

Proiect Explorarea Datelor Data Mining Analiza setului de date - Dermatology

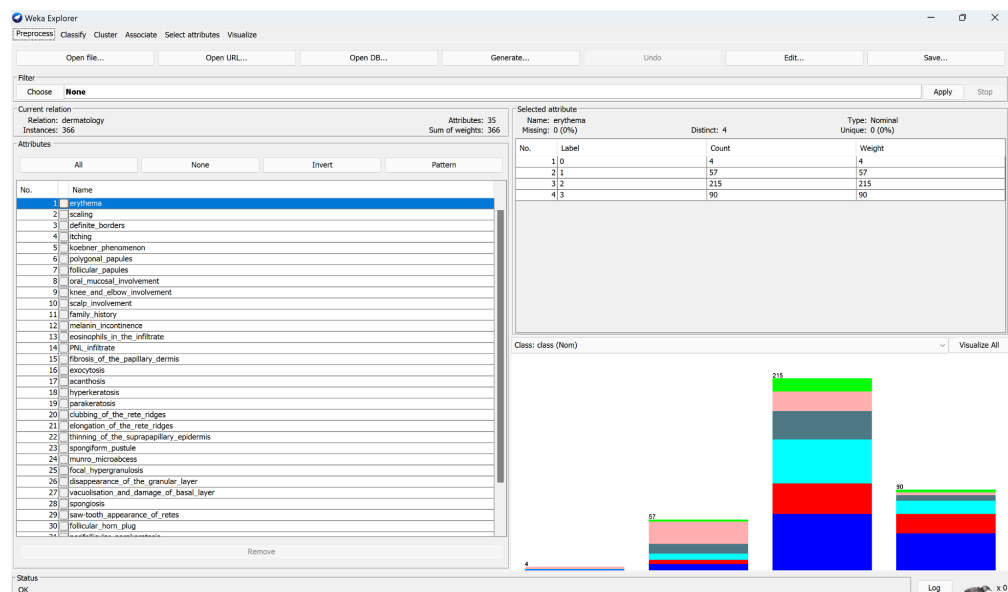
1. Descrierea setului de date

Am selectat un set de date denumit Dermatology, care furnizează informații esențiale despre diverse afecțiuni ale pielii. Acest set de date conține 34 de caracteristici, dintre care 33 sunt valori numerice continue, iar una este nominală. Numărul total de instanțe este de 366.

În cadrul acestui set de date, caracteristica de istoric familial este reprezentată printr-o valoare de 1 dacă afecțiuni cutanate au fost observate în familie și de 0 în caz contrar. Vârsta pacientului este înregistrată direct ca un număr întreg.

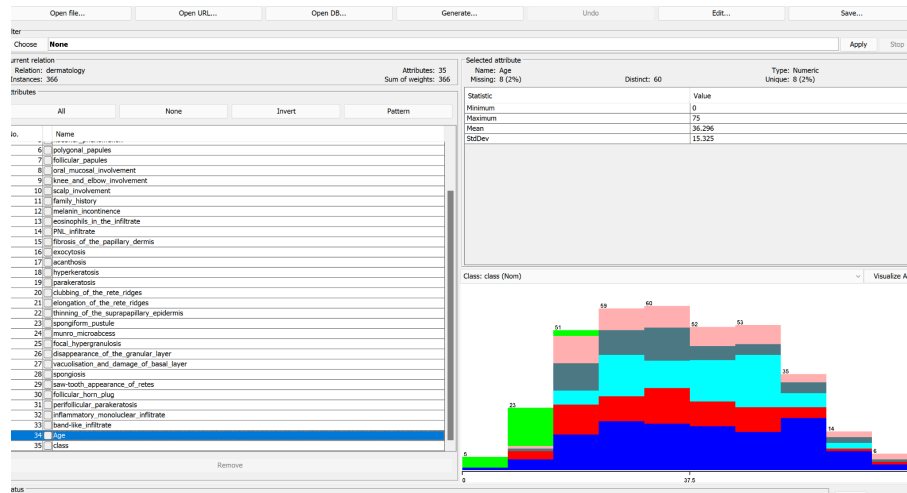
Celelalte caracteristici, fie clinice, fie histopatologice, sunt evaluate pe o scară de la 0 la 3. Valorile acestei scale indică absența caracteristicii (0), prezența caracteristicii într-o măsură scăzută (1 sau 2), sau prezența caracteristicii într-o măsură semnificativă (3).

