

UNIVERZA V LJUBLJANI  
FAKULTETA ZA RAČUNALNIŠTVO IN INFORMATIKO

Nejc Mušič

**Razlaga klasifikatorjev na podlagi  
podkonceptov**

DIPLOMSKO DELO

UNIVERZITETNI ŠTUDIJSKI PROGRAM  
PRVE STOPNJE  
RAČUNALNIŠTVO IN INFORMATIKA

MENTOR: prof. dr. Marko Robnik Šikonja

Ljubljana, 2023

To delo je ponujeno pod licenco *Creative Commons Priznanje avtorstva-Deljenje pod enakimi pogoji 2.5 Slovenija* (ali novejšo različico). To pomeni, da se tako besedilo, slike, grafi in druge sestavine dela kot tudi rezultati diplomskega dela lahko prosto distribuirajo, reproducirajo, uporabljajo, priobčujejo javnosti in predelujejo, pod pogojem, da se jasno in vidno navede avtorja in naslov tega dela in da se v primeru spremembe, preoblikovanja ali uporabe tega dela v svojem delu, lahko distribuira predelava le pod licenco, ki je enaka tej. Podrobnosti licence so dostopne na spletni strani [creativecommons.si](http://creativecommons.si) ali na Inštitutu za intelektualno lastnino, Streliška 1, 1000 Ljubljana.



Izvorna koda diplomskega dela, njeni rezultati in v ta namen razvita programska oprema je ponujena pod licenco GNU General Public License, različica 3 (ali novejša). To pomeni, da se lahko prosto distribuira in/ali predeluje pod njenimi pogoji. Podrobnosti licence so dostopne na spletni strani <http://www.gnu.org/licenses/>.

*Besedilo je oblikovano z urejevalnikom besedil L<sup>A</sup>T<sub>E</sub>X.*

**Kandidat:** Nejc Mušič

**Naslov:** Razlaga klasifikacijskih problemov na podlagi podkonceptov

**Vrsta naloge:** Diplomaska naloga na univerzitetnem programu prve stopnje  
Računalništvo in informatika

**Mentor:** prof. dr. Marko Robnik Šikonja

**Opis:**

Na nekaterih področjih je uporaba napovednih modelov strojnega učenja pogojena z možnostjo razlage napovedi. Večina obstoječih pristopov razlage tvori razlage posameznih napovedi v obliki faktorjev pomembnosti vhodnih atributov. Ta način ni vedno primeren, saj je prostor vhodnih atributov mnogokrat zelo velik ali človeku nerazumljiv. Preizkusite alternativno idejo razlage, ki iz učnih primerov danega razreda najprej tvori gruče, ki jih nato skuša interpretirati kot podkoncepte danega razreda z uporabo pravil in prototipov. Pristop preizkusite na nekaj umetnih in realnem problemu.

**Title:** Explanation of Classifiers Based on Subconcepts

**Description:**



# Kazalo

Povzetek

Abstract

<b>1</b>	<b>Uvod</b>	<b>1</b>
<b>2</b>	<b>Predstavitev uporabljenih tehnologij</b>	<b>5</b>
2.1	Algoritmi za gručenje . . . . .	5
2.2	Hevristike za oceno gručenj . . . . .	7
2.3	Klasifikatorji . . . . .	8
2.4	Metode razlage . . . . .	9
2.5	Zmanjšanje dimenzionalnosti . . . . .	10
<b>3</b>	<b>Razlaga s podskupinami</b>	<b>11</b>
<b>4</b>	<b>Podatkovne množice</b>	<b>17</b>
4.1	Umetne podatkovne množice . . . . .	17
4.2	Realne podatkovne množice . . . . .	20
<b>5</b>	<b>Primerjava algoritmov za gručenje</b>	<b>23</b>
5.1	Krožne gruče . . . . .	24
5.2	Trakovi 4-3 . . . . .	26
5.3	MNIST . . . . .	28
5.4	Diskusija . . . . .	31

<b>6</b>	<b>Razlaga realne podatkovne množice</b>	<b>33</b>
<b>7</b>	<b>Diskusija</b>	<b>43</b>
7.1	Število gruč . . . . .	43
7.2	Smiselnost prototipov . . . . .	44
7.3	Smiselnost pravil . . . . .	44
<b>8</b>	<b>Zaključki</b>	<b>47</b>
<b>A</b>	<b>Celotna primerjava gručenj</b>	<b>49</b>
A.1	Krožne gruče . . . . .	49
A.2	Trakovi 4-3 . . . . .	50
A.3	MNIST . . . . .	51
<b>B</b>	<b>Vrednosti SHAP preostalih gruč iz poglavja 6</b>	<b>55</b>
	<b>Literatura</b>	<b>59</b>



# Seznam uporabljenih kratic

kratica	angleško	slovensko
<b>MDEC</b>	Multidiversified Ensemble Clustering	Večdimenzionalno multidiversificirano združevanje v gruče
<b>XAI</b>	explainable artificial intelligence	razložljiva umetna inteligenca
<b>SHAP</b>	SHapley Additive exPlanations	Shapleyjeve aditivne razlage
<b>DBSCAN</b>	Density-Based Spatial Clustering of Applications with Noise	Gručenje prostorskih podatkov s šumom glede na gostoto
<b>HDBSCAN</b>	Hierarchical Density-Based Spatial Clustering of Applications with Noise	Hierarhično gručenje prostorskih podatkov s šumom glede na gostoto
<b>XGBoost</b>	Extreme Gradient Boosting	Ekstremni gradientni boosting
<b>DBCV</b>	Density-Based Clustering Validation	Preverjanje gručenja glede na gostoto
<b>t-SNE</b>	t-distributed Stochastic Neighbor Embedding	Stohastična vložitev sosedov s t-porazdelitvijo
<b>PCA</b>	Principal component analysis	Analiza glavnih komponent
<b>UMAP</b>	Uniform Manifold Approximation and Projection	Uniformna aproksimacija in projekcija
<b>MNIST</b>	Modified National Institute of Standards and Technology	Modificirana zbirka podatkov Nacionalnega inštituta za standarde in tehnologijo



# Povzetek

**Naslov:** Razlaga klasifikatorjev na podlagi podkonceptov

**Avtor:** Nejc Mušič

Pri nekaterih klasifikacijskih problemih je poleg natančnosti in zanesljivosti pomembna tudi razlaga odločitev (npr. v medicini). Ker dajejo kompleksni modeli, na primer nevronske mreže boljše rezultate, jih velikokrat raje uporabimo namesto preprostih, dobro razumljivih modelov. Diplomaska naloga poskuša kompleksne modele razložiti na podlagi gručenja podatkov v kombinaciji z odločitvenimi pravili, medoidi gruč in z vrednostmi SHAP. Preizkusimo nekaj algoritmov gručenja na nizko in visoko dimenzionalnih podatkovnih množicah, kjer pokažemo dobre in slabe lastnosti ter ustreznost glede na podatke. Predlagana metoda razlage potrebuje razdelitev primerov istega razreda na več gruč, kar se je pri realnih podatkovnih množicah izkazalo kot redko. Z eksperimenti ugotovimo, da je naša razlaga zahtevna za interpretacijo in ni povsem avtomatizirana. Za nekatere komponente naše razlage lahko vizualiziramo rezultate, kar pomaga pri interpretaciji.

**Ključne besede:** razložljiva umetna inteligenca, MDEC, DBSCAN, HDBSCAN, K-MEANS, SHAP, odločitvena pravila, medoid, prototip.



# Abstract

**Title:** Explanation of Classifiers based on Subconcepts

**Author:** Nejc Mušič

In some classification problems, in addition to accuracy and reliability, the explanation of decisions is also important (e.g., in medicine). Complex models, such as neural networks, often outperform simpler models like linear regression, which are inherently interpretable. The thesis addresses the challenge of explaining such complex models. Based on data clustering, we explain the classifier using decision rules, cluster medoids, and SHAP values. We test several clustering algorithms using low and high-dimensional datasets. The experiments aim to highlight the strengths and weaknesses of the algorithms, as well as their suitability with respect to specific datasets. Our explanation method relies on separability of instances from the same class into clean clusters, which is infrequent in real-world datasets. Through experiments, we discover that our explanation method may be challenging to interpret and lacks full automation. Some explanation components provide visualizations, which help the interpretation.

**Keywords:** XAI, MDEC, DBSCAN, HDBSCAN, K-MEANS, SHAP, decision rules, medoid, prototype.



# Poglavje 1

## Uvod

V zadnjem času smo priča izjemnim dosežkom na področju strojnega učenja. Kompleksni modeli pri klasifikacijskih problemih omogočajo dobro natančnost pri napovedovanju. Slaba lastnost t.i. "black box" modelov je nerazumljivost, predvsem kako in zakaj so sprejeli določene odločitve. Kot primer lahko vzamemo globoke nevronske mreže ali kompleksne ansamble (XGBoost, ekstremni gradientni boosting), ki imajo na tisoče ali celo milijone parametrov. To pomeni, da so njihovi odločitveni procesi izjemno zapleteni, kar otežuje razumevanje, kako določen vhod vpliva na izhod. Z diplomsko nalogo želimo prispevati k razvoju razlage kompleksnih modelov ter omogočiti boljše razumevanje njihovega delovanja, kar bo pripomoglo k dodatni izboljšavi modelov. Prav tako je razumevanje odločitev modelov ključno za zaupanje v njihove rezultate ter za širšo sprejetost uporabe umetne inteligence v različnih kritičnih domenah (npr. v medicini, varnosti računalniških omrežij, letalski in avtomobilski industriji).

V diplomski nalogi je preizkušen postopek razlage, kjer je klasifikator naučen na učni množici, na kateri poženemo tudi algoritme za gručenje. Za avtomatski izbor parametrov uporabimo hevrstiko silhuete in indeks DBCV (preverjanje gručenja glede na gostoto). Pri izboru najboljšega gručenja si pomagamo z vizualizacijami podatkov v dvodimezijskem prostoru, za kar uporabimo algoritem t-SNE (*slo.* stohastična vložitev sosedov s t-porazdelitvijo),

pri visokodimenzionalnih podatkih). Predpogoj metode razlage je ločenost podatkov na gruče, ki so medsebojno dobro ločljive in prisotnost posameznega razreda v dveh ali več različnih gručah. Nato posamezne gruče razložimo z odločitvenimi pravili po konceptu ena gruča proti vsem ostalim. Preverimo prisotnost razredov v gruči. V primeru, da je gruča homogena (primeri pripadajo enemu razreda), lahko gručo predstavimo z medoidom. V primeru, da je gruča heterogena, znotraj gruče uporabimo odločitvena pravila, ki ločijo primere različnih razredov. Dodatno izračunamo vrednosti SHAP za vsako gručo in jih predstavimo z grafom. Vrednosti SHAP nam razložijo pomembnost atributov za naš klasifikator, kar dodatno pomaga pri interpretaciji odločitvenih pravil in medoidov.

Nov primer z algoritmom za gručenje umestimo v najbližjo gručo in s tem razložimo z zgornjimi metodami. Dodatno izračunamo vrednosti SHAP za novi primer.

V zadnjih letih so bili razviti številni pristopi za razlago kompleksnih modelov. Pristop s CBR (sklepanje iz primerov) najde pomembne značilke (atributi) na podlagi lokalnih informacij in nato izbere primer, ki je bil pomemben za izgradnjo modela kot primer razlage [18]. Z odločitvenimi drevesi in odločitvenimi pravili, dobimo vpogled v odločitvene meje v podatkih in poenostavljeno predstavitev množice podatkov [1]. Vrednosti SHAP veljajo za dobro metodo razlage pomembnosti atributov, zato so bili tudi uporabljeni na različnih področjih, kot je na primer obdelava naravnega jezika (*ang.* Natural Language Processing) [14]. Še nekaj sorodnih del naštejemo v podpoglavju 4.6 pri primerjavi z obstoječimi razlagami.

Diplomsko nalogo razdelimo na osem poglavij in dva dodatka. Poleg že predstavljenega uvoda v poglavju 1 sledi predstavitev uporabljenih tehnologij in izbor parametrov za algoritme v poglavju 2. V poglavju 3 opišemo pristop razlage s podskupinami in predstavimo preprost ilustrativen primer razlage. V poglavju 4 predstavimo podatkovne množice, v poglavju 5 pa primerjavo algoritmov gručenja na umetnih podatkovnih množicah in realni podatkovni množici MNIST. V poglavju 6 sledi razlaga realne podatkovne množice

KDD99 in njene razlage. V poglavju 7 sledi diskusija o smiselnosti prototipov (medoidov), o pravilih, številu gruč in uspešnosti razlage. Diplomaska naloga se konča s poglavjem 8, kjer na kratko povzamemo rezultate. Izpostavimo dobre in slabe lastnosti pristopa ter predstavimo nekaj idej za nadaljne raziskave. V dodatku A podamo celotno primerjavo gruč, v dodatku B pa vrednosti SHAP preostalih gruč iz poglavja 6.





## Poglavje 2

# Predstavitev uporabljenih tehnologij

V tem poglavju predstavimo algoritme za gručenje (podpoglavje 2.1), heuristike za oceno gručenj (podpoglavje 2.2), uporabljene klasifikatorje (podpoglavje 2.3), metode razlage (podpoglavje 2.4) in tehniko zmanjšanja dimenzionalnosti (podpoglavje 2.5). Vse našete tehnologije uporabimo pri razlagah klasifikatorjev.

### 2.1 Algoritmi za gručenje

V tem podpoglavju predstavimo algoritme gručenja MDEC (Več dimenzionalno multidiversificirano združevanje v gručice), K-MEANS, DBSCAN (Gručenje prostorskih podatkov s šumom glede na gostoto) in HDBSCAN (Hierarhično gručenje prostorskih podatkov s šumom glede na gostoto). Izбира algoritma gručenja je odvisna od podatkov, zato smo izbrali algoritme z različnimi lastnostmi. Algoritem MDEC vrne dobre rezultate na visoko dimenzionalnih podatkih, K-MEANS na sferičnih podatkih, algoritem DBSCAN na poljubnih oblikah podatkov, kjer lahko nekatere točke izpustimo. Predstavimo tudi HDBSCAN, ki je nadgradnja algoritma DBSCAN, a je bolj stabilen pri gručenjih različnih gostot in samodejno določi nekatere parametre.

### 2.1.1 MDEC

Algoritem "Multidiversified Ensemble Clustering" [7] je zasnovan za gručenje visoko dimenzionalnih podatkov. Združuje več algoritmov gručenja za izboljšanje splošne uspešnosti gručenja. Algoritem MDEC ustvari veliko število raznolikih metrik z naključno skalirano eksponentno funkcijo podobnosti in jih poveže z naključnimi podprostor. S tem dobimo pare metrika-podprostor. Na osnovi matrik podobnosti, pridobljenih iz metrika-podprostor parov, nato sestavi ansambel raznolikih gručenj. Zatem uporabi kriterij, zasnovan na entropiji, za oceno raznolikosti združevanj glede na gruč, na osnovi katerega uporabi tri specifične algoritme ansambelskega združevanja (*ang.* consensus functions). Pri algoritmu smo opredelili le parameter za število skupin.

### 2.1.2 K-MEANS

K-MEANS algoritem [6] je iterativni algoritem, ki poskuša razdeliti nabor podatkov v  $K$  predhodno določenih neprekrivajočih se podskupin (gruč), pri čemer vsaka podatkovna točka pripada le eni skupini. Algoritem iterativno dodeljuje primere najbližjemu centroidu in zatem premika centroe. Pri algoritmu smo opredelili le parameter za število skupin.

### 2.1.3 DBSCAN

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) [8] je algoritem za gručenje, ki temelji na osnovi gostote podatkovnih točk v prostoru. Algoritem je sposoben zaznati gruč v podatkih, ki imajo različne oblike ter hkrati zaznati osamelce oziroma šum. Pri uporabi algoritma je ključna izbira dveh parametrov, t.i. epsilon in min samples. Parameter epsilon določa največjo razdaljo med dvema podatkovnima točkama, da se štejeta kot del iste skupine, celoštevilski parameter min samples pa določa najmanjše število podatkovnih točk, ki so vse znotraj razdalje epsilon druga od druge, in tvorijo osrednjo skupino (gruča mora imeti vsaj min samples osrednjih točk, da se šteje za veljavno). Vzamemo predlagano vrednost min

$\text{samples} = 2 * \text{dim}$ , kjer je  $\text{dim}$  dimenzija podatkovne množice. Za vrednost parametra  $\epsilon$  je predlagana tehnika, ki izračuna povprečno razdaljo med vsako točko in njenimi najbližjimi sosedami ( $k$ ), katerih vrednost je enaka  $\text{min samples}$ . Povprečne  $k$ -razdalje nato prikažemo v naraščajočem vrstnem redu na grafu. Optimalna vrednost za  $\epsilon$  se nahaja pri točki največje ukrivljenosti (tj. največji nagib na grafu) [16]. Največji nagib na grafu smo našli s pomočjo algoritma Kneed [9]. Pri obeh parametrih smo poskusili tudi okoliške vrednosti, ki se nahajajo okoli izbranih vrednosti parametrov  $\text{min samples}$  in  $\epsilon$ .

