

TARTU ÜLIKOOL
LOODUS- JA TÄPPISTEADUSTE VALDKOND
MATEMAATIKA JA STATISTIKA INSTITUUT

Johann Saavaste
**Isheemiliste südamehaiguste markerite leidmine
Tartu Ülikooli Eesti Geenivaramu
metaboloomika andmetele geneetilist algoritmi
rakendades**

Matemaatiline Statistika
Bakalaureusetöö (9 EAP)

Juhendajad:
PhD Jaanika Kronberg
PhD Krista Fischer

TARTU 2023

**ISHEEMILISTE SÜDAMEHAIGUSTE MARKERITE LEIDMINE
TARTU ÜLIKOOLI EESTI GEENIVARAMU METABOLOOMIKA
ANDMETELE GENEETILIST ALGORITMI RAKENDADES**

Bakalaureusetöö

Johann Saavaste

Lühikokkuvõte

Metaboloomika teadusharu uurib organismi ainevahetuse vahe- ja lõppsaadustena tekkinud väikesi molekule ehk metaboliite. Metaboliidi kontsentratsiooni kõrvalekalded võivad olla indikaatoriks haiguse tekkele. Igal aastal sureb maailmas enim inimesi isheemiliste südamehaiguste tagajärjel. Käesoleva bakalaureusetöö eesmärk on leida metaboliite, mis sobiksid isheemilisi südame-tõbesid kirjeldavateks markeriteks. Töös leitakse geneetilise algoritmiga kõige suurema seosega metaboliidid. Hinnatakse 2 logistilist regressioonimudelit, mille tunnustena võeti arvesse väljavalitud metaboliidid, sugu, vanus, kehamassiindeks ning elupaik.

CERCS teaduseriala: P160 Statistika, operatsioonianalüüs, programmeerimine, finants- ja kindlustusmatemaatika; B110 bioinformaatika, meditsiiniinformaatika, biomatemaatika, biomeetrika.

Märksõnad: Metaboloomika, geneetiline algoritm, lähim tsentroid, logistiline regressioon.

**FINDING MARKERS FOR ISCHEMIC HEART DISEASES USING
GENETIC ALGORITHM ON ESTONIAN GENOME CENTER'S
DATA**

Bachelor thesis

Johann Saavaste

Abstract

Metabolomics studies small molecules, also known as metabolites, which are intermediates or products of metabolism. Deviations of metabolite concentration could indicate development of a disease. Ischemic heart disease is the highest cause of death in the world each year. The aim of this bachelor's thesis is to find metabolites that would fit as descriptive markers of ischemic heart diseases. The metabolites with the largest association are found in the thesis using genetic algorithm. 2 logistic regression models are evaluated using chosen metabolites, sex, age, body mass index and place of residence as predictors.

CERCS research specialisation: P160 Statistics, operations research, programming, financial and actuarial mathematics; B110 bioinformatics, medical informatics, biomathematics, biometrics.

Key Words: Metabolomics, genetic algorithm, nearest centroid, logistic regression.

Sisukord

Sissejuhatus	5
1 Metaboliidid	6
1.1 Südamete isheemiatõved	7
2 Geneetiline algoritm	8
2.1 Rakendustarkvara R pakett GALGO	9
2.2 Klassifitseerimismeetodid	10
2.2.1 Lähim tsentroid	10
2.2.2 Tugivektor-masin	11
2.2.3 Juhuslik mets	12
2.3 Varasemad uuringud	12
3 Logistiline regressioon	13
4 Andmestik	14
5 Klassifikatsioonimeetodite võrdlus	16
6 Tulemused	18
6.1 Haigusjuhtude võrdlus kõiki kontrollgrupi isikuid sisaldava andmestikuga	19
6.1.1 Jälgimisaegseid juhte prognoosivad metaboliidid	19
6.1.2 Jälgimiseelseid juhte prognoosivad metaboliidid	21
6.2 Haigusjuhtude võrdlus haigusgrupiga sarnaseid kontrolle sisaldava andmestikuga	22

6.2.1	Jälgimisaegseid juhte prognoosivad metaboliidid	24
6.2.2	Jälgimiseelseid juhte prognoosivad metaboliidid	25
6.3	Logistiline regressioonimudel	27
7	Arutelu	30
	Kokkuvõte	32
	Kasutatud allikad	33
	Lisa 1. GALGO protsessiga leitud kõige sagedamad metaboliidid	37

Sissejuhatus

Metaboloomika on teadusharu, mis keskendub metaboliitide uurimisele organismis. Metaboliidid tekivad metabolismi ehk organismi ainevahetuse vahe- või lõppsaadusena. Haiguste ennetamiseks ja raviks uuritakse organismi metaboliitide kontsentratsioone tuvastamaks haiguse tekkeprotsessi juures olulisi metaboolseid muutuseid. Kui fikseerida metaboliitide kontsentratsioonide normid, on võimalik normilt kõrvalekaldumiste korral tuvastada muutusi organismis ning varakult haigusele jaole jõuda.

Kõige enam inimesi sureb maailmas südame isheemiatõbede tagajärjel. Maailma Terviseorganisatsiooni andmetel suri 2019. aastal haigusesse umbes 9 miljonit inimest, mis moodustas 16% kõikidest samal aastal toimunud surmajuhtumitest (Maailma Terviseorganisatsioon, 2020). Südame isheemiatõvesse haigestutakse, kui südamearterite kitsenemise tagajärjel jõuab südamelihastesse vähem verd ja hapnikku.

Käesoleva bakalaureusetöö eesmärk on leida Tartu Ülikooli Eesti Geenivaramu andmete põhjal leida metaboliidid, mis potentsiaalselt sobiksid isheemilisi südamehaigusi kirjeldavateks markeriteks.

Töö käigus rakendatakse metaboliitide leidmiseks geneetilist algoritmi 3 erineva klassifikatsioonimeetodiga. Haiguse hindamiseks logistilise regressioonimudeliga võetakse lisaks leitud metaboliitidele arvesse geenidoonori sugu, vanust, kehamassiindeksit ning elukohta.

Töö kirjutati tekstitöötlusprogrammiga $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ ning statistiline analüüs viidi läbi rakendustarkvaras R.

1 Metaboliidid

Metaboliidid on väikese molekulmassiga molekulid, mis saadakse organismi ainevahetuse vahe- või lõppsaadusena. Metabooloomikas tuvastatakse süstemaatiliselt organismi metaboliite ning mõõdetakse nende kontsentratsiooni. Kaks peamist tehnoloogiat metaboliitide tuvastamiseks on tuumamagnetresonantspektroskoopia (inglise keeles *nuclear magnetic resonance spectroscopy*, NMR) ning mass-spektromeetria (MS) (Idle ja Gonzalez, 2007).

Mõlemal tehnoloogial on omad eelised ja puudused. MS võimaldab meetodist sõltuvalt tuvastada üle 1000 metaboliidi, seevastu NMR-ga tuvastatavate metaboliitide arv jääb alla 100. MS on rahaliselt kulukam, kuna preparaadi ettevalmistus on keerukam ning analüüsiks kulub rohkem aega (Euroopa Bioinformaatika instituut, 2023).

Lipiidide alla kuuluvad erineva struktuuriga hüdrofoobsed biomolekulid nagu triglütseriidid ehk rasvad, fosfolipiidid ning steroidid. Triglütseriidid koosnevad glütseroolist ning 3 rasvhapest. Neis talletatakse organismi energiavaru ning lagundatakse rasvlahustuvaid vitamiine. Fosfolipiididel on üks rasvhappega side asendatud fosfaatgrupiga. Fosfaatgrupi hüdrofiilsuse tõttu leidub fosfolipiide rakumembraani koostises, kaitstes raku sisu teatud ehitusega molekulide eest. Steroidide hulgas on tähtsal kohal kolesterool, millest sünteesitakse teisi steroidhormoone näiteks östrogeen, testosteroon ja kortisool (Ahmed, Shah ja Ahmed, 2023).

Ksenobiootikumideks nimetatakse kemikaale, mida inimorganism ei sünteesi ja mis satuvad organismi väliskeskkonnast. Need jõuavad organismi üldjuhul toiduga. Peamiselt on tegu taimset päritolu metaboliitidega. Osa ksenobiootikume on tööstuslikul viisil sünteesitud ja reoveega levinud looduskeskkonda. Ksenobiootikumide hulka kuuluvad näiteks hügieenitarbed, ravimid, pestitsiidid ning tööstusjäätmed (Štefanac, Grgas ja Landeka Dragičević, 2021).

1.1 Südame isheemiatõved

RHK-10 süsteemi kuuluvaid I20-I25 haiguseid nimetatakse südame isheemiatõbedeks (edaspidi ISH). Sinna kuuluvad stenokardia ehk rinnaangiin, müokardiinfarkti erinevad juhud ning kroonilised südame isheemiatõved nagu ateroskleroosiline südamehaigus (Maailma Terviseorganisatsioon, 2016). ISH on põhjustatud südamearterite kitsenemisest, mille tõttu südamelihastesse jõuab vähem verd ning hapnikku. Rinnas tekkiv valu võib viia südameinfarktini (Ameerika Südameassotsiatsioon, 2022).

Südamearterid kitsenevad, kui arteri siseseintele ladestuvad rasvad ning kolesteroolijäägid, mis aja möödudes kõvastuvad. Varasem uuring on näidanud, et kõrgem kolesteroolijääkide kontsentratsioon suurendab ISH riski. Uuringus analüüsiti 74 000 Kopenhaageni elaniku kolesterooli taset veres, kellest 11 000 oli ISH diagnoositud ajavahemikus 1976–2010 (Varbo *et al.*, 2013).

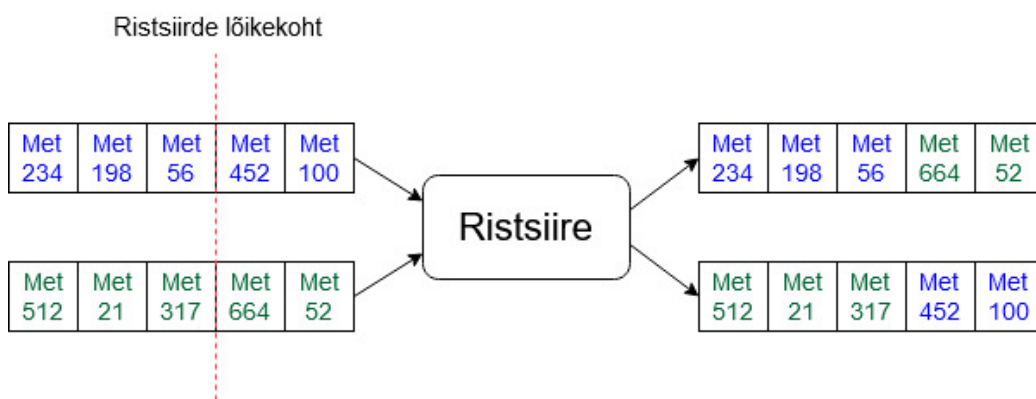
Oluliseks riskifaktoriks ISH esinemisele on leitud ka diabeedi olemasolu. Kopenhaageni linnaosa vähemalt 20-aastastest elanikest moodustati 13 000 pealine valim ning jagati teist tüüpi diabeedi järgi haigus-, soodumustega grupiks ning kontrollgrupiks. 20 aastat hiljem uuriti ISH esinemist valimis sh ka surmaga lõppenud juhtumid. Leiti, et teist tüüpi diabeetikutel on 2 kuni 3 korda suurem risk ISH-sse haigestuda. Diabeetikutest naiste suhteline haigestumise risk oli diabeedihaigete meeste omast oluliselt kõrgem (Almdal *et al.*, 2004).

