Helleur, R.J. (2008), “Simultaneous analysis of vitamins and caffeine in energy drinks by surfactant-mediated matrix-assisted laser desorption/ionization”, *Analytical and Bioanalytical Chemistry*, Vol. 391, pp. 2811–2818.
- Gutiérrez, S., Tardaguila, J., Fernández-Novales, J. and Diago, M.P. (2015), “Support Vector Machine and Artificial Neural Network Models for the Classification of Grapevine Varieties Using a Portable NIR Spectrophotometer”, edited by Scali, M. *PLOS ONE*, Vol. 10 No. 11, p. e0143197.
- Gütlein, M., Helma, C., Karwath, A. and Kramer, S. (2013), “A Large-Scale Empirical

- Evaluation of Cross-Validation and External Test Set Validation in (Q)SAR”, *Molecular Informatics*, Vol. 32, pp. 516–528.
- Halliwell, B. and Gutteridge, J.M.C. (1995), “The definition and measurement of antioxidants in biological systems”, *Free Radical Biology and Medicine*, Vol. 18 No. 1, pp. 125–126.
- Hanley, J.A. and McNeil, B.J. (1982), “The Meaning and Use of the Area under a Receiver Operating Characteristic (ROC) Curve”, *Radiology*, Vol. 143, pp. 29–36.
- Hargreaves, I.P. (2014), “Coenzyme Q10 as a therapy for mitochondrial disease”, *International Journal Of Biochemistry & Cell Biology*, Elsevier Science, Vol. 49, pp. 105–111.
- Hastie, T., Tibshirani, R. and Friedman, J. (2001a), “Model Assessment and Selection”, *The Elements of Statistical Learning; Data Mining, Inference, and Prediction*, Springer, New York, pp. 214–216.
- Hastie, T., Tibshirani, R. and Friedman, J. (2001b), “Unsupervised Learning”, *The Elements of Statistical Learning; Data Mining, Inference, and Prediction*, 1st ed., Springer, New York, pp. 472–475.
- Hastie, T., Tibshirani, R. and Friedman, J. (2001c), “Linear Methods for Classification”, *The Elements of Statistical Learning; Data Mining, Inference, and Prediction*, Springer, New York, pp. 84–90.
- Hastie, T., Tibshirani, R. and Friedman, J. (2001d), “Overview of Supervised Learning”, *The Elements of Statistical Learning; Data Mining, Inference, and Prediction*, Springer, New York, p. 31.
- Hastie, T., Tibshirani, R. and Friedman, J.H. (2009), “Boosting and additive trees”, *The Elements of Statistical Learning; Data Mining, Inference, and Prediction*, 2. ed., Springer, pp. 337–389.
- Héberger, K. (2010), “Sum of ranking differences compares methods or models fairly”, *TrAC Trends in Analytical Chemistry*, Vol. 29 No. 1, pp. 101–109.
- Héberger, K. and Rajkó, R. (2002), “Generalization of pair correlation method (PCM) for non-parametric variable selection”, *Journal of Chemometrics*, Vol. 16 No. 8-10, pp. 436–443.
- Heckman, M.A., Weil, J. and de Mejia, E.G. (2010), “Caffeine (1, 3, 7-trimethylxanthine) in foods: A comprehensive review on consumption, functionality, safety, and regulatory matters”, *Journal of Food Science*, Vol. 75 No. 3, pp. 77–87.
- Horovitz, W. (1975), “Sugar and Sugar products”, *Official Methods of Analysis of the Association of Official Analytical Chemists*, AOAC, Washington, USA, p. 573.
- Huang, D., Boxin, O.U. and Prior, R.L. (2005), “The chemistry behind antioxidant capacity assays”, *Journal of Agricultural and Food Chemistry*, Vol. 53 No. 6, pp. 1841–1856.
- Huang, H., Liu, L. and Ngadi, M.O. (2016), “Prediction of pork fat attributes using NIR Images of frozen and thawed pork”, *Meat Science*, Vol. 119, pp. 51–61.
- Huang, H., Yu, H., Xu, H. and Ying, Y. (2008a), “Near infrared spectroscopy for on/in-line monitoring of quality in foods and beverages: A review”, *Journal of Food Engineering*, Vol. 87, pp. 303–313.
- Huang, H., Yu, H., Xu, H. and Ying, Y. (2008b), “Near infrared spectroscopy for on/in-line monitoring of quality in foods and beverages: A review”, *Journal of Food Engineering*,

Vol. 87 No. 3, pp. 303–313.

- Huck, C.W., Guggenbichler, W. and Bonn, G.K. (2005), “Analysis of caffeine, theobromine and theophylline in coffee by near infrared spectroscopy (NIRS) compared to high-performance liquid chromatography (HPLC) coupled to mass spectrometry”, *Analytica Chimica Acta*, Vol. 538 No. 1-2, pp. 195–203.
- Internet. “Bruker Optik GmbH”, available at: [www.bruker.com](http://www.bruker.com), 2016.09.13.
- Kaffka, K.J. and Norris, K.H. (1976), “Rapid instrumental analysis of composition of wine”, *Acta Alimentaria*, Vol. 5 No. 3, pp. 267–279.
- Khasanov, V. V., Slizhov, Y.G. and Khasanov, V. V. (2013), “Energy drink analysis by capillary electrophoresis”, *Journal of Analytical Chemistry*, Vol. 68 No. 4, pp. 357–359.
- Kjeldahl, K. and Bro, R. (2010), “Some common misunderstandings in chemometrics”, *Journal of Chemometrics*, Vol. 24 No. 7-8, pp. 558–564.
- Kollár-Hunek, K. and Héberger, K. (2013), “Method and model comparison by sum of ranking differences in cases of repeated observations (ties)”, *Chemometrics and Intelligent Laboratory Systems*, Elsevier B.V., Vol. 127, pp. 139–146.
- Kotali, A., Nasiopoulou, D.A., Tsoleridis, C.A., Harris, P.A., Kontogiorgis, C.A. and Hadjipavlou-Litina, D.J. (2016), “Antioxidant activity of 3-[N-(acylhydrazono)ethyl]-4-hydroxy-coumarins”, *Molecules*, Vol. 21, p. 138.
- Kozik, T.M., Shah, S., Bhattacharyya, M., Franklin, T.T., Connolly, T.F., Chien, W., Charos, G.S., et al. (2016), “Cardiovascular responses to energy drinks in a healthy population: The C-energy study”, *W.B. Saunders*, Vol. 34 No. 7, pp. 1205–1209.
- Kraujalyte, V., Venskutonis, P.R., Pukalskas, A., Česonienė, L. and Daubaras, R. (2015), “Antioxidant properties, phenolic composition and potentiometric sensor array evaluation of commercial and new blueberry (*Vaccinium corymbosum*) and bog blueberry (*Vaccinium uliginosum*) genotypes”, *Food Chemistry*, Vol. 188, pp. 583–590.
- Kristensen, M., Savorani, F., Ravn-Haren, G., Poulsen, M., Markowski, J., Larsen, F.H., Dragsted, L.O., et al. (2010), “NMR and interval PLS as reliable methods for determination of cholesterol in rodent lipoprotein fractions”, *Metabolomics*, Vol. 6 No. 1, pp. 129–136.
- Ku, H.H. (1966), “Notes on the use of propagation of error formulas”, *Journal of Research of the National Bureau of Standards, Section C: Engineering and Instrumentation*, Vol. 70C No. 4, pp. 263–273.
- Lang, J.K. and Packer, L. (1987), “Quantitative determination of vitamin E and oxidized and reduced coenzyme Q by high-performance liquid chromatography with in-line ultraviolet and electrochemical detection”, *Journal of Chromatography*, Vol. 385, pp. 109–117.
- Lásztity, R. and Törley, D. (1987), *Az élelmiszer Analitika Elméleti Alapjai 1.*, Mezőgazdasági Kiadó.
- Leardi, R. (2007), “Genetic algorithms in chemistry.”, *Journal of Chromatography A*, Vol. 1158, pp. 226–233.
- Lee, B.-J., Huang, Y.-C., Chen, S.-J. and Lin, P.-T. (2012), “Coenzyme Q10 supplementation reduces oxidative stress and increases antioxidant enzyme activity in patients with coronary artery disease”, *Nutrition*, Elsevier, Vol. 28, pp. 250–255.