#### 2.1.4 HDBSCAN

HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) [13] izvaja algoritem DBSCAN pri različnih vrednostih  $\epsilon$  in združi rezultate, da najde gručenje, ki zagotavlja najboljšo stabilnost pri različnih vrednostih  $\epsilon$ . To omogoča HDBSCAN-u, da najde gručice z različnimi gostotami (za razliko od DBSCAN-a) in je bolj odporen na izbiro parametrov [13]. Pri HDBSCAN algoritmu sta pomembna parametra  $\text{min cluster size}$  in  $\text{min samples}$  [4]. Določili smo ju kot:  $\text{min cluster size} = (\text{ix1} * 5)$  in  $\text{min samples} = (\text{ix2} + 2) * 5$ , kjer  $\text{ix1} \in [2, 10]$  in  $\text{ix2} \in [0, 8]$  in  $\text{ix1}, \text{ix2} \in \mathbb{N}$ .

## 2.2 Hevristike za oceno gručenj

V tem razdelku predstavimo hevristiki za oceno gručenja koeficient silhuete in indeks DBCV. Hevristiko silhuete uporabimo za ocene gručenja pri algoritmih, ki temeljijo na razdalji med podatki, indeks DBCV pa pri algoritmih, ki temeljijo na gostoti podatkov.

### 2.2.1 Koeficient silhuete

Koeficient silhuete [19] je metoda za ocenjevanje kakovosti gručenja, ki se osredotoča na merjenje podobnosti točk v isti gruči v primerjavi z drugimi

gručami. Njegova glavna naloga je ponuditi kvantitativno merilo za oceno, kako dobro so točke znotraj iste gruče povezane in kako dobro so razmejene od drugih gruč. Za oceno smo uporabili hevrstiko povprečne silhuete, ki se ne osredotoča na oceno kakovosti posameznih točk, kot to počne koeficient silhuete za vsako točko, temveč daje celostno sliko o rezultatu gručenja.

## 2.2.2 Indeks DBCV

Indeks DBCV (Density-Based Clustering Validation) [15] je mera kakovosti gručenja, ki se uporablja za ocenjevanje učinkovitosti algoritmov za gručenje. Mera kot je koeficient silhuete je primerna za ocenjevanje gruč s sferičnimi oblikami, medtem ko algoritmi, ki temeljijo na gostoti podatkovnih točk, velikokrat zaznajo nepravilne oblike, kjer se bolje odreže indeks DBCV. Izračun indeksa DBCV tudi temelji na gostoti gruč.

## 2.3 Klasifikatorji

V tem razdelku predstavimo uporabljena klasifikatorja XGBoost in logistično regresijo.

### 2.3.1 XGBoost

XGBoost (Extreme Gradient Boosting) [17] je zmogljiva in priljubljena metoda strojnega učenja. Gre za algoritem, ki kombinira moč odločitvenih dreves in tehnike gradientnega spusta za dosego visoke točnosti. XGBoost razložimo s konceptoma bagging in gradient boosting. Pri bagging-u so osnovni klasifikatorji naučeni na naključnih podmnožicah originalnih podatkov. Končni rezultat je pridobljen z glasovanjem ali povprečenjem. Boosting je tehnika, kjer zgradimo močan klasifikator s pomočjo množice šibkih. Pri tehniki gradient boost klasifikator, ki je naslednji v vrsti klasifikatorjev, popravlja napake svojega prednika. Pri XGBoost se uteži dodelijo vsem atributom, ki se nato vstavijo v odločitveno drevo. Nato odločitveno drevo napove

rezultat. Utež spremenljivk, ki jih drevo napove napačno, se poveča. Te spremenljivke se zatem vstavijo v naslednje odločitveno drevo. Posamezne klasifikatorje na koncu združimo, da dobimo točnejši model.

### 2.3.2 Logistična regresija

Logistična regresija [3] se pogosto uporablja v strojnem učenju za reševanje klasifikacijskih problemov. Metoda napoveduje verjetnosti, da bo določen vhodni primer spadal v enega izmed razredov. Model logistične regresije temelji na linearni kombinaciji neodvisnih spremenljivk z utežmi (koeficienti) ter uporabi logistične funkcije za pretvorbo rezultata v verjetnost. Med učenjem se koeficienti prilagajajo tako, da optimiziramo logaritem verjetja z optimizacijsko metodo, kot je gradientni sestop.

## 2.4 Metode razlage

V tem razdelku predstavimo uporabljene metode razlage, to so odločitvena pravila, medoidi, vrednosti SHAP in tehnike za zmanjšanje dimenzionalnosti.

### 2.4.1 Odločitvena pravila

Kljub temu da so odločitvena pravila metoda strojnega učenja, smo jih uporabili za razlago gruč, saj so lahko razložljiva z obliko "IF conditions THEN response". Metoda, ki smo jo uporabili [5], iz ansambla odločitvenih dreves izvleče pravila z največjima vrednostima "Precision" in "Recall". Najboljša pravila smo nato uporabili za opis ene gruče proti ostalim. Pravila so uporabna tudi za ločitev primerov, ki pripadajo različnim razredom v eni gruči (medoid težko pojasni tako gručo).

### 2.4.2 Medoid

Za razlago posamezne gruče smo uporabili medoid, ki je izračunan na atributih primerov iz gruče. Gre za tipičen primer, ki enostavno opiše gručo. Pri

uporabi razlage z medoidom moramo biti pazljivi na sestavo gruče. Gruča, ki je sestavljena iz dveh ali več razredov, ne bo dobro razložena z medoidom (tu si pomagamo z odločitvenimi pravili znotraj gruče).

### 2.4.3 Algoritem SHAP

SHAP (SHapley Additive exPlanations) [11] je metoda, ki se uporablja za razumevanje in razlaganje pomembnosti posameznih značilk v modelih strojnega učenja. Njeno ime izhaja iz koncepta Shapleyevih vrednosti v teoriji iger, kjer se meri prispevek vsakega igralca k skupni vrednosti. Vrednosti SHAP bomo za vsako gručo posebej predstavili v obliki grafa. Z metodo SHAP bomo za izbran klasifikator izvedeli pomembnost atributov, kar nam bo pomagalo pri interpretaciji pravil in medoidov in s tem pri razumevanju posamezne gruče.

## 2.5 Zmanjšanje dimenzionalnosti

Za zmanjšanje dimenzionalnosti podatkov smo uporabili t-SNE (t-distributed Stochastic Neighbor Embedding) [22], ki se pogosto uporablja za vizualizacijo visoko dimenzionalnih podatkov v nižje dimenzionalnem prostoru. Nesistemično sta bili preizkušeni še tehniki UMAP, Uniform Manifold Approximation and Projection (*slo.* Uniformna aproksimacija in projekcija) in PCA, Principal component analysis (*slo.* analiza glavnih komponent), a je t-SNE vrnil boljše ločljive gruče, zato smo uporabili le t-SNE. Algoritem t-SNE v visoko dimenzionalnem prostoru izračuna verjetnostne porazdelitve, ki predstavljajo podobnost med primeri. Nato ustvari podobno verjetnostno porazdelitev v nižje dimenzionalnem prostoru in poskuša minimizirati razliko med tema dvema porazdelitvama. t-SNE smo uporabili tudi za predprocesiranje visoko dimenzionalnih podatkov, s čimer smo dobili boljši rezultat gručenja pri uporabi algoritmov K-MEANS, DBSCAN in HDBSCAN.

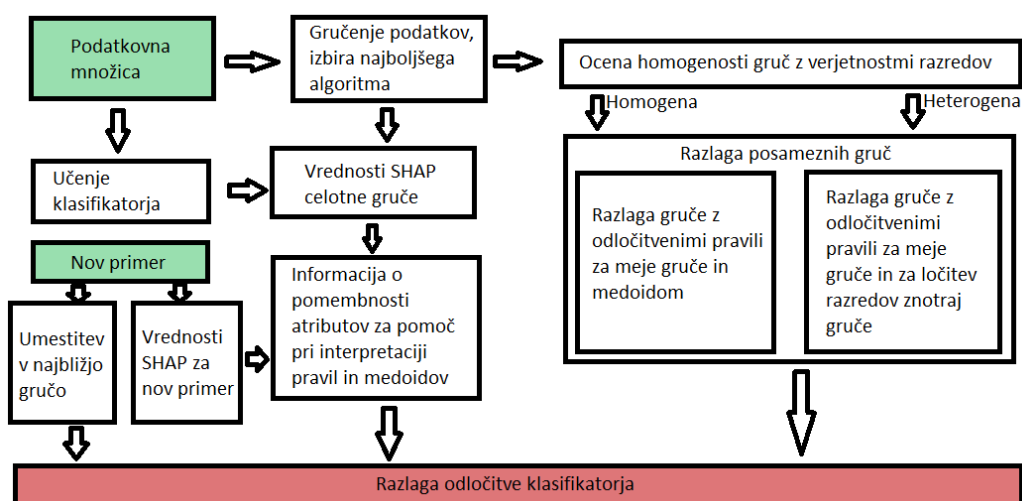
## Poglavje 3

# Razlaga s podskupinami

V tem poglavju predstavimo potek razlage, pri čemer predstavimo komponente. Na sliki 3.1 naše metode podamo diagram poteka in končamo s primerom razlage na umetni podatkovni množici.

Potek razlage opišemo po dveh poteh. Začnemo pri podatkovni množici, na kateri naučimo klasifikator in izvedemo različne algoritme gručenja. Izberemo najboljše gručenje za dano podatkovno množico. Za vsako gručo ocenimo ali je homogena ali heterogena z verjetnostmi razredov v gruči. Če je gruča homogena, uporabimo za razlago gručice odločitvena pravila za meje gručice in medoid kot razlago. V nasprotnem primeru uporabimo odločitvena pravila za meje gručice in odločitvena pravila za ločitev razredov znotraj gručice. S tem dobimo komponento, ki jo potrebujemo za interpretacijo odločitve klasifikatorja. Razlago posamezne gručice razširimo z grafi vrednosti SHAP, ki nam dajejo informacijo o pomembnosti atributov. To informacijo uporabimo pri interpretaciji odločitvenih pravil in medoidov, tako da upoštevamo attribute, ki najbolj vplivajo na klasifikator.

Za razlago novega primera, ga umestimo v najbližjo gručo in izračunamo njegove vrednosti SHAP. Če se vrednosti SHAP primera skladajo z vrednostmi SHAP za gručo, je interpretacija odločitve skladna s celotno gručo. Ko se odločimo o pomembnosti atributov, sledi razlaga odločitve z interpretacijo

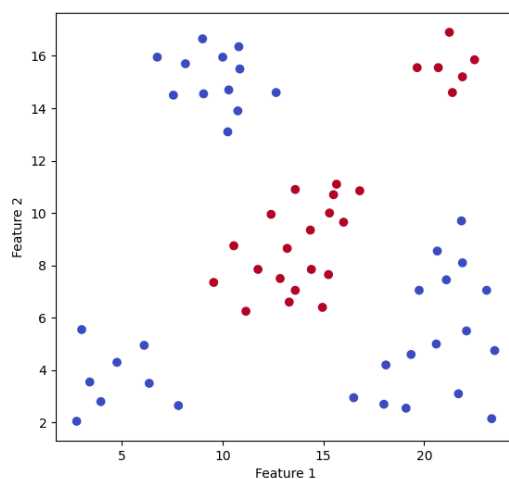


Slika 3.1: Diagram poteka razlage. Začetni stanji sta označeni z zeleno, končno z rdečo.

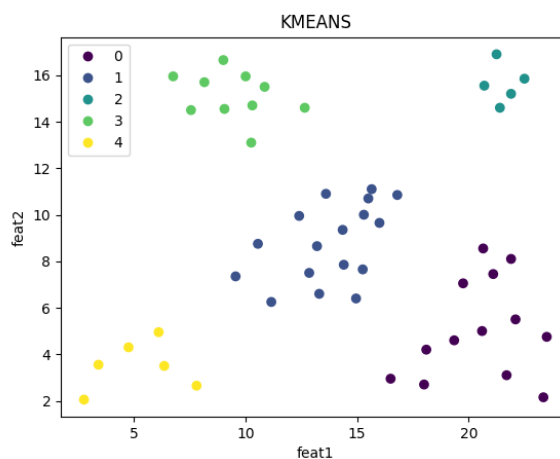
odločitvenih pravil in medoidov (če je gruča homogena). Diagram poteka vidimo na sliki 3.1.

Za lažje razumevanje podamo preprost primer na osnovi umetne množice podatkov, ki je prikazana na sliki 3.2 in jo imenujemo Homogene gruče. Vidimo pet gruč in vsaka od gruč v celoti vsebuje primere enega od dveh razredov, predstavljenih z modro in rdečo barvo. Podatke smo razdelili na učno in testno množico. Klasifikator, naučen na učni množici, je XGBoost in ima klasifikacijsko točnost 0.9 in F1-score 0.9 na testnih podatkih. Sledi gručenje podatkov, s katerim pridobimo podskupine podatkov. Testirali smo štiri algoritme gručenja. Njihova izbira je odvisna od posameznih podatkov. Za naš primer je bil uporabljen algoritem gručenja K-MEANS. Algoritem kot parameter zahteva število skupin. Za avtomatizacijo izbora parametra smo uporabili mero silhuete, ki pove, kako dobro so bili gručeni podatki. Algoritem K-MEANS smo pognali za različna števila skupin in največjo vrednost silhuete (0,601) dobili pri petih skupinah. Rezultat gručenja je viden na sliki 3.3.





Slika 3.2: Ilustracija umetne množice podatkov Homogene gruče, ki vsebuje več podkonceptov. Vidimo podatke dveh razredov predstavljenih z barvami (modra, rdeča) in pet gruč. Imamo dve značilki, ki se razprostirata v dvodimenzionalnem prostoru.



Slika 3.3: Rezultat gručenja na podatkovni množici Homogene gruče z algoritmom K-MEANS. Največja vrednost heuristike silhuete je bila pri petih skupinah.

Na tem mestu je potrebno omeniti, da smo za algoritme gručenja, ki temeljijo na gostoti podatkov uporabili hevrstiko indeks DBCV, ki je analogna meri silhuete, a je namenjena za oceno gručenja teh algoritmov. V dveh dimenzijah je že na videz jasno, kje so meje gruč, a v višjih dimenzijah temu ni tako. Vizualizacije ne povedo vsega, zato uporabimo odločitvena pravila, ki povedo, kako se izbrana gruča loči od preostalih. S tem pridobimo meje posamezne gruče, ki nam kasneje, pri razlagi primerov, pomagajo pri umestitvi primerov v gruče. Vidimo tudi, katere značilke so bolj pomembne za umestitev v gručo, saj so tiste, ki najbolj ustrezajo skupini, med pogoji odločitvenih pravil. Pravila pridobljena na našem primeru so prikazana v tabeli 3.1.

Tabela 3.1: Odločitvena pravila za meje posameznih gruč pri podatkovni množici Homogene gruče z gručenjem K-MEANS.