Pentru a facilita analiza, am transformat inițial setul de date din formatul său original în format .csv folosind Microsoft Office Excel. Acest format este recunoscut automat de către Weka, însă pentru a utiliza setul de date în Weka, acesta va fi convertit în formatul arff (Attribute-Relation File Format) prin intermediul pachetului Weka.

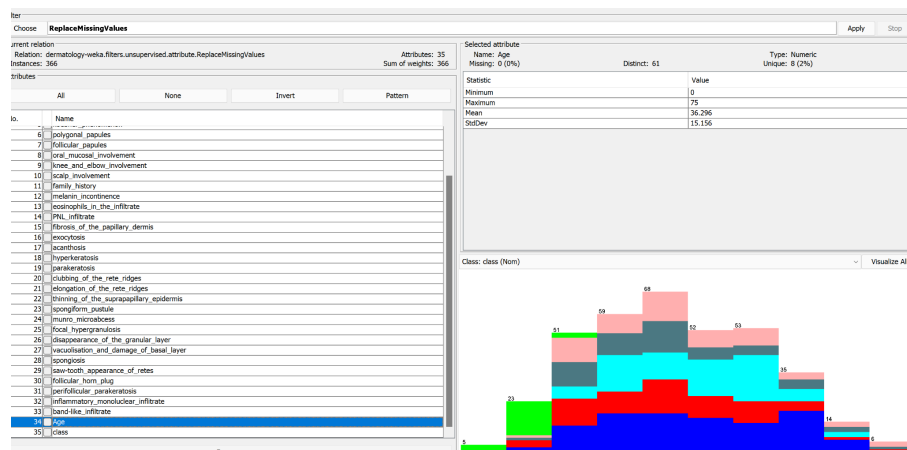


2. Analiza setului de atribute

În cadrul setului de date selectat, am identificat prezența unor valori lipsă, pe care trebuie să le înlocuim corespunzător.



Pentru aceasta, vom utiliza opțiunea "ReplaceMissingValues" din cadrul filtrului disponibil în Weka (Choose -> filters -> unsupervised -> attribute -> ReplaceMissingValues). După aplicarea acestui filtru, toate valorile lipsă vor fi înlocuite cu valoarea medie sau moda corespunzătoare atributului în cauză.



Rezultatul va fi salvat într-un nou fișier numit "dermatologyReplacedMissingValues.arff", care va fi utilizat în continuare pentru analiză.

Următorul pas constă în identificarea celor mai importante atribute din setul de date. Aceasta se va realiza folosind funcționalitatea "Select attributes". După selectarea metodei de căutare și a evaluatorului potrivit, vom iniția procesul. Rezultatele obținute vor evidenția cele 19 cele mai importante atribute, plus atributul "class", totalizând 20 de atribute relevante.

reprocess Classify Cluster Associate Select attributes Visualize

Attribute Evaluator

Choose **CfsSubsetEval** -P 1 -E 1

Search Method

Choose **BestFirst** -D 1 -N 5

Attribute Selection Mode

☒ Use full training set

☐ Cross-validation Folds Seed

to class

Start Stop

Result list (right-click for options)

7:38:32 - BestFirst + CfsSubsetEval

Attribute selection output

```

=== Attribute Selection on all input data ===

Search Method:
  Best first.
  Start set: no attributes
  Search direction: forward
  Stale search after 5 node expansions
  Total number of subsets evaluated: 506
  Merit of best subset found: 0.769

Attribute Subset Evaluator (supervised, Class (nominal): 35 class):
  CFS Subset Evaluator
  Including locally predictive attributes

Selected attributes: 2,3,4,5,7,9,13,14,15,16,20,21,22,25,26,28,29,31,33 : 19
scaling
definite_borders
itching
koebner_phenomenon
follicular_papules
knee_and_elbow_involvement
eosinophils_in_the_infiltrate
PNL_infiltrate
fibrosis_of_the_papillary_dermis
exocytosis
clubbing_of_the_rete_ridges
elongation_of_the_rete_ridges
thinning_of_the_suprapapillary_epidermis
focal_hypergranulosis
disappearance_of_the_granular_layer
spongiosis
saw-tooth_appearance_of_retes
perifollicular_parakeratosis
band-like_infiltrate

```

În continuare, vom salva acest nou set de date, eliminând toate celelalte atribute care nu se află printre cele mai importante. Pentru aceasta, vom utiliza metoda "BestFirst" pentru căutare, iar opțiunea "CfsSubsetEval" pentru evaluarea valorii fiecărui subset de atribute, luând în considerare atât capacitatea lor individuală de predictibilitate, cât și gradul de redundanță între ele.

