

Rīgas tehniskā universitāte

Mākslīgie intelekta pamati

2.praktiskais darbs

<https://www.kaggle.com/datasets/teertha/ushealthinsurancedataset?resource=download>

<https://github.com/DZUDMENS/PD2.git>

Izstrādāja: Gvido Urniežus

211RDB436 6.grupa

Saturs

1 daļa – Datu pirmapstrāde/izpēte.....	3
2.daļa-Nepārraudzītā mašīnmācīšanās.....	11
3.daļa-Pārraudzītā mašīnmācīšanās	19
Orange datu shemas attēlojums.....	23
Izmantotā avoti.....	23

1 daļa – Datu pirmapstrāde/izpēte

Lai izpildītu šī darba daļu, studentiem ir jāveic šādas darbības:

1. Ir jāizvēlas un jāapraksta datu kopa, pamatojoties uz informāciju, kas sniegta krātuvē, kurā datu kopa ir pieejama.

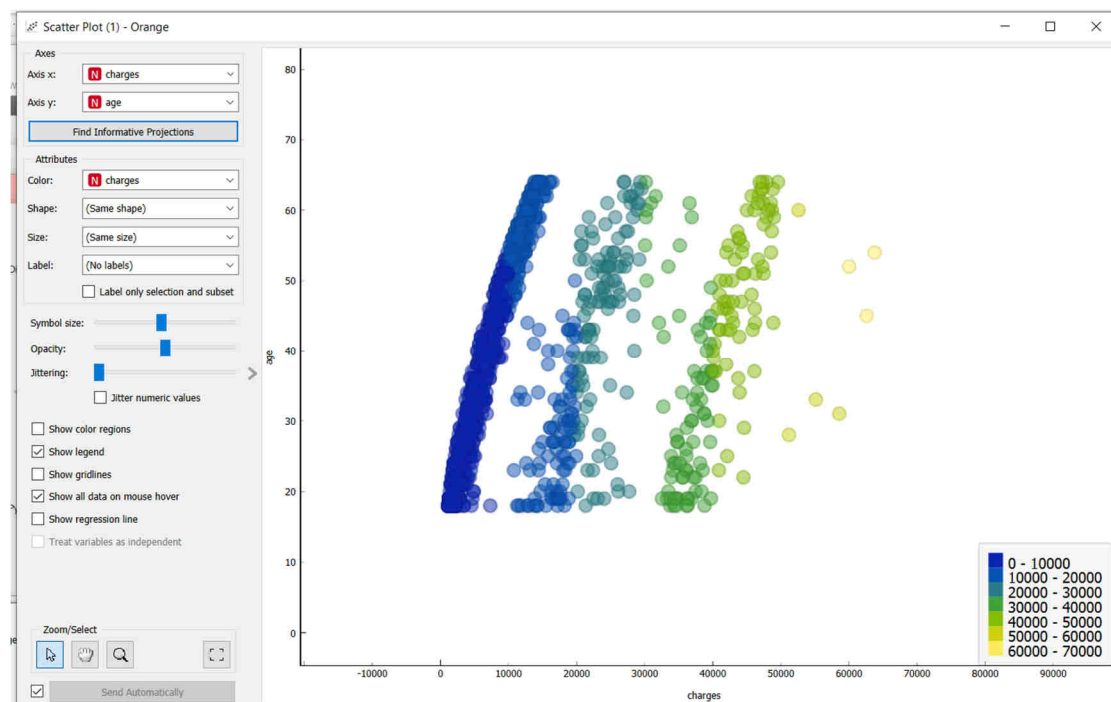
2. Ja no krātuves iegūtā datu kopa nav formātā, ar kuru ir viegli strādāt (piemēram, komatatzīmītas vērtības vai .csv fails), ir jāveic tās transformācija vajadzīgajā formātā.

3. Ja kādu pazīmju (atribūtu) vērtības ir tekstveida vērtības (piemēram, yes/no, positive/neutral/negative, u.c.), tās ir jātransformē skaitliskās vērtībās.

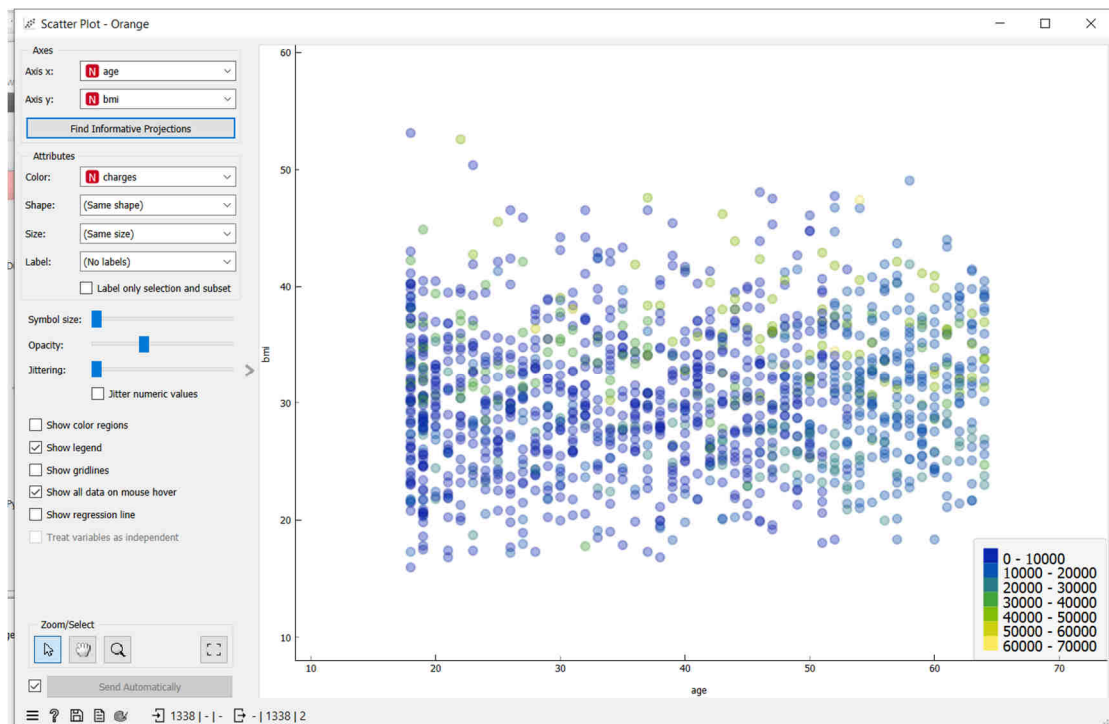
4. Ja kādiem datu objektiem trūkst atsevišķu pazīmju (atribūtu) vērtības, ir jāatrod veids, kā tās iegūt, studējot papildu informācijas avotus.