Gruča	Odločitvena pravila	Precision	Recall
0	feat1 > 16.07 in feat2 <= 8.95	1.0	1.0
1	feat1 <= 19.35 in feat2 <= 12.10 in feat2 > 5.88	1.0	1.0
2	feat1 > 20.68 in feat2 > 11.35	1.0	1.0
3	feat1 <= 16.68 in feat2 > 12.80	1.0	1.0
4	feat1 <= 7.97 in feat2 <= 9.72	1.0	1.0

Tabela 3.2: Medoidi posameznih gruč

Gruča	Feat1	Feat2	Razred
0	20.6	5	0
1	14.35	9.35	1
2	21.9	15.2	1
3	9.05	14.55	0
4	4.75	4.3	0

Gruče razložimo tudi z medoidom, ki predstavlja tipičen primer gruče. Prav

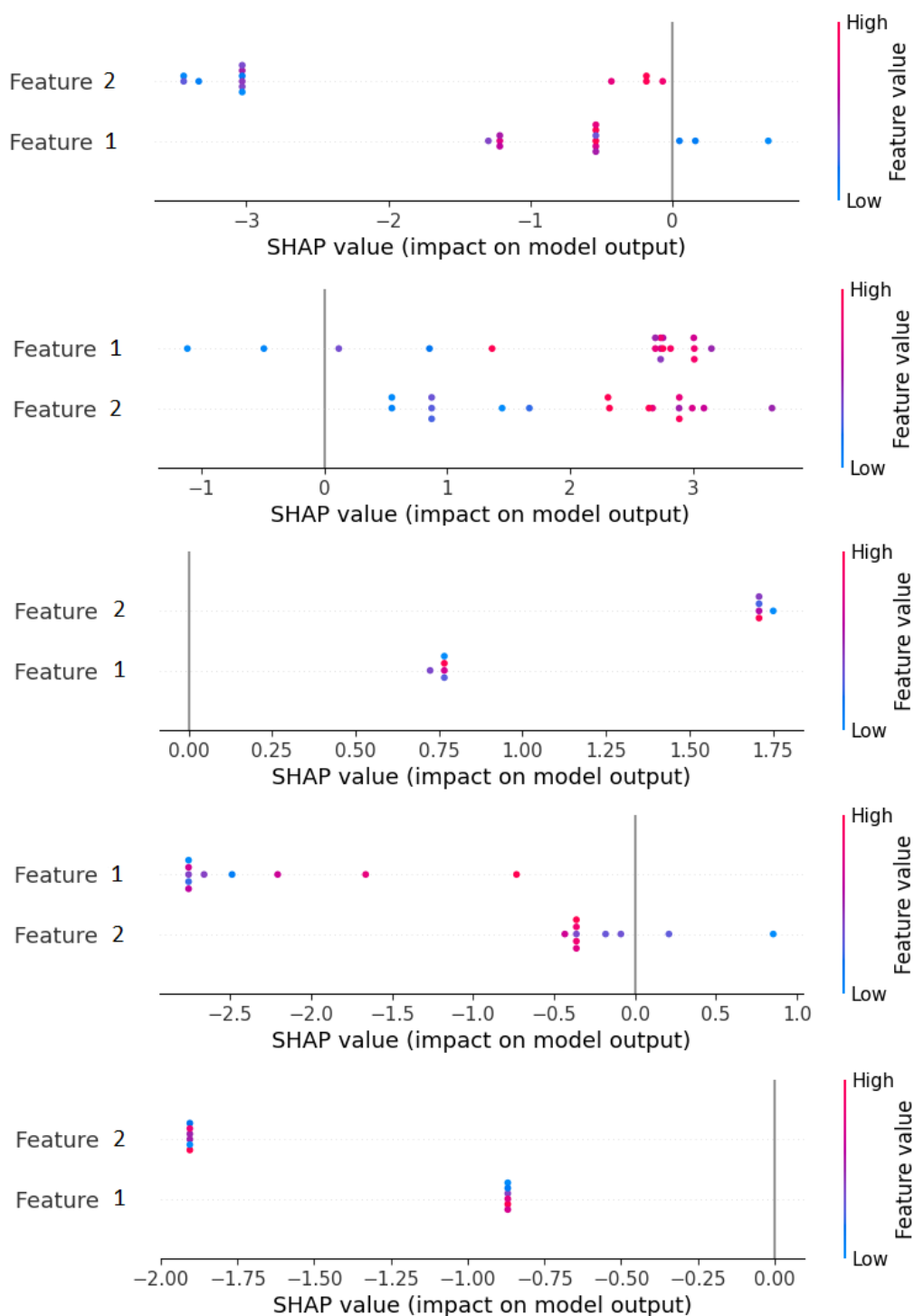
tako si bomo pogledali vrednosti SHAP celotne gruče in novega primera, ki ga vidimo v tabeli 3.3.

Tabela 3.3: Nov, še neviden primer v domeni Homogene gruče.

Feat1	Feat2	Predikcija klasifikatorja
13	9	1

Klasifikator ga je uvrstil v razred 1. Primer lahko po pravilih za meje gruče razvrstimo v gručo 1. Tudi knjižnica za gručenje lahko napove gručo novega primera, ki ga razvrsti v gručo 1. Najprej nas zanimajo verjetnosti razredov v gruči 1. Vidimo, da je gruča v celoti homogena in vsebuje le razred 1. V takem primeru nas zanima medoid, saj dobro opiše gručo. V tabeli 3.2 imamo podane medoide za gruče. Vidimo, da medoid spada v razred 1, kar podpre odločitev klasifikatorja. Pri razlagi nam velikokrat pomaga podatek o pomembnosti značilk, ki ga bomo pridobili z vrednostmi SHAP. Na sliki 3.4 vidimo vrednosti SHAP za posamezne gruče. Nov primer spada v gručo 1, zato analiziramo drugi graf na sliki (od zgoraj navzdol). Vidimo, da oba atributa pripomoreta h klasifikaciji v razred 1 (vse vrednosti so pozitivne, razen dveh nizkih v originalnem prostoru pri atributu 1). Opazimo tudi, da višje vrednosti obeh atributov v originalnem prostoru bolj pripomorejo h klasifikaciji v razred 1, kot nižje vrednosti. Nov primer ima vrednosti SHAP 3.008 in 3.091, kar pomeni, da sta obe značilki pomembni za klasifikacijo primera v razred 1.

Naša podatkovna množica (Homogene gruče) je sestavljena iz samih homogenih gruč. Če je kakšna gruča sestavljena iz več razredov, moramo uporabiti odločitvena pravila za ločitev razredov znotraj gruče.



Slika 3.4: Na grafu so prikazane vrednosti SHAP za vse gruče (0-4) od zgoraj navzdol. Vidimo, da sta za gručo 1 pomembni obe značilki in da so vrednosti SHAP bolj razpršene kot pri drugih gručah. Ta gruča je sredinska in je bolj kompleksna za klasifikator kot druge.

## Poglavje 4

# Podatkovne množice

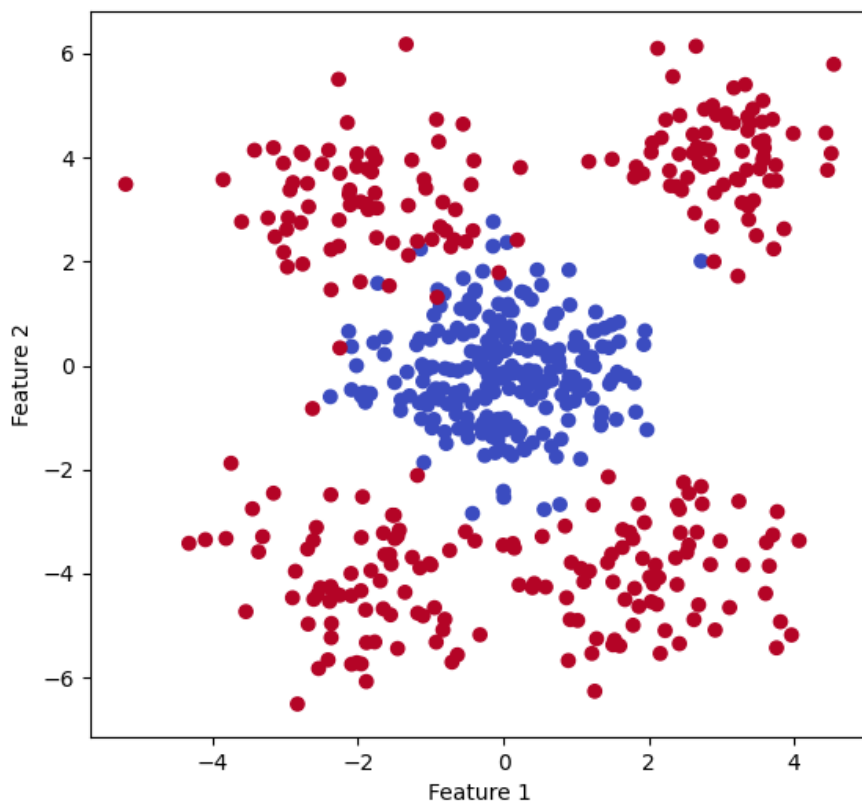
V tem poglavju predstavimo uporabljene podatkovne množice. Najprej predstavimo umetne podatkovne množice in nato še realne podatkovne množice. Pri realnih podatkovnih množicah opišemo tudi predprocesiranje.

### 4.1 Umetne podatkovne množice

V tem razdelku predstavimo dve umetni podatkovni množici, ki ju uporabimo za testiranje algoritmov gručenja in eno umetno podatkovno množico, ki jo uporabimo za razlago. Podatkovna množica Krožne gruče vsebuje gruče sferične oblike, kar poudari dobro lastnost K-MEANS algoritma. Podatkovna množica Trakovi 4-3 vsebuje gruče podolgovate oblike, kjer boljše rezultate vrnete algoritma DBSCAN in HDBSCAN. Sledi še podatkovna množica Homogene gruče, ki je bila zgrajena za enostavno predstavitev naše metode razlage.

#### 4.1.1 Umetna podatkovna množica Krožne gruče

Prva umetna podatkovna množica (slika 4.1) je dvodimenzionalna in ima 550 primerov. Podatkovna množica ima dva razreda, opisana sta z dvema atributom. Zanima nas, kako dobro algoritmi najdejo pet gruč. Gruče so sferične oblike in se na mejah rahlo prekrivajo. Primeri v posameznih gručah so po-

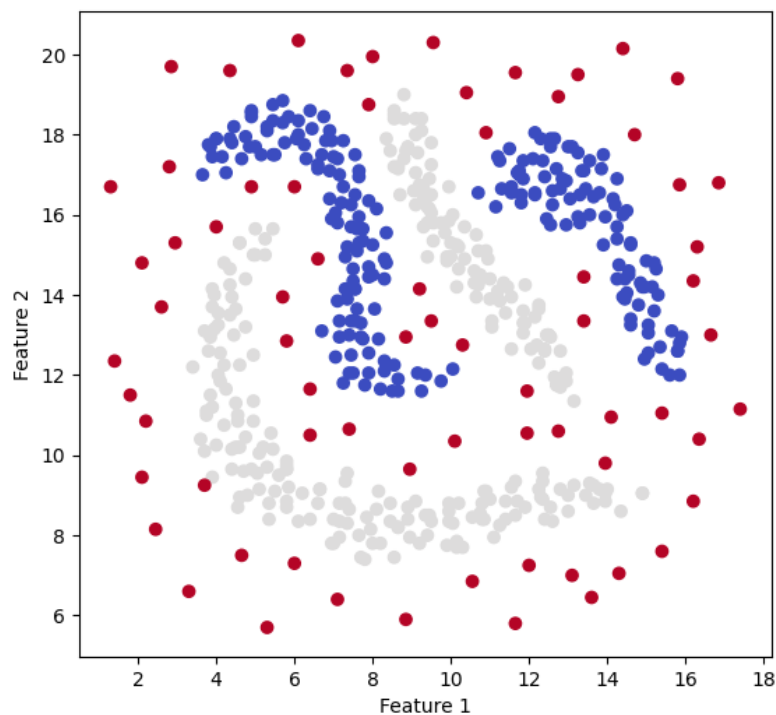


Slika 4.1: Dvodimenzionalna umetna podatkovna množica Krožne gruče s petimi gručami in dvema razredoma. Razred 0 je označen z modro barvo, razred 1 z rdečo.

razdeljeni normalno v dveh dimenzijah z različnimi standardnimi odkloni. V modri gruči je 250 primerov, v vsaki od rdečih pa 75. Podatkovna množica je bila namenjena primerjavi algoritmov in preizkusu hevrstike silhuete in indeksa DBCV.

#### 4.1.2 Umetna podatkovna množica Trakovi 4-3

Druga umetna podatkovna množica (slika 4.2) je dvodimenzionalna in ima 545 primerov, porazdeljenih v tri razrede in opisanih z dvema atributoma.



Slika 4.2: Dvodimenzionalna umetna podatkovna množica Trakovi 4-3. Ima 4 gruče in 3 razrede. Razreda 0 in 1 (modra in siva barva) predstavljata sestavo gruč, medtem ko razred 2 (rdeča barva) predstavlja šum.

Primeri so porazdeljeni v štiri zavite in podolgovate gruče, pozicionirane druga poleg druge. Primeri v posameznih gručah so zelo gosti do meje gruč, med in okoli gruč najdemo primere z nizko gostoto, ki predstavljajo šum. Podatkovna množica je bila namenjena primerjavi algoritmov in preizkusu hevrstike silhuete in indeksa DBCV.

### 4.1.3 Umetna podatkovna množica Homogene gruče

Tretjo umetno podatkovno množico Homogene gruče (slika 3.2) opišemo v poglavju 3. Podatkovna množica je bila namenjena opisu postopka razlage klasifikatorja.

## 4.2 Realne podatkovne množice

V tem podpoglavju predstavimo dve realni podatkovni množici, ki smo ju uporabili v naši analizi. Podatkovno množico MNIST uporabimo za testiranje algoritmov gručenja, množico KDD99 pa za razlago klasifikatorja z našo metodo. Pri obeh množicah lahko s tehniko zmanjšanja dimenzionalnosti na pridobljenem grafu razberemo ločene gručice.

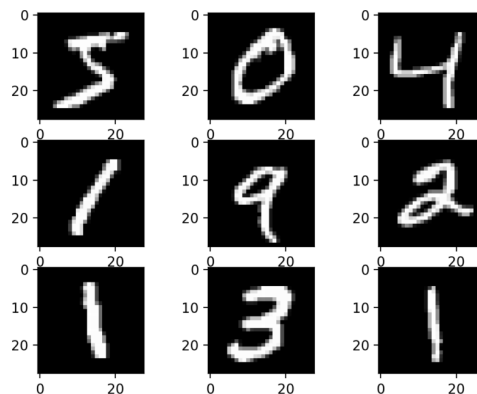
### 4.2.1 Podatkovna množica MNIST

Podatkovna množica MNIST (Modified National Institute of Standards and Technology) je široko uporabljena na področju strojnega učenja in računalniškega vida. Vsebuje zbirko ročno napisanih števil (slik), ki se pogosto uporabljajo za učenje različnih algoritmov za klasifikacijo (vsebuje 10 razredov: številke 0-9). Vsaka slika je sivinska in ima resolucijo  $28 \times 28 = 784$  pikslov. Podatkovna množica ima 70.000 slik, ki se ponavadi delijo na učno množico (60.000 slik) in testno množico (10.000 slik). Podatkovno množico smo uporabili za testiranje algoritmov gručenja. Zaradi časovne kompleksnosti algoritmov smo vzeli le prvih 10.000 primerov iz učne množice in vsako sliko, ki je 2D matrika, razvili v vektor. Primere slik vidimo na sliki 4.3 [2].

### 4.2.2 Podatkovna množica KDD99

Podatkovna množica KDD99 [21] je široko poznana in pogosto uporabljena na področju varnosti in odkrivanja anomalij v računalniških omrežjih. Podatkovna množica je uporabljena za razvoj in ocenjevanje sistemov, ki odkrivajo vdore v računalniška omrežja, in vsebuje različne vrste omrežnega prometa, vključno z običajnimi omrežnimi aktivnostmi in potencialnimi vdori. Vsebuje 23 različnih vrst napadov (ciljna spremenljivka). V našem primeru smo jih razdelili na škodljive in normalne povezave. Ima 41 atributov, kjer so tri kategorični (protocol, service, flag), ostali pa zvezni. Podatkovno množico smo (poleg ciljne spremenljivke) predprocesirali z naslednjimi koraki: odstranili smo primere, ki se ponavljajo, zatem smo uporabili tehniko eničnega





Slika 4.3: Primeri slik števk iz podatkovne množice MNIST, ki jih uporabljamo za strojno učenje in testiranje modelov.

kodiranja, kjer smo za vsako vrednost kategorične spremenljivke dodali nov atribut (dobimo 118 atributov) in vzeli naključnih 5.000 primerov zaradi časovne zahtevnosti algoritmov. V tabeli 4.1 vidimo imena vseh atributov (z enakimi kraticami so atributi poimenovani tudi v poglavju 6).