- Lee, J.H., Hoang, N.H., Huong, N.L., Shrestha, A. and Park, J.W. (2014), "Ultra-Performance Liquid Chromatography with Electrospray Ionization Mass Spectrometry for the Determination of Coenzyme Q10 as an Anti-Aging Ingredient in Edible Cosmetics", *Analytical Letters*, Vol. 47, pp. 367–376.
- León, L., Kelly, J.D. and Downey, G. (2005), "Detection of apple juice adulteration using near-infrared transfectance spectroscopy", *Applied Spectroscopy*, Vol. 59, pp. 593–599.
- Lin, L.I.-K. (1989), "A concordance correlation coefficient to evaluate reproducibility", *Biometrics*, Vol. 45 No. 1, pp. 255–68.
- Lin, L.I.-K. (1992), "Assay Validation Using the Concordance Correlation Coefficient", *Biometrics*, Vol. 48 No. 2, p. 599.
- Lockwood, K., Moesgaard, S., Yamamoto, T. and Folkers, K. (1995), "Progress on therapy of breast cancer with vitamin Q10 and the regression of metastases", *Biochemical and Biophysical Research Communications*, Vol. 212, pp. 172–177.
- López-Lluch, G., Rodríguez-Aguilera, J.C., Santos-Ocaña, C. and Navas, P. (2010), "Is coenzyme Q a key factor in aging?", *Mechanisms of Ageing and Development*, Elsevier Science, Vol. 131, pp. 225–235.
- Lopo, M., Páscoa, R.N.M.J., Graça, A.R. and Lopes, J.A. (2016), "Classification of Vineyard Soils Using Portable and Benchtop Near-Infrared Spectrometers: A Comparative Study", *Soil Science Society of America Journal*, Vol. 80 No. 3, p. 652.
- Lucena, R., Cárdenas, S., Gallego, M. and Valcárcel, M. (2005), "Continuous flow autoanalyzer for the sequential determination of total sugars, colorant and caffeine contents in soft drinks", *Analytica Chimica Acta*, Vol. 530 No. 2, pp. 283–289.
- Lunetta, S., Roman, M., Chandrah, A., Edamura, T., Honda, T., LeVanseler, K., Nagarajan, M., et al. (2008), "Determination of coenzyme Q10 content in raw materials and dietary supplements by high-performance liquid chromatography-UV: Collaborative study", *Journal of AOAC International*, Vol. 91, pp. 702–708.
- Magalhães, L.M., Machado, S., Segundo, M.A., Lopes, J.A. and Páscoa, R.N.M.J. (2016), "Rapid assessment of bioactive phenolics and methylxanthines in spent coffee grounds by FT-NIR spectroscopy", *Talanta*, Elsevier, Vol. 147, pp. 460–467.
- Magwaza, L.S., Opara, U.L., Terry, L. a., Landahl, S., Cronje, P.J.R., Nieuwoudt, H.H., Hanssens, A., et al. (2013), "Evaluation of Fourier transform-NIR spectroscopy for integrated external and internal quality assessment of Valencia oranges", *Journal of Food Composition and Analysis*, Vol. 31 No. 1, pp. 144–154.
- Malik, V.S., Schulze, M.B. and Hu, F.B. (2006), "Intake of sugar-sweetened beverages and weight gain: A systematic review", *American Journal of Clinical Nutrition*.
- Malinauskas, B.M., Aeby, V.G., Overton, R.F., Carpenter-Aeby, T. and Barber-Heidal, K. (2007), "A survey of energy drink consumption patterns among college students.", *Nutrition Journal*, Vol. 6, p. 35.
- Martinčič, R., Kuzmanovski, I., Wagner, A. and Novič, M. (2015), "Development of models for prediction of the antioxidant activity of derivatives of natural compounds", *Analytica Chimica Acta*, Vol. 868, pp. 23–35.
- Mattila, P. and Kumpulainen, J. (2001), "Coenzymes Q9 and Q10: Contents in Foods and