Trebuie să notăm că trei dintre cele mai importante atribute identificate sunt: "Scaling", "Define_borders" și "Itching". Acestea sunt toate atribute nominale care fac parte din categoria "Clinical Attributes", având valori în intervalul de la 0 la 3 (dacă nu este altfel specificat).

3. Analiza datelor

Tehnică:

- **Clasificarea**

Clasificarea este o tehnică predictivă folosită în învățarea supervizată, unde se prezic clasele pentru instanțele noi pe baza unor predictorii. Printre metodele comune de clasificare se numără algoritmul OneR și arborii de decizie, rețelele neuronale sau algoritmi genetici.

OneR

Algoritmul OneR este o metodă simplă care creează o singură regulă de clasificare bazată pe o singură trăsătură. Alegerea trăsăturii se face pe baza maximizării probabilității de clasificare. Pentru fiecare atribut, se construiește un clasificator separat, rezultând în mai multe clasificatoare, iar cel mai semnificativ este ales.

În practică, pentru a aplica algoritmul OneR în Weka, se folosește opțiunea Percentage Split pentru a diviza setul de date într-un set de antrenare și unul de testare (Classify -> Choose -> classifiers -> rules -> OneR). De obicei, se alege un procentaj mai mare pentru setul de antrenare pentru a asigura o cantitate suficientă de date pentru învățare. Noi am ales 75%.

Test options

- ☐ Use training set
- ☐ Supplied test set
- ☐ Cross-validation Folds: 10
- ☒ Percentage split %: 75

(Nom) class

Start Stop

Result list (right-click for options)

- 17:50:12 - rules.OneR

Classifier output

Time taken to build model: 0.02 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0.01 seconds

=== Summary ===

Correctly Classified Instances	44	48.3516 %
Incorrectly Classified Instances	47	51.6484 %
Kappa statistic	0.3286	
Mean absolute error	0.1722	
Root mean squared error	0.4149	
Relative absolute error	64.3867 %	
Root relative squared error	113.2415 %	
Total Number of Instances	91	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
1.000	0.161	0.744	1.000	0.853	0.790	0.919	0.744	1	
0.000	0.000	?	0.000	?	?	0.500	0.154	2	
1.000	0.487	0.288	1.000	0.448	0.385	0.757	0.288	3	
0.000	0.000	?	0.000	?	?	0.500	0.187	4	
0.000	0.000	?	0.000	?	?	0.500	0.110	5	
0.000	0.000	?	0.000	?	?	0.500	0.066	6	
Weighted Avg.	0.484	0.132	?	0.484	?	?	0.676	0.360	

=== Confusion Matrix ===

a	b	c	d	e	f	<-- classified as
29	0	0	0	0	0	a = 1
1	0	13	0	0	0	b = 2
0	0	15	0	0	0	c = 3
0	0	17	0	0	0	d = 4
9	0	1	0	0	0	e = 5
0	0	6	0	0	0	f = 6

Astfel, se obține o regulă de clasificare bazată pe o trăsătură, iar în urma testării, se evaluează performanța acesteia. Într-un exemplu, clasificarea se face în funcție de atributul "elongation_of_the_rete_ridges". Acuratețea clasificării, fiind în cazul nostru 48,3516%, și statistica Kappa, fiind în cazul nostru 0,3286 (nu foarte apropiată de valoarea 1, valoarea 1 semnificând acord total între clasificarea realizată prin această tehnică și valorile observate), sunt utilizate pentru a evalua cât de bine se potrivește modelul datelor și cât de apropiată este clasificarea de cea corectă. 47 de instanțe sunt clasificate greșit dintr-un total de 91 de instanțe din setul de test.

Matricea de confuzie oferă o comparație între clasificarea făcută de model și valorile reale. Valorile de pe coloane reprezintă predicțiile modelului, în timp ce cele de pe rânduri sunt valorile reale.

J48

Algoritmul C4.5, implementat în Weka sub numele de J48, este unul dintre cei mai utilizați algoritmi de inducție a arborilor de decizie. Este o extensie a algoritmului ID3, abordând probleme precum supra-antrenarea, gestionarea datelor continue și a celor cu valori lipsă, și îmbunătățirea eficienței computaționale. Acest algoritm generează un arbore de decizie prin împărțirea recursivă a setului de date, folosind o strategie de parcurgere în adâncime.