2 Geneetiline algoritm

Geneetiline algoritm (edaspidi GA) on optimeerimisalgoritm, mis tugineb Charles Darwini kirjeldatud looduslikule valikule. Sellega leitakse suure populatsiooniga kogumitest kõige parema sobivusega kombinatsioone ilma kõikide variantide ajamahuka läbivaatamiseta. Algoritmis kasutatavad mõisted on üle võetud bioloogiast ja geneetikast ning järgnevates lõikudes ei tasu neid seostada nimetatud valdkondadega (Mitchell, 1998).

Järgnevalt kirjeldatakse põlvkonna protsessi.

Kõigepealt luuakse juhusliku protsessi tulemusel metaboliitide kogumite ehk kromosoomide populatsioon. Defineeritakse sobivusfunktsioon ning arvutatakse iga kromosoomi sobivus. Seejärel valitakse kaks kõige suurema sobivusega kromosoomi. Valitud kromosoomidel teostatakse ristisiire – mõlemal jagatakse metaboliitide järjestused i -nda liikme järel kaheks osaks ning kromosoomid vahetavad omavahel ühe osa. Saadud järglastel vahetatakse osa metaboliite mõne teise vastu ehk toimub mutatsioon. Uue populatsiooniga arvutatakse taas iga kromosoomi sobivus (Katobh, Chauhan ja Kumar, 2021).



Joonis 1: Ristsiirde skeem

2.1 Rakendustarkvara R pakett GALGO

Alapeatükk on kirjutatud (Trevino ja Falciani, 2006) ning (Trevino ja Falciani, 2018) põhjal.

Programmeerimiskeeles R on GA rakendamiseks loodud pakett GALGO, millega on võimalik luua statistilisi mudeleid suuremõõtmeliste andmetele. GALGO-s luuakse *BigBang* objekt, milles rakendatakse GA protsessi mitme kromosoomi populatsiooniga ehk nišiga. Nišid võivad omavahel juhuslikult kromosoomi vahetada. Kromosoomi sobivuse hindamiseks luuakse treeningandmestik, mille põhjal leitakse prognoos testandmestikule sõltuvalt klassifitseerimismeetodist. Prognoosi täpsuse põhjal arvutatakse kromosoomi sobivus. Kui leitakse kromosoom, mille sobivus ületab eesmärgile vastavat sobivust, siis see valitakse välja ning lõpetatakse uute põlvkondade loomine. Väljavalitud soovitud sobivusega kromosoomi nimetatakse lahendiks. Perioodi esimese põlvkonna ning lahendi leidmise vahel nimetatakse evolutsiooniks. Vaikimisi sisaldab evolutsioon 200 põlvkonda. Kui evolutsiooni kestel lahendini ei jõuta, siis valitakse välja viimase põlvkonna parim kromosoom.

Kui programm on jookstanud läbi sadu evolutsioone, siis viiakse läbi ettesuunatud valik (inglise keeles *forward selection*). Selleks järjestatakse metaboliidid kahanevas järjekorras vastavalt sellele, mitmes kromosoomis see evolutsiooni lõpus sisaldus. Seejärel, alustades esialgu ainult järjekorras esimesest metaboliidist, luuakse mudel prognoosimaks haigusgruppi. Leitakse mudeli klassifitseerimise täpsus ning mudelisse lisatakse järjekorras järgmine metaboliit. Protsessi korratakse senikaua, kuni mudelisse on lisatud kõik metaboliidid. Parimaks mudeliks valitakse mudel, milles haigusgruppi klassifitseerimise täpsus on maksimaalne ning metaboliitide arv minimaalne.

2.2 Klassifitseerimismeetodid

Töös rakendatakse andmestikule geneetilist algoritmi programmeerismiskeeles R paketi GALGO abil 3 erineva klassifitseerimismeetodiga ja võrreldakse saadud tulemusi ning nende efektiivsust. Olgu treeningandmestik X , mis koosneb n paarist $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$, kus \vec{x}_i tähistab ühele vaatlusele vastavat m -mõõtmelist vektorit ning y_i vaatluse klassi. Olgu $y_i \in Y$ võimalike klasside hulk.

2.2.1 Lähim tsentroid

Alapeatükk on kirjutatud (Levner, 2005) põhjal.

Lähima tsentroidi (inglise keeles *nearest centroid*) meetodil arvutatakse treeningandmetes igale klassile selle vektorite aritmeetiline keskmine ehk tsentroid. Testandmete abil leitakse tundmatu klassiga vektori kaugused kõikidele klassidele vastavatest tsentroididest. Prognoositud klassiks valitakse vähima kaugusega tsentroidile vastav klass.

Klassile k vastav tsentroid leitakse valemiga

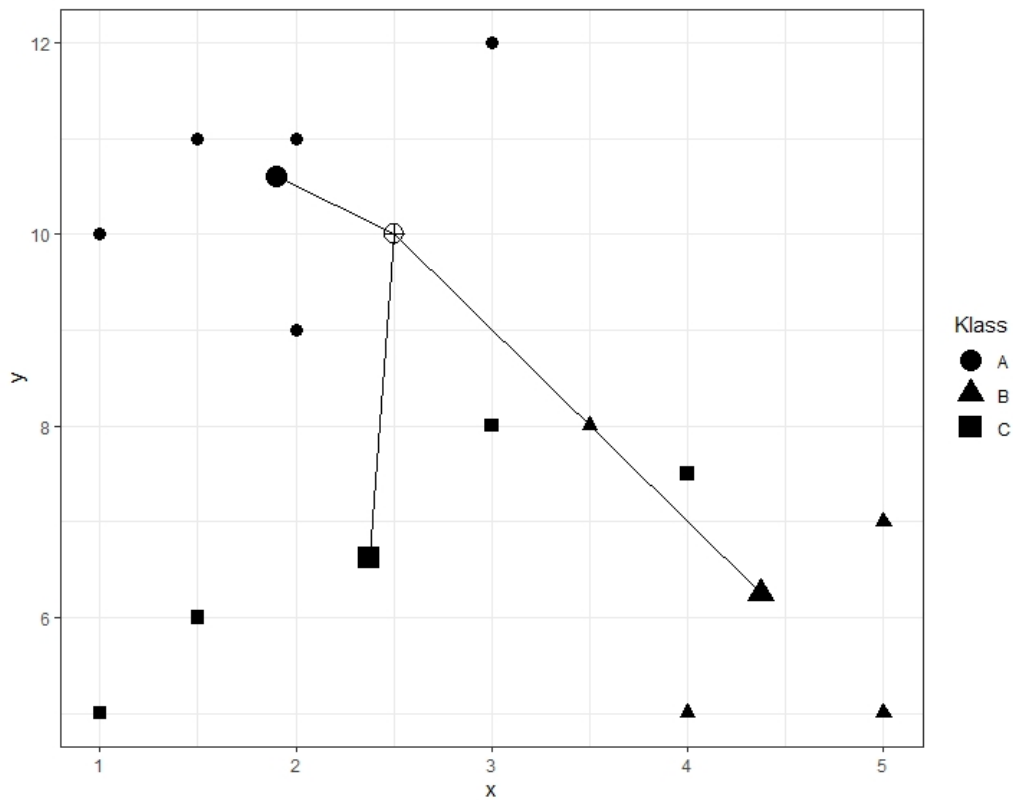
$$\vec{\mu}_k = \frac{1}{|D_k|} \sum_{y_i \in D_k} \vec{x}_i,$$

kus D_k on klassile k vastavate andmete hulk.

Testandmetes vaatlusega \vec{x}_j ning tundmatu klassiga y_j määratakse klass valemiga

$$\hat{y}_j = \arg \min_{\ell \in Y} \|\vec{\mu}_\ell - \vec{x}_j\| = \arg \min_{\ell \in Y} \sqrt{\sum_{i=1}^m (\mu_{\ell i} - x_{ji})^2}.$$

Joonisel 2 on kahemõõtmelised vaatlused kujutatud väiksemate kujunditena. Silmnähtavalt suurte kujunditega on kujutatud vaatluste klassidele vastavate tsentroidide asukohad. Kui hinnata tundmatu vaatluse klassi, mida joonisel on kujutatud ringjoonega ümbritsetud ristina, siis sellele kõige lähemal asub klassi A tsentroid.



Joonis 2: Lähim tsentroid kahemõõtmelises ruumis. Ringjoonega ümbritsetud ristiga tähistatud tundmatu klassiga vaatluse klassifitseerimisel valitakse klass A.

2.2.2 Tugivektor-masin

Tugivektor-masina (inglise keeles *support vector machine*, edaspidi TVM) meetodil paigutatakse treeningandmed m -mõõtmelisse ruumi punktidenä. Sobiva tuumafunktsiooni (inglise keeles *kernel function*) korral asetsevad punktid nii, et klasse saab ruumis eraldada hüperpinnaga. Hüperpind valitakse selline, et mõlema klassi punktid asetseksid sellest maksimaalsel kaugusel. Testandmed paigutatakse valitud hüperpinnaga samasse ruumi ning klassi prognoosimisel vaadeldakse punkti asukohta hüperpinna suhtes (Cervantes *et al.*, 2020).

2.2.3 Juhuslik mets

Juhusliku metsa (inglise keeles *random forest*) meetodi korral luuakse treeningandmetest mitu juhuslikku tagasipanekuga valimit. Iga valimi põhjal rajatakse otsustuspuu, mis muutujate omadusi arvestades määrab sellele vastava klassi. Testandmetega klassi prognoosimisel valitakse maksimaalne klassi määratud otsustuspuude arv. Juhusliku metsa eeliseks on tema madal tundlikkus üleprognoosimisele (Belgiu ja Drăguț, 2016).

2.3 Varasemad uuringud

Metaboloomika andmetest püütakse leida biomarkereid erinevatele haigustele. Näiteks on leitud GALGO abil biomarkerid teist tüüpi diabeedile. Uuringus analüüsiti 41–63-aastaste inimeste plasmaproovide mass-spektromeetria metaboloomika andmeid. 80 inimest jaotati võrdselt 4 gruppi, millest esimesse kuulusid diabeedi eelsoodumusega isikud, teise teist tüüpi diagnoosiga isikud, kolmandasse diabeetilise nefropaasi haigusega isikud ning neljandaks oli kontrollgrupp. Metaboliitidest biomarkereid prognoosivate mudelite koostamiseks kasutati järgnevaid klassifitseerimismeetodeid: k -lähimad naabrid, lähim tsentroid ning TVM meetod. Ettesuunatud valiku ning ristvalideerimise käigus leiti teist tüüpi diabeedi biomarkeriteks 5 metaboliiti (Morgan-Benita *et al.*, 2022). Lisaks on GALGO-t kasutatud rottide maksahaigustega seotud oluliste metaboliitide leidmiseks (Lin *et al.*, 2011).