Tabela 4.1: Imena in kratice atributov podatkovne množice KDD99

Kratika	Ime	Kratika	Ime	Kratika	Ime
F0	duration	F40	protocol_type_b'udp'	F80	service_b'ntp_u'
F1	src_bytes	F41	service_b'IRC'	F81	service_b'other'
F2	dst_bytes	F42	service_b'X11'	F82	service_b'pm_dump'
F3	land	F43	service_b'Z39_50'	F83	service_b'pop_2'
F4	wrong_fragment	F44	service_b'auth'	F84	service_b'pop_3'
F5	urgent	F45	service_b'bgp'	F85	service_b'printer'
F6	hot	F46	service_b'courier'	F86	service_b'private'
F7	num_failed_logins	F47	service_b'csnet_ns'	F87	service_b'red_i'
F8	logged_in	F48	service_b'ctf'	F88	service_b'remote_job'
F9	num_compromised	F49	service_b'daytime'	F89	service_b'rje'
F10	root_shell	F50	service_b'discard'	F90	service_b'shell'
F11	su_attempted	F51	service_b'domain'	F91	service_b'smtp'
F12	num_root	F52	service_b'domain_u'	F92	service_b'sql_net'
F13	num_file_creations	F53	service_b'echo'	F93	service_b'ssh'
F14	num_shells	F54	service_b'eco_i'	F94	service_b'sunrpc'
F15	num_access_files	F55	service_b'ecr_i'	F95	service_b'supdup'
F16	num_outbound_cmds	F56	service_b'efs'	F96	service_b'systat'
F17	is_host_login	F57	service_b'exec'	F97	service_b'telnet'
F18	is_guest_login	F58	service_b'finger'	F98	service_b'tftp_u'
F19	count	F59	service_b'ftp'	F99	service_b'tim_i'
F20	srv_count	F60	service_b'ftp_data'	F100	service_b'time'
F21	error_rate	F61	service_b'gopher'	F101	service_b'urh_i'
F22	srv_error_rate	F62	service_b'hostnames'	F102	service_b'urp_i'
F23	rerror_rate	F63	service_b'http'	F103	service_b'uucp'
F24	srv_error_rate	F64	service_b'http_443'	F104	service_b'uucp_path'
F25	same_srv_rate	F65	service_b'imap4'	F105	service_b'vmnet'
F26	diff_srv_rate	F66	service_b'iso_tsap'	F106	service_b'whois'
F27	srv_diff_host_rate	F67	service_b'klogin'	F107	flag_b'OTH'
F28	dst_host_count	F68	service_b'kshell'	F108	flag_b'REJ'
F29	dst_host_srv_count	F69	service_b'ldap'	F109	flag_b'RSTO'
F30	dst_host_same_srv_rate	F70	service_b'link'	F110	flag_b'RSTOS0'
F31	dst_host_diff_srv_rate	F71	service_b'login'	F111	flag_b'RSTR'
F32	dst_host_same_src_port_rate	F72	service_b'mtp'	F112	flag_b'S0'
F33	dst_host_srv_diff_host_rate	F73	service_b'name'	F113	flag_b'S1'
F34	dst_host_error_rate	F74	service_b'netbios_dgm'	F114	flag_b'S2'
F35	dst_host_srv_error_rate	F75	service_b'netbios_ns'	F115	flag_b'S3'
F36	dst_host_rerror_rate	F76	service_b'netbios_ssn'	F116	flag_b'SF'
F37	dst_host_srv_rerror_rate	F77	service_b'netstat'	F117	flag_b'SH'
F38	protocol_type_b'icmp'	F78	service_b'nns'	N/A	N/A
F39	protocol_type_b'tcp'	F79	service_b'nntp'	N/A	N/A

## Poglavje 5

# Primerjava algoritmov za gručenje

Razlaga s podkoncepti je odvisna od uspešnosti gručenja podatkovne množice. Za gručenje ne obstaja en sam algoritem, ki je v vseh primerih najboljši. V tem poglavju primerjamo štiri algoritme gručenja. Eksperimente izvedemo na dveh umetnih množicah (Krožne gruče in Trakovi 4-3) in eni realni podatkovni množici (MNIST). Cilj poglavja je analiza podatkov, ki dobro delujejo s posameznimi algoritmi gručenja in predstavitev dobrih in slabih lastnosti gručenj. Za ocenjevanje algoritmov gručenja bomo uporabili dve heuristiki: koeficient silhuete in indeks DBCV. Indeks silhuete je boljši za ocenjevanje algoritmov, ki delujejo na osnovi razdalj, indeks DBCV pa je boljši za algoritme, ki delujejo na osnovi gostote podatkov in del podatkov označijo kot šum [4]. DBSCAN in HDBSCAN ne razdelita vseh primerov v gruče, zato zapišemo tudi delež označenih primerov, ostali so označeni kot osamelci oziroma šum. Ker heuristike silhuete in indeksa DBCV ne moremo neposredno primerjati (poleg tega ne povesta vseh informacij o gručenju), rezultate gručenja preučimo tudi na grafih, kjer si v visoko dimenzionalnem prostoru pomagamo s tehniko t-SNE za zmanjšanje dimenzij. V nadaljevanju pri vsaki podatkovni množici pokažemo slike le najboljših gručenj. Gručenja ostalih algoritmov najdemo v dodatku A.

## 5.1 Krožne gruče

Tabela 5.1: Tabela ocen in deleža označenih primerov na podatkovni množici Krožne gruče.

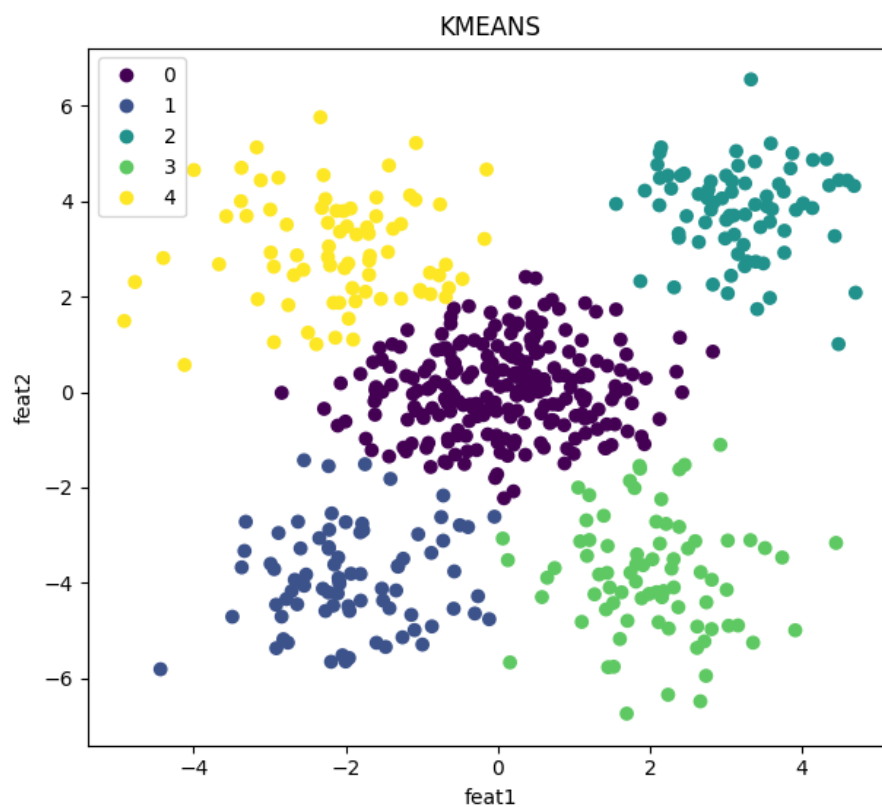
Algoritem	Silhueta	DBCV	%
MDEC	0.40	N/A	100
K-MEANS	0.52	N/A	100
DBSCAN	N/A	0.34	65
HDBSCAN	N/A	0.38	69

### 5.1.1 Algoritem MDEC

Algoritem MDEC vrne tri rezultate gručenja. Najboljša vrednost silhuete (tabela 5.1) je bila pri parametru za število gruč  $k = 3$ , kar ni najboljši rezultat, saj je razvidnih 5 skupin. Vredno je omeniti, da MDEC vsebuje naključne elemente, tako da se vrednosti in ocene lahko vsakič rahlo spremenijo, a so bile v večini podobne. Vzet je bil najboljši rezultat izmed izvedenih poskusov. Na sliki A.1 je bil to rezultat na skrajno desni (MDEC BG). Algoritem MDEC je pri dvodimenzionalnih podatkih ustvaril ravne pasove, kar ni bila dobra rešitev. Algoritem MDEC je namenjen uporabi v visoko dimenzionalnih prostorih, kjer se bolje odreže.

### 5.1.2 Algoritem K-MEANS

Algoritem K-MEANS je vrnil zelo dober rezultat, saj so oblike gruč sferične in dovolj ločene. Najboljša vrednost silhuete (tabela 5.1) je bila pri parametru za število gruč  $k = 5$ , kar je dejansko število skupin. Silhueta K-MEANS je v tem primeru boljša od silhuete MDEC algoritma, kar nam potrjuje boljše grupiranje K-MEANS algoritma. Na sliki 5.1 je razviden rezultat gručenja.



Slika 5.1: Rezultat K-MEANS algoritma pri številu gru $\check{c}$   $k = 5$  na podatkovni množici Kro $\check{z}$ ne gru $\check{c}$ e.

### 5.1.3 Algoritem DBSCAN

Algoritem DBSCAN nima parametra za nastavitev števila gruč, ima pa dva druga parametra - min samples in epsilon. Izbira parametrov je opisana v razdelku 2, tukaj predstavimo le najboljši rezultat. Algoritem je veliko primerov označil kot osamelce in jih ni dodelil v nobeno gručo. Oznako je dobilo 65 % primerov pri najboljšem rezultatu. Omeniti je potrebno, da DBCV ocene ni smiselno primerjati z silhueto, saj delujeta na drugačnih principih. DBSCAN je našel devet gruč, kar pomeni, da je našel štiri dodatne gruče, ki niso razvidne iz podatkov. Rezultati gručenja za DBSCAN so na sliki A.2 in v tabeli 5.1.

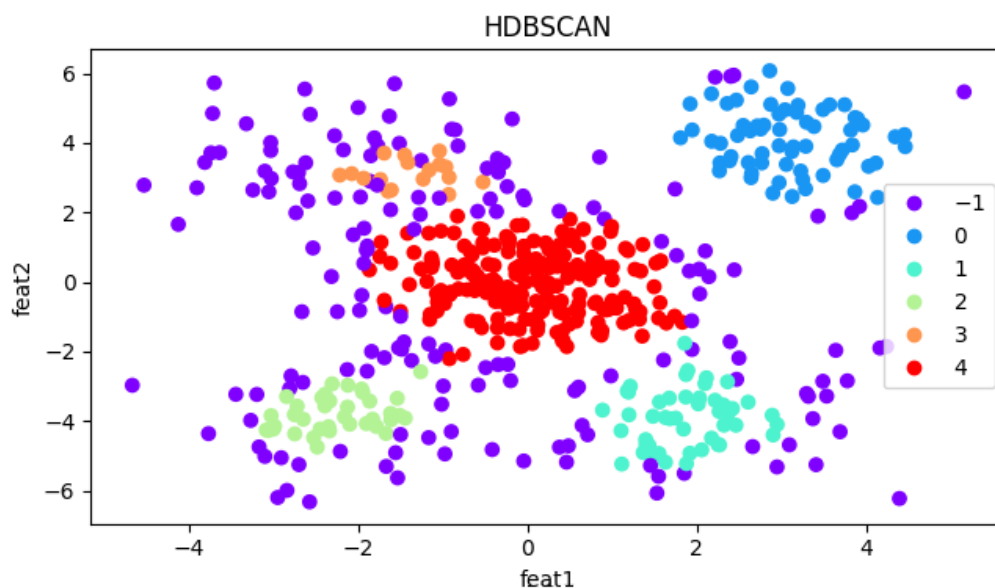
### 5.1.4 Algoritem HDBSCAN

Algoritem HDBSCAN samodejno določi epsilon. Namesto njega smo izbrali parametra min cluster size in min samples. Oznako je dobilo 69 % primerov pri najboljšem rezultatu. HDBSCAN je našel 5 gruč, kar pomeni, da je našel pravilno število gruč, enako kot K-MEANS. Vidimo, da je indeks DBCV boljši kot pri DBSCAN algoritmu, kar dodatno potrjuje boljše gručenje. Rezultati gručenja za HDBSCAN so na sliki 5.2 in v tabeli 5.1.

## 5.2 Trakovi 4-3

Tabela 5.2: Tabela ocen in deleža označenih primerov na podatkovni množici Trakovi 4-3.

Algoritem	Silhueta	DBCV	%
MDEC	0.25	N/A	100
K-MEANS	0.44	N/A	100
DBSCAN	N/A	0.0034	87
HDBSCAN	N/A	0.095	81



Slika 5.2: Rezultat HDBSCAN algoritma na podatkovni množici Krožne gruče.

### 5.2.1 Algoritem MDEC

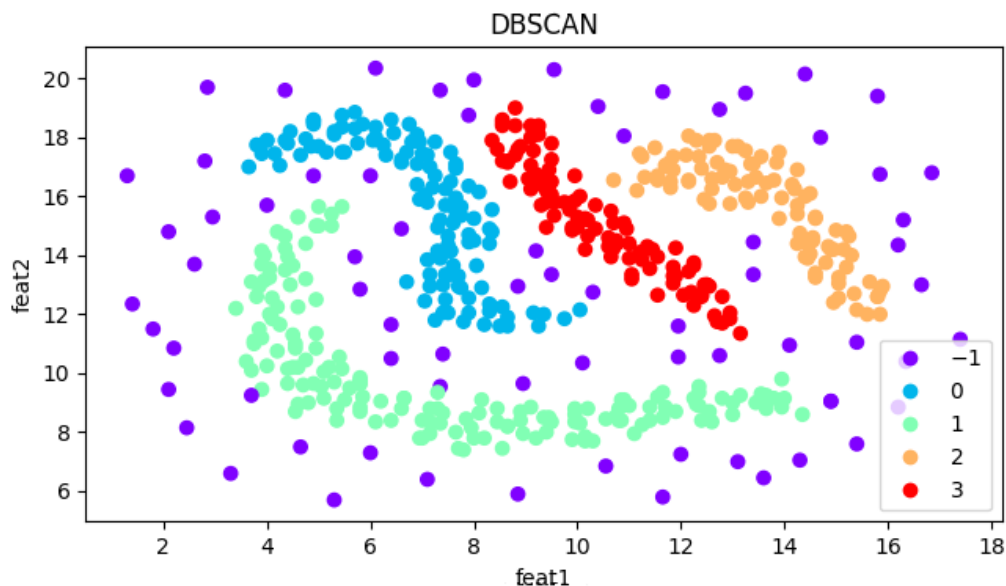
Ponovno vidimo, da je algoritem MDEC prostor razdelil na pasove. Rezultat je na sliki A.3. Najboljši rezultat (tabela 5.2) je dosegel MDEC\_HC (skrajno levo) pri parametru za število gruč  $k = 3$ .

### 5.2.2 Algoritem K-MEANS

Algoritem K-MEANS je vrnil najboljši rezultat pri številu gruč  $k = 10$ . Rezultat gručenja je viden v tabeli 5.2 in na sliki A.4. Vidimo, da se vsak centroid umesti v prostor in zavzame približno sferično obliko.

### 5.2.3 Algoritem DBSCAN

Algoritem DBSCAN je uporaben pri gručenju naravnih podolgovatih oblik, kar je razvidno tudi iz rezultata gručenja na sliki 5.3 in v tabeli 5.2. DBSCAN je našel štiri gruče kar pomeni, da je našel vse gruče in hkrati izločil šum iz



Slika 5.3: Rezultat DBSCAN algoritma na podatkovni množici Trakovi 4-3. Na sliki se vidi označbe vseh štirih skupin in izločitev šuma iz podatkov.

podatkov.

