



# Klasterizavimo algoritmai didelės apimties medicinos duomenims

Ataskaita už III-ąją doktorantūros kursą (2021/2022 metų I  
pusmetis)

Doktorantas: Roma Puronaitė  
Vadovas: prof. dr. Audronė Jakaitienė

Doktorantūros pradžios ir pabaigos metai: 2017 - 2023

2022 m. kovas



# Tyrimo objektas



- Klasterizavimo algoritmai
- Didelės apimties medicinos duomenys



# Tyrimo tikslas



- Pasiūlyti metodą didelės apimties medicinos duomenims klasterizuoti, atsižvelgiant į duomenų dinamikos laike savybes.



# Uždaviniai



- 1 Iširti didelės apimties medicinos duomenų klasterizavimui dažniausiai taikomus klasterizavimo metodus.
- 2 Pasiūlyti klasterizavimo algoritmą ar esamo metodo patobulinimą, kuris atsižvelgtų į dinamiką laike.
- 3 Pritaikyti atrinktus algoritmus ir pasiūlytą sprendimą realiems medicinos duomenų rinkiniams.
- 4 Pasiūlyti algoritmo integravimo į sveikatos priežiūros įstaigos informacinę sistemą modelį.

# Doktorantūros planas

Studijų metai	Egzaminai <sup>1</sup>		Dalyvavimas konferencijose <sup>2</sup>		Publikacijos <sup>3</sup>		
	Planas	Įvykdyta	Planas	Įvykdyta	Planas	Įvykdyta	Būklė <sup>4</sup>
I (2017/2018)	2	2		1		1	1 publikuota
II (2018/2019)	2	0	1	5		2	2 publikuota
III (2021/2022)	0	1 (skola iš II metų) + Antrasis bus laikomas 2022-03-30	1	1 + (1)**	1	2	2 publikuota
IV (2022/2023)	0	0			1		
Iš viso:	4	2	2	6	2	3+1	

\*Iki 2021 m. rugsėjo 30 d. akademinės atostogos

**pav.:** Visų studijų planas

# Ataskaita už III mokslo metus

## 1 Egzaminai:

Egzaminai		
Planas	Ivykdyta	Būklė
Netiesiniai statistikos modeliai masinių duomenų analizėje	Netiesiniai statistikos modeliai masinių duomenų analizėje, 2022 m. kovo 18 d.	Išlaikytas
Daugiamačių duomenų vizualizavimo metodai	Daugiamačių duomenų vizualizavimo metodai, 2022 m. kovo 30 d.	Nelaikytas

## 2 Konferencijos:

Dalyvavimas konferencijose		
<p><i>13th International Conference BIOMDLORE 2021</i>, 2022 m. spalio mėn. 21-23 d, Vilnius, Lietuva</p>	<p>Bliudzius, Antanas, <b>Puronaite, Roma</b>, Trinkunas, Justas, Jakaitiene, Audrone, Kasiulevicius, Vytautas, „Research on physical activity variability and changes of metabolic profile in patients with prediabetes using Fitbit activity trackers data“, <i>13th International Conference BIOMDLORE 2021</i>, 2022 m. spalio mėn. 21-23 d, Vilnius, Lietuva</p> <p>Note: [1] Antanas Bliudzius and Roma Puronaite contributed equally to this work.</p>	Tarptautinė
<p><i>31st International Biometric Conference</i>, 2022 m. liepos 10-15 d., Ryga, Latvija</p> <p>Priimtas žodinis pranešimas.</p>	<p><b>Puronaite, Roma</b>, Ramanauskaitė, Dovilė, Burneikaitė, Greta, Švaikevičienė, Kristina, Šavareikaitė, Alicija, Vaitkute, Samanta, Jakaitienė, Audronė, Dambrauskas, Laimis, Jurevičienė, Elena, Trinkūnas, Justas, Kasiulevičius, Vytautas, Kazėnaitė, Edita, „Challenges of modeling depression and anxiety risk using data from large healthcare databases: systematic review and situation analysis“, <i>31st International Biometric Conference</i>, 2022 m. liepos 10-15 d., Ryga, Latvija</p>	Tarptautinė