GA on varasemalt kasutatud mitmetes protsessides, kus traditsioonilise algoritmiga jäädi hätta. Näiteks hinnati GA abil matemaatiline mudel glükoosi metabolismile, mille täpsus oli oluliselt parem traditsioonilisel meetodil saadud mudelist (Morbiducci, Andrea Tura ja Grigioni, 2005).

3 Logistiline regressioon

Peatükk on kirjutatud (Nick ja Campbell, 2007) põhjal.

Olgu uuritav tunnus Y binaarne tunnus, mille väärtused on kodeeritud järgmiselt: 1 – uuritav sündmus toimus, 0 – uuritav sündmus ei toimunud. Logistiline regressiooni meetodil mudeldatakse sündmuse toimumise šansi logaritmi. Mudeli üldkuju on järgmine:

$$\ln \left(\frac{P(Y = 1 | X_1, \dots, X_k)}{1 - P(Y = 1 | X_1, \dots, X_k)} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k,$$

kus X_1, \dots, X_k tähistavad mudeli argumenttunnuseid, β_0 mudeli vabaliiget ning β_1, \dots, β_k argumenttunnuste kordajaid.

Mudeli valemist on võimalik avaldada valem sündmuse toimumise tõenäosuse arvutamiseks.

$$P(Y = 1 | X_1, \dots, X_k) = \frac{\exp(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}{1 + \exp(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}$$

4 Andmestik

Töös kasutatakse Eesti Geenivaramu aastatel 2002-2019 kogutud plasmaproovide mass-spektromeetria metabooloomika andmeid. Andmete kasutamine oli Eesti bioetika ja inimuuringute nõukogu loa 234T-12 “*Omic*s for Health” raames. Andmed väljastati Eesti Geenivaramust andmeväljastuse nr R26 alusel. Iga rida vastas ühele mõõdetud metaboliidi kontsentratsioonile. Veergudeks olid 999 doonorit. Metaboliidid on eristatud järjekorranumbri ning HMDB koodi kujul, mille kohta on võimalik teha päring Inimese Metaboliitide Andmebaasis (inglise keeles *the Human Metabolome Database*).

Töös vaadeldakse haigusgrupina doonoreid, kellel on pärast proovi andmist RHK-10 süsteemis diagnoositud vähemalt üks isheemilistest südamehaigustest. Töö andmebaasidega teostas juhendaja Jaanika Kronberg, kes andis autorile geenidoonorite klassifikatsiooni kujul 0/1/2, kus 0 on kontrollgrupp, 1 on isikud, kellel haigus tekkis pärast proovi andmist ja 2 on isikud, kellel oli haigus proovi andmise ajal diagnoositud. Doonoritest oli 69-l haigus diagnoositud peale proovi andmist, 138 doonoril enne proovi andmist ning ülejäänud 762 olid kontrollgrupis.

Esialguses andmestikus oli kokku mõõtmisi 1505 metaboliidi kohta. Moodustati 2 andmestikku, millest esimeses jäeti välja isikud, kel oli proovi andmisel haigus diagnoositud, ning teises isikud, kes haigestusid peale proovi andmist. Enam kui 1000 metaboliidil puudus vähemalt ühe isiku kohta mõõtmistulemus. Kuna TVM ja juhusliku metsa klassifitseerimismeetodid eeldavad, et andmed ei sisalda puuduvaid väärtusi, siis jäeti andmestikest välja metaboliidid, milles vähemalt 60% vaatlustest puudusid väärtused. Ülejäänud tühjad väljad asendati poolega vastava metaboliidi mõõtmistulemuste miinimumist. Andmestikus oli metaboliite, millel puudus HMDB kood, mistõttu eemaldati need samuti. Peale andmestiku puhastamist jäid alles 765 metaboliidi mõõtmistulemused.

Metaboliitidele lisaks kasutati andmeid doonorite tausta kohta nagu sugu ja vanus proovi andmise hetkel. Nende tunnuste põhjal leiti andmestikust igale isikule

haigusrühmast 2 temale vastavat kontrollgrupi liiget, millega saadud tulemusi töös eraldi vaadeldakse. Lisaks sisaldas andmestik teavet doonorite elupaiga (linn, väikelinn, maa) ja kehamassiindeksi kohta.

5 Klassifikatsioonimeetodite võrdlus

Haigusgrupina vaadeldi ainult isikuid, kellel diagnoositi ISH peale proovi andmist. Haiguse biomarkerite leidmiseks rakendati GA haigus- ja kontrollgrupi metabooloomika andmetele lähima tsentroidi, TVM ning juhusliku metsa klassifikatsioonimeetodiga. Protsess viidi läbi, kasutades rakendustarkvara R paketti GALGO (Trevino ja Falciani, 2018). Iga klassifikatsioonimeetodi peal prooviti rakendada samu parameetreid. Enamik parameetritest olid vaikeväärtustega, kuid sobivuseks määrati 80%, evolutsioonide arvuks 1000 ning kromosoomi suuruseks 10 metabooliiti. Tabelis 1 on välja toodud klassifikatsioonimeetodite peamised erinevused töös kasutatavate andmete peal rakendamisel.

Tabel 1: Klassifikatsioonimeetodite võrdlus töös kasutatavate andmete põhjal

	Lähim tsentroid	TVM	Juhuslik mets
Implementatsioon	Programmeerimiskeeles C	Autori modifitseeritud R-i kood	Autori modifitseeritud R-i kood
Põlvkonna kestus	0,2 s	1 s	20–30 s
Gruppide kaalud	kaalud haigus- ja kontrollgrupi osakaalude põhjal	haigusgrupil rohkem kaalu	haigusgrupil rohkem kaalu
Parim sobivus evolutsiooni lõpus (kuni 200 põlvkonda)	80%	20–40% (sõltub tuumameetodist)	10–25%

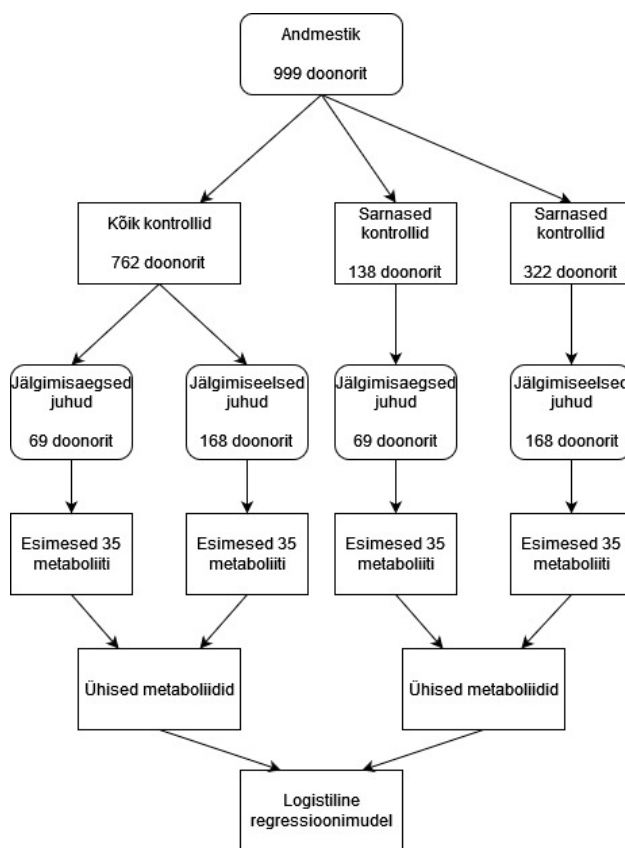
Kuna haigusgrupi osakaal andmestikus oli madal, siis igale klassifikatsioonimeetodile prooviti treening- ja testandmestiku loomisel arvestada ka haigusgrupi kaaluga. Vastasel juhul oleks prognoositud õigesti valdavalt kontrollgruppi. GALGO-s on lähima tsentroidi meetod kiiruse huvides implementeeritud programmeerimiskeeles

C, mistõttu autoril ei õnnestunud kirjutada koodi R-s, mis arvestaks haigus- ja kontrollgrupi kaale (Trevino ja Falciani, 2018). TVM ja juhusliku metsa meetodil õnnestus muuta koodi selliselt, et treening- ja testandmestikel oleks piisavalt haigusgrupi isikuid.

Paraku osutusid TVM ja juhusliku metsa meetodi rakendamine väga ajamahukaks. Suur ajakulu oli tingitud kromosoomide madalast sobivusest ning põlvkonna kestusest. Põlvkonna kestuses mängib rolli andmete analüüsiks kasutatav riistvara, milleks antud juhul oli autori isiklik sülearvuti. Kui üks põlvkond kestab TVM meetodil keskmiselt 1 sekundi, siis järgmise evolutsioonini jõutakse maksimaalselt 200 sekundiga. 1000 evolutsiooni korral oleks protsessi maksimaalne ajakulu olnud üle 50 tunni. Juhusliku metsa meetodil kestaks protsess kauemini, kuna põlvkonna ajakulu oli 20–30 sekundit. Protsessi kiirendamiseks prooviti alandada sobivust 50% peale ning vähendada evolutsioonis põlvkondade arvu 100 peale. Kuigi protsessi ajakulu vähenes, ei suurenud sellegipoolest lahendite arv. Lahendite saamiseks prooviti muuta tuumafunktsiooni parameetreid, kuid sellest hoolimata ei õnnestunud leida oluliselt rohkem lahendeid.

6 Tulemused

Töös vaadeldakse haigusgruppina doonoreid, kel oli haigus diagnoositud peale proovi andmist ehk jälgimisaegseid juhte, ja doonoreid, kellel oli haigus diagnoositud proovi andmise ajaks ehk jälgimiseelseid juhte. Mõlemad haigusgrupid sisalduvad eraldi andmestikes. Analüüs koosneb kahest osast. Esimeses osas analüüsitakse 2 andmestikku, mis sisaldab kõiki kontrollid. Teises osas valitakse haigusgrupi isikutele sarnased kontrollid ning moodustatud 2 andmestikku analüüsitakse sarnaselt esimese osaga. Analüüsitavad metaboliidid olid kõikides andmestikes ühesugused. Analüüs teostati rakendustarkvaras R versiooniga 4.1.1, kasutades paketti GALGO versiooniga 1.4 (Trevino ja Falciani, 2018). Joonisel 3 on visualiseeritud töö käik.