- Dietary Intake”, *Journal of Food Composition and Analysis*, Vol. 14, pp. 409–417.
- Melgarejo-Sánchez, P., Martínez, J.J., Legua, P., Martínez, R., Hernández, F. and Melgarejo, P. (2015), “Quality, antioxidant activity and total phenols of six Spanish pomegranates clones”, *Scientia Horticulturae*, Elsevier B.V., Vol. 182, pp. 65–72.
- Metz, C.E. (1978), “Basic principles of ROC analysis”, *Seminars in Nuclear Medicine*, Vol. 8 No. 4, pp. 283–298.
- Miller, N.J., Rice-Evans, C., Davies, M.J., Gopinathan, V. and Milner, A. (1993), “A novel method for measuring antioxidant capacity and its application to monitoring the antioxidant status in premature neonates”, *Clinical Science*, Vol. 84 No. 4, pp. 407–412.
- Monakhova, Y.B., Ruge, I., Kuballa, T., Lerch, C. and Lachenmeier, D.W. (2013), “Rapid determination of coenzyme Q10 in food supplements using 1H NMR spectroscopy”, *International Journal for Vitamin and Nutrition Research*, Vol. 83, pp. 67–72.
- Moo-Huchin, V.M., Estrada-Mota, I., Estrada-León, R., Cuevas-Glory, L., Ortiz-Vázquez, E., De Lourdes Vargas Y Vargas, M., Betancur-Ancona, D., et al. (2014), “Determination of some physicochemical characteristics, bioactive compounds and antioxidant activity of tropical fruits from Yucatan, Mexico”, *Food Chemistry*, Vol. 152, pp. 508–515.
- Moorthy, N.S.H.N., Cerqueira, N.M.F.S. a., Ramos, M.J. and Fernandes, P. a. (2015), “Ligand based analysis on HMG-CoA reductase inhibitors”, *Chemometrics and Intelligent Laboratory Systems*, Elsevier B.V., Vol. 140, pp. 102–116.
- Naes, T., Isaksson, T., Fearn, T. and Davies, T. (2002), “Validation”, *A User-Friendly Guide to Multivariate Calibration and Classification*, NIR Publications, Chicester, pp. 155–177.
- Nicholls, A. (2014), “Confidence limits, error bars and method comparison in molecular modeling. Part 1: The calculation of confidence intervals”, *Journal of Computer-Aided Molecular Design*, Kluwer Academic Publishers, Vol. 28 No. 9, pp. 887–918.
- Nørgaard, L., Saudland, A., Wagner, J., Nielsen, J.P., Munck, L. and Engelsen, S.B. (2000), “Interval partial least-squares regression (iPLS): A comparative chemometric study with an example from near-infrared spectroscopy”, *Applied Spectroscopy*, Vol. 54 No. 3, pp. 413–419.
- Nowik, W., Héron, S., Bonose, M. and Tchaplá, A. (2013), “Separation system suitability (3S): a new criterion of chromatogram classification in HPLC based on cross-evaluation of separation capacity/peak symmetry and its application to complex mixtures of anthraquinones.”, *The Analyst*, Vol. 138, pp. 5801–5810.
- O’Brien, M.C., McCoy, T.P., Rhodes, S.D., Wagoner, A. and Wolfson, M. (2008), “Caffeinated cocktails: Energy drink consumption, high-risk drinking, and alcohol-related consequences among college students”, *Academic Emergency Medicine*, Vol. 15, pp. 453–460.
- Ojha, P.K., Mitra, I., Das, R.N. and Roy, K. (2011), “Further exploring rm2 metrics for validation of QSPR models”, *Chemometrics and Intelligent Laboratory Systems*, Vol. 107 No. 1, pp. 194–205.
- Otto, M. (1999), “Pattern recognition and classification”, *Chemometrics*, 1st ed., Wiley–VCH, Weinheim, Germany, pp. 148–156.
- Pieszko, C., Baranowska, I. and Flores, A. (2010), “Determination of energizers in energy drinks”, *Journal of Analytical Chemistry*, SP MAIK Nauka/Interperiodica, Vol. 65 No. 12,

pp. 1228–1234.

- Pokol, G., Gyurcsányi, E.R., Simon, A., Bezúr, L., Horvai, G., Horváth, V. and Dudás, K.M. (2011), “Infravörös (IR) spektroszkópia”, *Analitikai Kémia*, Typotex Kiadó, pp. 270–280.
- Popov, I.N. and Lewin, G. (1994), “Photochemiluminescent detection of antiradical activity: II. Testing of nonenzymic water-soluble antioxidants”, *Free Radical Biology and Medicine*, Vol. 17 No. 3, pp. 267–271.
- Popov, I.N. and Lewin, G. (1996), “Photochemiluminescent detection of antiradical activity; IV: Testing of lipid-soluble antioxidants”, *Journal of Biochemical and Biophysical Methods*, Vol. 31 No. 1-2, pp. 1–8.
- Provost, F. and Domingos, P. (2000), *Well-Trained PETs: Improving Probability Estimation Trees*, No. IS-00-04, *CeDER Working Paper*, available at:<http://doi.org/10.1.1.33.309>.
- Rajalahti, T., Arneberg, R., Berven, F.S., Myhr, K.-M., Ulvik, R.J. and Kvalheim, O.M. (2009), “Biomarker discovery in mass spectral profiles by means of selectivity ratio plot”, *Chemometrics and Intelligent Laboratory Systems*, Elsevier B.V., Vol. 95 No. 1, pp. 35–48.
- Rajkó, R. and Héberger, K. (2001), “Conditional Fisher’s exact test as a selection criterion for pair-correlation method. Type I and Type II errors”, *Chemometrics and Intelligent Laboratory Systems*, Vol. 57 No. 1, pp. 1–14.
- Ratnam, D. V., Bhardwaj, V. and Kumar, M.N.V.R. (2006), “Simultaneous analysis of ellagic acid and coenzyme Q10 by derivative spectroscopy and HPLC”, *Talanta*, Vol. 70, pp. 387–391.
- Reissig, C.J., Strain, E.C. and Griffiths, R.R. (2009), “Caffeinated energy drinks-A growing problem”, *Drug and Alcohol Dependence*, Vol. 99, pp. 1–10.
- Ribeiro, J.S., Ferreira, M.M.C. and Salva, T.J.G. (2011), “Chemometric models for the quantitative descriptive sensory analysis of Arabica coffee beverages using near infrared spectroscopy”, *Talanta*, Elsevier B.V., Vol. 83 No. 5, pp. 1352–1358.
- Rodríguez-Acuña, R., Brenne, E. and Lacoste, F. (2008), “Determination of Coenzyme Q10 and Q9 in Vegetable Oils”, *Journal of Agricultural and Food Chemistry*, American Chemical Society, Vol. 56, pp. 6241–6245.
- Rodriguez-Saona, L.E., Fry, F.S., McLaughlin, M.A. and Calvey, E.M. (2001), “Rapid analysis of sugars in fruit juices by FT-NIR spectroscopy”, *Carbohydrate Research*, Vol. 336 No. 1, available at:[http://doi.org/10.1016/S0008-6215\(01\)00244-0](http://doi.org/10.1016/S0008-6215(01)00244-0).
- Roggo, Y., Chalus, P., Maurer, L., Lema-Martinez, C., Edmond, A. and Jent, N. (2007), “A review of near infrared spectroscopy and chemometrics in pharmaceutical technologies”, *Journal of Pharmaceutical and Biomedical Analysis*, Vol. 44 No. 3 SPEC. ISS., pp. 683–700.
- Roy, P.P. and Roy, K. (2008), “On Some Aspects of Variable Selection for Partial Least Squares Regression Models”, *QSAR & Combinatorial Science*, Vol. 27 No. 3, pp. 302–313.
- Rücker, C., Rücker, G. and Meringer, M. (2007), “ $\gamma$ -Randomization and its variants in QSPR/QSAR.”, *Journal of Chemical Information and Modeling*, American Chemical Society, Vol. 47 No. 6, pp. 2345–57.
- Salgó, A., Nagy, J., Mikó, É. and Boros, I. (1998), “Application of near infrared spectroscopy in