5. Ir jāatspoguļo datu kopa vizuāli un jāaprēķina statistiskie rādītāji:

a) ir jāizveido vismaz divas 2- vai 3-dimensiju izkliedes diagrammas (scatter plot), kas ilustrē klases atdalāmību, balstoties uz dažādām pazīmēm (atribūtiem); studentam ir jāizvairās izmantot datu objekta ID vai klases iezīmi kā mainīgo izkliedes diagrammā;

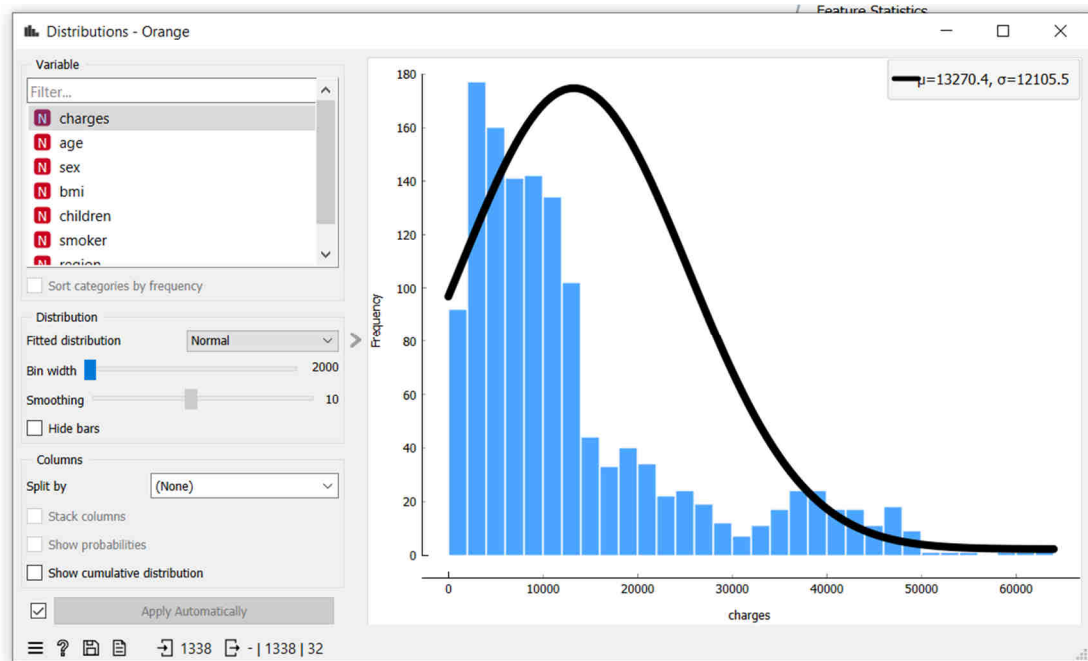


(1.att)

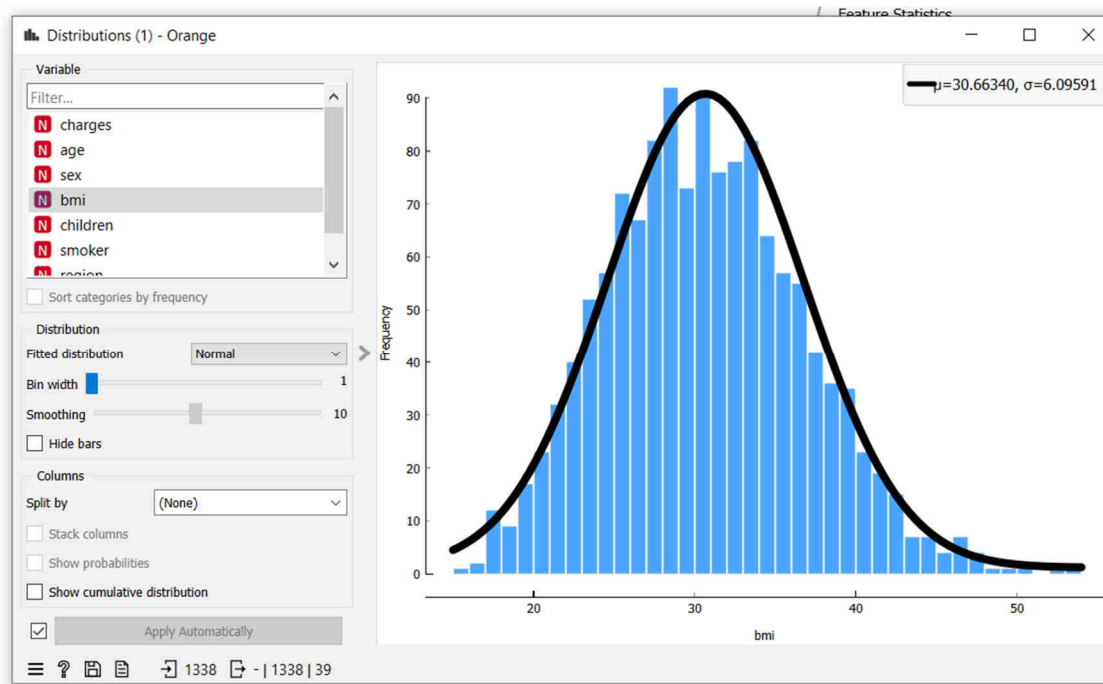


(2.att)

b) ir jāizveido vismaz 2 histogrammas, kas parāda klašu atdalīšanu, pamatojoties uz interesējošām pazīmēm (atribūtiem);



(3.att)



(4.att)

c) ir jāatspoguļo 2 interesējošo pazīmju (atribūtu) sadalījums;

The screenshot shows the 'Correlations - Orange' window. The 'Pearson correlation' method is selected, and '(All combinations)' is chosen for the filter. A list of 17 correlations is shown, each with a value, a variable, and another variable. The first correlation is +0.787 between 'charges' and 'smoker'. The rest of the correlations are listed in the table below.

Rank	Correlation	Variable 1	Variable 2
1	+0.787	charges	smoker
2	+0.299	age	charges
3	+0.198	bmi	charges
4	+0.158	bmi	region
5	+0.109	age	bmi
6	+0.076	sex	smoker
7	+0.068	charges	children
8	+0.057	charges	sex
9	+0.046	bmi	sex
10	+0.042	age	children
11	-0.025	age	smoker
12	-0.021	age	sex
13	+0.017	children	sex
14	+0.017	children	region
15	+0.013	bmi	children
16	+0.008	children	smoker
17	-0.006	charges	region

(5.att)

d) ir jāaprēķina statistiskie rādītāji (vismaz vidējās vērtības un dispersiju).



(6.att)paskaidrojumu skatīties pie 9.att.

1.Daļas atskaite

Datu kopas apraksts (sniedzot arī atsaucis izmantotajiem informācijas avotiem)

- Datu kopas nosaukums, avots, izveidotājs un/vai īpašnieks;

Nosaukums-US Health Insurance Dataset

Avots - <https://www.kaggle.com/datasets/teertha/ushealthinsurancedataset?resource=download>

Izveidotājs-Anirban Datta

- Datu kopas problēmsfēras apraksts;

Šī datu kopa var būt noderīga vienkāršā, bet izgaismojošā pētījumā, lai izprastu riska parakstīšanu veselības apdrošināšanā, dažādu apdrošinātā atribūtu mijiedarbību un noskaidrotu, kā tie ietekmē apdrošināšanas prēmiju.

- Datu kopas licencēšanas nosaucījumi (jā tādi ir);

<https://creativecommons.org/publicdomain/zero/1.0/>

- Veids, kā datu kopa tika savākta;

Nav norādīts.

Datu kopas satura apraksts (sniedzot arī atsauces uz izmantotajiem informācijas avotiem)

- Datu objektu skaits datu kopā;

1338 kopējie ieraksti.

- Datu kopas pazīmju (atribūtu) atspoguļojums kopā ar to lomām Orange rīkā;

Info				
1338 instances				
7 features (no missing values)				
Data has no target variable.				
0 meta attributes				
Columns (Double click to edit)				
	Name	Type	Role	Values
1	age	N numeric	feature	
2	sex	C categorical	feature	female, male
3	bmi	N numeric	feature	
4	children	N numeric	feature	
5	smoker	C categorical	feature	no, yes
6	region	C categorical	feature	northeast, northwest, southeast, southwest
7	charges	N numeric	target	

(7.att)Šajā attēlā ir norādīts cik daudz parametri ir doti Dataset.