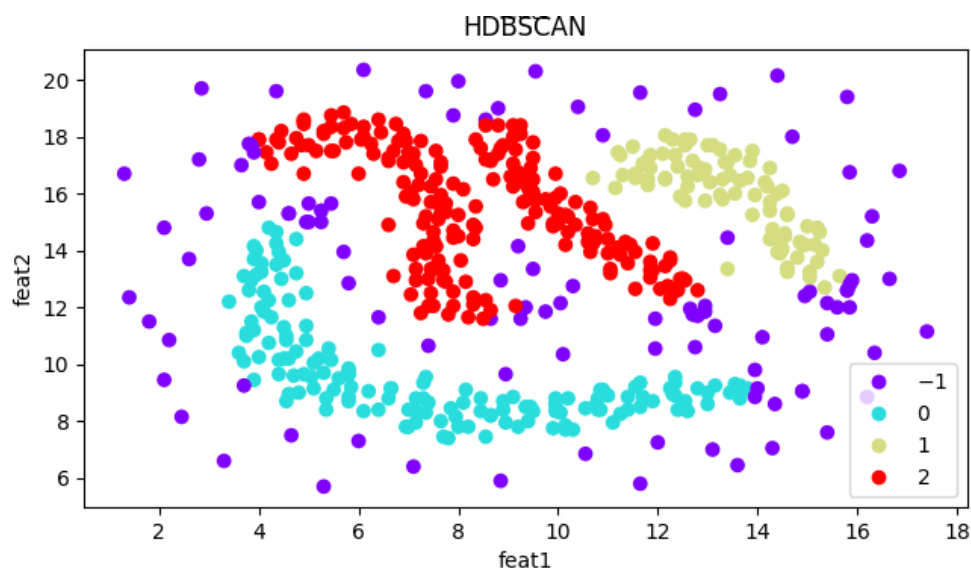
#### 5.2.4 Algoritem HDBSCAN

Algoritem HDBSCAN je hierarhična različica DBSCAN algoritma, a vrne slabši rezultat kot DBSCAN, kot je razvidno iz slike 5.4 in tabele 5.2. HDBSCAN je našel tri gruče, kar pomeni, da je dve ločeni gruči združil v eno. Najboljša DBCV vrednost je 0.095 kar je presenetljivo boljši indeks DBCV od algoritma DBSCAN glede na to, da vidimo slabše gručenje algoritma HDBSCAN. Šum je dobro ločen od skupin.

### 5.3 MNIST

Podatkovna množica MNIST ima 784 atributov, zato smo za prikaz v dveh dimenzijah uporabili vložitev t-SNE. Zaradi časovne zahtevnosti algoritmov smo uporabili le prvih 10.000 primerov. Prav tako smo pri K-MEANS, DB-



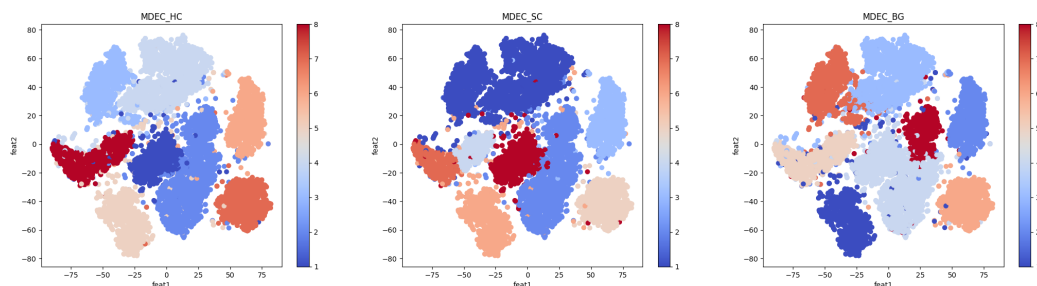


Slika 5.4: Rezultat HDBSCAN algoritma na podatkovni množici Trakovi 4-3. Na sliki se vidi označbe treh skupin in izločitev šuma iz podatkov.

SCAN in HDBSCAN algoritmi za gručenje vhodne podatke najprej pretvorili v dve dimenziji. Izjema je algoritem MDEC, ki je bolje deloval v originalnem, visoko dimenzionalnem prostoru. V tabeli 5.3 vidimo rezultate. Hevristika silhuete pri K-MEANS algoritmu je višja kot pri MDEC, vendar na slikah vidimo, da je bolje gručil algoritem MDEC (pri MDEC algoritmu je hevristika silhuete izračunana v višjih dimenzijah). Enako velja za HDBSCAN in DBSCAN, kjer ima HDBSCAN nižji indeks DBCV, je gručil bolje.

Tabela 5.3: Rezultati algoritmov za gručenje.

Algoritem	Silhueta	DBCV	%	Vizualizacija
MDEC	0.0471	N/A	100	slika 5.5
K-MEANS	0.4738	N/A	100	slika A.5
DBSCAN	N/A	0.282	63	slika A.6
HDBSCAN	N/A	-0.162	95	slika A.7



Slika 5.5: Rezultat MDEC algoritma. Algoritem MDEC\_HC (najbolj leva slika) je dosegel najboljšo hevristiko silhuete pri številu gruč  $k = 8$ . Našel je dve gruči manj, najdene gruče so smiselne.

Tabela 5.4: Najboljši algoritmi za posamezno podatkovno množico in algoritmi, ki vrnejo vsaj smiselne rezultate.

Podatkovna množica	Najboljši algoritem	Smiselni rezultati
Krožne gruče	K-MEANS	HDBSCAN
Trakovi 4-3	DBSCAN	HDBSCAN
MNIST	MDEC	HDBSCAN

## 5.4 Diskusija

Eksperimenti so pokazali nekatere lastnosti algoritmov. Algoritem K-MEANS vrne dobre rezultate s sferičnimi oblikami, kar se je izkazalo za neuporabno pri gručenju realne podatkovne množice MNIST. Če algoritem t-SNE vrne podolgovate oblike gruč, jih algoritem K-MEANS razdeli na več delov. Na umetni podatkovni množici Krožne gruče je vrnil najboljši rezultat, saj so bile gruče krožne oblike. Algoritem MDEC je v nizkih dimenzijah vrnil nesmiselne rezultate, saj je prostor delil na ravne pasove. Pri visoko dimenzionalni množici MNIST je vrnil najboljši rezultat brez uporabe tehnike t-SNE. Algoritem DBSCAN je odvisen od izbire pravih parametrov. Pri umetni množici Trakovi 4-3 je vrnil odličen rezultat, pri drugih dveh podatkovnih množicah je gručil slabše (še posebej pri podatkovni množici MNIST, kjer je vrnil nesmiseln rezultat). Za bolj stabilnega se je izkazal algoritem HDBSCAN, ki je pri vseh treh podatkovnih množicah vrnil smiseln rezultat.

S testiranjem algoritmov in hevristik lahko ugotovimo, da je pomemben del ocene gručenja pregled grafa gručenja, saj hevristike v nekaterih primerih odpovejo. V tabeli 5.4 vidimo algoritme z najboljšim rezultatom za vsako podatkovno množico, ki smo jo preizkusili in algoritme, ki so vrnili vsaj smiseln rezultat.



## Poglavje 6

# Razlaga realne podatkovne množice

V tem poglavju bomo s predlagano metodo razlage s podkoncepti razložili klasifikacije na realni podatkovni množici KDD99, za katero že obstajajo nekatere razlage, kar omogoča primerjavo. Najprej podatke gručimo in izberemo najprimernejši algoritem gručenja. Zatem ocenimo homogenost gruč s frekvenco razredov v posamezni gruči. Če je gruča homogena, jo razložimo z odločitvenimi pravili za meje gruče in medoidom, v nasprotnem primeru z odločitvenimi pravili za meje gruče in odločitvenimi pravili za ločitev razredov znotraj gruče. Nato izračunamo vrednosti SHAP za posamezno gručo; z njimi si pomagamo pri interpretaciji odločitvenih pravil in medoidov. Nov naključni primer umestimo v najbližjo gručo, izračunamo vrednosti SHAP zanj in interpretiramo odločitev klasifikatorja.

Podatkovna množica KDD99 ima strukturo, ki jo pričakuje naša metoda. Večina preostalih podatkovnih množic, ki smo jih pregledali, nima tako jasno primerov enega razreda razčlenjenega v več gruč ali pa ima slabo ločene gruče. Predprocesiranje podatkovne množice in imena atributov so na voljo v poglavju 4.2.2. Pri gručenju smo si pomagali s tehniko t-SNE, saj je pri algoritmih na osnovi gostote (v našem primeru na sliki 6.1 je to HDBSCAN) pomagala pri primerih, ki so bili označeni za šum zaradi visoke dimenzio-

nalnosti. S tem smo pridobili sedem gruč, ki so na prvi pogled smiselno opredeljene, saj vidimo goste gruče, ki so dobro ločene med seboj. Vrednost indeksa DBCV je 0.27, kar dodatno potrjuje smiselno gručenje. V tabeli 6.1 vidimo verjetnosti razredov za vsako gručo. Gruči 1 in 2 vsebujeta škodljive povezave, ostale normalne. Opazimo, da sta gruči 2 in 5 bolj mešani kot preostale gruče, kar bomo kasneje omenili pri razlagi z medoidi, kjer si bomo pomagali še z odločitvenimi pravili znotraj gruče. Za klasifikacijo povezav smo uporabili logistično regresijo, ki je imela na testnih podatkih klasifikacijsko točnost 0.9587 in F1 score 0.9587. V tabeli 6.2 imamo podana odločitvena pravila, ki med seboj ločijo gruče. Vidimo, da imamo za vsako gručo več pravil, saj so podatki visoko dimenzionalni in se jih lahko dobro loči po različnih atributih.

Vsi atributi niso enako pomembni za klasifikator, kar pomeni, da ni smiselno obravnavati vseh pravil enako. Pri večini pravil ene gruče vidimo majhno podmnožico atributov, ki se ponavljajo in so s tem verjetno bolj pomembni za meje kot drugi. Z odločitvenimi pravili odkrijemo dodatno znanje iz podatkov, ki nam pomaga predvsem, ko je gruča homogena, saj del razreda ločimo od preostalih primerov.

Pri interpretaciji pravil je pomembno pogledati meri Precision in Recall, saj so nekatera pravila manj točna ali pa ne zajamejo vseh primerov. Ko imamo občutek, kako so skupine medsebojno ločene, vsako lahko opišemo z medoidom, ki je tipičen primer gruče. Medoide vidimo v tabeli 6.3. Vidimo, da sta medoida M1 in M2 normalni mrežni povezavi, ostali medoidi predstavljajo škodljive. Opazimo, da imata edina atribut F1 (`src_bytes`) enak 0. Močno se razlikujeta tudi po atributu F19 (`count`), ki ima večjo vrednost od preostalih medoidov. Pri atributu F25 (`same_srv_rate`) imata nizke vrednosti, ostali imata vrednost 1. Prav tako je razlika pri F32 (`dst_host_same_src_port_rate`), kjer imata edina vrednost 0. Pri atributih F34 (`dst_host_serror_rate`), F35 (`dst_host_srv_serror_rate`), F86 (`service_private`) in F112 (`flag_S0`) imata vrednost 1, ostali imajo 0 (binarno kodirane vrednosti). Enako analizo lahko naredimo za vsak medoid posebj v primerjavi z ostalimi.

Zaradi velikega števila atributov smo obdržali v tabeli le tiste, ki se razlikujejo vsaj pri enem izmed medoidov. Poleg medoidov nas zanima tudi, kako so ločeni razredi znotaj gruče 2 in 5, kar lahko vidimo v tabeli 6.4. S tabelo 6.4 si pomagamo v primeru, da nov primer umestimo v gruči 2 ali 5, saj je medoid lahko zavajajoč. Zanima nas tudi, kateri atributi so pomembni za naš klasifikator (logistično regresijo), kar vidimo na sliki 6.2 za gručo 1, saj si bomo s tem pomagali pri interpretaciji odločitve. Vrednosti SHAP za preostale gruče najdemo v dodatku B.

Na sliki 6.2 vidimo, da so pomembni atributi F19 (count), F29 (dst\_host\_srv\_count), F20 (srv\_count) in F28 (dst\_host\_count), ostali so skoraj ničelni. Atributa F19 (count) in F29 (dst\_host\_srv\_count) pripomoreta h klasifikaciji v razred 1 (normalna povezava). Atribut 20 (srv\_count) z visokimi vrednostmi pripomore h klasifikaciji v razred 0 (škodljive povezave), nizke vrednosti pa h klasifikaciji v razred 1. Atribut 28 (dst\_host\_count) je nepomemben v primerjavi s prejšnjimi tremi in pripomore h klasifikaciji v razred 0.

Vidimo, da je veliko atributov nepomembnih. Pomembnost atributov nato upoštevamo pri odločitvenih pravilih in medoidih, kjer nam pravila in medoidi sestavljeni iz pomembnih atributov podajo informacijo o ločitvenih mejah in o sami gruči. S to informacijo lahko pogledamo nov primer in na podlagi njegovih atributov razložimo, zakaj spada v gručo (odločitvena pravila) in v kakšne vrste gručo spada (medoid).

Razložimo še naključen nov primer iz testne množice. Zaradi visoke dimenzionalnosti niso podane vse vrednosti atributov (pomembni atributi glede na vrednosti SHAP: F19 (count) = 218, F20 (srv\_count) = 6, F28 (dst\_host\_count) = 255, F29 (dst\_host\_srv\_count) = 6). Primer je logistična regresija opredelila v razred 1, algoritem za gručenje pa ga je razvrstil v gručo 1. Pogledamo tabelo frekvence razredov in vidimo, da je gruča 1 v celoti homogena s 100% primerov, ki pripadajo razredu 1 (škodljiva povezava).

Nov primer je dobil SHAP vrednosti F19 (count): 13.9, F29 (dst\_host\_srv\_count): 2.4, F20 (srv\_count): 0.34, F28 (dst\_host\_count): -0.23, vrednosti ostalih atributov so zanemarljivo majhne. Vrednosti SHAP novega primera potrjujejo

pomembnost zgornjih atributov v skupini 1. V tabeli z odločitvenimi pravili 6.2 vidimo devet pravil, ki z visoko natančnostjo opišejo gručo z majhno podmnožico atributov. Iz odločitvenih pravil za meje gruče razberemo, da primer zadostuje vsem pogojem, kar primer z gotovostjo postavi v gručo 1. Poleg tega, imamo še več pogojev z atributi, ki nas v tem trenutku ne zanimajo, saj niso pomembni za odločitev klasifikatorja. Skupina je homogena, zato jo dobro opiše medoid M1 v tabeli 6.3.

Medoid je iz razreda 1, kar pomeni, da je gruča sestavljena iz normalnih povezav (100% povezav). Medoid ima izmed vseh največjo vrednost atributa F19 (count), kar nam iz zgornje razlage vrednosti SHAP za gručo potrjuje, da je bil klasificiran v razred 1. Analizo za medoide M1 in M2 proti ostalim smo naredili zgoraj na strani 38, sedaj si pogledajmo v čem se razlikujeta M1 in M2, saj oba pripadata razredu 1. Očitna razlika se vidi pri atributu 29 (dst\_host\_srv\_count), kjer je medoid M1 bolj podoben M0, čeprav je M0 iz razreda 0. Enako velja za atribut F19 (count). Pri atributu F116 (flag\_SF) se medoida razlikujeta, saj ima le medoid M2 vrednost 0. V primeru, da bi bil nov primer iz gruče 2 ali 5 (nečisti gruči), bi si ogledali še pravila znotraj gruče, ki ločijo razrede.

Na podatkovni množici KDD99 že obstajajo poskusi razlage, [12], kjer je bila podatkovna množica razložena z odločitvenim drevesom. Atributi so bili glede na pomembnost ocenjeni z mero entropije. Pri odločitvenem drevesu sta bila dva izmed pomembnih atributov enaka (najpomembnejši: Count in Dst\_host\_count). Ta pristop omogoča razlago celotne podatkovne množice z enim drevesom, kar olajša interpretacijo, medtem ko se naš pristop bolj poglobi v prostorsko povezanost primerov in iz tega pridobi informacije na podlagi odločitvenih pravil in medoidov.