- the sugar industry”, *Journal of Near Infrared Spectroscopy*, Vol. 6 No. 1-4.
- Schüürmann, G., Ebert, R.-U., Chen, J., Wang, B. and Kühne, R. (2008), “External validation and prediction employing the predictive squared correlation coefficient test set activity mean vs training set activity mean.”, *Journal of Chemical Information and Modeling*, American Chemical Society, Vol. 48 No. 11, pp. 2140–5.
- Seifert, S.M., Schaechter, J.L., Hershorin, E.R. and Lipshultz, S.E. (2011), “Health effects of energy drinks on children, adolescents, and young adults.”, *Pediatrics*, Vol. 127 No. 3, pp. 511–528.
- Sereshti, H. and Samadi, S. (2014), “A rapid and simple determination of caffeine in teas, coffees and eight beverages.”, *Food Chemistry*, Elsevier Ltd, Vol. 158, pp. 8–13.
- Shi, L.M., Fang, H., Tong, W., Wu, J., Perkins, R., Blair, R.M., Branham, W.S., et al. (2001), “QSAR Models Using a Large Diverse Set of Estrogens”, *Journal of Chemical Information and Modeling*, American Chemical Society, Vol. 41 No. 1, pp. 186–195.
- Singleton, V.L., Orthofer, R. and Lamuela-Raventós, R.M. (1998), “Analysis of total phenols and other oxidation substrates and antioxidants by means of folin-ciocalteu reagent”, *Methods in Enzymology*.
- Singleton, V.L. and Rossi, J.A. (1965), “Colorimetry of total phenolics [in grapes and wine] with phosphomolybdic-phosphotungstic acid reagents”, *American Journal of Enology and Viticulture*, Vol. 16, pp. 144–158.
- Sinija, V.R. and Mishra, H.N. (2009), “FT-NIR spectroscopy for caffeine estimation in instant green tea powder and granules”, *LWT - Food Science and Technology*, Vol. 42 No. 5, available at:<http://doi.org/10.1016/j.lwt.2008.12.013>.
- Slavova-Kazakova, A., Karamać, M., Kancheva, V. and Amarowicz, R. (2015), “Antioxidant Activity of Flaxseed Extracts in Lipid Systems”, *Molecules*, Vol. 21, p. 17.
- Subramanian, A. and Rodriguez-Saona, L. (2009), “Fourier Transform Infrared (FTIR) Spectroscopy”, in Sun, D.-W. (Ed.), *Infrared Spectroscopy for Food Quality Analysis and Control*, 1st ed., Vol. 1, Academic Press, pp. 145–174.
- Szigedi, T., Dernovics, M. and Fodor, M. (2011), “Determination of protein, lipid and sugar contents of bakery products by using fourier-transform near infrared spectroscopy”, *Acta Alimentaria*, Vol. 40 No. SUPPL. 21, available at:<http://doi.org/10.1556/AAlim.40.2011.Suppl.21>.
- Tahir, H.E., Xiaobo, Z., Tinting, S., Jiyong, S. and Mariod, A.A. (2016), “Near-Infrared (NIR) Spectroscopy for Rapid Measurement of Antioxidant Properties and Discrimination of Sudanese Honeys from Different Botanical Origin”, *Food Analytical Methods*, Vol. 9 No. 9, pp. 2631–2641.
- Tang, P.H., Miles, M. V, DeGrauw, A., Hershey, A. and Pesce, A. (2001), “HPLC analysis of reduced and oxidized coenzyme Q10 in human plasma”, *Clinical Chemistry*, Vol. 47, pp. 256–265.
- Todeschini, R., Consonni, V. and Maiocchi, A. (1999), “The K correlation index: theory development and its application in chemometrics”, *Chemometrics and Intelligent Laboratory Systems*, Vol. 46, pp. 13–29.
- Tóth, G., Bodai, Z. and Héberger, K. (2013), “Estimation of influential points in any data set

from coefficient of determination and its leave-one-out cross-validated counterpart”, *Journal of Computer-Aided Molecular Design*, Vol. 27 No. 10, pp. 837–844.

- Vass, A., Deák, E. and Dernovics, M. (2014), “Quantification of the Reduced Form of Coenzyme Q10, Ubiquinol, in Dietary Supplements with HPLC-ESI-MS/MS”, *Food Analytical Methods*, Food Analytical Methods, available at:<http://doi.org/10.1007/s12161-014-9911-x>.
- Vochyánová, B., Opekar, F., Tůma, P. and Štulík, K. (2012), “Rapid determinations of saccharides in high-energy drinks by short-capillary electrophoresis with contactless conductivity detection.”, *Analytical and Bioanalytical Chemistry*, Vol. 404 No. 5, pp. 1549–54.
- Wiklund, S., Nilsson, D., Eriksson, L., Sjöström, M., Wold, S. and Faber, K. (2007), “A randomization test for PLS component selection”, *Journal of Chemometrics*, Vol. 21 No. 10-11, pp. 427–439.
- Williams, P.C. (2001), “Implementation of near infrared spectroscopy”, in Williams, P. and Norris, K. (Eds.), *Near-Infrared Technology in the Agricultural and Food Industries*, American Association of Cereal Chemists, pp. 145–169.
- Wold, S. (1976), “Pattern recognition by means of disjoint principal components models”, *Pattern Recognition*, Vol. 8 No. 3, pp. 127–139.
- Wold, S., Esbensen, K. and Geladi, P. (1987), “Principal component analysis”, *Chemometrics and Intelligent Laboratory Systems*, Vol. 2 No. 1-3, pp. 37–52.
- Workman, J.J. (2000), “Functional groupings and calculated locations in wavenumbers (cm<sup>-1</sup>) for IR spectroscopy”, *The Handbook of Organic Compounds: NIR, IR, Raman, and UV-Vis Spectra Featuring Polymers and Surfactants*, Academic Press, San Diego, pp. 229–236.
- Workman, J.J. and Weyer, L. (2007), *Practical Guide to Interpretive Near-Infrared Spectroscopy*, CRC Press, available at:  
<http://books.google.com/books?id=TJZYWo7gUN8C&pgis=1> (accessed 9 July 2014).
- Wu, Z., Xu, E., Long, J., Pan, X., Xu, X., Jin, Z. and Jiao, A. (2015), “Comparison between ATR-IR, Raman, concatenated ATR-IR and Raman spectroscopy for the determination of total antioxidant capacity and total phenolic content of Chinese rice wine”, *Food Chemistry*, Elsevier Ltd, Vol. 194, pp. 671–679.
- Wu, Z., Xu, E., Long, J., Wang, F., Xu, X., Jin, Z. and Jiao, A. (2015), “Rapid Measurement of Antioxidant Activity and  $\gamma$ -Aminobutyric Acid Content of Chinese Rice Wine by Fourier-Transform Near Infrared Spectroscopy”, *Food Analytical Methods*, Springer New York LLC, Vol. 8 No. 10, pp. 2541–2553.
- Yang, H., Irudayaraj, J. and Paradkar, M.M. (2005), “Discriminant analysis of edible oils and fats by FTIR, FT-NIR and FT-Raman spectroscopy”, *Food Chemistry*, Vol. 93 No. 1, pp. 25–32.
- Ying, Y. and Liu, Y. (2008), “Nondestructive measurement of internal quality in pear using genetic algorithms and FT-NIR spectroscopy”, *Journal of Food Engineering*, Vol. 84 No. 2, pp. 206–213.
- Youden, W.J. (1975), *Statistical Manual of the Association of Official Analytical Chemists: Statistical Techniques for Collaborative Test*, AOAC International, Gaithersburg.