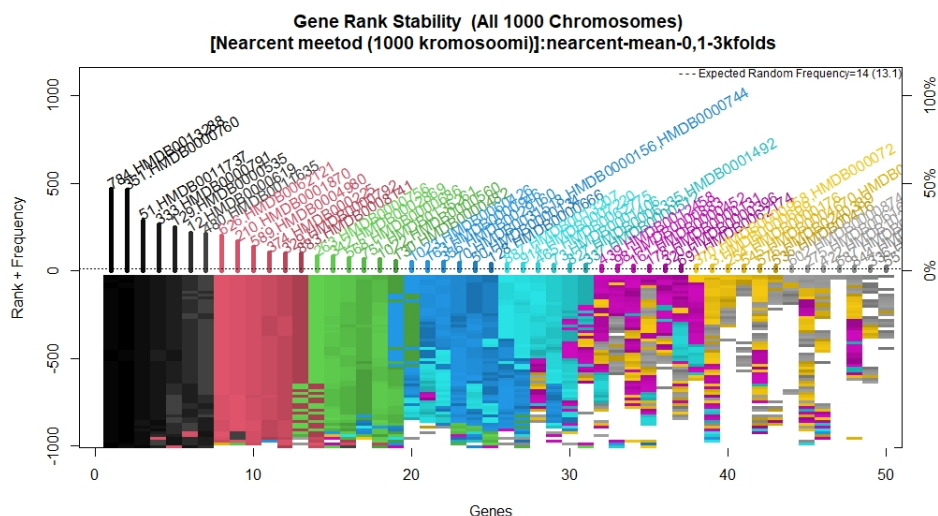
Joonis 3: Lihtsustatud analüüsi käik

Iga andmestiku tulemustega saab tutvuda [lisa 1](#) tabelitest.

6.1 Haigusjuhtude võrdlus kõiki kontrollgrupi isikuid sisaldava andmestikuga

Mõlema andmestiku korral rakendati GALGO *BigBang* objektis ühesuguseid parameetreid, kus kromosoomi suuruseks määrati 10 metaboliiti, evolutsiooni pikkuseks 200 põlvkonda ning kokku viidi läbi 1000 evolutsiooni (Trevino ja Falciani, 2018). Ainukese parameetrina erines soovitud sobivus, kuna jälgimiselsetel juhtudel ei õnnestunud enamiku evolutsioonide käigus saavutada sobivuseks 80%, mistõttu alandati seda 70% peale. Sellegipoolest jäeti jälgimisaegsete juhtude korral soovitud sobivus 80% peale.

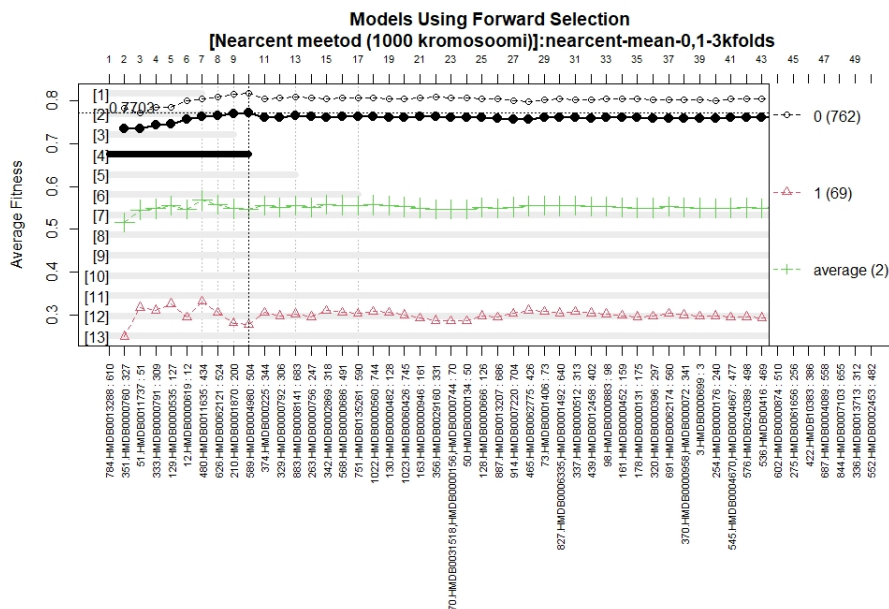
6.1.1 Jälgimisaegseid juhte prognoosivad metaboliidid



Joonis 4: Metaboliitide üldine astakute stabiilsus ning sagedus jälgimisaegseid juhte ja kõiki kontrole sisaldava andmestiku põhjal

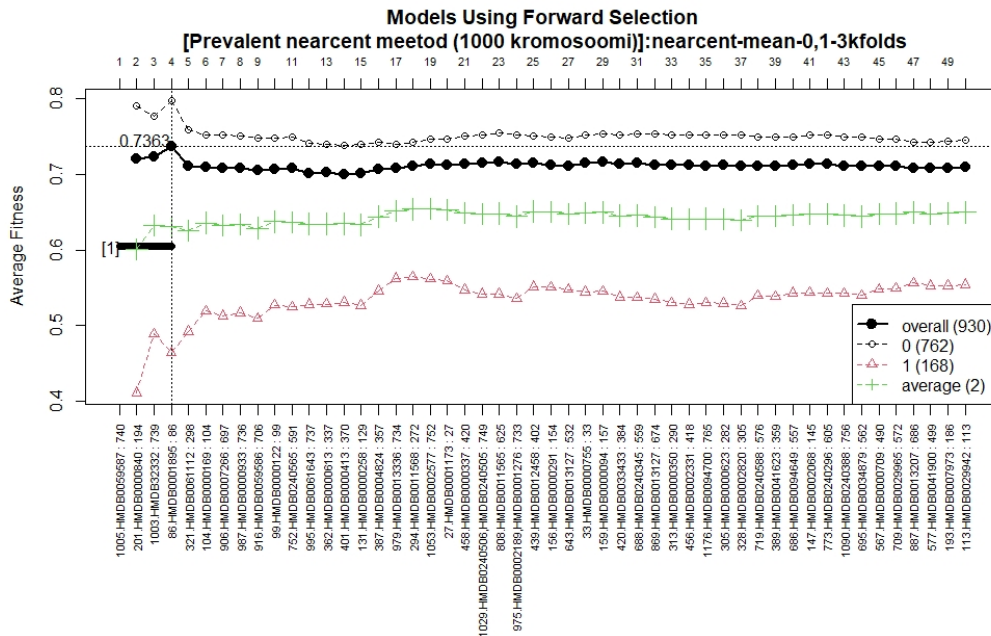
Tulemused andmestikus, kus haigusgruppi kuulusid jälgimisaegsete juhtudega isikud, on esitatud joonisel 4. Geenide nimetusena on kujutatud 50 kõige sagedamat metaboliiti kahanevas järjekorras, mis sisaldasid lahenditeks valitud kromosoomides. Vertikaalse telje ülemine osa kujutab metaboliitide sagedust. Värvidega on näidatud metaboliidi asukohta järjekorras eelneva evolutsiooni lõpus. Värvid kõige

alumises reas näitavad metaboliitide asukohta järjestuses peale esimest evolutsiooni. Mida ühevärvilisemalt on metaboliidid kujutatud, seda stabiilsemana on nende järjestus evolutsioonide käigus püsinud. (Trevino ja Falciani, 2018) On näha, et esimesel 10 metaboliidil on püsinud järjestus ligikaudu viimased 900 evolutsiooni muutumatusena. Kõige rohkem leidus lahendites 464 korda metaboliit HMDB0013288 ehk nonaüülkarnitiin, millele järgnes 462 korda esinenud HMDB0000760 ehk hüokoolhape (inglise keeles *Hyochoolic acid*).



Joonis 5: Jälgimisaegsete juhtude metaboliitide järjestuse ettesuunatud valik

Joonisel 5 on abstsistteljel taas metaboliidid järjestatud sageduse järgi. Punaste kolmnurkadena on kujutatud jälgimisaegsete grupi ning musta ringjoonena kontrollgrupi klassifitseerimise täpsust. Mustade ringidena on kujutatud mudeli üldist täpsust ehk mõlema grupi õigesti klassifitseeritud vaatluste osakaalu kõikidest vaatlustest. Kõige parema täpsusega mudeli saaks GALGO hinnangul, kui mudelisse lisada kõik metaboliidid alates vasakult kuni katkendliku vertikaalse jooneni, kuna mustade ringidena tähistatud üldine täpsus on sel juhul kõige kõrgem (Trevino ja



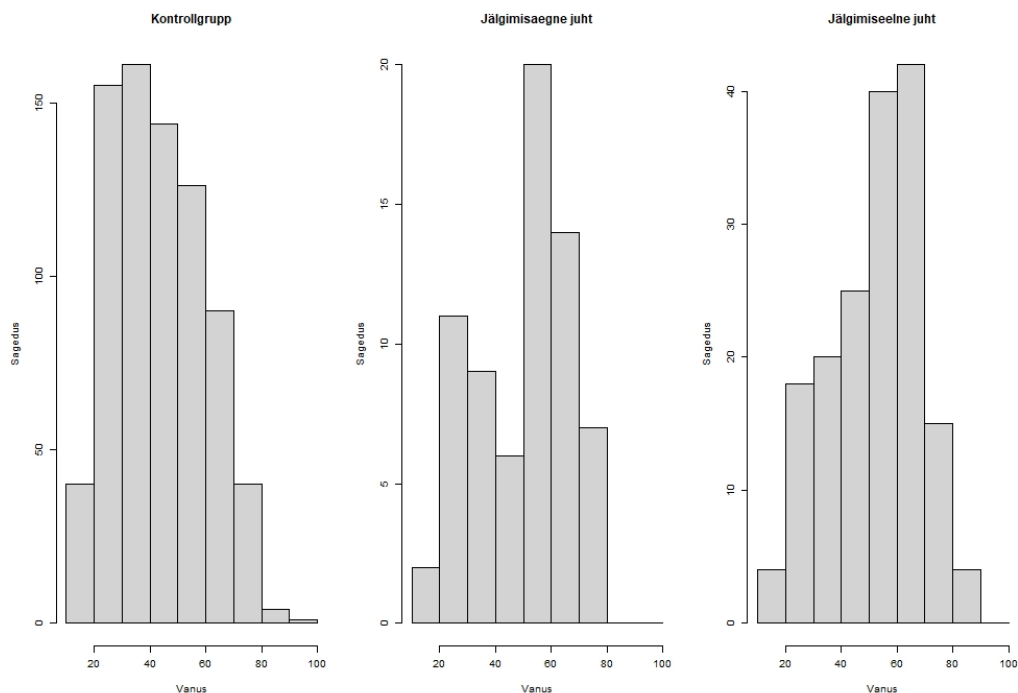
Joonis 7: Jälgimiseelsete juhtude metaboliitide järjestuse ettesuunatud valik

74%. Sellegipoolest on haigusgrupi prognoosimise täpsus selles mudelis madalam kui suuremal osal teistest mudelitest. Kõige parema täpsusega haigusgruppi prognoosiva mudeli saaks, kui võtta mudelisse esimesed 18 metaboliiti. Sellise mudeli haigusgrupi prognoosimise täpsus oleks joonise põhjal 56% ning üldine täpsus 71%, mis ei erine palju GALGO valitud mudeli täpsusest.