- Klašu skaits datu kopā, katras klases nozīme un klašu atspoguļošanas veids (klasēm atbilstošo iezīmju skaidrojums); ja datu kopa nodrošina vairākas iespējamās datu klasifikācijas, tad atskaitē skaidri ir jāidentificē, kāda tieši klasifikācija tiek apskatīta darbā;

Sex-šī kategorija norāda kāds dzimums ir šai personai.

Smoker-šī kategorija norāda vai persona pīpe tabakas izstrādājumus vai nē

Region-šī kategorija norāda kādā reģionā US dzīvo.

- Datu objektu skaits, kas pieder katrai klasei;

Sex-male(vīrietis),female(sieviete)

Smoker-yes(jā),no(nē)

Region-

northeast(ziemeļaustrumi),northwest(ziemeļrietumi),southeast(dienvidaustrumos),southwest(dien
vidrietumos)

- Pazīmju (atribūtu) skaits un nozīme datu kopā, kā arī to vērtību tipi un diapazoni (ši informācija būtu jāatspoguļo tabulā, norādot pazīmes (atribūta) apzīmējumu, skaidrojumu, vērtību tipu un datu kopā pieejamo vērtību diapazonu);

Nr.	Name(nosaukums)	Type(tips)	Skaidrojums	Role(loma)	Values(vērtības)
1	Age	Skaitlisks	Vecums	feature	Skaitlis(18-64)
2	Sex	Kategorija	Dzimums	feature	(Female,male)(0,1)
3	Bmi(body mass index)	Skaitlisks	Ķermeņa masas indekss	feature	Skaitlis(16-53,1)
4	Children	Skaitlisks	Cik bērni ir cietušajam	feature	Skaitlis(0-5)
5	Smoker	Kategorija	Tabakas lietošana	feature	(Yes,no)(1,0)
6	Region	Kategorija	Reģions	feature	(northeast, northwest Southeast, southwest (0,1,2,3)
7	Charges	Skaitlisks	Izmaksas	target	Skaitlis(1120-63800)

- Datu faila struktūras fragments, kurā ir redzamas visas datu faila kolonnas un to vērtības vismaz dažiem datu objektiem;

	charges	age	sex	bmi	children	smoker	region
1	16884.9	19	0	27.900	0	1	3
2	1725.55	18	1	33.770	1	0	2
3	4449.46	28	1	33.000	3	0	2
4	21984.5	33	1	22.705	0	0	1
5	3866.86	32	1	28.880	0	0	1
6	3756.62	31	0	25.740	0	0	2
7	8240.59	46	0	33.440	1	0	2
8	7281.51	37	0	27.740	3	0	1
9	6406.41	37	1	29.830	2	0	0
10	28923.1	60	0	25.840	0	0	1
11	2721.32	25	1	26.220	0	0	0
12	27808.7	62	0	26.290	0	1	2
13	1826.84	23	1	34.400	0	0	3
14	11090.7	56	0	39.820	0	0	2
15	39611.8	27	1	42.130	0	1	2
16	1837.24	19	1	24.600	1	0	3
17	10797.3	52	0	30.780	1	0	0
18	2395.17	23	1	23.845	0	0	0
19	10602.4	56	1	40.300	0	0	3
20	36837.5	30	1	35.300	0	1	3
21	13228.8	60	0	36.005	0	0	0
22	4149.74	30	0	32.400	1	0	3

(8.att)Šajā attēlā ir norādīts data table kurā ir atspoguļoti dati jau pārveidoti uz numerācija kā piem. (yes ir 1, un no ir 0).

Secinājumi, kas izriet no izkliedes diagrammu, histogrammu un sadalījumu analīzes (sk. I daļas 5. solis) par datu kopas klašu atdalāmību. Studentiem ir jāatbild uz šādiem jautājumiem:

- Vai klases datu kopā ir līdzsvarotas, vai dominē viena klase (vai vairākas klases)? Tas tiek noteikts, spriežot pēc tā, cik daudz datu objektu pieder katrai kopai.

Nevarētu teikt ka visas klases ir līdzsvarotas jo tomēr ir bula tipa klases, šajā dataset dominē viena klase kas ir (charges) kas norāda cik daudz patērētu līdzekļu tika izmatoti uz personu veselības apdrošināšanā.

- Vai datu vizuālais atspoguļojums ļauj redzēt datu struktūru? Runa ir par to, vai datu objekti, kuri pieder dažādām klasēm, ir skaidri atdalāmi.

Šajā tabulā ir atspoguļojama datu struktūra, jo katru klasi var skaidri atdalīt.

- Cik datu grupējums ir iespējams identificēt, pētot datu vizuālo atspoguļojumu? Runa ir par to, vai ir kaut cik atdalāmi datu grupējumi, ja gadījumā dažādu klašu datu objekti saplūst kopā.

Ir iespēja atdalīt datu grupējumus un labāk parskatīt datus.

- Vai identificētie datu grupējumi atrodas tuvu viens otram vai tālu viens no otra?

Identificētie datu grupējumi atrodas viens otram tuvu.

Secinājumi, kas izriet no statistisko rādītāju (vidējo vērtību un dispersijas vērtību) analīzes.



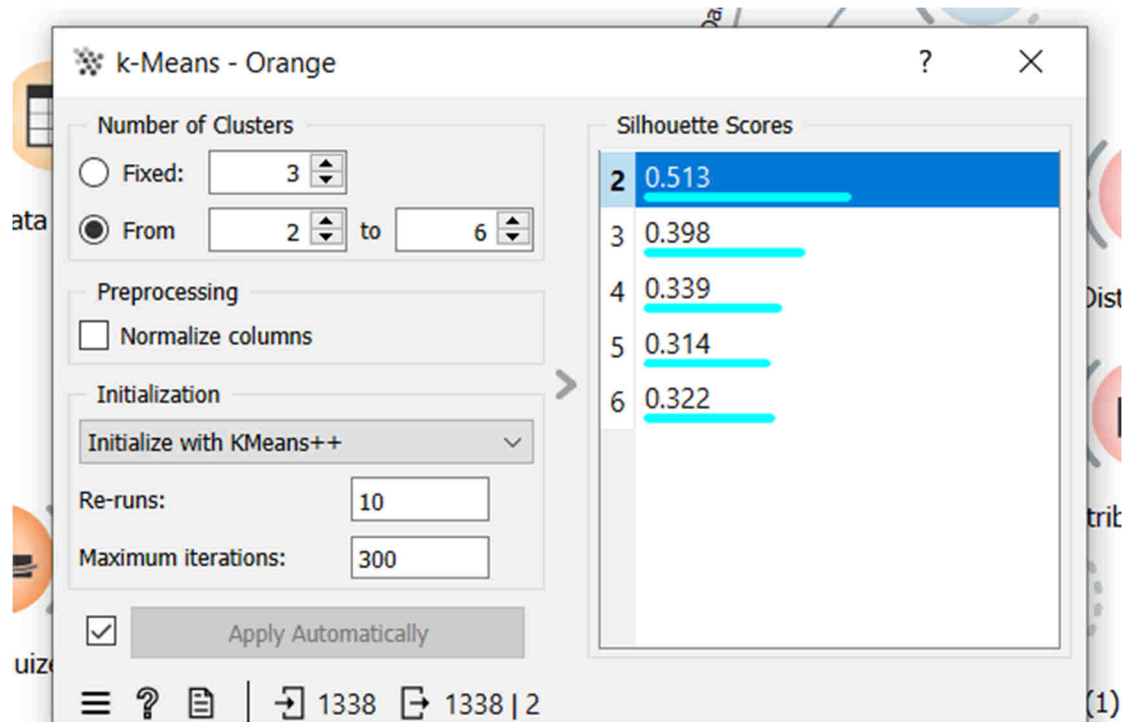
(9.att)

Šajā statistiskajā radījumā ir aprēķināta videja vērtība un dispersija. Vidējais vecums ir 39 gadi un vidējie patērētie līdzekļi ir 13270.4\$ kas ir uz pusi sievietēm un uz pusi vīriešiem 1 vienu bērnu ģimenē. Dispersija patērētajiem līdzekļiem ir 0.912216.