# Ataskaita už III mokslo metus

## 1 Publikacijos:

Publikacijos			
Planas	Įvykdyta	Būklė	Publikacijos tipas
Technology and health care	Bliūdžius, A., <b>Puronaitė, R.</b> , Trinkūnas, J., Jakaitienė, A., & Kasiulevičius, V. (2022). Research on physical activity variability and changes of metabolic profile in patients with prediabetes using Fitbit activity trackers data. <i>Technology and health care</i> , 30(1), 231-242. doi:10.3233/THC-219006 [DB: Scopus; Social Sciences Citation Index (Web of Science)] [Indėlis: 0,268] [Citav. rod.: 1.285 (2020, SCIE)] [M.kr.: M 001,N 009] <i>Note: Antanas Bliudzius and Roma Puronaitė contributed equally to this work.</i>	Publikuota	IF 1.285
International journal of environmental research and public health	Jurevičienė, E., Burneikaitė, G., Dambrauskas, L., Kasiulevičius, V., Kazėnaitė, E., Navickas, R., <b>Puronaitė, R.</b> , Smailytė, G., Visockienė, Ž., & Danila, E. (2022). Epidemiology of chronic obstructive pulmonary disease (COPD) comorbidities in Lithuanian national database: a cluster analysis. <i>International journal of environmental research and public health</i> , 19(2), 1-14. doi:10.3390/ijerph19020970 [DB: MEDLINE; Embase; Scopus; Social Sciences Citation Index (Web of Science); Science Citation Index Expanded (Web of Science)] [Indėlis: 0,067] [Citav. rod.: 3.390 (2020, SCIE); 3.390 (2020, SSCI)] [M.kr.: M 001]	Publikuota	IF 3.390

# Ataskaita už III mokslo metus

## Planas:

Darbo pavadinimas	Atlikimo terminai	Pastabos
<p>1 Mokslinių tyrimų disertacijos tema apžvalga ir analizė (Lietuvoje ir užsienyje):</p> <p>Anotuotos bibliografijos sudarymas. Mokslinės literatūros apžvalga. Egzistuojančių metodų taikymo medicinos duomenimis analizavimas.</p>	<p>2017 m. spalio mėn. – 2018 m. birželio mėn.</p>	<p>Naudojant sisteminės literatūros apžvalgos metodą, atlikta literatūros analizė šiomis temomis:</p> <p>Sisteminė literatūros apžvalga – poliligtumo modelių klasterizavimo metodai (išnagrinėta 151 publikacija, identifiukuoti šioje tematikoje naudojami klasterizavimo metodai) ir</p> <p><b>Sisteminė literatūros apžvalga – netiesiniai statistikos modeliai: poliligtumo analizės tyrimai naudojant didelės apimties medicinos duomenų masyvus (išnagrinėta 206 publikacijos, identifiukuoti šioje tematikoje naudojami klasterizavimo metodai)</b></p> <p><b>Apžvelgtos antriniu tikslu naudojamos administracinės duomenų bazės analizės strategijos.</b></p>
<p>2 Mokslinio tyrimo vykdymas:</p> <p>2.1. Tyrimo metodikos sudarymas: 1. Disertacijos tikslo formulavimas. 2. Disertacijos uždavinių formulavimas.</p>	<p>2018 m. birželio mėn. – 2018 m. lapkričio mėn.</p>	<p>Suformuluotas tikslas ir uždaviniai</p>
<p>2.2. Teorinis tyrimas: Matematinio modelio sudarymas. Algoritmų konstravimas ir tobulinimas.</p>	<p>2018 m. lapkričio mėn. – 2019 m. birželio mėn.</p>	<p>Nagrinėjami ir išbandomi nagrinėtoje literatūroje pasiūlyti metodai atsižvelgiantys į duomenų dinamiką.</p> <p>Empiriniui tyrimui atrinkta:</p> <ul style="list-style-type: none"> <li>• Skerspjūvio tipo analizė</li> <li>• Temporalinio fenotipavimo metodai (PARAFAC2, neneigiamas matricos faktorizavimas)</li> <li>• Laiko dimensijos transformavimas (siūloma transformacijai naudoti Poincare grafikus (indeksus), laiko eilučių išlyginimą, netiesinius statistinius modelius)</li> <li>• Kiti: gilusis mokymas, gilieji rekurentiniai neuroniniai tinklai, arba giliojo mokymo teksto analizės metodai.</li> </ul>