6.2 Haigusjuhtude võrdlus haigusgrupiga sarnaseid kontrole sisaldava andmestikuga

Kontrollgrupiga võrreldes on haigusgrupp oluliselt väiksema osakaaluga. Seda eriti jälgimisaegsete juhtudega andmestikus, kus haigusgrupp moodustab alla 10% vaadeldavatest geenidoonoritest. Kui uurida gruppide vanuselist jaotust proovi andmise ajal joonisel 8, on näha, et nende jaotused ei ole ühesugused. Kontrollgruppi kuuluvad peamiselt alla 50-aastased. Seevastu ISH diagnoosiga isikud olid peami-

selt üle 50 aasta vanused. Haigusgrupi isikutele vastavate kontrollidega saaks teha paremaid järeldusi.



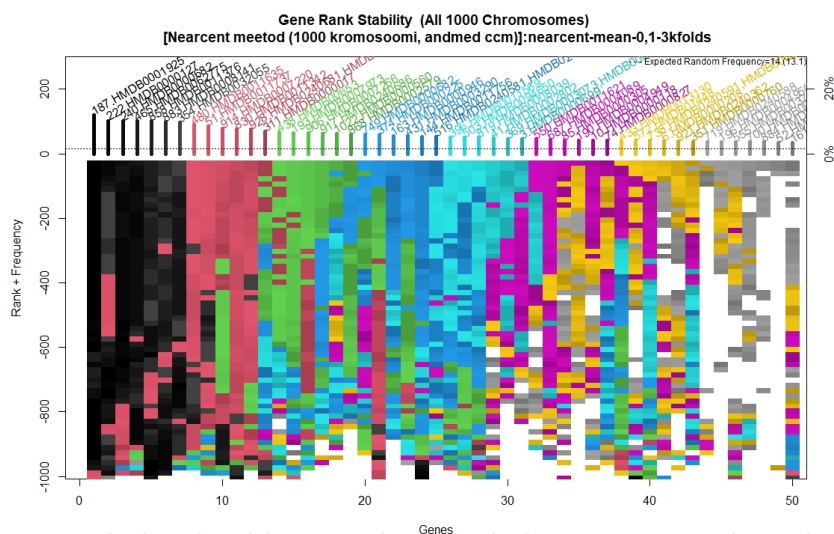
Joonis 8: Kontrollgrupi, jälgimisaegsete ja -eelsete juhtude vanuste jaotused proovi andmise ajal

Autori kirjutatud R-i skriptiga valiti igale haigusgrupi isikule kontrollgrupi isikute hulgast välja 2 temaga kõige sarnasemat. Kontroll pidi olema haigusgrupi isikuga samast soost ning mõlema vanusevahe ei tohtinud ületada 2 aastat. Kui ühele haigusgrupi isikule vastas rohkem kui 2 kontrollgrupi isikut, siis valiti nende hulgast 2 välja juhuslikult. Saadi 2 andmestikku, millest ühes olid jälgimisaegsete juhtudega isikud ning nendega sarnased kontrollid ja jälgimiseelsete juhtudega isikud ning nendega sarnased kontrollid. Kõigile 168 jälgimiseelsete juhtudega isikutele ei õnnestunud leida 2 sarnast kontrolli, mistõttu valiti välja vaid 322 kontrolli. Seega kontrollgrupi isikute osakaal jälgimisaegsete juhtudega andmestikus on suurem kui jälgimiseelsete juhtudega andmestikus.

Mõlemal andmestikul rakendati GALGO-ga GA protsessi ühesuguste parameetrite-

ga. Kokku tehti 1000 evolutsiooni, igas evolutsioonis maksimaalselt 200 põlvkonda. Iga kromosoom sisaldas 10 metaboliiti ning soovitud sobivuseks valiti 70%.

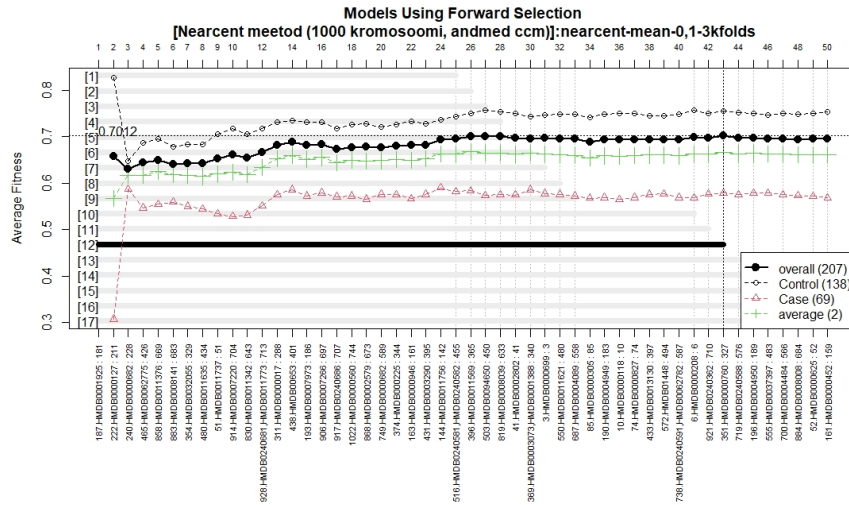
6.2.1 Jälgimisaegseid juhte prognoosivad metaboliidid



Joonis 9: Metaboliitide üldine astakute stabiilsus ning sagedus jälgimisaegseid juhte ning sarnaseid kontrollid sisaldava andmestiku põhjal

Joonisel 9 paistab, et esimese 10 metaboliidi järjestus on olnud kõigi 1000 evolutsiooni käigus pidevas muutumises. Graafikult paistab, et lõpliku järjestuseni on jõutud viimase 100 evolutsiooni käigus. Sellele viitab ka madal vahe metaboliitide sageduste arvude vahel. Kõige sagedamat metaboliiti HMDB0001925 ehk ibuprofeeni esines 120 korda, millele järgnes 102 korda esinenud HMDB0000127 ehk glükuroonhape.

Joonisel 10 on näha, et kõige parema üldise täpsusega 70% on mudel, millesse lisada esimesed 43 metaboliiti. Kõige parema täpsusega prognoosib haigusgruppi esimese 3 metaboliidiga mudel ligikaudu 60% täpsusega. Sellegipoolest mudeli üldine täpsus 63% on teiste mudelitega võrreldes madalaim. Jooniselt paistab, et 14 metaboliidiga mudel prognoosib haigusgruppi sarnase täpsusega, ent selle üldine täpsus on lähemal GALGO poolt valitud mudeli täpsusele.

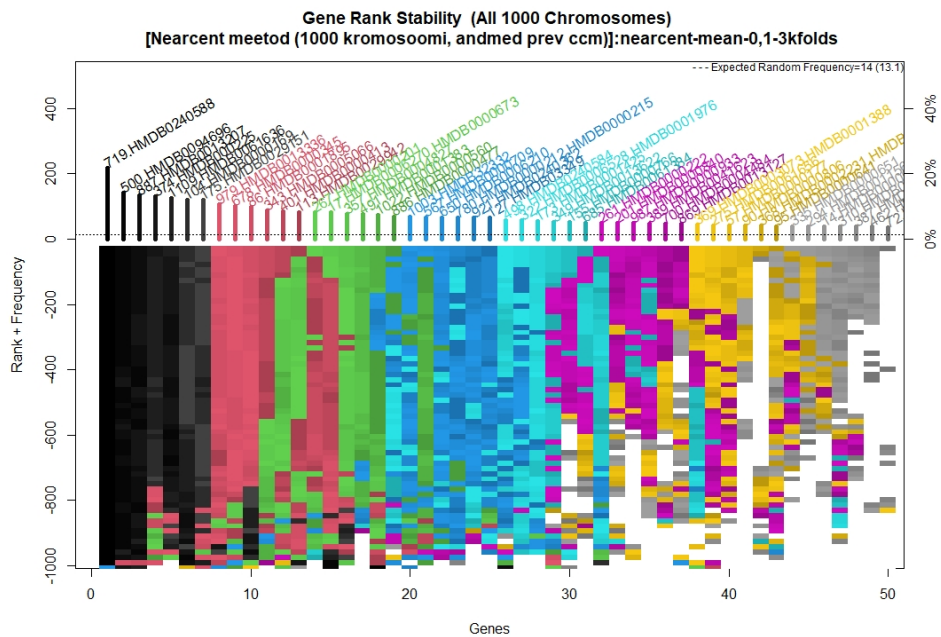


Joonis 10: Sarnaste kontrollidega jälgimisaegsete juhtude metaboliitide järjestuse ettesuunatud valik

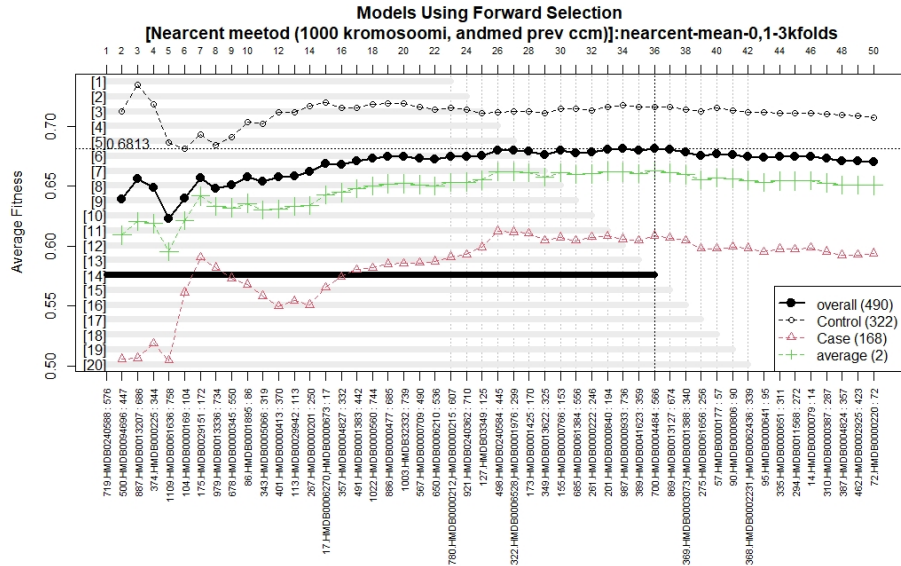
6.2.2 Jälgimiseelseid juhte prognoosivad metaboliidid

Joonisel 11 on näha, et esimese 10 metaboliidi järjestus on püsinud enam-vähem muutumatuna viimased 700 evolutsiooni. Kuigi need metaboliidid on mitmeid korra vahepeal omavahel järjestust vahetanud, ei paista, et mõni metaboliit tagapool oleks esikümne hulka jõudnud. Kirjud värvid esimeste metaboliitide all viitavad sellele, et metaboliidid asusid esimese 200 evolutsiooni järel tagapool. Kõige sagedamini valiti 217 korda metaboliiti HMDB0240588 ehk müristoleüülkarnitiin (inglise keeles *myristoleoylcarnitine*) ning teisel kohal on 141 korral esinenud HMDB0094696 ehk N-metüül-L-proliin (inglise keeles *N-methyl-L-proline*).

