

Naiwny klasyfikator Bayesa

Przemysław Klęsk & Joanna Kołodziejczyk

1 Problem do rozwiązania

Opis zadania przedstawia implementację naiwnego klasyfikatora Bayesa (NBC) w dwóch wariantach (dyskretnym i ciągłym) w języku Python.

2 NBC wersja dyskretna

Implementacja NBC w wersji dyskretniej powinna być wykonana dla zbioru „wine” z repozytorium UCI <https://archive.ics.uci.edu/ml/datasets/Wine>. Przed przystąpieniem do implementacji pobierz zbiór danych. Dotyczy on przydziału wina na podstawie składu chemicznego do trzech kategorii i zapoznaj się z nim. Zwróć uwagę, która ze zmiennych (kolumna) jest zmienną decyzyjną.

2.1 Kolejne kroki w NBC dla zmiennych dyskretnych

1. Wczytaj dane z pobranego pliku tekstowego `wine.data` do macierzy `numpy` (wykorzystaj funkcję `numpy.genfromtxt`) i rozdziel tę macierz na dwie macierze X (o wymiarze 178×13) i y (178×1 — etykiety klas).
2. Dyskretyzację danych „wine” można wykonać wykorzystując gotowy obiekt `KBinsDiscretizer` (z pakietu `sklearn.preprocessing`) lub samodzielnie na poziomie opracowywanej klasy NBC (liczba przedziałów, na którą dyskretyzujemy cechy oryginalnie ciągłe, powinna być parametrem nastawialnym przez użytkownika).
3. Podziel dane na część uczącą i testową (wykorzystaj funkcję `train_test_split` z pakietu `sklearn.model_selection`).
4. Napisz klasę reprezentującą naiwny klasyfikator Bayesa w wariantcie ze zmiennymi dyskretnymi. Klasę przygotuj zgodnie z ideą biblioteki `scikit-learn` — m.in.: wykonaj dziedziczenie po klasach `BaseEstimator` i `ClassifierMixin`, przygotuj metody `fit` (uczenie) i `predict` (klasyfikowanie - wskazanie identyfikatora klasy) oraz pomocniczo `predict_proba` (podanie wartości obliczonego prawdopodobieństwa). (Szczegóły w kolejnych podrozdziałach).
5. Zastanów się i zaplanuj wg własnego uznania wygodne struktury danych do przechowywania:

- rozkładu a priori klas $P(Y = y)$,
- rozkładów warunkowych $P(X_j = v | Y = y)$.

Mogą to być tablice, słowniki, listy lub odpowiednie połączenia / zagnieżdżenia tych struktur. Do tego celu potrzebne będzie także ustalenie dyskretnych dziedzin zmiennych, tj. wykrycie, jakie unikalne wartości poszczególne zmienne mogą przyjmować, np. z wykorzystaniem funkcji `numpy.unique`. Przemyśl, czy informacje o dziedzinach należy zdobywać na poziomie funkcji `fit` na podstawie danych uczących, czy też lepiej przekazać je klasyfikatorowi już podczas konstrukcji.

6. Wyznacz dokładność ($Accuracy = \frac{NP}{LZ} * 100\%$, gdzie NP — liczba próbek poprawnie sklasyfikowanych w zadanym zbiorze, LZ — licznosc zbioru) otrzymanego klasyfikatora na zbiorach uczącym i testowym.
7. Obliczenia powtórz uwzględniając poprawkę LaPlace'a (możesz do tego celu wprowadzić przełącznik w konstruktorze Twojej klasy). Powtórz obliczenia - uczenie i pomiary dokładności. Zwróć uwagę, czy poprawka LaPlace'a podnosi dokładność testową dla tego zbioru danych.

2.2 Uczenie NBC - metoda fit

Uczenie NBC w wariacie dyskretnym polega na wyznaczeniu i zapamiętaniu (w pewnej strukturze danych, np. w tablicy lub słowniku) wszystkich możliwych prawdopodobieństw, które są potrzebne we wzorze ((1)). Utożsamiamy prawdopodobieństwa $P(X = x | Y = y)$ i $P(Y = y)$ z częstościami występującymi w zbiorze uczącym.

2.3 Klasyfikacja - odpowiedź modelu metoda predict i predict_proba

2.3.1 Wyznaczanie klasy - metoda predict

W metodzie `predict` następuje obliczanie odpowiedzi klasyfikatora w wariacie dyskretnym może być realizowane zgodnie ze wzorem (1):

$$y^* = \arg \max_{y \in \{1, \dots, K\}} \prod_{j=1}^n P(X_j = x_j | Y = y) P(Y = y) \quad (1)$$

tj. jako iloczyn prawdopodobieństw (bez zabiegu logarytmowania).

2.3.2 Wyznaczanie wartości prawdopodobieństwa przynależności do każdej klasy - metoda predict_proba

Obliczanie prawdopodobieństwa w metodzie `predict_proba` sprowadza się do wyznaczenia wartości prawdopodobieństwa, przynależności próbki testowej do każdej z klas.

