

# Klasyfikacja i regresja: Implementacja naiwnego klasyfikatora bayesowskiego w `scikit-learn`

31 października 2018

## Opis pliku z zadaniami

Wszystkie zadania na zajęciach będą przekazywane w postaci plików `.pdf`, sformatowanych podobnie do tego dokumentu. Zadania będą różnego rodzaju. Za każdym razem będą one odpowiednio oznaczone:

- Zadania do wykonania na zajęciach oznaczone są symbolem  $\triangle$  – nie są one punktowane, ale należy je wykonać w czasie zajęć.
- Punktowane zadania do wykonania na zajęciach oznaczone są symbolem  $\diamond$  – należy je wykonać na zajęciach i zaprezentować prowadzącemu.
- Zadania do wykonania w domu oznaczone są symbolem  $\star$  – są one punktowane, należy je dostarczyć w sposób podany przez prowadzącego i w wyznaczonym terminie (zwykle przed kolejnymi zajęciami).

# 1 Transformacja danych do formatu wspieranego przez sklearn △

## Treść

Wykonaj poniższe kroki w celu odpowiedniego wczytania i przekształcenia danych:

1. Zapoznaj się z typami danych obsługiwanymi przez `sklearn` oraz proponowanymi metodami wczytywania danych (<http://scikit-learn.org/stable/datasets/index.html#external-datasets>).
2. Zapoznaj się z podstawowymi informacjami o pakietach `numpy` (<http://www.numpy.org/>) oraz `pandas` (<https://pandas.pydata.org/>).
3. Wczytaj dane z plików `.csv` dostępnych na stronie:
  - `data/grypa-train.csv`,
  - `data/grypa-test.csv`,

za pomocą funkcji z pakietu `pandas` ([https://pandas.pydata.org/pandas-docs/stable/generated/pandas.read\\_csv.html#pandas.read\\_csv](https://pandas.pydata.org/pandas-docs/stable/generated/pandas.read_csv.html#pandas.read_csv)).

Poniżej podany jest spodziewany wynik:

```
1 print(data_train)
2
3 # dreszcze katar bol_glowy goraczka grypa
4 #0      tak   nie   sredni      tak   nie
5 #1      tak   tak     nie      nie   tak
6 #2      tak   nie   duzy      tak   tak
7 #3      nie   tak   sredni     tak   tak
8 #4      nie   nie     nie     nie   nie
9 #5      nie   tak   duzy      tak   tak
10 #6      nie   tak   duzy     nie   nie
11 #7      tak   tak   sredni     tak   tak
12
13
14 print(data_test)
15
16 # dreszcze katar bol_glowy goraczka
17 #0      nie   tak     nie      tak
```

4. Przekształć wartości nominalne na liczby całkowite:

Poniżej podany jest spodziewany wynik:

```

1 print(data_train_converted)
2
3 #   dreszcze   katar bol_glowy   goraczka   grypa
4 #0         1     0   sredni     1         0
5 #1         1     1     nie       0         1
6 #2         1     0    duzy     1         1
7 #3         0     1   sredni     1         1
8 #4         0     0     nie       0         0
9 #5         0     1    duzy     1         1
10 #6         0     1    duzy     0         0
11 #7         1     1   sredni     1         1
12
13
14 print(data_test_converted)
15
16 #   dreszcze   katar bol_glowy   goraczka
17 #0         0     1     nie       1

```

5. Podziel wczytane dane na cechy oraz klasę, a następnie przekształć dane z obiektów pandas do macierzy numpy.

Poniżej podany jest spodziewany wynik:

```

1 print(X_train)
2
3 #[[1 0 1 1]
4 # [1 1 0 0]
5 # [1 0 2 1]
6 # [0 1 1 1]
7 # [0 0 0 0]
8 # [0 1 2 1]
9 # [0 1 2 0]
10 # [1 1 1 1]]
11
12 print(y_train)
13
14 #[0 1 1 1 0 1 0 1]
15
16 print(X_test)
17
18 #[[0 1 0 1]]

```

6. Zaimportuj dowolny klasyfikator z pakietu `sklearn` (może to być na przykład algorytm drzew decyzyjnych `DecisionTreeClassifier`). Wytrenuj model na danych treningowych, a następnie porównaj jego działanie na danych treningowych i testowych. Wypisz predykcje modelu dla obu zbiorów danych. Przeanalizuj odpowiedź klasyfikatora. W jakim stopniu predykcja jest poprawna?

Poniżej podany jest spodziewany wynik:

```
1 # Create and train a DecisionTreeClassifier
2 dt = DecisionTreeClassifier()
3 clr = dt.fit(X_train, y_train)
4
5 # Predict the test example
6 ypred = clr.predict(X_test)
7 print(ypred)
8 #[1]
9
10 # Predict the train examples
11 ypred = clr.predict(X_train)
12 print(ypred)
13 # [0 1 1 1 0 1 0 1]
14
15 # Average accuracy on the training set (known decisions)
16 accuracy_score(y_train, ypred)
17 #1.0
```

## 2 Implementacja naiwnego klasyfikatora bayesowskiego dla danych nominalnych

Op.\*

### Treść

Naiwny klasyfikator Bayesa jest prostym klasyfikatorem probabilistycznym, który został przedstawiony na wykładzie. Zadaniem jest zaimplementowanie tego klasyfikatora w `pythonie` przy wykorzystaniu biblioteki `scikit-learn`. Można skorzystać z szablonu, który jest do pobrania ze strony przedmiotu.

Szczegóły zadania:

- Zakładamy, że na dziedzinie cech są wartości nominalne kodowane za pomocą liczb całkowitych.
- Przetestuj rozwiązanie na zbiorze *grypa*. Poprawność można zweryfikować przeliczając prawdopodobieństwa *ręcznie*.

Kilka podpowiedzi:

- Zapoznaj się i skorzystaj z dostarczonego kodu.
- Przechowuj potrzebne informacje do obliczenia prawdopodobieństw w odpowiednich strukturach danych jako pola w klasie klasyfikatora.
- Metoda `fit` tworzy klasyfikator na podstawie danych wejściowych.
- Metoda `predict_proba` zwraca dla danych przykładów obliczone prawdopodobieństwa warunkowe klas w wektorze `numpy`.
- Metoda `predict` zwraca dla danych przykładów przewidywaną odpowiedź (tj. klasę).
- Folder `classifiers_students` zawiera szkielet klasy klasyfikatora.
- Folder `data` zawiera potrzebne dane.
- Folder `tasks_students` zawiera skrypty z prostymi testami napisanego klasyfikatora.
- Folder `utils` zawiera przydatne narzędzia np. do konwersji danych i ewaluacji klasyfikatorów.