Kui vaadata joonist 12, siis on näha, et kõige parema üldise täpsusega 68% mudeli saab, kui mudeli argumentideks võtta esimesed 36 metaboliiti. Võttes mudeli argumentideks 26 metaboliiti, oleks mudeli haigusgrupi prognoosi täpsus natuke suurem ning üldine täpsus sarnane GALGO valitud mudeliga võrreldes. Sarnaned täpsused oleksid ka juhul, kui võtta mudeli argumentideks esimesed 27 metaboliiti.



Joonis 11: Metaboliitide üldine astakute stabiilsus ning sagedus sarnaste kontrollidega jälgimiseelsete juhtude põhjal



Joonis 12: Sarnaste kontrollidega jälgimiseelsete juhtude metaboliitide järjestuse ettesuunatud valik

6.3 Logistiline regressioonimudel

Autor jooksutas GALGO-ga andmeid mitu korda läbi ning sooritas saadud metaboliitide järjestustele ettesuunatud valikuid. Ettesuunatud valikutes tulid parima täpsusega mudelid metaboliitide arvuga 1 kuni 43. Optimaalseks metaboliitide arvuks mudelisse valiti 35, mille põhjal leiti sõltuvalt kontrollgrupi isikute arvust esimese 35 metaboliidi hulgast ühised metaboliidid jälgimisaegsete ja -eelsete juhtude vahel. Kokku saadi 4 metaboliiti, millest 1 leiti kõikide kontrollide põhjal ja ülejäänud sarnaste kontrollidega andmestikust. Metaboliitide nimetused ja järjekorranumbrid vastava haigusgrupi andmestiku kohta on välja toodud tabelis 2.

Tabel 2: Jälgimisaegsete ja -eelsete juhtude ühised metaboliidid esimese 35 hulgast

HMDB kood	Ingliskeelne nimetus	Järjekorranumber järjestuses	
		Jälgimisaegne	Jälgimiseelne
HMDB0012458	<i>γ</i> alpha-Hydroxy-3-oxo-4-cholestenoate	32	25
HMDB0002250	<i>Dodecanoylcarnitine</i>	11	4
HMDB0000560	<i>Goshuyic acid</i>	18	18
HMDB0013207	<i>9-Hexadecenoylcarnitine</i>	26	3

Koos isiku taustinfoga nagu sugu, vanus, kehamassiindeks ning elukoha tüüp lisati saadud metaboliidid logistilise regressioonimudelitesse, mis prognoosivad jälgimisaegseid ja -eelseid juhte. Mudel hinnati andmestikuga, mis sisaldas kõiki kontrollgrupi isikuid. Kui tunnuse p-väärtus ületas olulisuse nivood $\alpha = 0,05$, jäeti see mudelist välja. Tunnuseid jäeti ükshaaval välja, alustades kõige suurema p-väärtusega tunnusest.

Jälgimisaegsete juhtude mudeldamisel osutusid statistiliselt oluliseks vaid 2 tun-

nust: isiku vanus ja metaboliidi dodekanoüülkarnitiini (inglise keeles *Dodecanoylcarnitine*, *Dodecan*) kontsentratsioon. Tabelis 3 on välja toodud lõpliku mudeli parameetrite hinnangud koos standardvigadega ja p-väärtustega.

Tabel 3: Jälgimisaegse juhu lõplik mudel

Tunnus	Parameetri hinnang	Standardviga	P-väärtus
Vabaliige	-3,60	0,41	$2 \cdot 10^{-16}$
Vanus	0,02	0,0078	0,012
<i>Dodecan</i>	0,23	0,12	0,045

Mudelis on mõlema tunnuse parameetrite hinnangud positiivse väärtusega, mis tähendab, et suurema väärtuse korral suureneb šanss peale proovi andmist ISH-sse haigestuda. Olgu vaatluse all 2 sama vanusega isikut, kellest ühel on dodekanoüülkarnitiini kontsentratsioon 1 ühiku võrra suurem kui teisel. Mudeli põhjal oleks suurema dodekanoüülkarnitiini kontsentratsiooniga isikul 1,26 korda suurem šanss haigestuda kui teisel.

Tabel 4: Jälgimiseelse juhu lõplik mudel

Tunnus	Parameetri hinnang	Standardviga	P-väärtus
Vabaliige	-3,39	0,51	$1,31 \cdot 10^{-14}$
Naine	-0,051	0,18	0,0058
Kehamassiindeks	0,04	0,016	0,011
Vanus	0,03	0,0057	$7,7 \cdot 10^{-6}$
<i>Gosh. acid</i>	0,3	0,086	0,00047

Jälgimiseelsete juhtude lõplikusse mudelisse jäi alles 4 statistiliselt olulist tunnust: sugu, kehamassiindeks, vanus ning metaboliidi (Z,Z)-5,8-tetradekadieenhappe (inglise keeles *Goshuyic acid*) kontsentratsioon. Soo puhul on mudelis baastunnuseks meessugu. Tabelis 4 on välja toodud mudeli parameetrite hinnangud, standardvead ning p-väärtused.

Mudelis on tunnuste parameetri hinnangud positiivsed, välja arvatud sool. On näha, et kõrgema kehamassiindeksiga, vanusega ja *Gosh. acid* kontsentratsiooniga meestel on suurem šans olla varem diagnoositud ISH-ga proovi andmisel. Olgu vaatluse all 2 sama kehamassiindeksiga ning vanusega meest, kellest ühel on ühiku võrra suurem metaboliidi *Gosh. acid* kontsentratsioon kui teisel. Mudeli põhjal on suurema metaboliidiga isikul 1,35 korda suurem šans olla varasema ISH diagnoosiga kui kui teisel.

7 Arutelu

Käesoleva töö tulemustena leitud metaboliite on kavas uurida lähemalt edaspidistes uuringutes, et leida südame isheemiatõbede tekkeprotsessidega seotud metaboliite.

Töös tutvustati geneetilist algoritmi ja rakendati seda Tartu Ülikooli Eesti Geenivaramu andmestikule. Tulemusteks saadi 4 andmestiku põhjal 4 erinevat metaboliitide järjestust vastavalt sellele, mitmes parima sobivusega kromosoomis metaboliit sisaldus. Kui igast järjestusest võtta esimesed 35 metaboliiti ehk kõige sagedamad, siis leiduvad vaid üksikud metaboliidid, mis sisalduvad vähemalt 2 järjestuses. Neist mudelisse lisati 4 metaboliiti, mis kattusid andmestikega, millel olid ühised kontrollid. Hinnati 2 logistilist regressioonimudelit. Kummagi mudeli tunnuste hulgas oli 1 metaboliit statistiliselt oluline.

Statiliselt oluliste metaboliitide hulka kuulus dodekanoüülkarnitiin. See atsüülkarnitiinide hulka lipiidide grupis, mille peamiseks ülesanneteks on transportida orgaanilisi happeid ja rasvhappeid mitokondrisse energia tootmiseks (Inimese Metaboliitide Andmebaas, 2023a). Kõrgema kontsentratsiooniga atsüülkarnitiinide metaboliite on seostatud teist tüüpi diabeetikutest südame- ja veresoonkonnahaiguse diagnoosiga patsientidel (Zhao *et al.*, 2020) ning südame vasaku vatsakese puudulikkusega (Zordoky *et al.*, 2015).

Teine statistiliselt oluline metaboliit oli (Z,Z)-5,8-tetradekadieenhape, mis on pika ahelaga rasvhape ja asub peamiselt rakumembraanis. See tekib vahesaadusena küllastumata rasvhapete oksüdeerumisel ning sisaldub rapsiõlis (Inimese Metaboliitide Andmebaas, 2023b).

Kõige sagedamate metaboliitide hulgas leidub enim lipiide, kuna andmestikus üle poole metaboliitidest on lipiidid. Neist enamik on rasvhapped, kuid eksisteerivad ka mõned steroidid.

Ksenobiootikume ülejäänud metaboliidi tüüpidega võrreldes esineb samuti teistest rohkem. Näiteks kõikide kontrollidega jälgimiseelsete juhtude andmestikus kuulus

esimese 4 metaboliidi hulka 3 ksenobiootikumi. Nende hulgas esikohal paiknev perfluorooktaanhape (PFOA) on väga levinud kemikaal, kuna organism ei sünteesi ega lagunda seda. PFOA-d peetakse mürgiseks ning on leitud seoseid kolesterooliga. (Steenland, Fletcher ja Savitz, 2010). On leitud ka, et PFOA tase on kõrgem vanemaeliste inimeste hulgas (Steenland *et al.*, 2009).

Jälgimisaegsete juhtude puhul töö autor ei olnud teadlik, mitu aastat peale proovi andmist haigus diagnoositi. Kui haigus avaldus 10 aastat hiljem, siis vanade metaboolomika andmete põhjal ei pruugi seda prognoosida. Ideaaltingimustes eeldatakse, et isik ei ole vahepeal oma käitumisharjumusi muutnud. Erinevalt haiglas läbiviidavatest sarnastest uuringutest, kus on kontrollitud keskkond, geenidoonoritele säärased korraldusi ei määrata. Seega võib haigus tekkida muu elustiili käigus, mille kohta andmed puuduvad.

Tulemusteni jõuti ainult ühe klassifikatsioonimeetodiga, kuigi plaanis oli võrrelda 3 meetodi tulemusi. Peamine takistus teistel meetoditel oli madal sobivus, mis põlvkondade jooksul tõusis vähehaaval. Et jõuda soovitud sobivuseni, oli vaja läbida rohkem põlvkondi, mistõttu kuluks protsessiks rohkem aega. Autori isikliku rüperaaliga võib selleks kuluda nädalaid. Kui tulevikus andmestik uuesti TVM ja juhusliku metsa klassifikatsioonimeetodiga GALGO-s läbi jooksutada, siis tarvitseb seda teha võimsama riistvaraga.

Töös leitud markeritest ei saa teha põhjuslikke järeldusi, kuna ajapuuduse tõttu ei olnud võimalik metaboliite täpsemalt uurida. Need tuleks tulevikus valideerida iseseisva andmestikuga. Geenivaramul on olemas veel metaboolomika andmestikke, millega saab tulemusi valideerida.

Kokkuvõte

Bakalaureusetöö eesmärk oli Tartu Ülikooli Geenivaramu metaboolika andmete põhjal leida metaboliidid, mis sobiksid potentsiaalseteks südame isheemiatõbesid kirjeldavateks markeriteks.

Tulemuste leidmiseks rakendati andmetele geneetilist algoritmi. Töös üritati võrrelda 3 klassifikatsioonimeetodil saadud tulemusi. Tulemused saadi lähima tsentroidi meetodil. Paraku protsessiks nõutava mahuka ajakulu tõttu tugivektor-masina ning juhusliku metsa meetodil tulemusteni ei jõutud.

Südame isheemiatõvega diagnoositud geenidoonorid jagati kahte gruppi: jälgimiseaegsed juhud ning jälgimiseelsed juhud. Kummagi haigusgrupiga eraldi moodustati andmestikud, millest esimesse kuulusid kõik kontrollgruppi kuulunud geenidoonorid ning teise haigusgrupi isikutega sarnased kontrollid. Kõikidel andmestikel rakendati geneetilist algoritmi lähima tsentroidi klassifikatsioonimeetodiga.