În cazul nostru, am aplicat algoritmul J48 în Weka, utilizând opțiunea "Percentage Split". Astfel, am construit regulile de clasificare pe baza a 75% din setul nostru de date, restul fiind rezervat pentru testare. Am obținut un arbore de decizie cu 40 de noduri și 30 de frunze. Acuratețea clasificării a fost de 91,2088%, cu doar 8 din cele 91 de instanțe fiind clasificate greșit, iar 83 corect. Statistica Kappa a fost de 0,8888, apropiindu-se de valoarea 1, ceea ce indică un grad ridicat de acord între clasificarea obținută și valorile observate.

Classifier output

```
Time taken to build model: 0.04 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances      83          91.2088 %
Incorrectly Classified Instances    8           8.7912 %
Kappa statistic                    0.8888
Mean absolute error                 0.0372
Root mean squared error             0.1736
Relative absolute error             13.9011 %
Root relative squared error         47.3716 %
Total Number of Instances          91

=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC     ROC Area  PRC Area  Class
      0.966   0.081   0.848     0.966   0.903     0.858   0.942    0.830    1
      0.786   0.026   0.846     0.786   0.815     0.783   0.927    0.710    2
      1.000   0.013   0.938     1.000   0.968     0.962   0.993    0.938    3
      0.765   0.000   1.000     0.765   0.867     0.852   0.901    0.837    4
      1.000   0.000   1.000     1.000   1.000     1.000   1.000    1.000    5
      1.000   0.000   1.000     1.000   1.000     1.000   1.000    1.000    6
Weighted Avg.   0.912   0.032   0.918     0.912   0.910     0.887   0.951    0.861

=== Confusion Matrix ===

  a  b  c  d  e  f  <-- classified as
28  0  1  0  0  0 | a = 1
 3 11  0  0  0  0 | b = 2
 0  0 15  0  0  0 | c = 3
 2  2  0 13  0  0 | d = 4
 0  0  0  0 10  0 | e = 5
 0  0  0  0  0  6 | f = 6
```

Naive Bayes

Se utilizează pentru a clasifica date neetichetate prin estimări, bazându-se pe datele de antrenament etichetate. Este un algoritm generativ bazat pe ipoteza "naivă" că predictele condiționate de apartenența la o clasă sunt independente între ele.

În cazul nostru, am aplicat algoritmul Naive Bayes în Weka, selectând opțiunea "Percentage Split". Astfel, regulile de clasificare au fost construite pe baza a 75% din setul nostru de date, iar restul instanțelor au fost folosite pentru testare. Acuratețea clasificării a fost de 98,9011%. Statistica Kappa a fost de 0.9862, indicând un acord între clasificarea obținută și valorile observate.

Classifier

Choose **NaiveBayes**

est options

☐ Use training set
☐ Supplied test set Set...
☐ Cross-validation Folds 10
☒ Percentage split % 75

More options...

Nom) class

Start Stop

result list (right-click for options)

7:50:12 - rules.OneR
8:26:33 - trees.J48
8:30:51 - bayes.NaiveBayes

Classifier output

Time taken to build model: 0 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0.01 seconds

=== Summary ===

Correctly Classified Instances	90	98.9011 %
Incorrectly Classified Instances	1	1.0989 %
Kappa statistic	0.9862	
Mean absolute error	0.0072	
Root mean squared error	0.0417	
Relative absolute error	2.707 %	
Root relative squared error	11.3726 %	
Total Number of Instances	91	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	1
	0.929	0.000	1.000	0.929	0.963	0.957	1.000	1.000	2
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	3
	1.000	0.014	0.944	1.000	0.971	0.965	1.000	1.000	4
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	5
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	6
Weighted Avg.	0.989	0.003	0.990	0.989	0.989	0.987	1.000	1.000	

=== Confusion Matrix ===


```

a  b  c  d  e  f  <-- classified as
29  0  0  0  0  0 | a = 1
0 13  0  1  0  0 | b = 2
0  0 15  0  0  0 | c = 3
0  0  0 17  0  0 | d = 4
0  0  0  0 10  0 | e = 5
0  0  0  0  0  6 | f = 6

```

Regresia logistică

Regresia logistică este utilizată pentru a investiga relația dintre multiple variabile independente și o variabilă dependentă binară. Această variabilă dependentă se referă adesea la apartenența la două categorii, cum ar fi prezența/absența sau da/nu.

În Weka, pentru a modela o variabilă binară folosind regresia logistică, selectăm opțiunea "SimpleLogistic" din panoul "Classify". De asemenea, bifăm opțiunea "Percentage split" pentru a împărți setul de date într-un set de antrenare și unul de testare, folosind 75% pentru antrenare și 25% pentru testare.