2.daļa-Nepārraudzītā mašīnmācīšanās

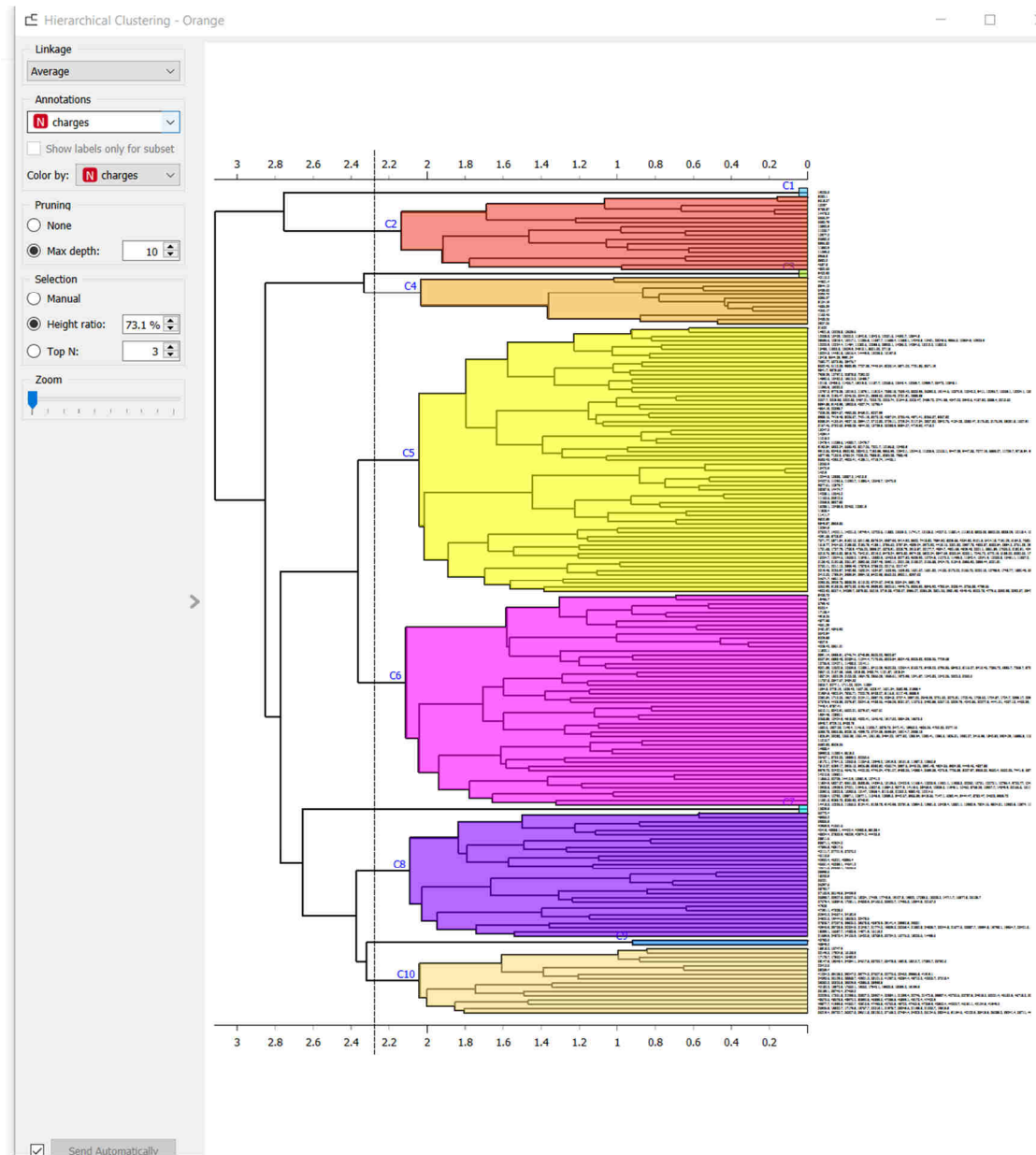
Lai izpildītu šo darba daļu, studentiem ir jāveic šādas darbības:

1. Jāpielieto divi studiju kursā apskatītie nepārraudzītās mašīnmācīšanās algoritmi: (1) hierarhiskā klasterizācija un (2) K-vidējo algoritms.