- Yubero, D., Montero, R., Ramos, M., Neergheen, V., Navas, P., Artuch, R. and Hargreaves, I. (2015), "Determination of urinary coenzyme Q10 by HPLC with electrochemical detection: Reference values for a paediatric population", *BioFactors*, Vol. 41 No. 6, pp. 424–430.
- Zhang, X., Li, W., Yin, B., Chen, W., Kelly, D.P., Wang, X., Zheng, K., et al. (2013), "Improvement of near infrared spectroscopic (NIRS) analysis of caffeine in roasted arabica coffee by variable selection method of stability competitive adaptive reweighted sampling (SCARS)", *Spectrochimica Acta - Part A: Molecular and Biomolecular Spectroscopy*, Elsevier B.V., Vol. 114, pp. 350–356.

## 12. MELLÉKLETEK

**M1** A teljesítményparaméterek összehasonlításához felhasznált paraméterek listája

Teljesítmény paraméter	Mikor számolható?	Formula <sup>b</sup>
$R^2, R^2_{ext}$	Tréning, belső és külső validálás	$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{RSS}{TSS}$
$R^2_{adj.}$	Tréning	$R^2_{adj.} = R^2 - (1 - R^2) \times \frac{p}{n - p - 1}$
$R^2 - R^2_{adj.}$	Tréning	Lásd. fent
LOF (Friedman, 1991b)	Tréning	$LOF = \frac{RSS}{n \left( 1 - \frac{M + d(M-1)/2}{n} \right)^2}$
$K_x$	Tréning	PCA alapján, lásd. (Todeschini et al., 1999)
$\Delta K$	Tréning	PCA alapján, lásd. (Todeschini et al., 1999)
RMSE	Tréning, belső és külső validálás	$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$
MAE	Tréning, belső és külső validálás	$MAE = \frac{\sum_{i=1}^n  y_i - \hat{y}_i }{n}$
RSS	Tréning	$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
CCC (Lin, 1989, 1992)	Tréning, belső és külső validálás	$CCC = \frac{2 \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \hat{y})}{\sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - \hat{y})^2 + n(\bar{y} - \hat{y})^2}$
$s$	Tréning	$s = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$
$F$	Tréning	$F = \left( \frac{\sum_{i=1}^N (\bar{y} - \hat{y}_i)^2}{p-1} \right) / \left( \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{n-p} \right)$
$Q^2_{LoO}$	Belső validálás	$Q^2_{LoO} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_{i/i})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{PRESS}{TSS}$
$R^2 - Q^2_{LoO}$	Belső validálás	Lásd. fent
PRESS	Belső és külső validálás	$PRESS = \sum_{i=1}^n (y_i - \hat{y}_{i/i})^2$

$Q^2_{LMO}$	Belső validálás	$Q^2_{LMO} = 1 - \frac{\sum_{j=1}^m \sum_{i=1}^n (y_i - \hat{y}_{i/j})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$
$R^2_{Y-SCRAMBLE}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$RMSE_{Avg, Y-SCRAMBLE}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$Q^2_{Y-SCRAMBLE}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$R^2_{RND-DESCR}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$Q^2_{RND-DESCR}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$R^2_{RND-RESP}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$Q^2_{RND-RESP}$ (Rücker et al., 2007)	Belső validálás	Lásd. fent
$Q^2_{F1}$ (Consonni et al., 2010; Schüürmann et al., 2008)	Külső validálás	$Q^2_{F1} = 1 - \frac{\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y}_{TR})^2}$
$Q^2_{F2}$ (Consonni et al., 2010; Shi et al., 2001)	Külső validálás	$Q^2_{F2} = 1 - \frac{\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y}_{EXT})^2}$
$Q^2_{F3}$ (Consonni et al., 2009, 2010)	Külső validálás	$Q^2_{F3} = 1 - \frac{\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i)^2 / n_{EXT}}{\sum_{i=1}^{n_{TR}} (y_i - \bar{y}_{TR})^2 / n_{TR}}$
$\overline{r^2}_m$ (Ojha et al., 2011; Roy and Roy, 2008)	Külső validálás	$\overline{r^2}_m = \frac{r_m^2 + r'^2_m}{2}$
$\Delta r^2_m$	Külső validálás	$\Delta r^2_m = r_m^2 - r'^2_m$


<sup>b</sup> A teljesítmény paraméterek egyenleteiben használt szöveges kifejezések:  $tr$  = tréning,  $cv$  = kereszt-ellenőrzés,  $ext$  = külső validálás. A jelölések magyarázata:  $y_i$ : referencia érték,  $\bar{y}$ : a mért értékek átlaga,  $\hat{y}_i$ : becsült érték,  $\bar{\hat{y}}$ : becsült értékek átlaga,  $\hat{y}_{i/i}$ : az  $i$ -dik minta becsült értéke, ha az  $i$ -dik mintát kihagyjuk a kalibrációs alcsoportból,  $\hat{y}_{i/j}$ : az  $i$ -dik minta becsült értéke, ha a  $j$ -edik részt kihagyjuk a kalibrációs alcsoportból (a teljes adatkészlet  $m$  részre van felosztva),  $n$ : mintaszám,  $i$ : minta indexelése,  $p$ : a változók száma a modellben.