În rezultatele obținute, observăm o acuratețe de clasificare de 98,9011%. Din cele 91 de instanțe, 90 au fost clasificate corect, iar 1 a fost clasificata greșit. Valoarea indicelui Kappa statistics este de 0,9862, apropiindu-se de valoarea maximă de 1, ceea ce indică o concordanță excelentă între clasificarea realizată prin regresia logistică și valorile observate.

Classifier

Choose **SimpleLogistic** -I 0 -M 500 -H 50 -W 0.0

est options

☐ Use training set
☐ Supplied test set Set...
☐ Cross-validation Folds
☒ Percentage split %

More options...

Vom) class

Start Stop

result list (right-click for options)

7:50:12 - rules.OneR
3:26:33 - trees.J48
3:30:51 - bayes.NaiveBayes
3:34:22 - functions.SimpleLogistic

Classifier output

Time taken to build model: 0.46 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances

90

98.9011 %

Incorrectly Classified Instances

1

1.0989 %

Kappa statistic

0.9862

Mean absolute error

0.0086

Root mean squared error

0.0525

Relative absolute error

3.2063 %

Root relative squared error

14.3294 %

Total Number of Instances

91

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	FRC Area	Class
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	1
	1.000	0.013	0.933	1.000	0.966	0.960	0.999	0.995	2
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	3
	0.941	0.000	1.000	0.941	0.970	0.964	0.999	0.997	4
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	5
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	6
Weighted Avg.	0.989	0.002	0.990	0.989	0.989	0.987	1.000	0.999	

=== Confusion Matrix ===

a b c d e f <-- classified as

29 0 0 0 0 0 | a = 1

0 14 0 0 0 0 | b = 2

0 0 15 0 0 0 | c = 3

0 1 0 16 0 0 | d = 4

0 0 0 0 10 0 | e = 5

0 0 0 0 0 6 | f = 6

- Gruparea

Gruparea, cunoscută și sub numele de analiză de tip cluster, este o tehnică utilizată pentru a împărți un set de date în clase sau grupuri, fără a avea clase prestabilite. Fiecare clasă este definită de un model, care poate fi un obiect abstract sau un exemplu reprezentativ al clasei.

Simple Kmeans

Algoritmul SimpleKMeans este unul dintre algoritmii utilizați în gruparea datelor. Acesta începe prin crearea unei partiționări inițiale a datelor și apoi încearcă să îmbunătățească această partiționare prin mutarea iterativă a obiectelor dintr-un grup în altul. Procesul se oprește atunci când se reușește împărțirea datelor în k grupuri, astfel încât fiecare grup să conțină cel puțin un obiect și fiecare obiect să aparțină unui singur grup.

Pentru a folosi algoritmul SimpleKMeans în Weka, trebuie să accesăm opțiunea "Cluster" și să selectăm "SimpleKMeans" din lista de opțiuni disponibile. Înainte de a începe procesul, bifăm opțiunea "Classes to clusters evaluation". De asemenea, putem alege distanța de evaluare a similarității între obiecte, cum ar fi distanța euclidiană sau Manhattan.

În rezultatele obținute, observăm instanțele incorecte de clustere sunt 51,0929%, iar algoritmul a efectuat 4 iterații. Interesant este că valorile obținute pentru distanța euclidiană sunt similare cu cele pentru distanța Manhattan. Precizia datelor este de 48,90% ($100\% - (187 / 366) * 100 = 48.09\%$).

Clusterer

Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation (Nom) class

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

18:40:59 - SimpleKMeans

Algorithm. More Capabilities

Memory 100

Density 2.0

PruningRate 10000

canopyT1 -1.25

canopyT2 -1.0

debug False

StdDevs False

Function Choose **EuclideanDistance -R first-last**

Capabilities False

LogValues False

DistanceCalc False

Clusterer output

vacuolisation_and_damage_of_basal_layer	0	0	0
spongiosis	0	2	0
saw-tooth_appearance_of_retes	0	0	0
follicular_horn_plug	0	0	0
perifollicular_parakeratosis	0	0	0
inflammatory_mononuclear_infiltrate	2	2	2
band-like_infiltrate	0	0	0
Age	36.2961	34.8238	39.6783

Time taken to build model (full training data) : 0.04 seconds

=== Model and evaluation on training set ===

Clustered Instances

0	255 (70%)
1	111 (30%)

Class attribute: class

Classes to Clusters:

0	1	<-- assigned to cluster
5	107	1
61	0	2
72	0	3
49	0	4
49	3	5
19	1	6

Cluster 0 <-- 3

Cluster 1 <-- 1

Incorrectly clustered instances : 187.0 51.0929 %

Clusterer

Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.ManhattanDistance -R first-last" -I 500 -num-slots 1 -S 10