Iga andmestikuga saadud tulemustest valiti välja esimesed 35 metaboliiti. Nende hulgast saadi 4 metaboliiti, mis kuulusid võrdse arvu kontrollidega andmestike ühisossa. Töös hinnati 2 logistilist regressioonimudelit prognoosimaks haigusgruppe, kus lisaks saadud metaboliitidele kasutati tunnustena doonori sugu, vanust proovi andmise hetkel, kehamassiindeksit ning elupaika. Mõlema mudeli peale sai kokku 2 statistiliselt olulist metaboliiti: dodekanoüülkarnitiin ja (Z,Z)-5,8-tetradekadieenhappe. Varasemad uuringud on leidnud seoseid dodekanoüülkarnitiiniga ja südamega seotud haigustega.

Kasutatud allikad

- Ahmed, S., P. Shah ja O. Ahmed (2023). *StatPearls [Internet]. Biochemistry, Lipids*. Treasure Island (FL): StatPearls Publishing. URL: <https://www.ncbi.nlm.nih.gov/books/NBK525952/> (vaadatud 06.05.2023).
- Almdal, Thomas, Henrik Scharling, Jan Skov Jensen ja Henrik Vestergaard (2004). “The independent effect of type 2 diabetes mellitus on ischemic heart disease, stroke, and death: a population-based study of 13 000 men and women with 20 years of follow-up”. *Archives of internal medicine* 164.13, lk. 1422–1426.
- Ameerika Südameassotsiatsioon (2022). *Silent Ischemia and Ischemic Heart Disease*. URL: <https://www.heart.org/en/health-topics/heart-attack/about-heart-attacks/silent-ischemia-and-ischemic-heart-disease> (vaadatud 17.04.2023).
- Belgiu, Mariana ja Lucian Drăguț (2016). “Random forest in remote sensing: A review of applications and future directions”. *ISPRS Journal of Photogrammetry and Remote Sensing* 114, lk. 24–31. ISSN: 0924-2716. DOI: <https://doi.org/10.1016/j.isprsjprs.2016.01.011>. URL: <https://www.sciencedirect.com/science/article/pii/S0924271616000265> (vaadatud 28.03.2023).
- Cervantes, Jair, Farid Garcia-Lamont, Lisbeth Rodríguez-Mazahua ja Asdrubal Lopez (2020). “A comprehensive survey on support vector machine classification: Applications, challenges and trends”. *Neurocomputing* 408, lk. 189–215. ISSN: 0925-2312. DOI: <https://doi.org/10.1016/j.neucom.2019.10.118>. URL: <https://www.sciencedirect.com/science/article/pii/S0925231220307153> (vaadatud 28.03.2023).
- Euroopa Bioinformaatika instituut (2023). *Comparison of NMR and MS*. URL: <https://www.ebi.ac.uk/training/online/courses/metabolomics->

- [introduction/designing-a-metabolomics-study/comparison-of-nmr-and-ms/](#) (vaadatud 21.04.2023).
- Idle, Jeffrey R ja Frank J Gonzalez (2007). “Metabolomics”. *Cell metabolism* 6.5, lk. 348–351.
- Inimese Metaboliitide Andmebaas (2023a). *Showing metabocard for Dodecanoylcarnitine (HMDB0002250)*. URL: <https://hmdb.ca/metabolites/HMDB0002250> (vaadatud 09.05.2023).
- (2023b). *Showing metabocard for Goshuyic acid (HMDB0000560)*. URL: <https://hmdb.ca/metabolites/HMDB0000560> (vaadatud 09.05.2023).
- Katobh, S., S. S. Chauhan ja V. Kumar (2021). “A review on genetic algorithm: past, present, and future”. *Multimedia Tools and Applications* 80.5, lk. 8091–8126. URL: <https://link.springer.com/article/10.1007/s11042-020-10139-6> (vaadatud 27.03.2023).
- Levner, Ilya (2005). “Feature selection and nearest centroid classification for protein mass spectrometry”. *BMC bioinformatics* 6.1, lk. 1–14.
- Lin, Xiaohui, Quancai Wang, Peiyuan Yin, Liang Tang, Yexiong Tan, Hong Li, Kang Yan ja Guowang Xu (2011). “A method for handling metabonomics data from liquid chromatography/mass spectrometry: combinational use of support vector machine recursive feature elimination, genetic algorithm and random forest for feature selection”. *Metabolomics* 7, lk. 549–558.
- Maailma Terviseorganisatsioon (2016). *International Statistical Classification of Diseases and Related Health Problems 10th Revision (ICD-10)*. URL: <https://icd.who.int/browse10/2016/en#/I20-I25> (vaadatud 28.03.2023).

- Maaailma Terviseorganisatsioon (2020). *The top 10 causes of death*. URL: <https://www.who.int/news-room/fact-sheets/detail/the-top-10-causes-of-death> (vaadatud 08.05.2023).
- Mitchell, Melanie (1998). *An introduction to genetic algorithms*. MIT press.
- Morbiducci, Umberto, Andrea Tura ja Mauro Grigioni (2005). “Genetic algorithms for parameter estimation in mathematical modeling of glucose metabolism”. *Computers in Biology and Medicine* 35.10, lk. 862–874. ISSN: 0010-4825. DOI: <https://doi.org/10.1016/j.combiomed.2004.07.005>. URL: <https://www.sciencedirect.com/science/article/pii/S0010482504000976>.
- Morgan-Benita, Jorge, Ana G Sánchez-Reyna, Carlos H Espino-Salinas, Juan José Oropeza-Valdez, Huizilopoztli Luna-García, Carlos E Galván-Tejada, Jorge I Galván-Tejada, Hamurabi Gamboa-Rosales, Jose Antonio Enciso-Moreno ja José Celaya-Padilla (2022). “Metabolomic Selection in the Progression of Type 2 Diabetes Mellitus: A Genetic Algorithm Approach”. *Diagnostics* 12.11, lk. 2803.
- Nick, Todd G ja Kathleen M Campbell (2007). “Logistic regression”. *Topics in biostatistics*, lk. 273–301.
- Steenland, Kyle, Tony Fletcher ja David A Savitz (2010). “Epidemiologic evidence on the health effects of perfluorooctanoic acid (PFOA)”. *Environmental health perspectives* 118.8, lk. 1100–1108.
- Steenland, Kyle, Chuangfang Jin, Jessica MacNeil, Cathy Lally, Alan Ducatman, Veronica Vieira ja Tony Fletcher (2009). “Predictors of PFOA levels in a community surrounding a chemical plant”. *Environmental health perspectives* 117.7, lk. 1083–1088.

- Štefanac, Tea, Dijana Grgas ja Tibela Landeka Dragičević (oktoober 2021). “Xenobiotics-division and methods of detection: A review”. en. *J. Xenobiot.* 11.4, lk. 130–141.
- Trevino, V. ja F. Falciani (2006). “GALGO: an R package for multivariate variable selection using genetic algorithms”. *Journal of the American Statistical Association* 22.9, lk. 1154–1156. URL: https://www.researchgate.net/publication/7266382_GALGO_An_R_package_for_multivariate_variable_selection_using_genetic_algorithms (vaadatud 26.03.2023).
- Trevino, Victor ja Francesco Falciani (2018). *galgo: Genetic Algorithms for Multivariate Statistical Models from Large-Scale Functional Genomics Data*. R package version 1.4. URL: <https://github.com/vtrevino/GALGO>.
- Varbo, Anette, Marianne Benn, Anne Tybjærg-Hansen, Anders B Jørgensen, Ruth Frikke-Schmidt ja Børge G Nordestgaard (2013). “Remnant cholesterol as a causal risk factor for ischemic heart disease”. *Journal of the American College of Cardiology* 61.4, lk. 427–436.
- Zhao, Shuo, Xiao-Fei Feng, Ting Huang, Hui-Huan Luo, Jian-Xin Chen, Jia Zeng, Muyu Gu, Jing Li, Xiao-Yu Sun, Dan Sun, Xilin Yang, Zhong-Ze Fang ja Yun-Feng Cao (mai 2020). “The association between acylcarnitine metabolites and cardiovascular disease in Chinese patients with type 2 diabetes mellitus”. en. *Front. Endocrinol. (Lausanne)* 11, lk. 212.
- Zordoky, Beshay N, Miranda M Sung, Justin Ezekowitz, Rupasri Mandal, Beomsoo Han, Trent C Bjorndahl, Souhaila Bouatra, Todd Anderson, Gavin Y Oudit, David S Wishart, Jason R B Dyck ja Alberta HEART (mai 2015). “Metabolomic fingerprint of heart failure with preserved ejection fraction”. en. *PLoS One* 10.5, e0124844.

Lisa 1. GALGO protsessiga leitud kõige sagedamad metaboliidid

Tabel 5: Kõige sagedamad 35 metaboliiti kõikide kontrollidega jälgimisaegsetel juhtudel

Jrk	HMDB kood	Metaboliidi ingliskeelne nimetus
1	HMDB0013288	<i>Nonanoylcarnitine</i>
2	HMDB0000760	<i>Hyocholic acid</i>
3	HMDB0011737	<i>gamma-Glutamylglutamic acid</i>
4	HMDB0000791	<i>Octanoylcarnitine</i>
5	HMDB0000535	<i>Caproic acid</i>
6	HMDB0000619	<i>Cholic acid</i>
7	HMDB0011635	<i>p-Cresol sulfat</i>
8	HMDB0062121	<i>Dihydroferulic acid</i>
9	HMDB0001870	<i>Benzoic acid</i>
10	HMDB0004980	<i>cis-4-Decenoic acid</i>
11	HMDB0002250	<i>Dodecanoylcarnitine</i>
12	HMDB0000792	<i>Sebacic acid</i>
13	HMDB0008141	PC(18:2(9Z,12Z)/18:3(9Z,12Z,15Z))
14	HMDB0000756	<i>Hexanoylcarnitine</i>
15	HMDB0002869	<i>Campesterol</i>

16	HMDB0000686	<i>Isoursodeoxycholic acid</i>
17	HMDB0135261	<i>Propyl 4-hydroxybenzoate sulfate</i>
18	HMDB0000560	<i>Goshuyic acid</i>
19	HMDB0000482	<i>Caprylic acid</i>
20	HMDB0060426	<i>8-Methoxykynurenate</i>
21	HMDB0000946	<i>Ursodeoxycholic acid</i>
22	HMDB0029160	<i>gamma-Glutamyltryptophan</i>
23	HMDB0031518 vôi HMDB0000156	<i>D-Malic acid vôi Malic acid</i>
24	HMDB0000134	<i>Fumaric acid</i>
25	HMDB0000666	<i>Heptanoic acid</i>
26	HMDB0013207	<i>9-Hexadecenoylcarnitine</i>
27	HMDB0007220	DG(18:1(9Z)/18:3(6Z,9Z,12Z)/0:0)
28	HMDB0062775	<i>4-Vinylphenol sulfate</i>
29	HMDB0001406	<i>Niacinamide</i>
30	HMDB0006335 vôi HMDB0001492	<i>beta-Tocopherol vôi gamma-Tocopherol</i>
31	HMDB0000512	<i>N-Acetyl-L-phenylalanine</i>
32	HMDB0012458	<i>7alpha-Hydroxy-3-oxo-4-cholestenoate</i>
33	HMDB0000883	<i>L-Valine</i>
34	HMDB0000452	<i>L-alpha-Aminobutyric acid</i>