## M2 A HPLC mérésekhez használt eluensek táblázata

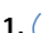
Név	Forgalmazó
Acetonitril (HPLC minőség)	Scharlau (Barcelona, Spanyolország)
Tetrahidrofurán (izokratikus HPLC minőség)	VWR (Radnor, PA, USA)
Standard Q10 koenzim ( $\geq 98\%$ )	Sigma–Aldrich csoport (Schnelldorf, Németország)
Metanol (HPLC minőség)	Scharlau (barcelona, Spanyolország)
Koffein standard ( $\geq 98\%$ )	Sigma–Aldrich csoport (Schnelldorf, Németország)
Ultranagy tisztaságú desztillált víz (Milli–Q rendszer)	Merck–Millipore (Milford, MA, USA)


### M3 A ROC görbék készítésének sémája

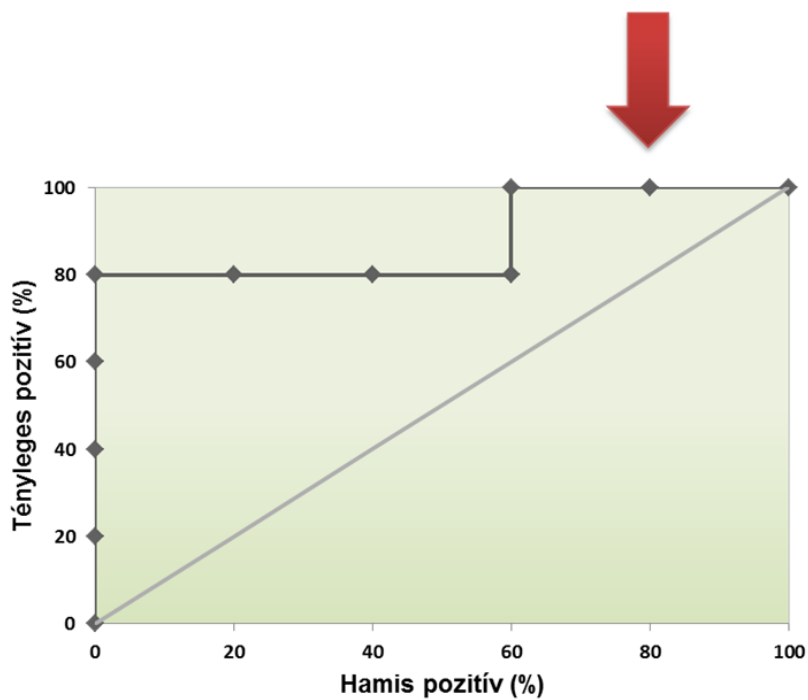
Minta	Csoport	Becsült valószínűség
1	+	0.987
2	+	0.876
3	-	0.758
4	+	0.547
5	-	0.384
6	-	0.765
7	-	0.642
8	+	0.896
9	+	0.812
10	-	0.423



Csoport	Rangsorolt
+	0.987
+	0.896
+	0.876
+	0.812
-	0.765
-	0.758
-	0.642
+	0.547
-	0.423
-	0.384

1. 

2. 



## M4 A többosztályos ROC görbék megalkotásának programkódja

```
#!/bin/bash
# ROC görbe generálás

sed -i 's;/ /g' $1
sed -i 's/.000//g' $1

FEJLEC=`head -1 $1`
tail -n +2 $1 > torzs$$tmp

AUC_OSSZEG=0
for i in `seq 1 $2`
do
    cp torzs$$tmp torzs$$tmp$i
    sed -i "s/ $i / 0 /g" torzs$$tmp$i

    UJFAJL=""`basename $1 .csv`_class"$i".csv
    echo $FEJLEC > $UJFAJL
    awk -v osztaly=$i '{print sqrt(($3-osztaly)^2),$1,$2}' torzs$$tmp$i | sort | awk '{print $2,$3,$1}'
    >> $UJFAJL

    roc_create.sh $UJFAJL
    cp $UJFAJL "`basename $UJFAJL .csv`".txt
    enrichment.sh "`basename $UJFAJL .csv`".txt

    EF_FAJL=""`basename EF_$UJFAJL .csv`.txt"
    AUC=`grep AUC $EF_FAJL | awk '{print $2}`"
    AUC_OSSZEG=`echo "scale=7; $AUC_OSSZEG + $AUC" | bc`
done

AUC_ATLAG=`echo "scale=7; $AUC_OSSZEG / $2" | bc`
echo $AUC_ATLAG

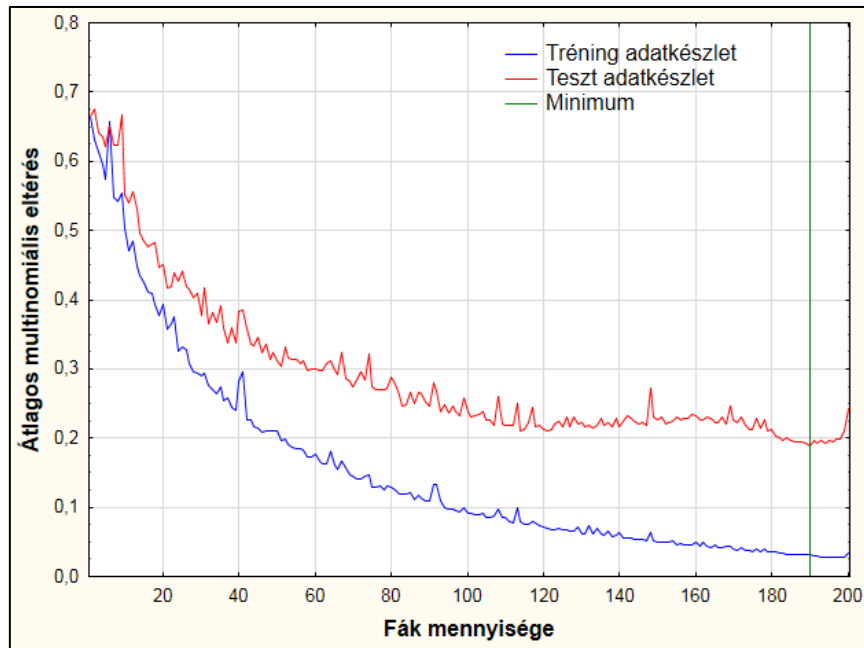
echo 0 0 0 > ROC_""`basename $1 .csv`_atlag.csv
for i in `seq 1 100`
do
    #y=`echo "scale=7; 100 * ( ( $i / 100 ) ^ ( ( 1 - $AUC_ATLAG ) / $AUC_ATLAG ) )" | bc -l`
    y=`echo "scale=7; 100 * e ( ( ( 1 - $AUC_ATLAG ) / $AUC_ATLAG ) * l ( $i / 100 ) )" | bc -l`
    echo $i $i $y >> ROC_""`basename $1 .csv`_atlag.csv
done

rm *$$tmp*

exit 0
```