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation (Nom) class

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

18:40:59 - SimpleKMeans

18:54:50 - SimpleKMeans

Algorithm. More Capabilities

Memory 100

Density 2.0

PruningRate 10000

canopyT1 -1.25

canopyT2 -1.0

debug False

StdDevs False

Function Choose **ManhattanDistance -R**

Capabilities False

LogValues False

DistanceCalc False

Clusterer output

vacuolisation_and_damage_of_basal_layer	0	0	0
spongiosis	0	2	0
saw-tooth_appearance_of_retes	0	0	0
follicular_horn_plug	0	0	0
perifollicular_parakeratosis	0	0	0
inflammatory_mononuclear_infiltrate	2	2	2
band-like_infiltrate	0	0	0
Age	36	35	39

Time taken to build model (full training data) : 0.03 seconds

=== Model and evaluation on training set ===

Clustered Instances

0	255 (70%)
1	111 (30%)

Class attribute: class

Classes to Clusters:

0	1	<-- assigned to cluster
5	107	1
61	0	2
72	0	3
49	0	4
49	3	5
19	1	6

Cluster 0 <-- 3

Cluster 1 <-- 1

Incorrectly clustered instances : 187.0 51.0929 %

Farthest First

Algoritmul Furthest First este o altă metodă de grupare în care clasificarea unei noi observații nu se face în funcție de valoarea medie a clasei, ci prin selectarea instanței reale cea mai aproape de medie. Această abordare permite luarea deciziei de clasificare în funcție de o observație concretă, ceea ce poate conduce la clustere ce reflectă datele mai realistic.

În rezultatele obținute, se observă o acuratețe a setului de date de 42,62% ($100\% - (213 / 366) * 100 = 42.62\%$). Comparativ cu algoritmul SimpleKMeans, această tehnică de grupare produce rezultate mai slabe cu câteva procente. Este important de menționat că performanța algoritmilor de grupare poate varia în funcție de natura și distribuția datelor, precum și de parametrii specifici utilizați în algoritm.

The screenshot displays the 'Clusterer' application window. The 'Choose' dropdown is set to 'FarthestFirst -N 2 -S 1'. Under 'Cluster mode', 'Classes to clusters evaluation' is selected with a percentage of 66, and 'Store clusters for visualization' is checked. The 'Ignore attributes' field is empty. The 'Start' button is visible. Below the settings, a 'result list' shows three entries: '8:40:59 - SimpleKMeans', '8:54:50 - SimpleKMeans', and '9:01:27 - FarthestFirst' (highlighted in blue).

The 'Clusterer output' pane shows the following text:

```
FarthestFirst
=====

Cluster centroids:

Cluster 0
 3 2 1 2 0 0 0 0 0 0 0 0 2 0 2 1 0 2 0 0 0 0 0 0 0 2 0 0 0 2 0 25.0
Cluster 1
 3 3 2 1 1 0 0 0 2 2 1 0 0 0 0 0 3 2 3 2 2 2 1 1 0 0 0 0 0 0 1 0 42.0

Time taken to build model (full training data) : 0.01 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      285 ( 78%)
1       81 ( 22%)

Class attribute: class
Classes to Clusters:

 0 1 <-- assigned to cluster
31 81 | 1
61 0 | 2
72 0 | 3
49 0 | 4
52 0 | 5
20 0 | 6

Cluster 0 <-- 3
Cluster 1 <-- 1

Incorrectly clustered instances :      213.0      58.1967 %
```

4. Obținere de noi cunoștințe

Pentru a obține noi cunoștințe din setul de date original, am împărțit inițial setul de date într-un set de training și un set de testare. Am folosit o tehnică de grupare pentru a extrage aceste noi cunoștințe.

Procesul de împărțire a setului de date a implicat utilizarea opțiunii "RemovePercentage" din meniul de preprocesare în Weka, unde am setat un procentaj de 50% și am aplicat această modificare.

The screenshot shows the Weka GUI with the 'RemovePercentage' filter selected. The filter is configured to remove 50.0% of the data. A dialog box titled 'weka.gui.GenericObjectEditor' is open, showing the filter's configuration. The 'percentage' field is set to 50.0. The background shows a list of attributes and a data visualization.