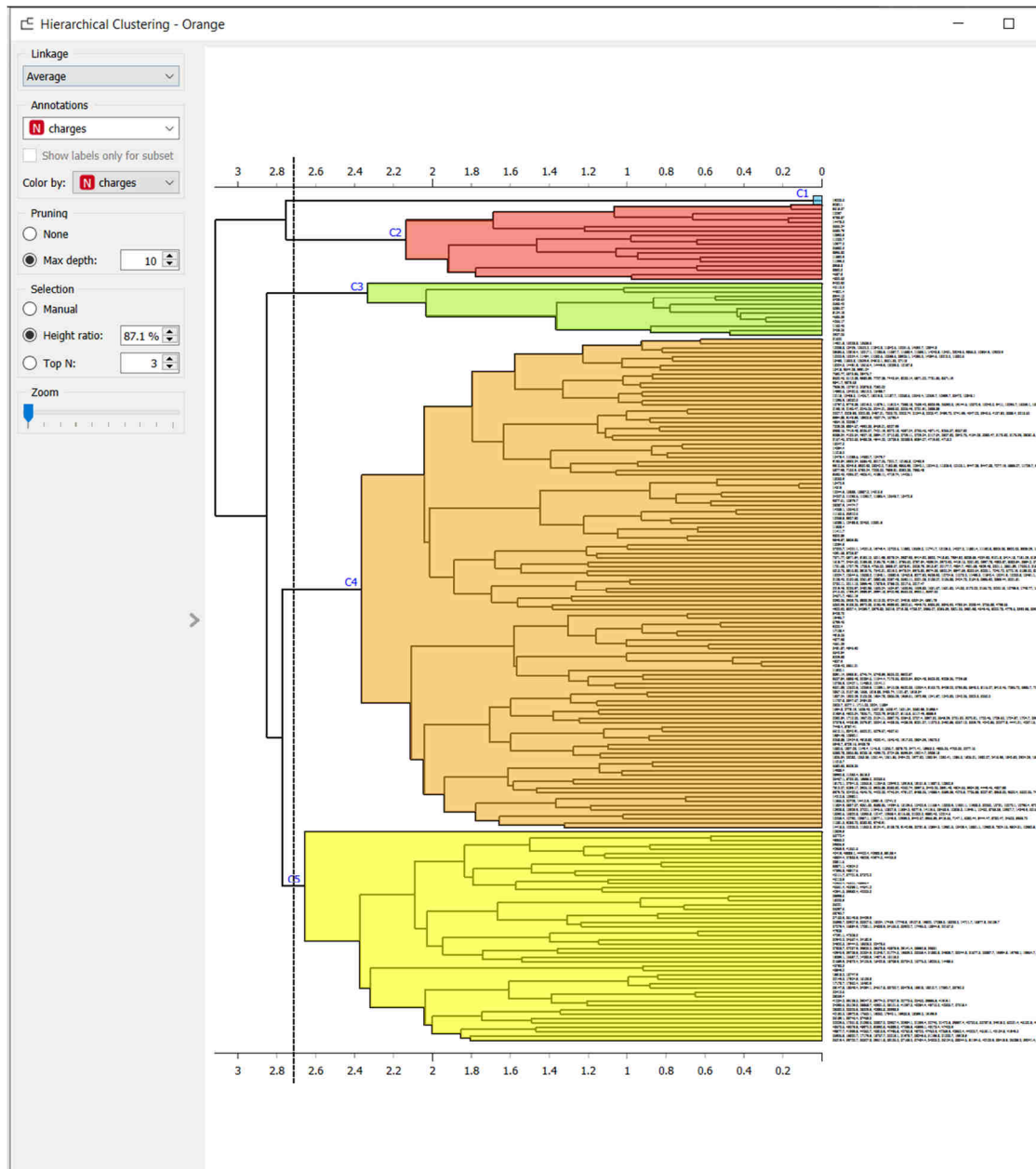


(10.att)

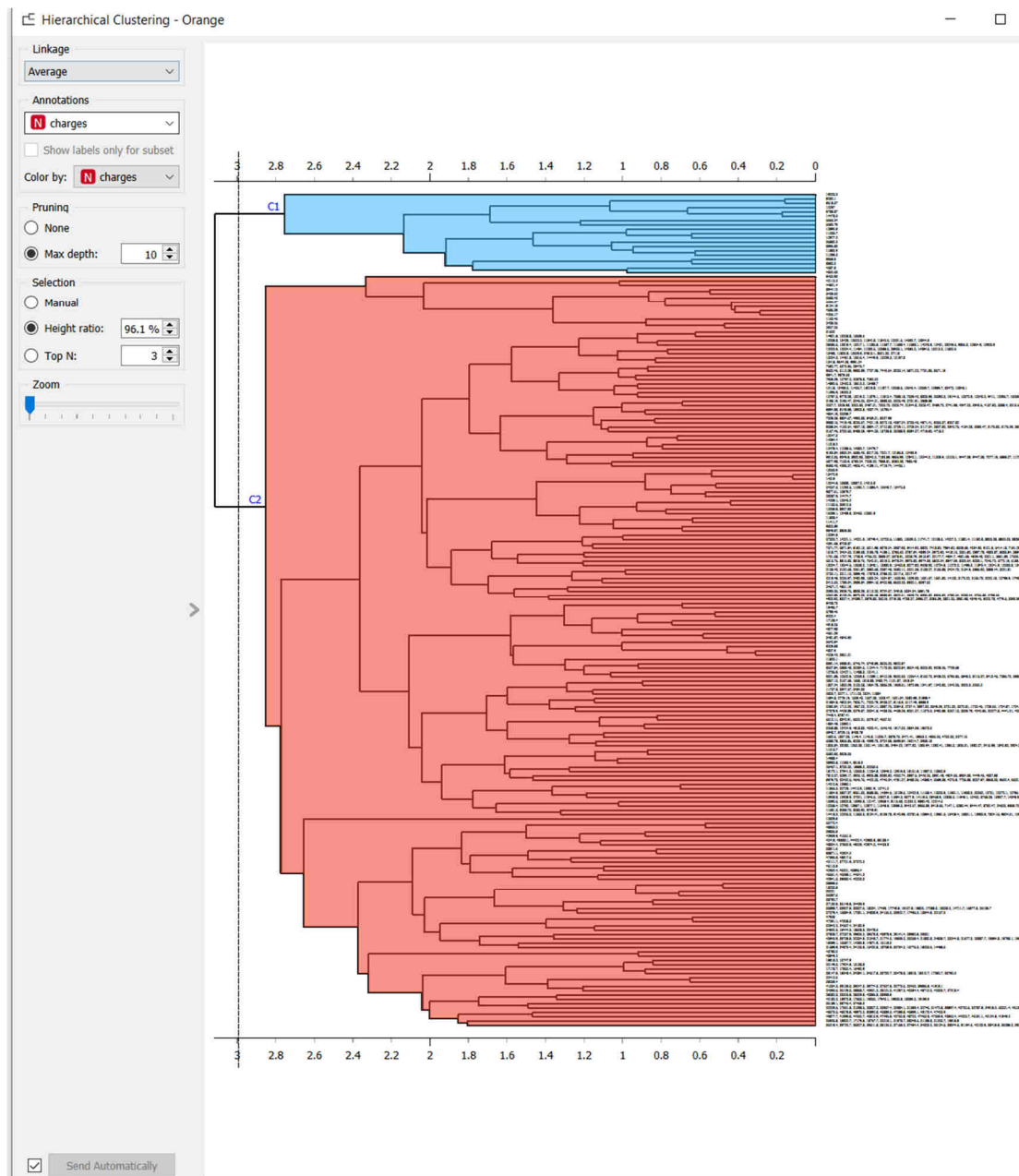
2. Hierarhiskās klasterizācijas algoritmam ir jāveic vismaz 3 eksperimenti, brīvi pārvietojot atdalāmo līniju un analizējot, kā mainās klasteru skaits un saturs;



(11.att)



(12.att)

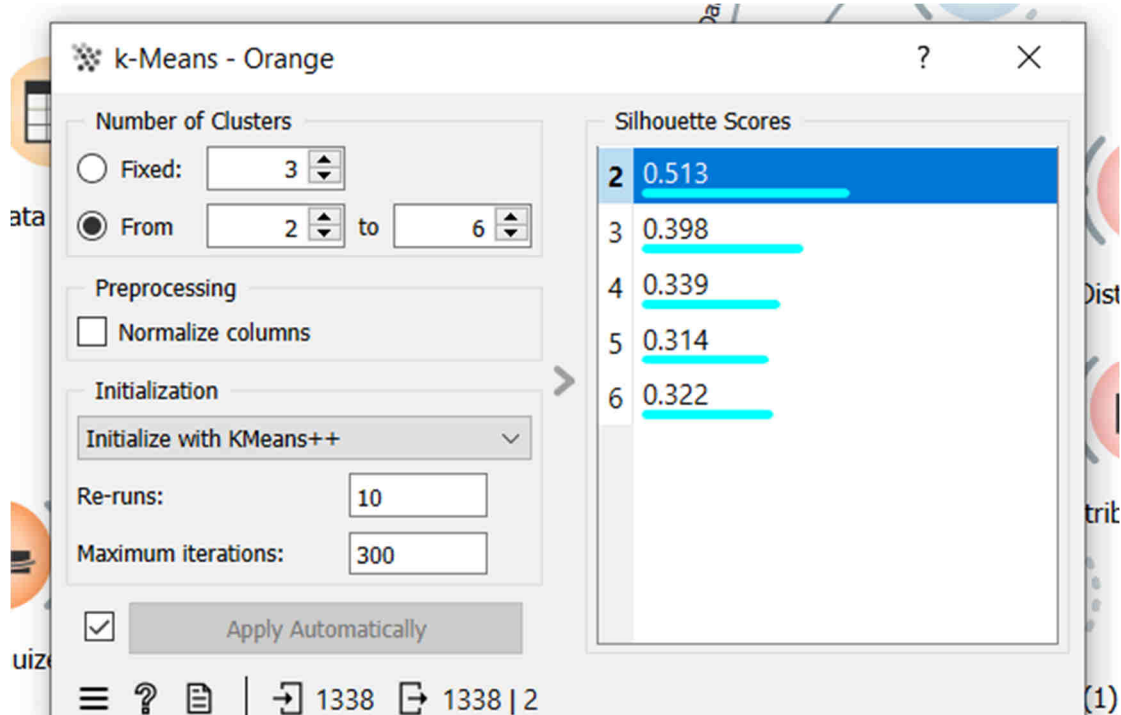


(13.att)