Remah Younis in sod. [23] so različne modele za strojno učenje razložili s SHAP vrednostmi in jih potrdili na podlagi ocenjevanja gostote verjetnosti z jedrom KDE (kernel density estimation) grafov. Večina klasifikatorjev je imela med najpomembnejšimi značilkami značilko count, dst\_host\_srv\_count

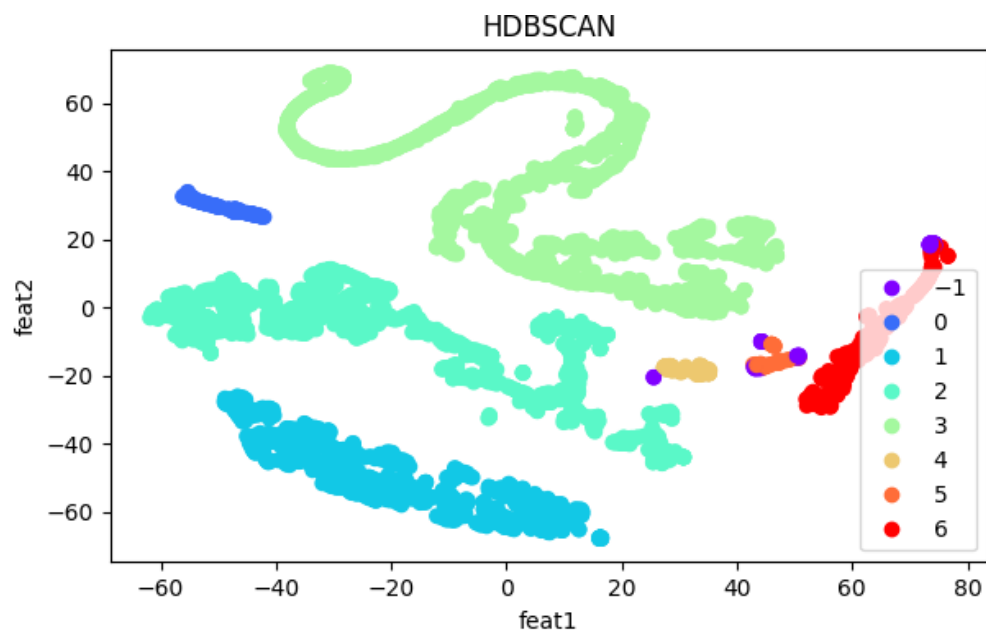


in `dst_host_count`. Razlika je bila pri porazdelitvi pomembnosti med atributi. Za pomembne je bilo označenih večje število značilk kot pri naši razlagi. Vrednosti SHAP so bile upravičene, kar pri nas ne drži, a poleg pomembnosti atributov nismo izvedeli ničesar o podatkovni množici.

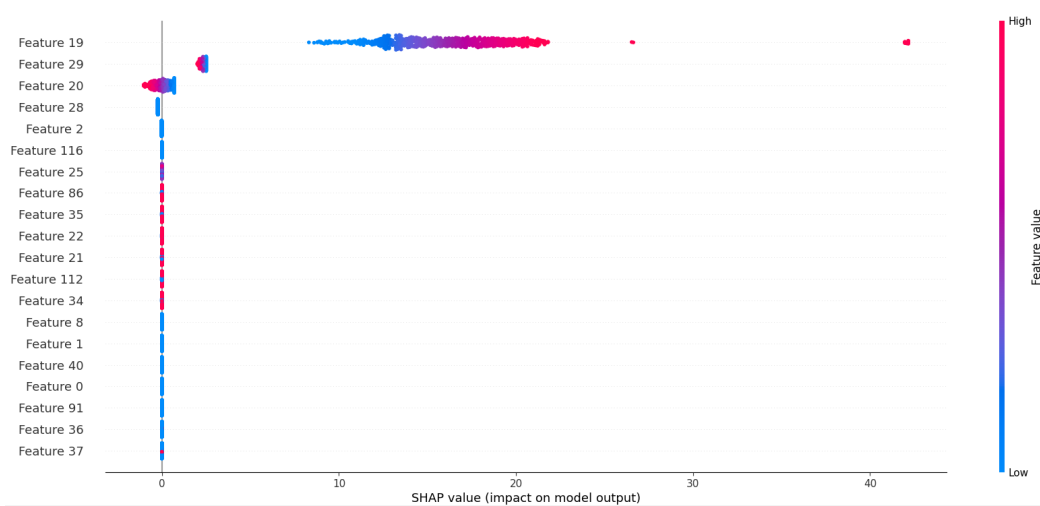
Z algoritmom K-MEANS so na podlagi gruč Mohammad Khubeb Siddiqui in sod. [20] iskali povezavo med tipom napada in uporabljenim protokolom. Osredotočili so se na analizo podatkov v posamezni gruči, enako kot pri naši razlagi, a niso pojasnili klasifikatorjev. Ugotovili so, da je protokol TCP najpogosteje deležen napadov. Naši medoidi imajo vsi vključen TCP protokol, razen medoida M0, zaradi česar je zanimiv za nadaljnjo analizo. Podatkovno množico so opisali tudi z odločitvenimi pravili [10]. Za razliko od naše razlage so odločitvena pravila uporabili na negručenih podatkih, kar je rezultiralo v drugačnih pravilih. Uporabili so odločitvena pravila za ločitev razredov, kar smo mi storili v nečistih gručah. V pogojih odločitev se pojavljajo atributi, ki tvorijo tudi naša odločitvena pravila, a z drugačnimi vrednostmi. Njihova razlaga, kot pri razlagi z odločitvenim drevesom, ne gruči podatkov pred uporabo odločitvenih pravil, zato je o podobnih zaključih težko sklepati.

Tabela 6.1: Frekvenca razredov v gručah.

Razred	G0	G1	G2	G3	G4	G5	G6
0	0.966	0	0.344	0.993	0.787	0.641	0.972
1	0.334	1	0.656	0.007	0.213	0.359	0.028



Slika 6.1: Rezultat t-SNE + HDBSCAN algoritma na učni množici. Algoritem je našel 7 gruč. DBSV indeks je 0.5617.



Slika 6.2: Vrednosti SHAP za gručo 1.

Tabela 6.2: Odločitvena pravila za meje posameznih gruĉ.

Gruĉa	Odloĉitvena pravila	Precision	Recall
0	$F0 > 640.5$ in $F2 > 52.5$	1.0	1.0
0	$F0 > 629.5$ in $F30 \leq 0.7755$	1.0	1.0
0	$F0 > 724.5$ in $F31 > 0.01$	1.0	1.0
0	$F0 > 640.5$ in $F1 \leq 2567830.0$	1.0	1.0
0	$F0 > 640.5$ in $F28 > 11.5$	1.0	1.0
0	$F0 > 640.5$ in $F33 \leq 0.02$	1.0	1.0
0	$F0 > 640.5$ in $F60 \leq 0.5$	1.0	1.0
1	$F116 \leq 0.5$ in $F19 > 162.0$ in $F32 \leq 0.005$	1.0	1.0
1	$F19 > 162.0$ in $F25 \leq 0.265$ in $F28 > 185.0$	1.0	1.0
1	$F19 > 162.0$ in $F30 \leq 0.115$ in $F32 \leq 0.005$	1.0	0.9991
1	$F19 > 162.0$ in $F28 > 185.0$ in $F30 \leq 0.115$	1.0	0.9987
1	$F19 > 162.0$ in $F20 \leq 49.0$ in $F32 \leq 0.005$	1.0	0.9985
1	$F19 > 162.0$ in $F28 > 185.0$ in $F29 \leq 48.5$	1.0	0.9985
1	$F19 > 162.0$ in $F32 \leq 0.005$ in $F39 > 0.5$	1.0	0.9985
1	$F1 \leq 14.0$ in $F19 > 162.0$ in $F32 \leq 0.005$	1.0	0.9985
1	$F116 \leq 0.5$ in $F19 > 162.0$ in $F28 > 185.0$	1.0	0.9985
2	$F1 \leq 137.0$ in $F19 \leq 162.0$ in $F2 \leq 1269.5$	1.0	0.991
3	$F1 \leq 27526.0$ in $F63 > 0.5$ in $F8 > 0.5$	0.999	0.984
3	$F1 \leq 27540.5$ in $F1 > 71.0$ in $F63 > 0.5$	1.0	0.981
3	$F1 \leq 27540.5$ in $F2 > 36.5$ in $F63 > 0.5$	1.0	0.981
3	$F23 \leq 0.135$ in $F63 > 0.5$ in $F9 \leq 0.5$	0.998	0.982
3	$F1 > 70.5$ in $F63 > 0.5$ in $F9 \leq 0.5$	0.999	0.979
3	$F63 > 0.5$ in $F8 > 0.5$ in $F9 \leq 0.5$	0.999	0.978
4	$F1 \leq 475.0$ in $F1 > 143.5$ in $F102 \leq 0.5$ in $F60 > 0.5$	1.0	0.66
5	$F1 \leq 1482.5$ in $F1 > 475.0$ in $F20 \leq 279.0$ in $F60 > 0.5$	1.0	0.77
6	$F1 > 598.0$ in $F91 > 0.5$	1.0	0.812
6	$F1 > 644.5$ in $F12 \leq 4.5$ in $F32 \leq 0.315$ in $F33 \leq 0.11$ in $F6 \leq 2.5$ in $F8 > 0.5$ in $F97 \leq 0.5$	1.0	0.812

Tabela 6.3: Medoidi z atributi, ki se razlikujejo.

Atribut	M0	M1	M2	M3	M4	M5	M6
F0	3058	0	0	0	0	0	5
F1	147	0	0	286	201	854	1536
F2	105	0	0	1768	0	0	328
F8	0	0	0	1	1	0	1
F19	1	240	109	24	6	2	1
F20	1	1	9	35	6	2	1
F21	0	1	1	0	0	0	0
F22	0	1	1	0	0	0	0
F25	1	0.04	0.08	1	1	1	1
F26	0	0.06	0.06	0	0	0	0
F27	0	0	0	0.11	0	0	0
F28	255	255	255	143	170	106	80
F29	1	1	15	255	28	100	154
F30	0	0.04	0.06	1	0.09	0.47	0.98
F31	0.6	0.07	0.07	0	0.02	0.04	0.03
F32	0.93	0	0	0.01	0.09	0.47	0.01
F33	0	0	0	0.01	0.07	0.02	0.01
F34	0	1	1	0	0	0	0
F35	0	1	1	0	0	0	0
F39	0	1	1	1	1	1	1
F40	1	0	0	0	0	0	0
F60	0	0	0	0	1	1	0
F63	0	0	0	1	0	0	0
F81	1	0	0	0	0	0	0
F86	0	1	1	0	0	0	0
F91	0	0	0	0	0	0	1
F112	0	1	1	0	0	0	0
F116	1	1	0	1	1	1	1
Razred	0	1	1	0	0	0	0

Tabela 6.4: Pravila za primere različnih razredov znotraj gruče. V oklepajih so zapisane vrednosti Precision in Recall.

Gruča	Pravila za razred 0	Pravila za razred 1
2	F111 $\leq$ 0.5 in F25 $>$ 0.49 in F29 $>$ 1.5 in F31 $\leq$ 0.845 in F33 $\leq$ 0.44 in F35 $\leq$ 0.03 in F4 $\leq$ 1.5 in F54 $\leq$ 0.5 (1,0.96)	F1 $\leq$ 74.5 in F25 $\leq$ 0.49 (0.9992, 0.93)
5	F30 $\leq$ 0.665 (1, 1)	F30 $>$ 0.665 (1, 1)



# Poglavje 7

## Diskusija

V tem poglavju podamo zaključke o izbiri števila gruč (podpoglavje 7.1), smiselnosti prototipov (podpoglavje 7.2) in smiselnosti odločitvenih pravil (podpoglavje 7.3). Zaključke podamo predvsem na podlagi razlag s podkoncepti na podatkovnih množicah Homogene gruče in KDD99. Na koncu poglavja sledi diskusija o lastnostih uspešnih razlag (preprostost, interpretabilnost, vizualizacija, prilagodljivost) in o prihodnji avtomatizaciji postopka.

### 7.1 Število gruč

Število gruč igra pomembno vlogo pri algoritmih, kot sta K-MEANS in MDEC, in zato tudi pri razlagi. Pri umetnih podatkih smo poznali dejansko število gruč, tako da je bila najboljša ocena gručenja takrat, ko smo za število gruč uporabili dejansko vrednost. Težje je ta parameter najti v realnih podatkih, ko ne vemo, koliko gruč se skriva v njih. Zato je podatke smiselno preučiti v nižje dimenzionalnih prostorih, preizkusiti različne vrednosti števila gruč in iskati najboljšo oceno gručenja. Veliko število gruč se pozna tudi pri razlagi, kjer se nestrjenost gruč kaže v slabo razločljivih prototipih in pravilih, ki bi morala podajati isto razlago. V nasprotju s tem premajhno število gruč preveč posploši razlago in poda razlage, ki bi morale biti ločene na več pravil in prototipov (dobimo nečiste skupine). V obeh

primerih je lahko razlaga neverodostojna.

## 7.2 Smiselnost prototipov

Prototipi, kot so medoidi, so smiselni takrat, ko je gruča homogena in jo lahko predstavimo z reprezentativnim primerom. Pomemben je tudi prostorski razpored primerov različnih razredov v gruči, saj lahko medoid predstavlja manjšinski razred gruče, če so primeri iz manjšinskega razreda razporejeni v središču sferične gruče. Prav tako medoid ne pove, katere značilke so pomembne. Če vzamemo podolgovato gručo v obliki linije v dvodimenzionalnem prostoru, je pri medoidu pomembna le ena dimenzija, kar ne bo razvidno iz prototipa. Pri takšni težavi si lahko pomagamo z vrednostmi SHAP. V višje dimenzionalnih prostorih so medoidi neintuitivni in nepregledni, kar lahko naslovimo z vrednostmi SHAP za pridobitev pomembnosti atributov.

## 7.3 Smiselnost pravil

Pri odločitvenih pravilih se lahko zgodi, da so gruče dobro ločene v eni dimenziji in slabše v drugih, zato je smiselno gledati samo pogoje v pravilu, ki dobro loči skupine. V visokih dimenzijah se zgodi, da več različnih podmnožic pravil dobro loči skupino. Tu si lahko pomagamo z vrednostmi SHAP za gručo ter za posamezen primer, saj nam te povedo, katere značilke so pomembne, in to upoštevamo v pogojih pravil. Treba je omeniti, da pravila delijo prostor s hiperpravokotniki, kar pomeni, da ne zajamejo v celoti kompleksnih oblik, kar vpliva tudi na razlago. Smiselnost pravil se da oceniti tudi z vrednostmi Precision in Recall, ki jih izračunamo za pravila.

Predstavljen postopek razlage je težek za interpretacijo in v nekaterih primerih nepopoln. Večina realnih podatkovnih množic nima pričakovane strukture in zato naše metode razlage ne moremo uporabiti. Rezultati posameznih metod zahtevajo človeško interpretacijo. Zaradi tega postopek ni v celoti av-



tomatiziran, saj je ključen izbor pravega algoritma gručenja, hevristike za oceno gručenja pa niso zadostne za izbor. Po izvedbi posameznih komponent metodologije razlage dobimo rezultate, ki jih je potrebno interpretirati, kar zahteva ekspertno znanje. Dodatno oviro predstavlja visoka dimenzionalnost podatkov, ki dodatno oteži razlago in tudi delne človeške ocene. Dobra lastnost naše metode razlage je razčlenitev podatkov na gruče, saj s tem pridobimo dodatne informacije o podkonceptih, ki so drugim metodam razlage, kot sta, na primer odločitvena pravila in medoidi, skrite. Rezultati nekaterih komponent naše metodologije (gručenje, SHAP) so lahko vizualno predstavljeni, kar olajša razumevanje in omogoča lažjo interpretacijo.



## Poglavje 8

### Zaključki

V diplomski nalogi smo predstavili metodo razlage klasifikatorjev na podlagi podkonceptov (gruče primerov istega razreda) z uporabo algoritmov gručenja, odločitvenih pravil, medoidov in vrednosti SHAP. Izvedli smo eksperimente s štirimi algoritmi za gručenje na dveh umetnih in dveh realnih podatkovnih množicah. Za izbor parametrov in oceno kvalitete gručenja smo uporabili heuristiki silhuete in indeks DBCV. Pri visoko dimenzionalnih podatkih je bila uporabljena tudi tehnika za zmanjšanje dimenzij t-SNE.