Attributes: 35
Sum of weights: 366

Selected attribute: Name: Age
Missing: 0 (0%)
Distinct: 61
Type: Numeric
Unique: 8 (2%)

Statistic Value
Minimum 0
Maximum 75
Mean 36.296
StdDev 15.156

weka.gui.GenericObjectEditor
weka.filters.unsupervised.instance.RemovePercentage
About
A filter that removes a given percentage of a dataset.
More
Capabilities

debug False
doNotCheckCapabilities False
invertSelection False
percentage 50.0
The percentage of the data to select
Open... Save... OK Cancel

Setul de date rezultat a fost salvat sub denumirea "dermatologyTrainSet.arff".

The screenshot shows the Weka GUI with the 'RemovePercentage' filter selected. The filter is configured to remove 50.0% of the data. A dialog box titled 'weka.gui.GenericObjectEditor' is open, showing the filter's configuration. The 'percentage' field is set to 50.0. The background shows a list of attributes and a data visualization.

Attributes: 35
Sum of weights: 183

Selected attribute: Name: erythema
Missing: 0 (0%)
Distinct: 3
Type: Nominal
Unique: 0 (0%)

No. Label Count Weight
1 0 0 0
2 1 20 20
3 2 118 118
4 3 45 45

Class: class (Nom)
Visualize All

Apoi, am revenit la starea inițială a setului de date folosind funcția "undo" și am repetat procedura, dar de data aceasta am setat opțiunea "invertSelection" la true. Astfel, am obținut setul de date de test.

Iter Choose **RemovePercentage -P 50.0 -V** Apply Stop

Current relation: dermatology-weka.filters.unsupervised.instance.RemovePercentage-P50.0-V
Instances: 183 Attributes: 35 Sum of weights: 183

Selected attribute: Name: erythema
Missing: 0 (0%) Distinct: 4 Type: Nominal Unique: 0 (0%)

No.	Label	Count	Weight
1	0	4	4
2	1	37	37
3	2	97	97
4	3	45	45

Class: class (Nom) Visualize All

Remove

Setul de date de test este similar celui original sau celui de training, având aceeași dimensiune. Singura diferență constă în faptul că valorile atributului clasă au fost înlocuite cu "?", reprezentând datele necunoscute.

```
@attribute focal_hypergranulosis {0,1,2,3}
@attribute disappearance_of_the_granular_layer {0,1,2,3}
@attribute spongiosis {0,1,2,3}
@attribute saw-tooth_appearance_of_retes {0,1,2,3}
@attribute perifollicular_parakeratosis {0,1,2,3}
@attribute band-like_infiltrate {0,1,2,3}
@attribute class {1,2,3,4,5,6}

@data
2,0,3,0,0,1,0,0,0,3,0,0,0,0,0,3,0,0,0,?
3,3,2,1,0,1,0,1,0,1,2,2,2,0,0,0,0,0,0,?
1,2,3,1,0,0,0,0,0,1,0,0,0,2,0,3,2,0,3,?
2,2,0,0,0,3,0,3,0,0,2,2,2,0,3,0,0,0,0,?
3,2,2,2,0,0,0,0,0,1,0,0,0,2,2,2,3,0,3,?
3,2,0,0,0,0,2,1,0,2,0,0,0,0,0,2,0,0,0,?
1,0,2,0,0,0,0,0,3,1,0,2,0,0,0,0,0,0,0,?
2,3,3,3,0,0,0,0,0,2,0,0,0,0,2,3,2,0,3,?
```

Pentru a aplica setul de date de test în Weka, am urmat următorii pași:

1. Am ales o tehnică de grupare, în acest caz am selectat FarthestFirst.
2. Am selectat opțiunea "Supplied test set" și am încărcat setul de date de test.
3. Am reevaluat gruparea inițială folosind setul de date de test.
4. Am obținut rezultatele analizei.

```
U 4 ~~~ assigned to CLUSTER
5 55 | 1
19 0 | 2
34 0 | 3
30 0 | 4
18 11 | 5
11 0 | 6

Cluster 0 <-- 3
Cluster 1 <-- 1

Incorrectly clustered instances : 64.0 51.3661 %

=== Detailed Accuracy By Class ===
Correctly Classified Instances 69 94.5205 %
Incorrectly Classified Instances 4 5.4795 %
Kappa statistic 0.9306
Mean absolute error 0.0377
Root mean squared error 0.1128
Relative absolute error 14.1294 %
Root relative squared error 30.8643 %
Total Number of Instances 73
```

Este important să remarcăm că setul de date de test a fost utilizat pentru a evalua performanța modelului de grupare pe date noi, neetichetate anterior. Această evaluare ne permite să înțelegem cât de bine generalizează modelul nostru pe datele pe care nu le-a văzut anterior.