3. K-vidējo algoritmam ir jāaprēķina Silhouette Score vismaz 5 dažādām k vērtībām, un jāanalizē algoritma darbība.

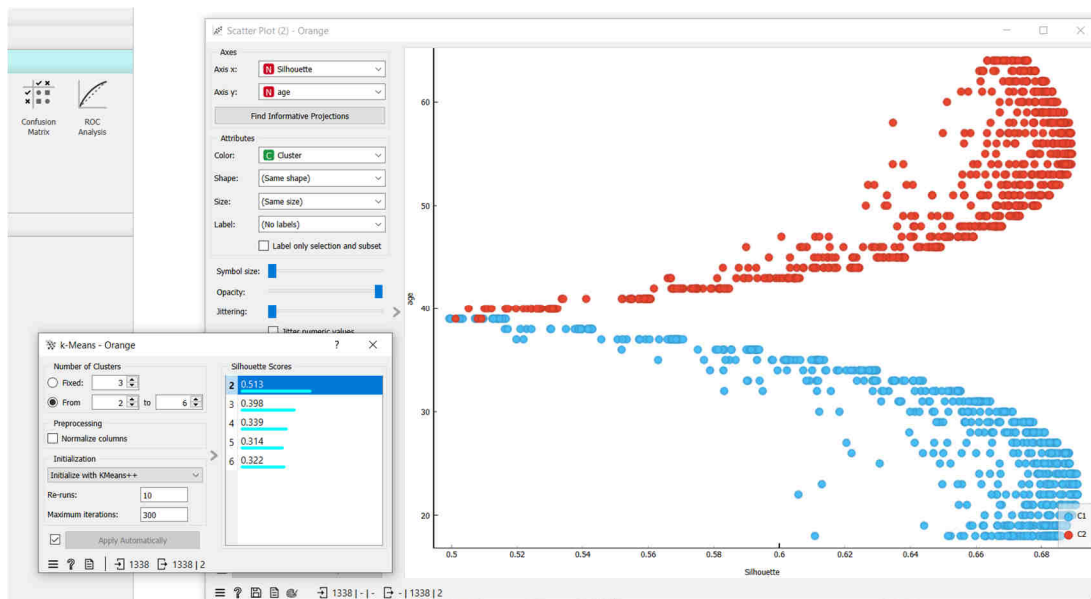
Darba atskaitē ir jāiekļauj šāda informācija par šo darba daļu:

- Katram algoritmam ir jāapraksta Orange rīkā pieejamie hiperparametri un to nozīme.

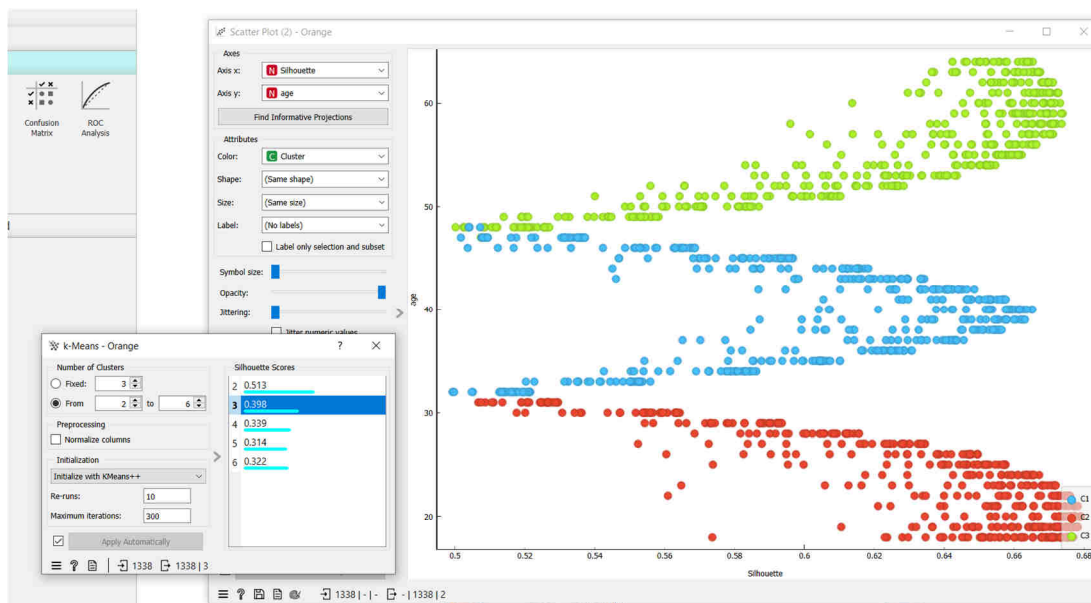


(14.att)

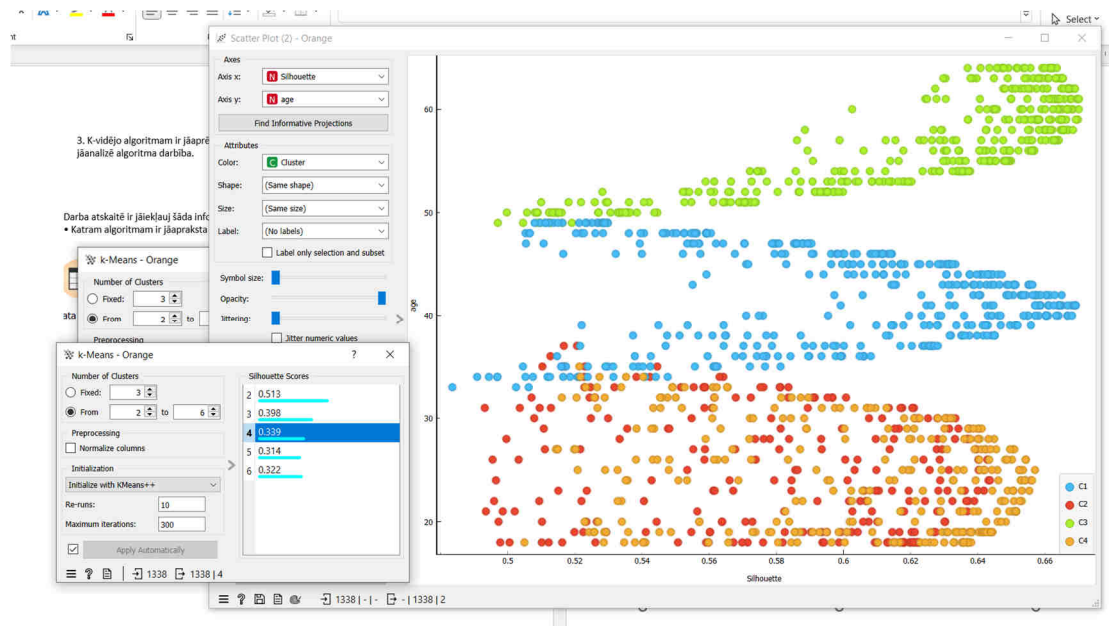
- Katram algoritmam ir jāapraksta veiktie eksperimenti, skaidri norādot izmantotās hiperparametru vērtības, un sniedzot secinājumus par algoritma darbību no tā viedokļa, cik iegūtie rezultāti atbilst zināmajam klašu skaitam datu kopā.



