

Zdolność kredytowa – zadanie z uczenia maszynowego (klasyfikacja)

Zdolność kredytowa wyznacza maksymalną kwotę kredytu, której bank może udzielić na dany cel (na przykład zakup mieszkania). Oprócz wysokości dochodów, wpływa na nią wiele innych czynników, takich jak miesięczny dochód, czy miesięczne wydatki na utrzymanie. W tym ćwiczeniu zostaną zbudowane modele pozwalające odpowiedzieć na pytanie, czy dany wniosek kredytowy powinien być przyjęty (tak), czy też odrzucony (nie) przez bank.

Zbiór danych

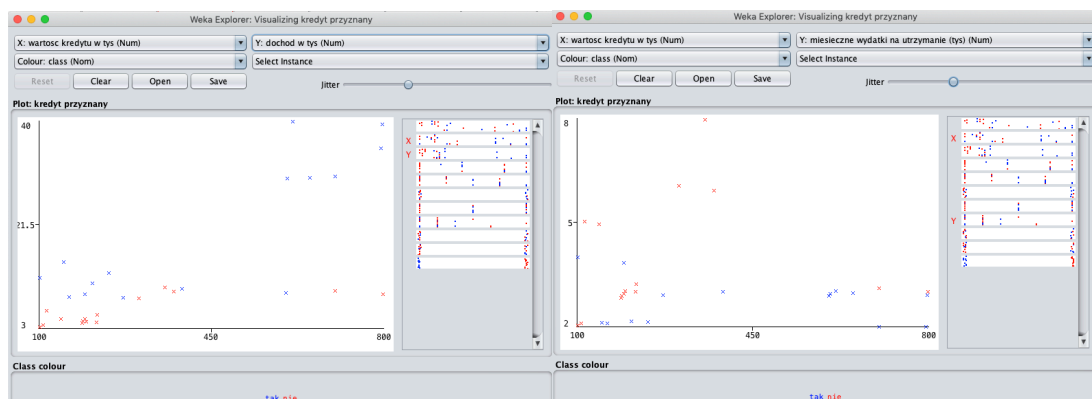
Przygotowany zbiór danych w formacie ARFF zawiera 30 obserwacji dla 10 atrybutów niezależnych (w tym trzy 3 kategorię i 7 numerycznych) oraz jeden atrybut zależny (kategoryczny). Zbiór jest doskonale zbalansowany (15 obserwacji na tak, 15 na nie).

The screenshot shows the Weka Explorer interface. The 'Current relation' panel displays 'Relation: kredyt przyznany' with 11 attributes and 30 instances. The 'Attributes' list includes: 1. wartosc nieruchomosci w tys, 2. wartosc kredytu w tys, 3. dochod w tys, 4. okres splaty w latach, 5. liczba osob w gospodarstwie domowym, 6. posiada majatek własny, 7. wiek kredytobiorcy, 8. miesieczne wydatki na utrzymanie (tys), 9. forma zatrudnienia, 10. inne kredyty i zobowiazania, and 11. class. The 'Selected attribute' panel shows 'Name: class' with 2 distinct values: 'tak' (15) and 'nie' (15). Below this, a bar chart visualizes the distribution of the 'class' attribute, showing 15 instances for 'tak' (blue bar) and 15 instances for 'nie' (red bar).

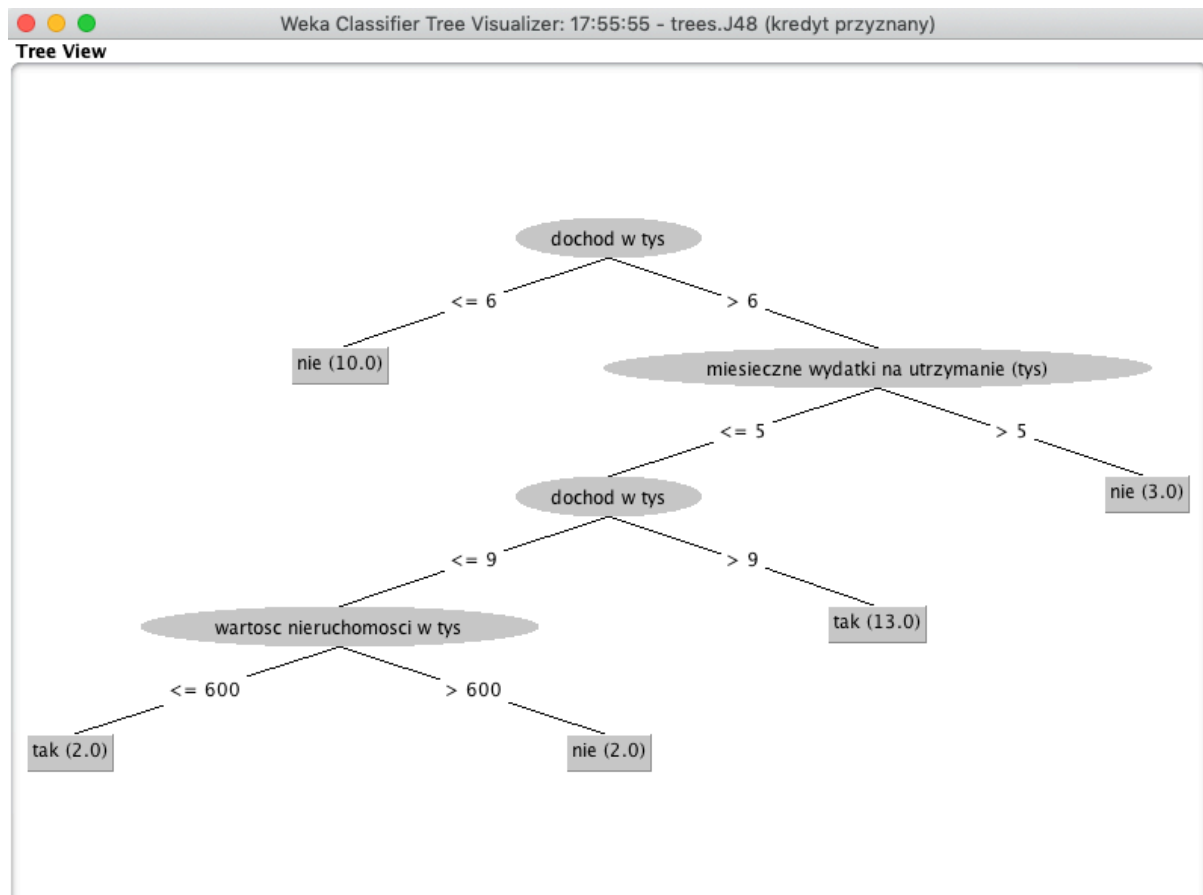
Celem ćwiczenia jest porównanie działania kilku klasyfikatorów i wybranie najlepszego z nich jako modelu decyzyjnego wspomagającego przydzielanie kredytów w przyszłości. Kryteriami oceny będą poprawność klasyfikacji oraz czytelność wiedzy.