**M5** Az energiatalok PCA adatkészlete alapján történő optimalás a fejlesztett fák módszerénél



## M6 Publikációs lista

	<b>Hatástényezővel (IF) rendelkező angol nyelvű folyóiratcikkek</b>
1.	<b>A. RÁCZ</b> , K. HÉBERGER, M. FODOR (2016): Quantitative determination and classification of energy drinks using near-infrared spectroscopy. <i>Analytical and Bioanalytical Chemistry</i> , Vol. 408. pp. 6403-6411. (IF=3.125)
2.	H. TIMA, <b>A. RÁCZ</b> , ZS. GULD, CS. MOHÁCSI-FARKAS, G. KISKÓ (2016): Deoxynivalenol, zearalenone and T-2 in grain based swine feed in Hungary. <i>Food Additives &amp; Contaminants: Part B</i> . Doi: <a href="http://dx.doi.org/10.1080/19393210.2016.1213318">http://dx.doi.org/10.1080/19393210.2016.1213318</a> (IF=1.467)
3.	F. ANDRIC, D. BAJUSZ, <b>A. RÁCZ</b> , S. SEGAN and K. HÉBERGER (2016): Multivariate assessment of lipophilicity scales – computational and reversed phase thin-layer chromatographic indices. <i>Journal of Pharmaceutical and Biomedical Analysis</i> , Vol. 127. pp. 81-93. (IF=2.867)
4.	<b>A. RÁCZ</b> , D. BAJUSZ, M. FODOR and K. HÉBERGER (2016): Comparison of classification methods with "n-class" receiver operating characteristic curves: A case study of energy drinks. <i>Chemometrics and Intelligent Laboratory Systems</i> , Vol. 151. pp. 34-43. (IF=2.321)
5.	D. BAJUSZ, <b>A. RÁCZ</b> and K. HÉBERGER (2015): Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations? <i>Journal of Cheminformatics</i> , Vol. 7:(1) paper 20. 13 p. (IF=4.547)
6.	<b>A. RÁCZ</b> , N. PAPP, E. BALOGH, M. FODOR and K. HÉBERGER (2015): Comparison of antioxidant capacity assays with chemometric methods. <i>Analytical Methods: Advancing Methods and Applications</i> , Vol. 7. pp. 4216-4224. (IF=1.821)
7.	<b>A. RÁCZ</b> , A. VASS, K. HÉBERGER and M. FODOR (2015): Quantitative determination of coenzyme Q10 from dietary supplements by FT-NIR spectroscopy and statistical analysis. <i>Analytical and Bioanalytical Chemistry</i> , Vol. 407:(10). pp. 2887-2898. (IF=3.436)
8.	<b>A. RÁCZ</b> , D. BAJUSZ, K. HÉBERGER (2015): Consistency of QSAR models: Correct split of training and test sets, ranking of models and performance parameters. <i>SAR and QSAR in Environmental Research</i> , Vol. 26:(7-9). pp. 683-700. (IF=1.596)
9.	O.H.J. CHRISTIE, <b>A. RÁCZ</b> , J. ELEK, K. HÉBERGER (2014): Classification and unscrambling a class-inside-class situation by object target rotation: Hungarian silver coins of the Árpád Dynasty, ad 997–1301. <i>Journal of Chemometrics</i> , Vol. 28:(4). pp. 287-292. (IF=1.803)
10.	<b>A. RÁCZ</b> , K. HÉBERGER, R. RAJKÓ, J. ELEK (2013): Classification of Hungarian medieval silver coins using X-ray fluorescent spectroscopy and multivariate data analysis. <i>Heritage Science</i> , Vol. 1:(2) paper 1/1/2. 9 p. (IF(Chemistry Central Journal)=1.663)

	<b>Magyar nyelvű publikációk</b>
11.	<b>RÁCZ, A.</b> (2015): Mire jó a kemometria? <i>Élet és Tudomány</i> , Vol. 70:(4). pp. 118-120.
12.	<b>RÁCZ, A., VASS, A., HÉBERGER, K., FODOR, M.</b> (2015): Q10 tartalmú étrendkiegészítők hatóanyagtartalmának mennyiségi meghatározása FT-NIR spektroszkópiával. <i>Élelmiszer - Tudomány Technológia</i> , Vol. 69:(2). pp. 25-32.

	<b>Angol nyelvű könyvfejezetek</b>
1.	<b>A. RÁCZ, D. BAJUSZ, K. HÉBERGER:</b> <i>Chapter title:</i> Cheminformatics/chemometrics in analytical chemistry In <b>Chemoinformatics – a textbook</b> (Editors: Johann Gasteiger, Thomas Engel) <b>WILEY</b> , 2017 (in press)
2.	<b>D. BAJUSZ, A. RÁCZ, K. HÉBERGER:</b> <i>Chapter 30010, title:</i> Chemical data formats, fingerprints and other molecular descriptions for database analysis and searching In <b>Comprehensive Medicinal Chemistry III.</b> (Editors: Andy Davis, Colin Edge) <b>ELSEVIER</b> , 2016 (in press)
3.	<b>D. BAJUSZ, A. RÁCZ, K. HÉBERGER:</b> <i>Chapter title:</i> Which performance parameters are best suited to assess the predictive ability of models? In <b>Advances in QSAR modeling with applications in Pharmaceutical, Chemical, Food, Agricultural and Environmental Sciences</b> (Editor: Kunal Roy) <b>SPRINGER</b> , 2017 (in press)