We wzorze ((1)) wykorzystuje się iloczyn, czyli wartości wiarygodności (likelihood):

$$Likelihood(Y = y) = \prod_{j=1}^n P(X_j = x_j | Y = y) P(Y = y).$$

Aby z wiarygodności obliczyć prawdopodobieństwo, że dana próbka należy do klasy y należy podzielić obliczoną dla danej wartości klasy wiarygodność, przez sumę wiarygodności dla wszystkich klas:

$$Probability(Y = y) = \frac{Likelihood(Y = y)}{\sum_{y=1}^K Likelihood(Y = y)},$$

2.4 Poprawka LaPlace'a

Przypuśćmy, że w m próbach zaobserwowaliśmy k wystąpień pewnego zdarzenia A dotyczącego zmiennej o q unikalnych wartościach. Szacując prawdopodobieństwo na podstawie częstości, powinniśmy napisać $P(A) \approx k/m$. Stosując poprawkę LaPlace'a, oszacowanie przybiera postać:

$$P(A) \approx \frac{k + 1}{m + q}. \quad (2)$$

3 NBC wersja ciągła

Nie wszystkie zbiory danych zawierają dane dyskretne. Istnieje wariant NBC dla atrybutów ciągłych. Jako rozszerzenie opracuj nowy klasyfikator bayesowski realizujący klasyfikację danych z winem w wariancie ciągłym, czyli bez wykonywania dyskretyzacji danych.

3.1 Kolejne kroki w NBC dla zmiennych ciągłych

1. Wczytaj dane z pobranego pliku tekstowego `wine.data` do macierzy `numpy` tak samo jak w wariancie dyskretnym.
2. Pomiń dyskretyzację!
3. Podziel dane na część uczącą i testową jak w wariancie dyskretnym.
4. Napisz klasę reprezentującą naiwny klasyfikator Bayesa w wariancie ze zmiennymi ciągłymi. Przyjmując takie same założenia jak w zadaniu z danymi dyskretnymi.
5. Zastosuj estymaty funkcji gęstości oparte na rozkładach normalnych. W szczególności zaplanuj odpowiednie struktury danych do przechowywania średnich i odchyłeń standardowych dla poszczególnych gęstości warunkowych. (Szczegóły w podrozdziale 3.3).
6. Wyznacz dokładność jak w wariancie dyskretnym i porównaj oba warianty NBC.
7. Porównaj zgodność działania otrzymanego klasyfikatora (Twojej implementacji) z gotową implementacją `GaussianNB` dostępną w pakiecie `sklearn.naive_bayes`.

3.2 Uzyskanie odpowiedzi z klasyfikatora

W ramach tego ćwiczenia obliczanie odpowiedzi klasyfikatora (w metodach `predict_proba`, `predict`) może być realizowane zgodnie ze wzorem:

$$y^* = \arg \max_{y \in \{1, \dots, K\}} \prod_{j=1}^n p_j(x|Y = y) \quad (3)$$

3.3 Funkcje gęstości

Przypuśćmy że wszystkie rozpatrywane zmienne wejściowe są ciągłe, i przypomnijmy notację dla zbioru danych postaci: $D = \{(\mathbf{x}_i, y_i)\}_{i=1, \dots, m}$, gdzie $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{in}) \in \mathbb{R}^n$ są wektorami cech rzeczywistoliczbowych, zaś y_i etykietami klas. A zatem chcąc przygotować NBC w wariancie ciągłym gaussowskim, musimy wyznaczyć $2 \cdot n \cdot K$ parametrów. Oznaczmy je z użyciem pary indeksów (gdzie $j = 1, \dots, n$, $y = 1, \dots, K$):

- μ_{jy} – to średnie
- σ_{jy} — to odchylenia standardowe

dla wszystkich warunkowych rozkładów zmiennych $X_j|Y = y$.

Gęstość wybranego rozkładu to:

$$p_j(X_j = x|Y = y) = \frac{1}{\sigma_{jy}\sqrt{2\pi}} e^{-\frac{(x-\mu_{jy})^2}{2\sigma_{jy}^2}}, \quad (4)$$

stosuje się poniższe wzory do wyznaczenia estymat odpowiednio średniej i odchylenia standardowego:

$$\mu_{jy} = \frac{1}{m} \sum_{\substack{i=1 \\ y_i=y}}^m x_{ij}, \quad (5)$$

$$\sigma_{jy} = \sqrt{\frac{1}{m-1} \sum_{\substack{i=1 \\ y_i=y}}^m (x_{ij} - \mu_{jy})^2}. \quad (6)$$

Uwaga — czynnik normalizujący $\frac{1}{m-1}$ widoczny w drugim wzorze nie jest pomyłką, a wynika z posłużenia się tzw. *estymatorem nieobciążonym*.

4 Wymagania

Zadanie na 2 tygodnie na dwie oceny 5.0 :-)

1. Implementacja wariantu dyskretnego. (4 pkt)
2. Dodanie poprawki LaPlace'a. (1 pkt)

3. Implementacja wariantu ciągłego i porównanie z klasyfikatorem z biblioteki. (4 pkt)
4. Przetestowanie na wybranym zbiorze z repozytorium UCI innym niż „wine” odpowiedniego wariantu (dla zmiennych ciągłych lub dyskretnych adekwatnie do danych). (1 pkt)
5. Dla chętnych: wykorzystać zabieg logarytmowania, zmodyfikować zmodyfikować wariant dyskretny zgodnie ze wzorem (6.26) ze skryptu a ciągły zgodnie ze wzorem 6.28.

5 Przekazanie zadań

Kod z rozwiązaniem proszę podpiąć w Teams. Proszę w nazwach plików źródłowych zawierać swoje nazwisko celem łatwiejszej identyfikacji.