35	HMDB0000131	<i>Glycerol</i>
----	-------------	-----------------

Tabel 6: Kõige sagedamad 35 metaboliiti kõikide kontrollidega jälgimiseelse-
tel juhtudel

Jrk	HMDB kood	Metaboliidi ingliskeelne nimetus
1	HMDB0059587	<i>Perfluorooctanoic acid (PFOA)</i>
2	HMDB0000840	<i>Salicyluric acid</i>
3	HMDB32332	<i>Hydroxylated lecithin</i>
4	HMDB0001895	<i>Salicylic acid</i>
5	HMDB0061112	<i>3-Carboxy-4-methyl-5-propyl-2-furanpropionic acid</i>
6	HMDB0000169	<i>D-Mannose</i>
7	HMDB0007266	DG(18:2(9Z,12Z)/22:6(4Z,7Z,10Z,13Z,16Z,19Z)/0:0)
8	HMDB0000933	<i>Traumatic acid</i>
9	HMDB0059586	<i>Perfluorooctanesulfonate (PFOS)</i>
10	HMDB0000122	<i>D-Glucose</i>
11	HMDB0240565	<i>Umbelliferone sulfate</i>
12	HMDB0061643	<i>3-carboxy-4-methyl-5-pentyl-2-furanpropanoic acid</i>
13	HMDB0000613	<i>Erythronic acid</i>

14	HMDB0000413	<i>3-Hydroxydodecanedioic acid</i>
15	HMDB0000258	<i>Sucrose</i>
16	HMDB0004824	<i>N2,N2-Dimethylguanosine</i>
17	HMDB0013336	<i>3-Hydroxyhexadecanoylcarnitine</i>
18	HMDB0011568	MG(18:2(9Z,12Z)/0:0/0:0)
19	HMDB0002577	<i>Cholic acid glucuronide</i>
20	HMDB0001173	<i>5'-Methylthioadenosine</i>
21	HMDB0000337	<i>(S)-3,4-Dihydroxybutyric acid</i>
22	HMDB0240506 vôi	<i>Enterolactone 3''-sulfate vôi</i>
	HMDB0240505	<i>Enterolactone 3'-sulfate</i>
23	HMDB0011565	MG(16:1(9Z)/0:0/0:0)
24	HMDB0002189 vôi	<i>N8-Acetylspermidine vôi</i>
	HMDB0001276	<i>N1-Acetylspermidine</i>
25	HMDB0012458	<i>7alpha-Hydroxy-3-oxo-4-cholestenoate</i>
26	HMDB0000291	<i>Vanillylmandelic acid</i>
27	HMDB0013127	<i>3-Hydroxybutyrylcarnitine</i>
28	HMDB0000755	<i>Hydroxyphenyllactic acid</i>
29	HMDB0000094	<i>Citric acid</i>
30	HMDB0033433	<i>(S)-Homostachydrine</i>
31	HMDB0240345	<i>N2,N5-Diacetylornithine</i>
32	HMDB0013127	<i>3-Hydroxybutyrylcarnitine</i>

33	HMDB0000350	<i>3-Hydroxysebacic acid</i>
34	HMDB0002331	<i>Imidazoleacetic acid riboside</i>
35	HMDB0094700	<i>Perfluorohexanesulfonate (PFHxS)</i>

Tabel 7: Kõige sagedamad 35 metaboliiti sarnaste kontrollidega jälgimisaegsetel juhtudel

Jrk	HMDB kood	Metaboliidi ingliskeelne nimetus
1	HMDB0001925	<i>Ibuprofen</i>
2	HMDB0000127	<i>D-Glucuronic acid</i>
3	HMDB0000682	<i>Indoxyl sulfate</i>
4	HMDB0062775	<i>4-Vinylphenol sulfate</i>
5	HMDB0011376	PE(P-18:0/18:2(9Z,12Z))
6	HMDB0008141	PC(18:2(9Z,12Z)/18:3(9Z,12Z,15Z))
7	HMDB0032055	<i>N-Acetylhistidine</i>
8	HMDB0011635	<i>p-Cresol sulfate</i>
9	HMDB0011737	<i>gamma-Glutamylglutamic acid</i>
10	HMDB0007220	DG(18:1(9Z)/18:3(6Z,9Z,12Z)/0:0)
11	HMDB0011342	PE(P-16:0/18:1(9Z))
12	HMDB0240681 või HMDB0011773	Cer(d16:1/16:0) või Cer(d18:1/14:0)
13	HMDB0000017	<i>4-Pyridoxic acid</i>

14	HMDB00653	<i>Cholesterol sulfate</i>
15	HMDB0007973	PC(16:0/18:2(9Z,12Z))
16	HMDB0007266	DG(18:2(9Z,12Z)/22:6(4Z,7Z,10Z,13Z, 16Z,19Z)/0:0)
17	HMDB0240686	Cer(d18:2(4E,14Z)/16:0)
18	HMDB0000560	<i>Goshuyic acid</i>
19	HMDB0002579	<i>Glycochenodeoxycholic acid 3-glucuronide</i>
20	HMDB0000682	<i>Indoxyl sulfate</i>
21	HMDB0002250	<i>Dodecanoylcarnitine</i>
22	HMDB0000946	<i>Ursodeoxycholic acid</i>
23	HMDB0003290	<i>Gulonic acid</i>
24	HMDB0011756	<i>N-Acetylleucine</i>
25	HMDB0240581 vôi HMDB0240582	<i>5alpha-Pregnan-3beta,20beta-diol disulfate vôi 5alpha-Pregnan-3alpha,20beta-diol disulfate</i>
26	HMDB0011569	MG(18:3(6Z,9Z,12Z)/0:0/0:0)
27	HMDB0094650	<i>5alpha-Pregnan-3beta,20alpha-diol disulfate</i>
28	HMDB0008039	PC(18:0/18:2(9Z,12Z))
29	HMDB0002802	<i>Cortisone</i>

30	HMDB0003073 või HMDB0001388	<i>gamma-Linolenic acid</i> või <i>alpha-Linolenic acid</i>
31	HMDB0000699	<i>1-Methylnicotinamide</i>
32	HMDB0011621	<i>Cinnamoylglycine</i>
33	HMDB0004089	<i>Formylanthranilic acid</i>
34	HMDB0000305	<i>Vitamin A</i>
35	HMDB0004949	<i>Cer(d18:1/16:0)</i>

Tabel 8: Kõige sagedamad 35 metaboliiti sarnaste kontrollidega jälgimiseelsetel juhtudel

Jrk	HMDB kood	Metaboliidi ingliskeelne nimetus
1	HMDB0240588	<i>Myristoleoylcarnitine</i>
2	HMDB0094696	<i>N-Methyl-L-proline</i>
3	HMDB0013207	<i>9-Hexadecenoylcarnitine</i>
4	HMDB0002250	<i>Dodecanoylcarnitine</i>
5	HMDB0061636	<i>3-hydroxydecanoyl carnitine</i>
6	HMDB0000169	<i>D-Mannose</i>
7	HMDB0029151	<i>gamma-Glutamylhistidine</i>
8	HMDB0013336	<i>3-Hydroxyhexadecanoylcarnitine</i>
9	HMDB0000345	<i>3-Hydroxyadipic acid</i>
10	HMDB0001895	<i>Salicylic acid</i>

11	HMDB0005066	<i>Tetradecanoylcarnitine</i>
12	HMDB0000413	<i>3-Hydroxydodecanedioic acid</i>
13	HMDB0029942	<i>D-Arabinose</i>
14	HMDB0000201	<i>L-Acetylcarnitine</i>
15	HMDB0006270 vői HMDB0000673	<i>Linoelaidic acid vői Linoleic acid</i>
16	HMDB0004827	<i>Proline betaine</i>
17	HMDB0001383	<i>Sphinganine 1-phosphate</i>
18	HMDB0000560	<i>Goshuyic acid</i>
19	HMDB0000477	<i>7Z,10Z-Hexadecadienoic acid</i>
20	HMDB32332	<i>Hydroxylated lecithin</i>
21	HMDB0000709	<i>L-Cysteinylglycine disulfide</i>
22	HMDB0006210	<i>Heptadecanoyl carnitine</i>
23	HMDB0000212 vői HMDB0000215	<i>N-Acetylgalactosamine vői N-Acetyl-D-glucosamine</i>
24	HMDB0240362	<i>N-stearoyl-sphingadienine</i>
25	HMDB03349	<i>L-Dihydroorotic acid</i>
26	HMDB0240584	<i>(4E,15E)-Bilirubin</i>
27	HMDB0006528 vői HMDB0001976	<i>Docosapentaenoic acid (22n-3) vői Docosapentaenoic acid (22n-6)</i>
28	HMDB0001425	<i>Estrone sulfate</i>

29	HMDB0013622	<i>10Z-Nonadecenoic acid</i>
30	HMDB0000766	<i>N-Acetyl-L-alanine</i>
31	HMDB0061384	<i>N-(3-acetamidopropyl)pyrrolidin-2-one</i>
32	HMDB0000222	<i>Palmitoylcarnitine</i>
33	HMDB0000840	<i>Salicyluric acid</i>
34	HMDB0000933	<i>Traumatic acid</i>
35	HMDB0041623	<i>N6-Carbamoyl-L-threonyladenosine</i>

Lihtlitsents lõputöö reprodutseerimiseks ja üldsusele kättesaadavaks tegemiseks

Mina, Johann Saavaste,

1. annan Tartu Ülikoolile tasuta loa (lihtlitsentsi) minu loodud teose "Isheemiliste südamehaiguste markerite leidmine Tartu Ülikooli Eesti Geenivaramu metaboolomika andmetele geneetilist algoritmi rakendades", mille juhendajad on Jaanika Kronberg ja Krista Fischer, reprodutseerimiseks eesmärgiga seda säilitada, sealhulgas lisada digitaalarhiivi DSpace kuni autoriõiguse kehtivuse lõppemiseni.
2. Annan Tartu Ülikoolile loa teha punktis 1 nimetatud teos üldsusele kättesaadavaks Tartu Ülikooli veebikeskkonna, sealhulgas digitaalarhiivi DSpace kaudu Creative Commons'i litsentsiga CC BY NC ND 4.0, mis lubab autorile viidates teost reprodutseerida, levitada ja üldsusele suunata ning keelab luua tuletatud teost ja kasutada teost ärieesmärgil, kuni autoriõiguse kehtivuse lõppemiseni.
3. Olen teadlik, et punktides 1 ja 2 nimetatud õigused jäävad alles ka autorile.
4. Kinnitan, et lihtlitsentsi andmisega ei riku ma teiste isikute intellektuaalomandi ega isikuandmete kaitse õigusaktidest tulenevaid õigusi.

Johann Saavaste

09.05.2023