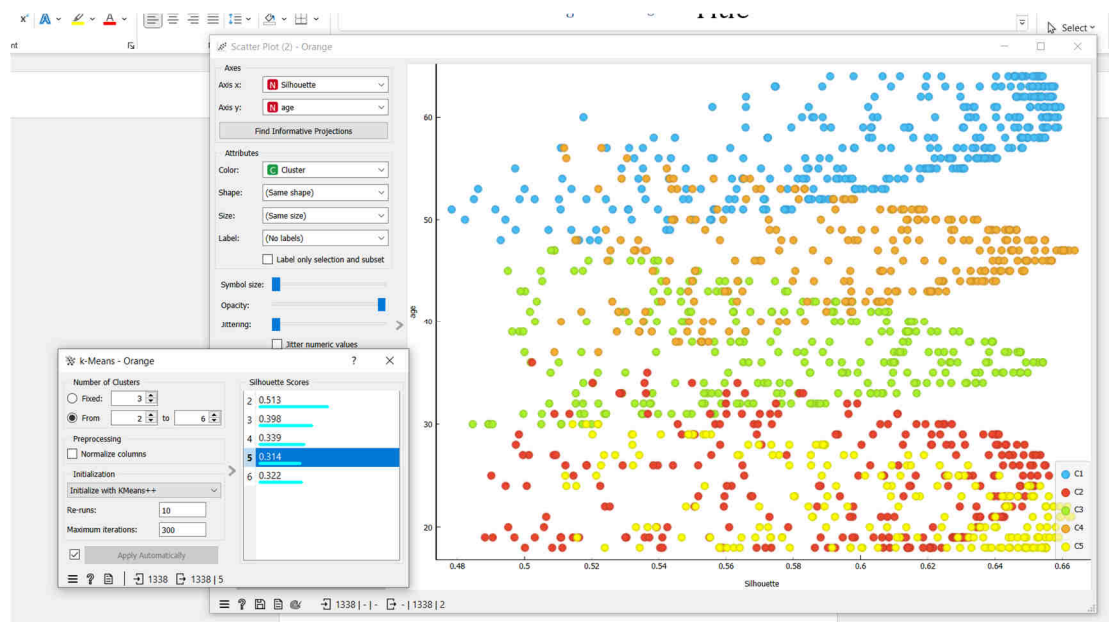
(15.att)



(16.att)



(17.att)



(18.att)skatoties uz šo eksperimentu un 19.att eksperimentu var saprast cik būtiski ir tomēr K-Means score, jo atšķiroties tikai par 0.008 grafiks izmainās.



(19.att)

• Balstoties uz abu algoritmu darbības analīzi, ir jādod studenta secinājumi par to, vai datu kopā esošās klases ir labi vai slikti atdalāmas.

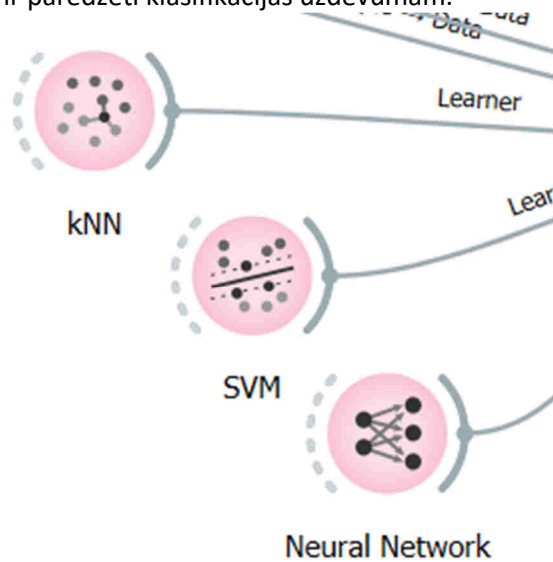
Balstoties uz abu algoritmiem ir slikti atdalāmas, jo parak daudz cluster.

3.daļa-Pārraudzītā mašīnmācīšanās

Lai izpildītu šo darba daļu, studentiem ir jāveic šādas darbības:

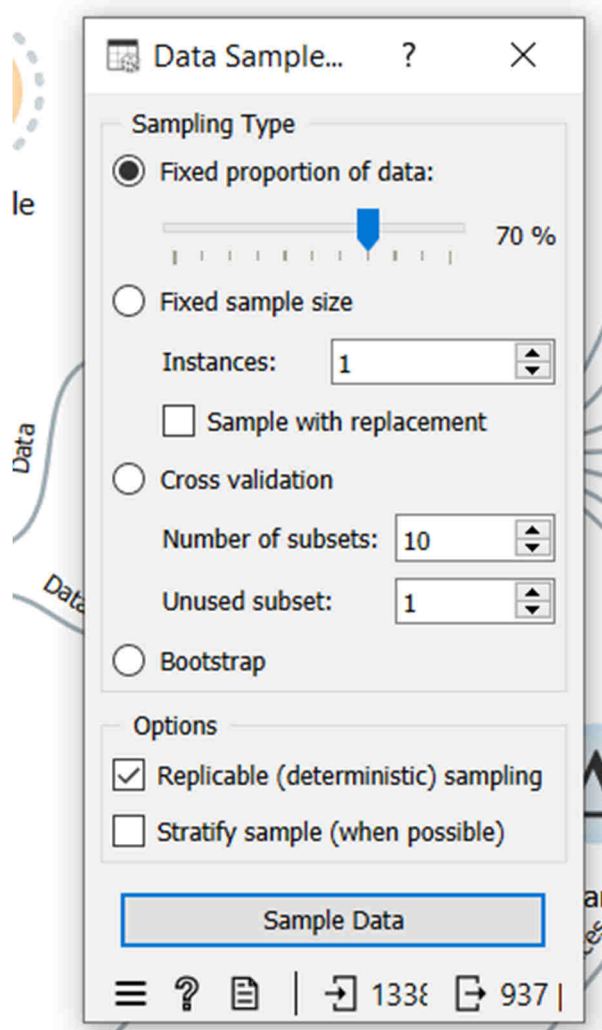
1. Ir jāizvēlas vismaz divi pārraudzītās mašīnmācīšanās algoritmi, kas ir paredzēti klasifikācijas uzdevumam. Studenti drīkst izmantot studiju kursā aplūkotos algoritmus vai arī jebkurus citus algoritmus, kuri ir paredzēti klasifikācijas uzdevumam.

Data Sampler



(20.att)

2. Ir jāsadala datu kopa apmācību un testa datu kopās.



(21.att)

3. Katram algoritmam, lietojot apmācību datu kopu, ir jāveic vismaz 3 eksperimenti, mainot algoritma hiperparametru vērtības un analizējot algoritmu veikspējas metrikas;
4. Katram algoritmam ir jāizvēlas tas apmācītais modelis, kas nodrošina labāko algoritma veikspēju;
5. Katra algoritma apmācītais modelis ir jāpielieto testa datu kopai.