<b>Magyar nyelvű konferencia (összefoglaló)</b>	
	<b>RÁCZ, A., FILIP, A., BAJUSZ, D., HÉBERGER, K.</b> (2016): Számításos és vékonyréteg-kromatográfiás lipofilicitási indexek összehasonlítása kemometriai módszerekkel. KeMoMo–QSAR 2016 szimpózium, Szeged, 2016. május <a href="http://www.chemicro.hu/QSAR/kovetkezo.html">http://www.chemicro.hu/QSAR/kovetkezo.html</a>
	<b>RÁCZ, A., HÉBERGER, K., FODOR, M.</b> (2016): Energiaitalok minőségi és mennyiségi elemzése FT-NIR spektroszkópiával. Aktualitások a táplálkozástudományi kutatásokban – VI. PhD Konferencia, Budapest, 2016. február ( <a href="http://www.mttt.hu/index.php?content=57">http://www.mttt.hu/index.php?content=57</a> )
	<b>RÁCZ, A., FODOR, M., HÉBERGER, K.</b> (2015): Energiaitalok egy kemometrikus szemével. KeMoMo–QSAR 2016 szimpózium, Szeged, 2015. május <a href="http://www.chemicro.hu/QSAR/20150514.html">http://www.chemicro.hu/QSAR/20150514.html</a>
	<b>RÁCZ, A., VASS, A., HÉBERGER, K., FODOR, M.</b> (2015): Q10 tartalmú étrendkiegészítők minőségellenőrzése FT-NIR módszerrel. Aktualitások a táplálkozástudományi kutatásokban – V. PhD Konferencia, Budapest, 2015. január. p. 32 (ISBN 978-963-88108-8-5)

<p><b>RÁCZ, A., FODOR, M., HÉBERGER, K. (2014):</b> Változó kiválasztás – avagy legális út regressziós modellek javítására. KeMoMo–QSAR 2014 szimpózium, Szeged, 2014. május <a href="http://www.chemicro.hu/QSAR/kivonatok19/kivonat1902.html">http://www.chemicro.hu/QSAR/kivonatok19/kivonat1902.html</a></p>
<p><b>RÁCZ, A., PAPP, N., FODOR, M., HÉBERGER, K. (2014):</b> Antioxidáns kapacitás meghatározási technikák összehasonlítása kemometriai módszerek segítségével. Aktualitások a táplálkozástudományi kutatásokban – IV. PhD Konferencia, Budapest, 2014. január. p. 18 (ISBN 978-963-88108-7-8)</p>
<p><b>RÁCZ, A., ELEK, J., PAPP, G. (2013):</b> Égetett szeszesitalok tömegspektrometriás és infravörös spektroszkópiás vizsgálatának többváltozós elemzése KeMoMo–QSAR 2014 szimpózium, Szeged, 2013. április</p>
<p><b>RÁCZ A., ELEK J., NEMES Z. (2012):</b> Hogyan segíthet a modern műszeres analitika történelmi kérdések tisztázásában, avagy interdiszciplináris kutatások egy római kori sírkő körül. XXXV. KEN, Szeged, 2012. október. p. 246 (ISBN: 978-963-315-099-3)</p>
<p><b>RÁCZ, A., ELEK, J., NEMES, Z. (2013):</b> Feltárhat-e rejtett összefüggéseket egy római kori sírkőről röntgenfluoreszcenciás elemzés adatok főkomponens-elemzése? KeMoMo–QSAR 2014 szimpózium, Szeged, 2012. szeptember</p>
<p><b>RÁCZ, A., ELEK, J., HÉBERGER K., RAJKÓ, R., LENGYEL, A. (2011):</b> Régészeti leletek XRF vizsgálata és értékelése, avagy mennyi adatra van szükségünk főkomponens elemzéshez MKE 1. Nemzeti konferencia, Sopron, 2011. május (<a href="http://www.mkenk2011.mke.org.hu/hu/tudomanyos-program.html">http://www.mkenk2011.mke.org.hu/hu/tudomanyos-program.html</a>)</p>
<p><b>Nemzetközi konferencia (összefoglaló)</b></p>
<p><b>D. BAJUSZ, A. RÁCZ, K. HÉBERGER (2015):</b> Revival of an old debate: Cross- vs. External validation in QSAR modeling. Conferentia Chemometrica 2015, Budapest (Hungary), 2015, September. p. L23 (ISBN: 978-963-7067-31-0)</p>
<p><b>A. RÁCZ, D. BAJUSZ, K. HÉBERGER (2015):</b> Large scale statistical comparison of similarity metrics for fingerprint-based calculations. Conferentia Chemometrica 2015, Budapest (Hungary), 2015, September. p. P04 (ISBN: 978-963-7067-31-0)</p>
<p><b>A. RÁCZ, D. BAJUSZ, K. HÉBERGER (2015):</b> <i>n</i>-class ROC curves as novel, intuitive tools for method comparison. Conferentia Chemometrica 2015, Budapest (Hungary), 2015, September. p. P24 (ISBN: 978-963-7067-31-0)</p>
<p><b>K. HÉBERGER, D. BAJUSZ, A. RÁCZ (2015):</b> Consistency of QSAR models: correct split of training and test sets, ranking of models and performance parameters. 8<sup>th</sup> International symposium on computational methods in toxicology and pharmacology integrating internet resources, Chios (Greece), 2015, June. p. 26</p>
<p><b>J. ELEK, A. MEZŐSI, A. RÁCZ (2013):</b> Estimation of active ingredient content by multivariate calibration: NIR examination of tablets used in ED treatment. Conferentia Chemometrica 2013, Sopron (Hungary), 2013, September. p. L26 (ISBN: 978-963-9970-38-0)</p>
<p><b>A. RÁCZ, J. ELEK, G. PAPP (2013):</b> Multivariate data analysis of Hungarian spirit drinks' IR spectroscopic data. Conferentia Chemometrica 2013, Sopron (Hungary), 2013, September. p. P22 (ISBN: 978-963-9970-38-0)</p>

O.H.J. CHRISTIE, **A. RÁCZ**, J. ELEK, K. HÉBERGER (2012): Object target rotation for principal components analysis: metal content data of hungarian silver coins from the Árpád dynasty.  
XIII. Chemometrics in Analytical Chemistry, Budapest (Hungary), 2012, June. p. 72 (ISBN: 978-963-9970-24-3)

**A. RÁCZ**, J. ELEK, K. HÉBERGER, R. RAJKÓ, A. LENGYEL (2012): Principal component analysis of an XI-XV century silver coins' XRF dataset and verification by additional multilinear methods.  
WSC 8, Drakino, Russia, 2012, February.  
(<http://wsc.chemometrics.ru/wsc8/presentations/?page=2>)

**A. RÁCZ**, J. ELEK, K. HÉBERGER, R. RAJKÓ, A. LENGYEL (2011): Classification of medieval silver coins using X-ray fluorescent spectroscopy and multivariate data analysis.  
Conferentia Chemometrica 2011, Sümeg (Hungary), 2011, September. p. P33 (ISBN: 978-963-9970-15-1)

**A. RÁCZ**, J. ELEK, K. HÉBERGER, R. RAJKÓ, A. LENGYEL (2011): Principal Component Analysis of an XI-XV century silver coins' XRF dataset and verification by additional multilinear methods.  
Conferentia Chemometrica 2011, Sümeg (Hungary), 2011, September. p. L26 (ISBN: 978-963-9970-15-1)