Primerjava algoritmov gručenja pokaže, da se algoritem MDEC slabo odreže v nizko dimenzionalnih prostorih (2D), medtem ko se je dobro odrezal v visoko dimenzionalnem prostoru. Algoritem K-MEANS se je odrezal dobro na sferičnih gručah, tudi če so si bile gruče blizu. V visoko dimenzionalnem prostoru se je bolje odrezal v kombinaciji s t-SNE, a vseeno ne tako dobro kot algoritem HDBSCAN. Algoritma DBSCAN in HDBSCAN sta brez zmanjšanja dimenzionalnosti s t-SNE veliko primerov klasificirala kot šum. Algoritem DBSCAN je v kombinaciji s t-SNE slabo gručil, razen na umetni podatkovni množici Trakovi 4-3. Algoritem HDBSCAN se je izkazal za bolj stabilnega (od DBSCAN) v kombinaciji s t-SNE, zato je bil uporabljen tudi pri razlagi realne podatkovne množice KDD99.

Naš postopek razlage s podkoncepti se je v večini primerov realnih podatkovnih množic izkazal za neuporabnega, saj imajo le redke podatkovne množice

primere enega razreda razdeljene v dve ali več gruči (premajhna razčlenjenost primerov istega razreda na gruče).

Predstavljena metoda razlage razdeli podatkovno množico na manjše dele, kar naredi razlago bolj razumljivo in bolj podrobno, kot če bi metode uporabili na celotni podatkovni množici. Pri razlagi z medoidi so se pojavile težave v visoko dimenzionalnih prostorih, saj je medoid z veliko atributi dokaj ne-intuitiven in težko razumljiv; za take primere smo si pomagali z vrednostmi SHAP. Ogledali smo si le pomembne attribute in prikazali le attribute, v katerih so se medoidih različnih gruč razlikovali. Odločitvena pravila so v nekaterih primerih vsebovala nepomembne attribute, a so vseeno vsebovala dodatno informacijo o gruči. Težavo predstavljajo gruče sestavljene iz dveh ali več razredov, saj je medoid lahko zavaajoč. Uporabili smo tudi odločitvena pravila, tokrat znotraj gruče, da so pomagala ločiti primere enega razreda od preostalih.

Razvoj metodologije za razlago klasifikatorjev se lahko nadaljujejo v smeri dodatnih razlag za posamezne gruče. V to je lahko vključena vizualizacija dodatnih primerov (ne samo medoida) iz različnih delov gruče in dodatne analize nečistih gruč, predvsem, kako naj bo gruča predstavljena. Attribute bi lahko opisali z statističnimi metodami, saj o njih naše razlage ne povedo dosti. Uporabno bi bilo tudi avtomatizirati primernost podatkovnih množic za metodo s podkoncepti, za kar bi lahko razvili ustrezne hevristike. Izvorno kodo diplomske naloge najdemo na repozitoriju Github.<sup>1</sup>

---

<sup>1</sup><https://github.com/musicn/Diploma>

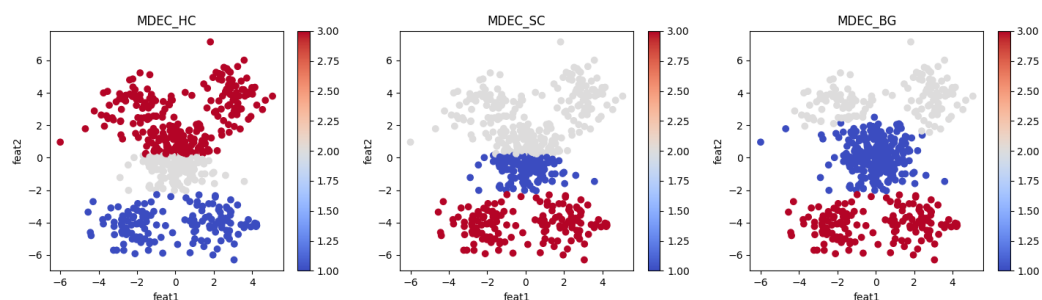
# Dodatek A

## Celotna primerjava gručenj

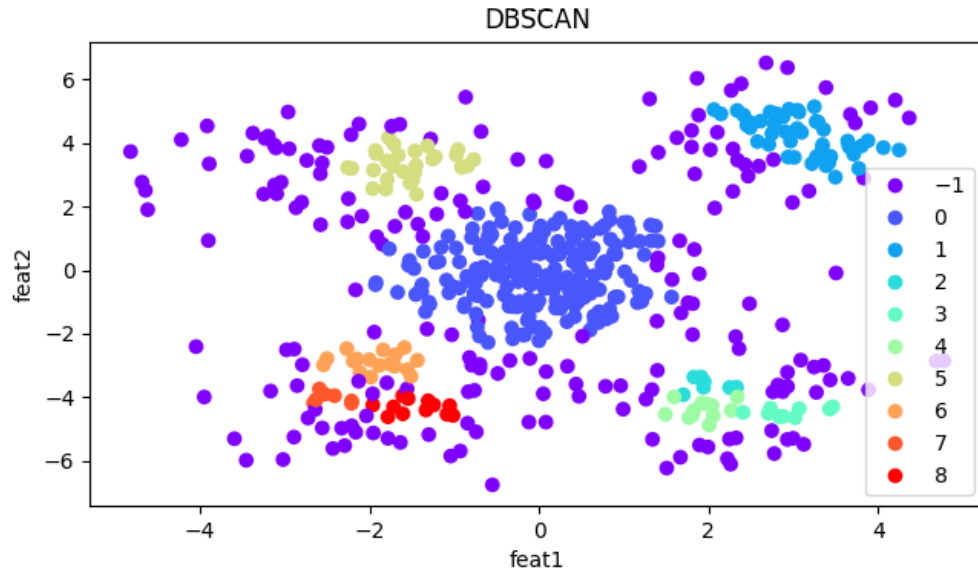
V tem poglavju predstavimo slike gručenj, ki niso bila med boljšimi v poglavju 5. V podpoglavju A.1 predstavimo preostala gručenja na podatkovni množici Krožne gruče, v podpoglavju A.2 preostala gručenja na podatkovni množici Trakovi 4-3 in v podpoglavju A.3 preostala gručenja na podatkovni množici MNIST.

### A.1 Krožne gruče

V tem podpoglavju predstavimo na sliki A.1 gručenje algoritma MDEC in na sliki A.2 gručenje algoritma DBSCAN na podatkovni množici Krožne gruče.



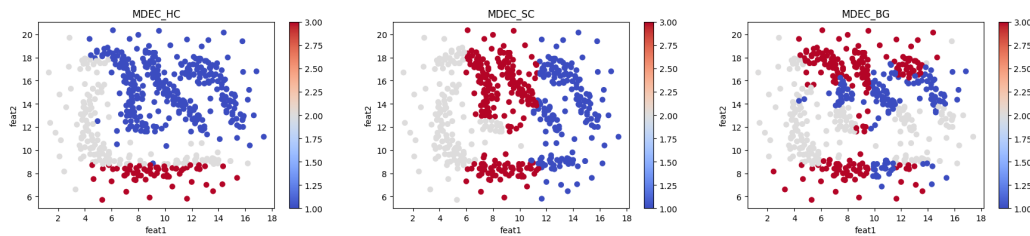
Slika A.1: Tri rezultati MDEC algoritma pri številu gruč  $k = 3$  na podatkovni množici Krožne gruče. Zanima nas najbolj desni rezultat z najboljšo hevristiko.



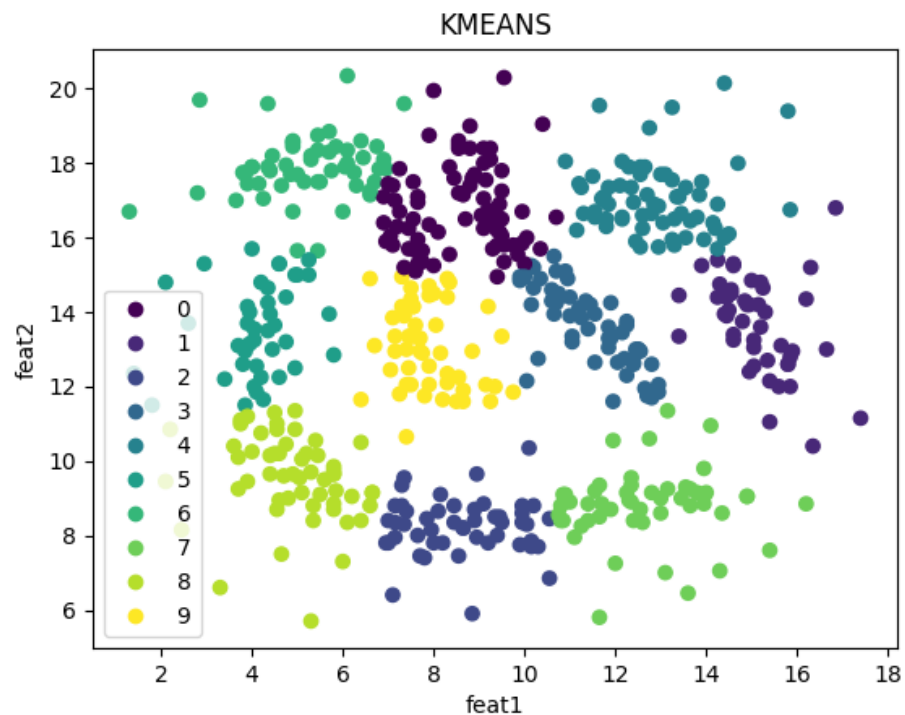
Slika A.2: Rezultat DBSCAN algoritma na podatkovni množici Krožne gruč.

## A.2 Trakovi 4-3

V tem podpoglavju predstavimo na sliki A.3 gručenje algoritma MDEC in na sliki A.4 gručenje algoritma K-MEANS na podatkovni množici Trakovi 4-3.



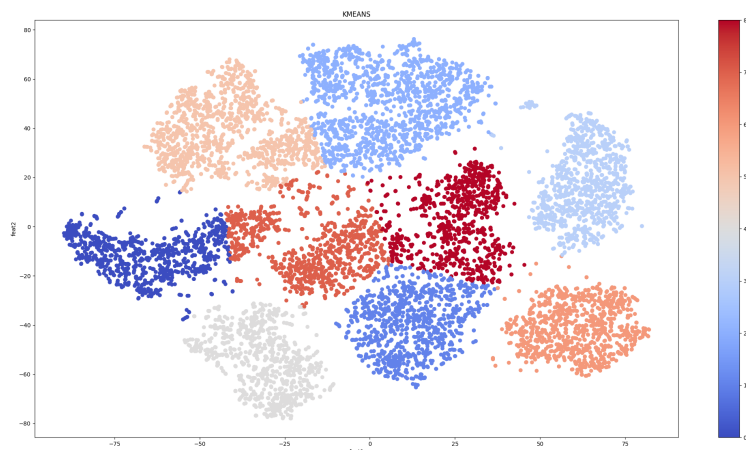
Slika A.3: Tri rezultati MDEC algoritma pri številu gruč  $k = 3$  na podatkovni množici Trakovi 4-3. Najbolj levi je dobil najboljšo oceno.



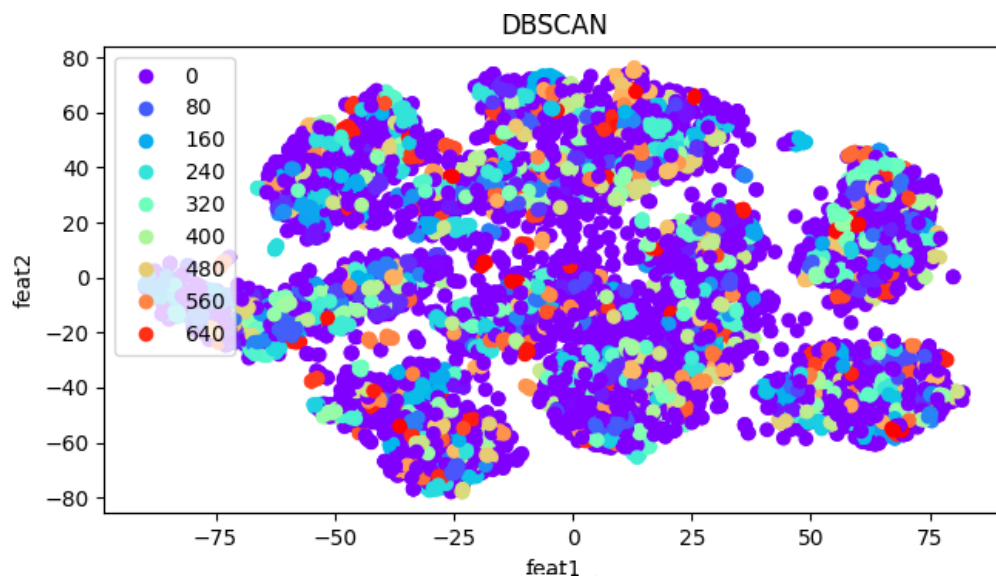
Slika A.4: Rezultat K-MEANS algoritma pri številu gruč  $k = 10$  na podatkovni množici Trakovi 4-3.

### A.3 MNIST

V tem podpoglavju predstavimo na sliki A.5 gručenje algoritma K-MEANS, na sliki A.6 gručenje algoritma DBSCAN in na sliki A.7 gručenje algoritma HDBSCAN na podatkovni množici MNIST.

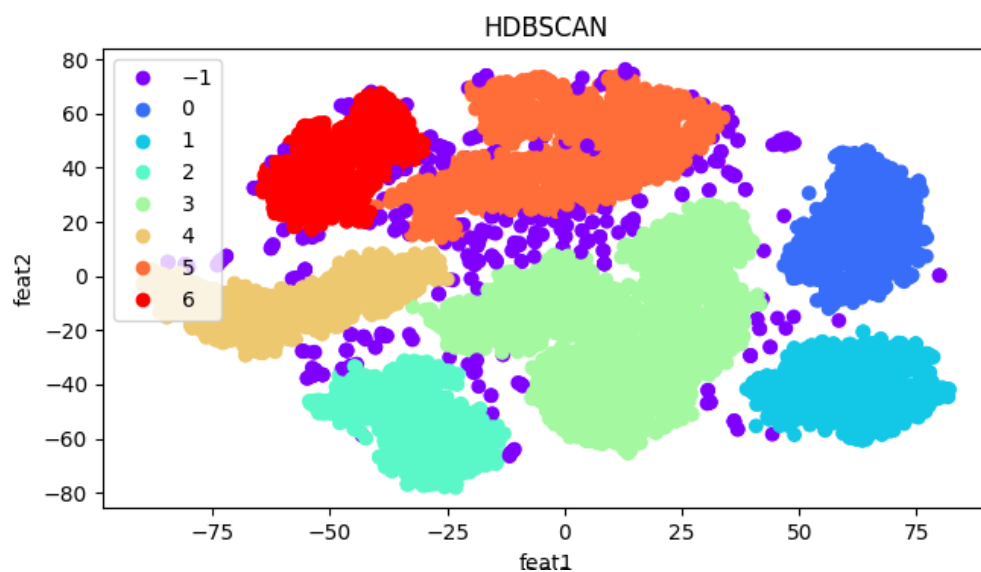


Slika A.5: Rezultat K-MEANS algoritma na podatkovni množici MNIST. Vidimo, da je algoritem dosegel najboljšo hevristiko silhuete pri številu gruč  $k = 9$ . Vse gruče niso smiselne, kar se vidi pri temno modri gruči.



Slika A.6: Rezultat DBSCAN algoritma na podatkovni množici MNIST. Algoritem ni deloval po pričakovanjih, kar se vidi po nerazumnih oznakah. Ločene gruče so neobstoječe. Naša izbira parametrov se v tem primeru ne obrestuje, saj je algoritem našel 640 gruč.