Eksploracja danych

Pobieżna analiza wizualna pozwoliła na zidentyfikowanie silnych predyktorów najlepiej nadających się do klasyfikacji. Pierwszym z nich jest dochód. Im wyższy dochód, tym większa szansa na przyznanie kredytu. Dochód powyżej 12 tysięcy gwarantuje przyznanie kredytu. Z kolei miesięczne wydatki powyżej 5 tysięcy złotych oznaczają nieprzyznanie kredytu.



Algorytm J48 (F-Measure 0.778) wygenerował drzewo decyzyjne (poniżej), z którego wynika, iż nie ma co marzyć o kredycie, jeśli zarabia się mniej niż 6 (tys.) oraz miesięczne wydatki ma się na poziomie powyżej 5 (tys.). W przypadku wartości nieruchomości powyżej 600 (tys.) trzeba zarabiać powyżej 9 (tys.) aby otrzymać kredyt.



Podobne informacje otrzymamy z JRip (F-Measure 0.778), chociaż tutaj próg akceptowalnych miesięcznych wydatków został podniesiony do 6 (tys.), a przy wartościach kredytu powyżej 700 (tys.) i wartości nieruchomości poniżej 850 (tys.) również otrzymamy odmowę.

JRIP rules:

=====

```

(dochod w tys <= 6) => class=nie (10.0/0.0)
(miesieczne wydatki na utrzymanie tys) >= 6 => class=nie (3.0/0.0)
(wartosc kredytu w tys >= 700) and (wartosc nieruchomosci w tys <= 850) => class=nie (2.0/0.0)
=> class=tak (15.0/0.0)
  
```

Number of Rules : 4

Na koniec użyto modułu Experimenter do uzyskania wyników porównawczych. Metoda podziału na zbiór uczący oraz testowy była taka sama, wybór obserwacji do poszczególnych zbiorów losowy (stad nieco inne wyniki, chociaż co do ostatecznej konkluzji zgodne). Klasyfikator oparty o algorytm AdaBoostM1 okazał się najlepszy oraz statystycznie (poziom

istotności 5%) lepszy od referencyjnego klasyfikatora zbudowanego za pomocą algorytmu BayesNet.

```

Test output

Tester:      weka.experiment.PairedTTester -G 3,4,5 -D 1 -R 2 -S 0.05 -result-matrix "we
Analysing:    F_measure
Datasets:     1
Resultsets:   6
Confidence:   0.05 (two tailed)
Sorted by:    -
Date:         15/03/2020, 18:26

Dataset              (1) bayes.B | (2) baye (3) tree (4) rule (5) meta (6) func
-----
'kredyt przyznany'   (10)  0.70 |  0.80    0.76    0.71    0.80 v  0.77
-----
                    (v/ /*) | (0/1/0) (0/1/0) (0/1/0) (1/0/0) (0/1/0)

Key:
(1) bayes.BayesNet '-D -Q bayes.net.search.local.K2 -- -P 1 -S BAYES -E bayes.net.estim
(2) bayes.NaiveBayes '' 5995231201785697655
(3) trees.J48 '-C 0.25 -M 2' -217733168393644444
(4) rules.JRip '-F 3 -N 2.0 -O 2 -S 1' -6589312996832147161
(5) meta.AdaBoostM1 '-P 100 -S 1 -I 10 -W trees.DecisionStump' -1178107808933117974
(6) functions.Logistic '-R 1.0E-8 -M -1 -num-decimal-places 4' 3932117032546553727

```

Z punktu widzenia czytelności wiedzy (w tym przypadku bardzo istotnej, gdyż bank musi być w stanie uzasadnić decyzję odmowną), algorytmy J48 lub JRip byłyby bardziej preferowane.

Wynik

Ostatecznie J48 osiągnął najlepszy wynik w ogólnym porównaniu (poprawność klasyfikacji oraz czytelność wiedzy).