# Ataskaita už III mokslo metus

## Planas:

<p>1. Siūlomų algoritmų pritaikymas medicinos duomenims. 2. Siūlomų algoritmų tobulinimas, atsižvelgiant į gautus rezultatus.</p>	<p><b>2018 m. birželio mėn. – 2022 m. vasario mėn.*</b></p>	<ul style="list-style-type: none"> <li>• Diabetas</li> <li>• Lėtinė obstrukcinė plaučių liga</li> <li>• Depresija ir nerimas</li> </ul> <p>Kiekvieno paciento atveju buvo sudaryti dvinariai kintamieji atitinkantys ligos (išskirta 30 lėtinių ligų) buvimą arba nebuvimą atitinkamo paciento ligos istorijoje per tyrimo laikotarpį. Šiuose pogrupuose pritaikius hierarchinį klasterizavimą (angl. HCA) gauti ligų klasteriai, kurie buvo vertinti gydytojų specialistų, dalis šių klasterių įvertinti kaip turintys klinikinį paaiškinimą. Taip pat pritaikyta tiriamoji faktorinė analizė (angl. EFA). Koreliacijos matrica sudaryta naudojant tetrachorinę koreliaciją tarp dvinarių kintamųjų porų.</p> <p>Atlikta vizuali daugiamačių polilogotumo duomenų analizė</p> <p>Pradėta:</p> <ul style="list-style-type: none"> <li>• Temporalinio fenotipavimo metodai (PARAFAC2, neneigiamas matricos faktorizavimas)</li> </ul> <p>Numatoma:</p> <ul style="list-style-type: none"> <li>• Laiko dimensijos transformavimas (siūloma transformacijai naudoti Poincare grafikus (indeksus), laiko eilučių išlyginimą, netiesinius statistinius modelius)</li> <li>• Kiti: gilusis mokymas, gilieji rekurentiniai neuroniniai tinklai, arba giliojo mokymo teksto analizės metodai.</li> </ul>
<p>2.4. Gautų duomenų analizė, apibendrinimas, išvadų parengimas:</p> <p>1. Algoritmų tikslumo įvertinimas, palyginimas su kitų autorių metodais, atrinktais remiantis išanalizuota moksline literatūra. 2. Gautų rezultatų apibendrinimas. 3. Išvadų parengimas.</p>	<p>2020 m. vasario mėn. – 2020 m. liepos mėn.</p> <p><b>2022 m. vasario mėn. – 2022 m. liepos mėn.*</b></p>	<p><b>2021/2022 metais I pusmetis</b></p> <p>Planuojamos publikacijos (pradėta rengti): Depresijos ir nerimo prognozavimas remiantis administracinio pobūdžio duomenimis, pritaikant Poincare plot metodą ir apskaičiuotus indeksus (Health Informatics tematikos žurnale)</p> <p>Pranešimas konferencijoje (pimta): <b>Puronaite, Roma</b>, Ramanauskaitė, Dovilė, Burneikaitė, Greta, Švaikevičienė, Kristina, Švareikaitė, Alicija, Vaitkute, Samanta, Jakaitienė, Audronė, Dambrauskas, Laimis, Jurevičienė, Elena, Trinkūnas, Justas, Kasilevičius, Vytautas, Kazėnaitė, Edita, „Challenges of modeling depression and anxiety risk using data from large healthcare databases: systematic review and situation analysis“, <i>31st International Biometric Conference</i>, 2022 m. liepos 10-15 d., Ryga, Latvija</p> <p>Pranešimas pristatomas 2021/2022 metais II pusmetis</p>