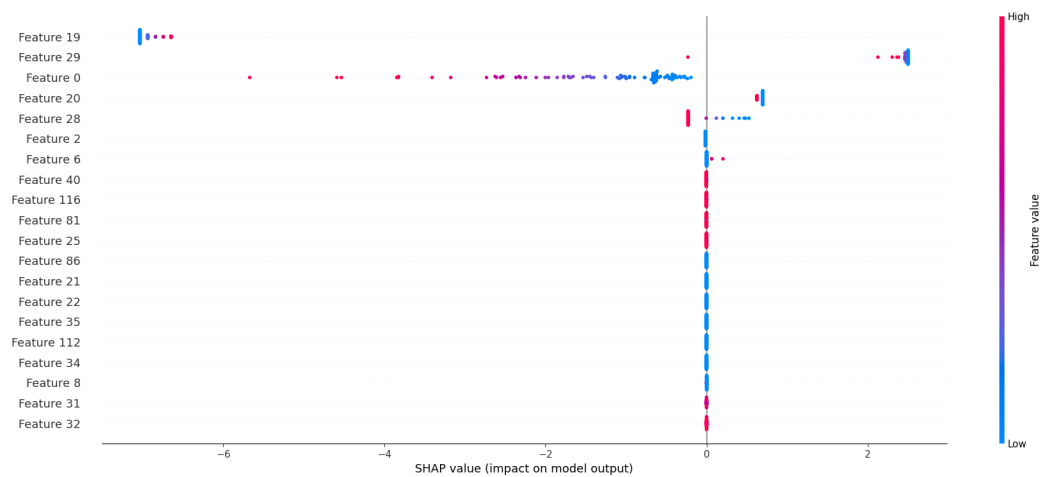
Slika A.7: Rezultat HDBSCAN algoritma na podatkovni množici MNIST. Algoritem je našel smiselne gruče (7), a ne vseh. Na sliki se vidi, da bi morali biti oranžna in zelena gruča razdeljeni na dodatne podgruče.



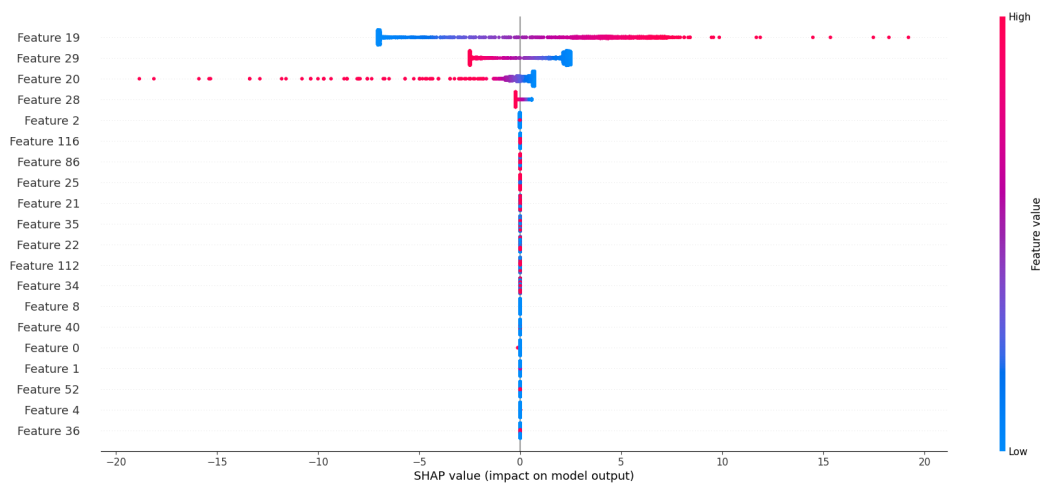
## Dodatek B

# Vrednosti SHAP preostalih gruč iz poglavja 6

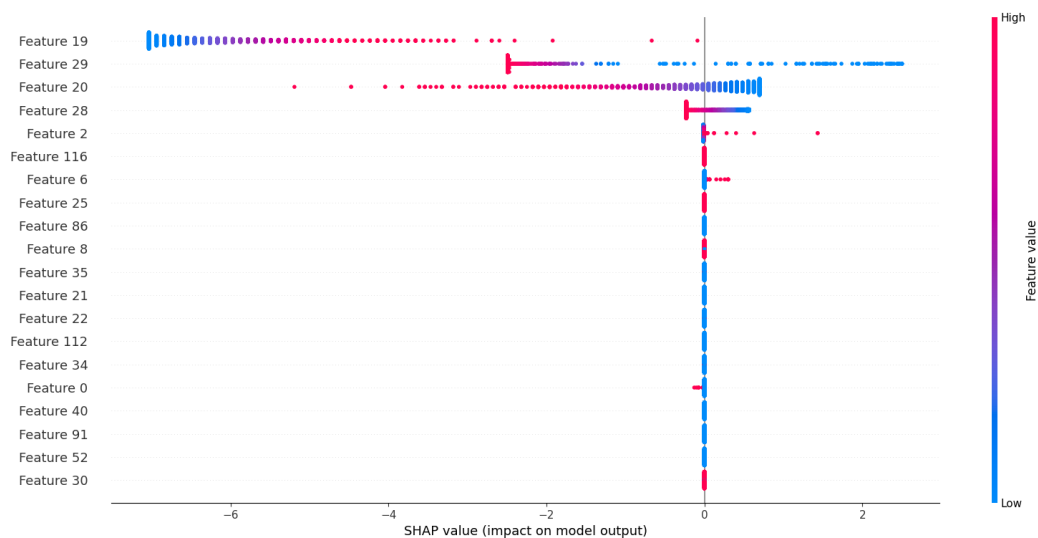
V tem poglavju predstavimo vrednosti SHAP preostalih skupin iz razlage podatkovne množice KDD99. Teh vrednosti SHAP nismo potrebovali za razlago novega primera.



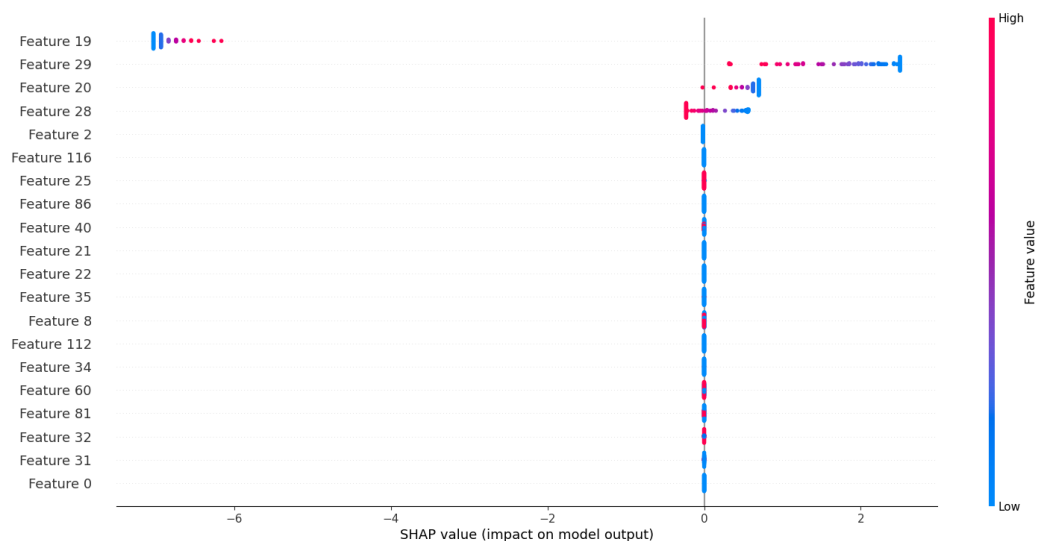
Slika B.1: Vrednosti SHAP za gručo 0 na podatkovni množici KDD99.



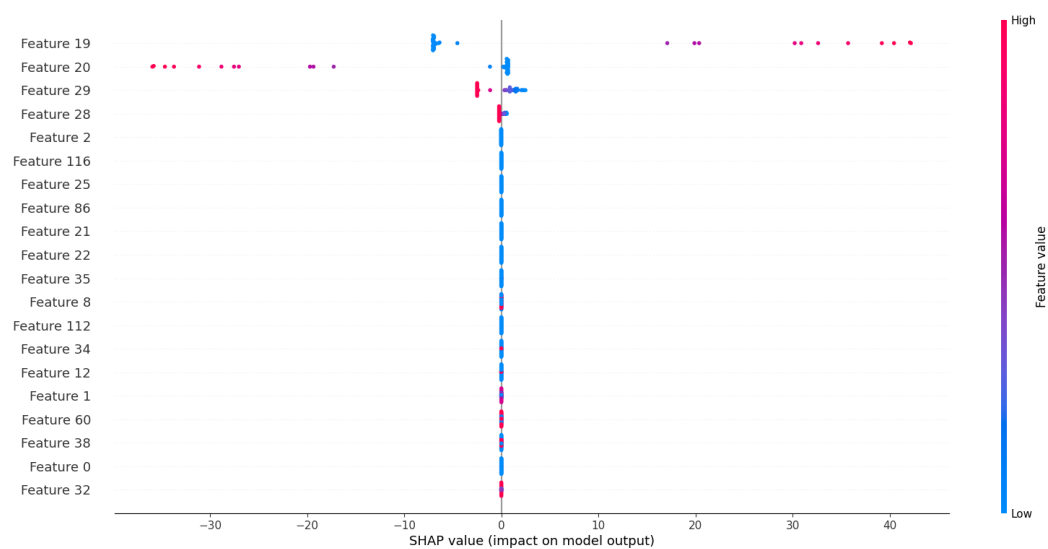
Slika B.2: Vrednosti SHAP za gručo 2 na podatkovni množici KDD99.



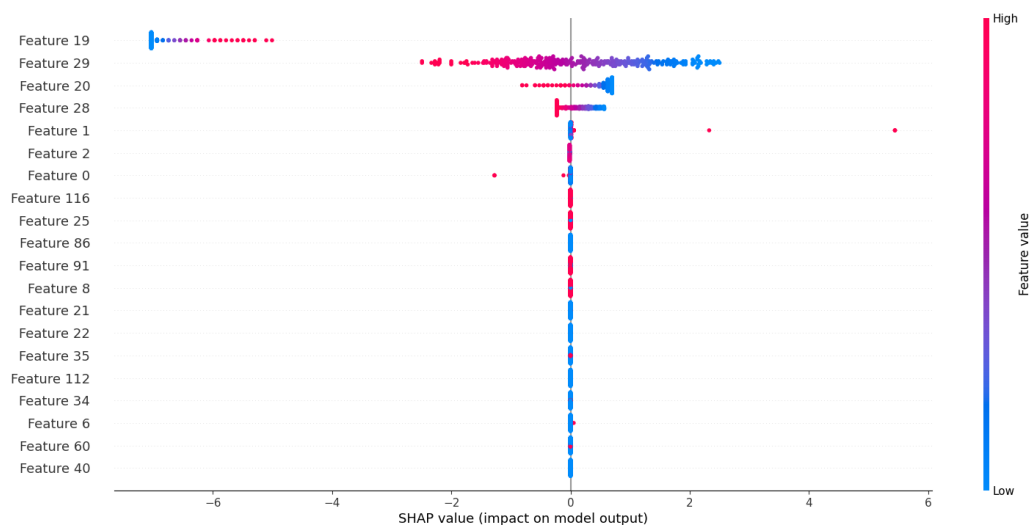
Slika B.3: Vrednosti SHAP za gručo 3 na podatkovni množici KDD99.



Slika B.4: Vrednosti SHAP za gručo 4 na podatkovni množici KDD99.



Slika B.5: Vrednosti SHAP za gručo 5 na podatkovni množici KDD99.



Slika B.6: Vrednosti SHAP za gručo 6 na podatkovni množici KDD99.

# Literatura

- [1] Chidanand Apté in Sholom Weiss. “Data mining with decision trees and decision rules”. V: *Future generation computer systems* 13.2-3 (1997), str. 197–210.
- [2] Jason Brownlee. “How to Develop a CNN for MNIST Handwritten Digit Classification”. V: *Machine Learning Mastery* (2019). Datum dosega 11.8.2023. URL: <https://machinelearningmastery.com/how-to-develop-a-convolutional-neural-network-from-scratch-for-mnist-handwritten-digit-classification/>.
- [3] Alfred DeMaris. “A tutorial in logistic regression”. V: *Journal of Marriage and the Family* (1995), str. 956–968.
- [4] Charles Frenzel. “How To Tune HDBSCAN”. V: *Towards Data Science* (2021). Datum dosega 11.8.2023. URL: <https://towardsdatascience.com/tuning-with-hdbscan-149865ac2970>.
- [5] Jerome H. Friedman in Bogdan E. Popescu. “Predictive learning via rule ensembles”. V: *The Annals of Applied Statistics* 2.3 (2008), str. 916–954. DOI: 10.1214/07-AOAS148.
- [6] John A Hartigan in Manchek A Wong. “Algorithm AS 136: A k-means clustering algorithm”. V: *Journal of the royal statistical society. series c (applied statistics)* 28.1 (1979), str. 100–108.

- [7] Dong Huang, Chang-Dong Wang, Jian-Huang Lai in Chee-Keong Kwoh. “Toward multidiversified ensemble clustering of high-dimensional data: From subspaces to metrics and beyond”. V: *IEEE Transactions on Cybernetics* 52.11 (2021), str. 12231–12244.
- [8] Kamran Khan, Saif Ur Rehman, Kamran Aziz, Simon Fong in Sababady Sarasvady. “DBSCAN: Past, present and future”. V: *The fifth international conference on the applications of digital information and web technologies (ICADIWT 2014)*. IEEE. 2014, str. 232–238.
- [9] Daniel Kleine. “Detecting knee- / elbow points in a graph of a function”. V: *Towards Data Science* (2021). Datum dosega 20.7.2023. URL: <https://towardsdatascience.com/detecting-knee-elbow-points-in-a-graph-d13fc517a63c>.
- [10] Shih-Wei Lin, Kuo-Ching Ying, Chou-Yuan Lee in Zne-Jung Lee. “An intelligent algorithm with feature selection and decision rules applied to anomaly intrusion detection”. V: *Applied Soft Computing* 12.10 (2012), str. 3285–3290.
- [11] Scott M Lundberg in Su-In Lee. “A unified approach to interpreting model predictions”. V: *Advances in neural information processing systems* 30 (2017).
- [12] Basim Mahbooba, Mohan Timilsina, Radhya Sahal in Martin Serrano. “Explainable artificial intelligence (XAI) to enhance trust management in intrusion detection systems using decision tree model”. V: *Complexity* 2021 (2021), str. 1–11.
- [13] Leland McInnes, John Healy in Steve Astels. “hdbscan: Hierarchical density based clustering.” V: *J. Open Source Softw.* 2.11 (2017), str. 205. DOI: 10.21105/joss.00205.
- [14] Edoardo Mosca, Ferenc Sziget, Stella Tragianni, Daniel Gallagher in Georg Groh. “SHAP-based explanation methods: a review for NLP interpretability”. V: *Proceedings of the 29th International Conference on Computational Linguistics*. 2022, str. 4593–4603.



- [15] Davoud Moulavi, Pablo A Jaskowiak, Ricardo JGB Campello, Arthur Zimek in Jörg Sander. “Density-based clustering validation”. V: *Proceedings of the 2014 SIAM international conference on data mining*. SIAM. 2014, str. 839–847.
- [16] Tara Mullin. “DBSCAN Parameter Estimation Using Python”. V: *Medium* (2020). Datum dosega 8.7.2023. URL: <https://medium.com/@tarammullin/dbscan-parameter-estimation-ff8330e3a3bd>.
- [17] Ashutosh Nayak. “XGBoost: An Intuitive Explanation”. V: *Towards Data Science* (2019). Datum dosega 15.8.2023. URL: <https://towardsdatascience.com/xgboost-an-intuitive-explanation-88eb32a48eff>.
- [18] Conor Nugent in Pádraig Cunningham. “A case-based explanation system for black-box systems”. V: *Artificial Intelligence Review* 24 (2005), str. 163–178.
- [19] Peter J. Rousseeuw. “Silhouettes: A graphical aid to the interpretation and validation of cluster analysis”. V: *Journal of Computational and Applied Mathematics* 20 (1987), str. 53–65. DOI: [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7).
- [20] Mohammad Khubeb Siddiqui in Shams Naahid. “Analysis of KDD CUP 99 dataset using clustering based data mining”. V: *International Journal of Database Theory and Application* 6.5 (2013), str. 23–34.
- [21] Saurabh Singh. “Building an Intrusion Detection System using KDD Cup’99 Dataset”. V: *Medium* (2020). Datum dosega 16.8.2023. URL: <https://medium.com/analytics-vidhya/building-an-intrusion-detection-model-using-kdd-cup99-dataset-fb4cba4189ed>.
- [22] Laurens Van der Maaten in Geoffrey Hinton. “Visualizing data using t-SNE.” V: *Journal of machine learning research* 9.11 (2008).

- [23] Remah Younisse, Ashraf Ahmad in Qasem Abu Al-Haija. “Explaining Intrusion Detection-Based Convolutional Neural Networks Using Shapley Additive Explanations (SHAP)”. V: *Big Data and Cognitive Computing* 6.4 (2022), str. 126.