The screenshot displays the Orange3 data mining software interface. The 'Test and Score' window is the primary focus, showing the results of training and testing three models: kNN, Neural Network, and SVM. The 'Cross validation' section is active, with 'Number of folds' set to 5 and 'Stratified' checked. The 'Test on test data' option is also selected. The results table shows the following metrics:

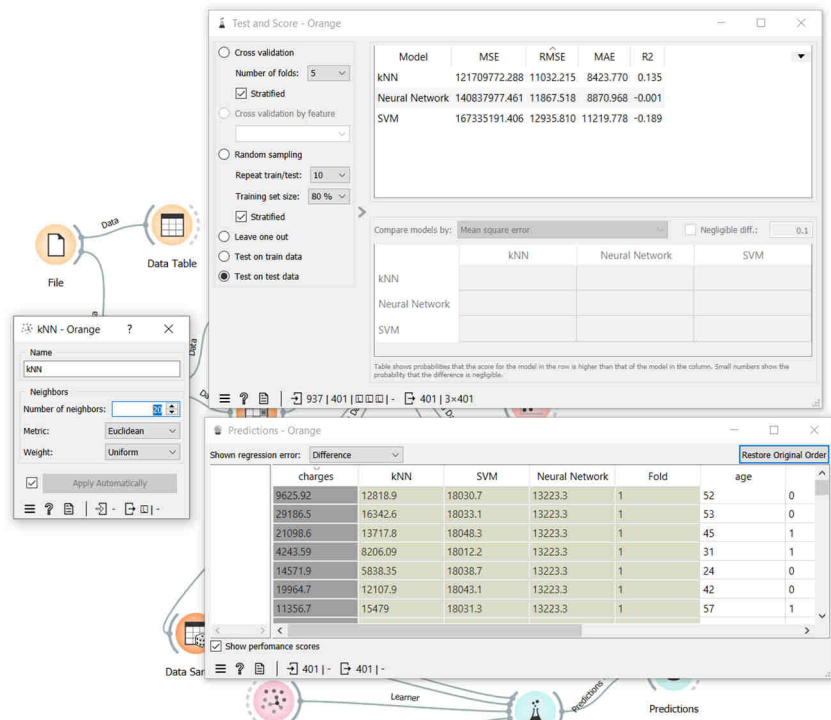
Model	MSE	RMSE	MAE	R2
kNN	122678825.575	11076.047	8168.992	0.128
Neural Network	140837977.461	11867.518	8870.968	-0.001
SVM	167335191.406	12935.810	11219.778	-0.189

The 'Compare models by' section shows a comparison of the models based on the 'Mean square error' metric. The 'Predictions' window shows the results of applying the trained models to the test data. The 'Shown regression error' is set to 'Difference'. The table shows the predicted values for the 'age' variable for each model and the 'Fold' used for training.

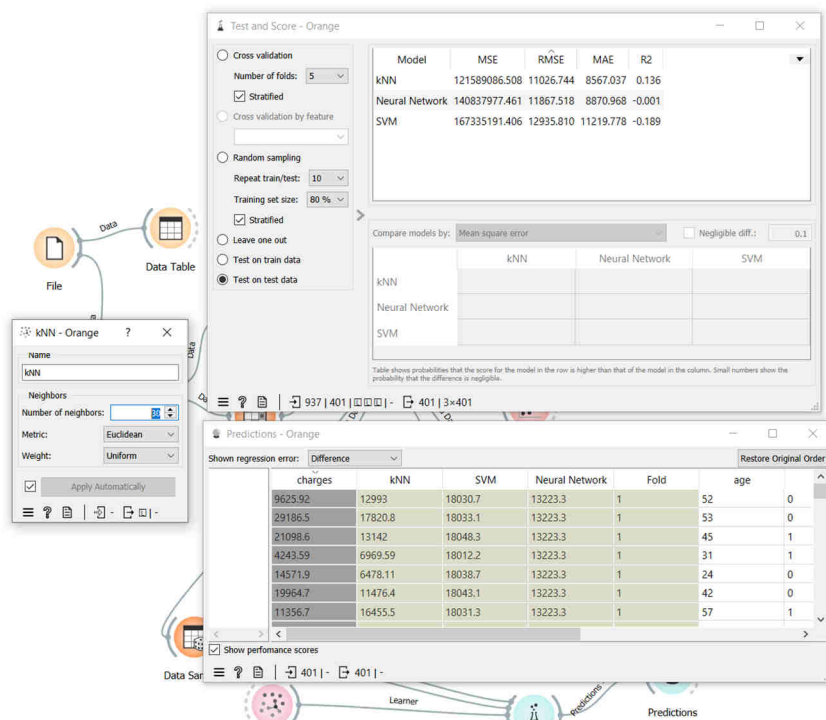
charges	kNN	SVM	Neural Network	Fold	age
9625.92	13119.4	18030.7	13223.3	1	52
29186.5	13986.1	18033.1	13223.3	1	53
21098.6	12293.6	18048.3	13223.3	1	45
4243.59	7454.31	18012.2	13223.3	1	31
14571.9	6198.83	18038.7	13223.3	1	24
19964.7	13135.9	18043.1	13223.3	1	42
11356.7	15805.8	18031.3	13223.3	1	57

The 'Show performance scores' checkbox is checked. The 'Predictions' window also shows the 'age' variable and the 'Fold' used for training.

(22.att)



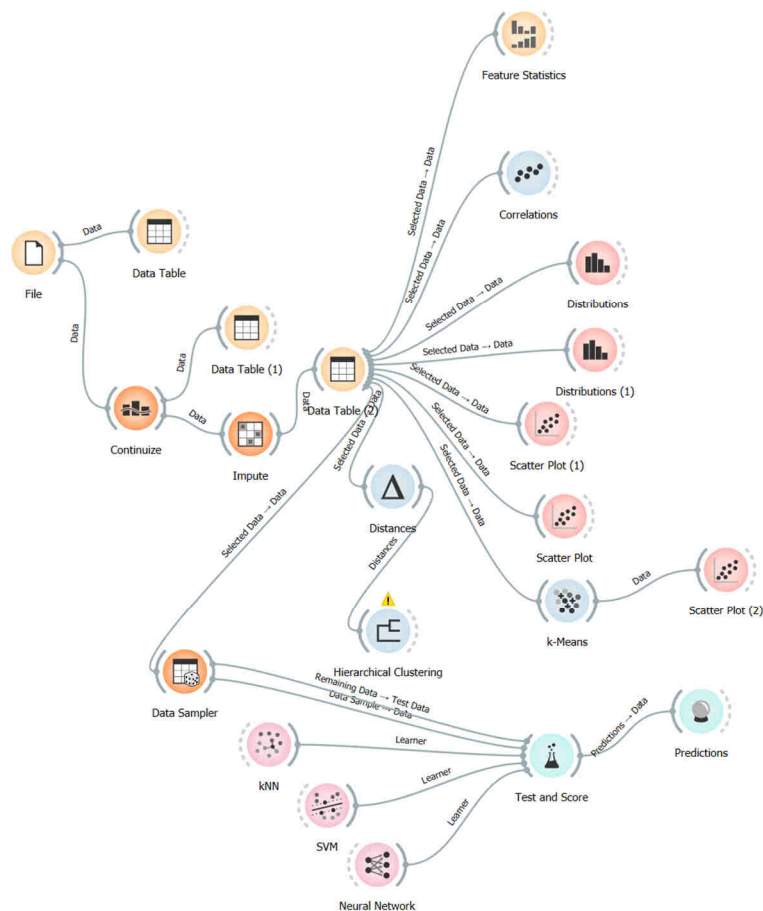
(23.att)



(24.att)

6. Ir jānovērtē un jāsalīdzina apmācīto modeļu veiktspēja. Skatoties uz 22 23 un 24 attēlu un mainot knn datus mainījās beigu prediction dati un testa dati.

Orange datu shemas attēlojums



Izmantotā avoti

<https://estudijas.rtu.lv/course/view.php?id=252548>

<https://www.kaggle.com/datasets/teertha/ushealthinsurancedataset?resource=download>

<https://orangedatamining.com/>