Graficul acestui tip de grupare este:



Predictii:

Pentru a face predicții, am utilizat fișierul de test și am reaplicat evaluarea setului de date pentru a observa diferențele și rezultatele.

- **OneR:** Toate instanțele sunt prezise cu o șansă de 100% ca fiind negative.
- **J48:** Există o tendință de a prezice anumite clase datorită structurii arborelui de decizie. Acest lucru se datorează probabil prezenței unor caracteristici semnificative care sunt asociate puternic cu anumite clase dermatologice. De exemplu, anumite atribute precum "melanin incontinence" sau "eosinophils in the infiltrate" pot avea o influență semnificativă asupra predicțiilor clasificatorului.
- **NaiveBayes:** Predicțiile diferă în funcție de instanță, însă majoritatea sunt apropiate de 1, indicând o probabilitate mare de a fi negative. Am creat un model Naive Bayes pentru fiecare pereche posibilă de clase. De exemplu, avem un model pentru a distinge între psoriazis și dermatită seboreică, altul pentru a distinge între psoriazis și lichen plan, și tot așa, până când acoperim toate posibilele perechi de clase. Pentru fiecare pereche de clase, am antrenat modelul Naive Bayes folosind doar datele corespunzătoare celor două clase. Aceasta înseamnă că am creat subseturi de date care conțin doar exemplele pentru cele două clase respective. După antrenarea modelelor pentru fiecare pereche de clase, vom clasifica un nou exemplu prin aplicarea tuturor modelelor antrenate și votarea pentru clasa care primește cele mai multe voturi.

=== Predictions on user test set ===

inst#	actual	predicted	error	prediction	inst#	actual	predicted	error	prediction
1	1:?	3:3	1		1	1:?	2:2	0.909	
2	1:?	1:1	1		2	1:?	1:1	0.991	
3	1:?	3:3	1		3	1:?	3:3	0.987	
4	1:?	1:1	1		4	1:?	1:1	0.999	
5	1:?	3:3	1		5	1:?	3:3	0.995	
6	1:?	3:3	1		6	1:?	2:2	0.926	
7	1:?	1:1	1		7	1:?	5:5	0.997	
8	1:?	3:3	1		8	1:?	3:3	0.975	
9	1:?	3:3	1		9	1:?	4:4	0.928	
10	1:?	3:3	1		10	1:?	4:4	0.928	
11	1:?	1:1	1		11	1:?	1:1	0.999	

Concluzii:

- Setul de date este extrem de complex, cu un număr mare de instanțe și date. Folosirea unor algoritmi simpli, cum ar fi OneR, duce la o clasificare cu un procentaj mai mare de instanțe incorect clasificate.
- Surprinzător, clasificatorul NaiveBayes a avut o precizie de 98,9011%, apropiată de 100%, indicând că este cel mai potrivit pentru acest set de date, cu 1 instanță clasificată incorect.
- Regresia liniară a avut o acuratețe de 98,9011%, egală cu cea de la NaiveBayes. J48 a avut o acuratețe de 91,2088%.
- Algoritmii mai complexi, precum NaiveBayes, Regresia Liniară și J48, sunt mai buni în procesarea unor astfel de date.
- Problema se repetă și în cazul tehnicilor de grupare, unde numărul mare de atribute duce la o eroare ridicată. Pentru SimpleKMeans, acuratețea este de 48,90%, iar pentru FarthestFirst este de doar 42,62%, ambele cu erori considerabile.
- Impactul numărului mare de caracteristici în comparație cu dimensiunea setului de date este evident în performanța algoritmilor. Acest aspect conduce la variații semnificative în ritmul de învățare între diversele tehnici. De exemplu, Naive Bayes se evidențiază prin viteză și precizie.
- Metodele de clusterizare nu ating nivelul de precizie dorit pe acest set de date, datorită numărului mare de caracteristici implicate. Interpretarea rezultatelor este puternic influențată de aceste caracteristici, ceea ce poate genera clusterizări extrem de polarizate. Această problemă devine evidentă în cazul abordării FarthestFirst.

Bibliografie:

<http://weka.sourceforge.net/doc.dev/weka/classifiers/bayes/NaiveBayes.html>

<https://www.saedsayad.com/oner.html>

<http://inf.ucv.ro/documents/rstoean/7.%20Evaluarea%20performantei.pdf>

<http://weka.sourceforge.net/doc.dev/weka/classifiers/trees/J48.html>

http://ionut.mironica.ro/teaching/TACAI_lab1.pdf

Laboratoarele facute in cadrul semestrului, insotite de lucrarile de laborator