# Ataskaita už III mokslo metus

## 1 Planas:

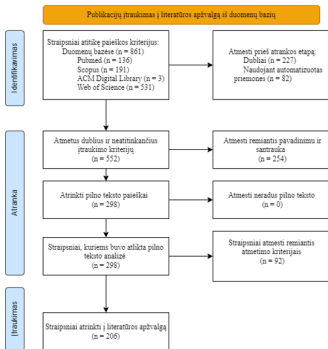
		Planuojama: Pranešimas tarptautinėje statistikos konferencijoje: ISCB2023 (Rezultatų apibendrinimas) 2022/2023 metais II pusmetis
3	<p>Atskirų daktaro disertacijos dalių (tyrimo metodikos, rezultatų, ginamų teiginių, išvadų, ir kt.) parengimas:</p> <ol style="list-style-type: none"> <li>1. Tyrimų apžvalga ir analizė.</li> <li>2. Tyrimo metodikos sudarymas.</li> <li>3. Teorinis tyrimas.</li> <li>4. Empirinis tyrimas.</li> <li>5. Gautų duomenų analizė, apibendrinimas.</li> <li>6. Išvados, įvadas, literatūros sąrašas.</li> </ol>	<p>2020 m. rugsėjo mėn. – 2021 m. gegužės mėn.</p> <p><b>2022 m. rugsėjo mėn. – 2023 m. gegužės mėn.*</b></p> <p>2021/2022 metais I pusmetis Pradėtos rengti dalys: mokslinės literatūros apžvalga, tyrimo metodika, duomenų analizė.</p>
4	Daktaro disertacijos parengimas ir svarstymas padalinyje	<p>2021 m. birželio mėn.</p> <p><b>2023 m. birželio mėn.*</b></p>
5	Daktaro disertacijos gynimas	<p>2021 m. rugsėjo mėn.</p> <p><b>2023 m. rugsėjo mėn.*</b></p>

# Ataskaita už III mokslo metus

Literatūros apžvalga papildyta nauja dalimi - Netiesiniai statistikos modeliai: poliligotumo analizės tyrimai naudojant didelės apimties medicinos duomenų masyvus

Sisteminės apžvalgos pagrindinės komponentės:

- 1 Poliligotumas – kelios vienu metu egzistuojančios lėtinės ligos.
- 2 Netiesiniai modeliai.
- 3 Antrinis (pakartotinis) sveikatos duomenų naudojimas.



## Ataskaita už III mokslo metus

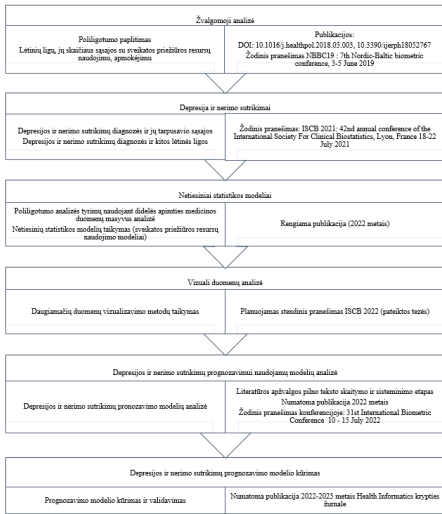
- 1 Dažniausiai naudojami pakartotinai naudojamų (ne pagal pirminę paskirtį) duomenų šaltiniai, kai analizei buvo naudoti netiesiniai statistiniai modeliai: administraciniai, elektroniniai sveikatos įrašai, sveikatos tyrimų (statistinių) duomenys.
- 2 Dažniausiai naudoti netiesiniai statistikos modeliai: logistinė regresija, Puasono, neigiama binominė regresijos.
- 3 Nėra vieningo poliligtumo apibrėžimo: poliligtumas vertinamas pagal ligų skaičių, pagal ligų grupes, naudojant poliligtumo indeksus.
- 4 Dažniausios temos: poliligtumas ir sveikatos priežiūros resursų naudojimas (išlaidos, paslaugų dažnis) - 15 proc., poliligtumo tyrimai (poliligtumo struktūra, klasteriai) - 10,7 proc., poliligtumas ir socialiniai netolygumai (pagal gyvenamąją vietą, socialinį statusą) - 9,2 proc.

# Ataskaita už III mokslo metus

## Duomenų parengimas analizei ir klasterizavimo algoritmai

Skerspjuvio tipo	Temporalinis fenotipavimas (fenotipavimas laike)	Laiko dimensijos transformacija	Kiti būdai
<ul style="list-style-type: none"> <li>Būsena kintamiesiems įvertinama apimant visą laikotarpį, pvz. pacientas turėjo x diagnozę tarp 2012-2014 metų.</li> <li>Pirminis objektas matrica:               <ul style="list-style-type: none"> <li>Pacientai x Diagnozės</li> </ul> </li> <li>Pvz. hierarchinis klasterizavimas, tiriamoji faktorinė analizė.</li> </ul>	<ul style="list-style-type: none"> <li>Būsena gali kisti, atsižvelgiama į laiko žymę, eilės tvarką, pvz. pacientas turėjo x diagnozę 2012, 2014 metais, bet neturėjo 2013.</li> <li>Pirminis objektas tenzorius:               <ul style="list-style-type: none"> <li>Pacientai x Diagnozės x Vizitai (laiko žymė)</li> </ul> </li> <li>Pvz. tenzorių faktorizacija (PARAFAC2 algoritmas ir kt.)</li> </ul>	<ul style="list-style-type: none"> <li>Būsenos kaita įvertinama atliekant transformaciją</li> <li>Pirminis objektas tenzorius:               <ul style="list-style-type: none"> <li>Pacientai</li> <li>Diagnozės</li> <li>Vizitas (laiko žymė)</li> </ul> </li> <li>Po transformacijos klasterizavimas atliekamas matricoms</li> <li>Matrica:               <ul style="list-style-type: none"> <li>pvz. Pacientai x Diagnozės</li> </ul> </li> <li>Laiko eilučių išlyginimo ir panašumo vertinimas</li> <li>Pvz. dinaminis laiko skalės kraipymas -&gt; klasterizavimo algoritmai (HCA, EFA)</li> </ul>	<ul style="list-style-type: none"> <li>Gilieji neuroniniai tinklai</li> <li>Gilūs mokymai</li> <li>Teksto analizė</li> <li>Kryptiniai grafai</li> <li>Jungtinis klasterizavimas</li> </ul>

# Mokslinio tyrimas "Lėtinių neinfekcinių ligų paplitimo tarpusavio sąveikos, sveikatos priežiūros paslaugų bei vaistų vartojimo bei klinikinių baigčių vertinimas Lietuvoje".





# III (2021 / 2022) mokslo metų II pusmečio planas



- 1 Egzaminas: Daugiamačių duomenų vizualizavimo metodai (2022 m. kovo 30 d.)
- 2 Publikacija sisteminės apžvalgos "Netiesiniai statistikos modeliai: poliligtumo analizės tyrimai naudojant didelės apimties medicinos duomenų masyvus" pagrindu
- 3 Publikacija (Darbinis pavadinimas: "Multimorbidity, depression and anxiety patterns: temporal phenotyping approach"- statusas parengta)