

WSI Modele bayesowskie sprawozdanie

Małgorzata Grzanka

10.06.2024

Spis treści

1	Wstęp	1
2	Zasada działania naiwnego klasyfikatora Bayes’a	1
2.1	Sposób klasyfikacji	1
2.2	Traning modelu i dane ciągłe	2
3	Hiperparametry	2
4	Badanie jakości działania klasyfikatora	2
4.1	Ostateczny wynik	2
4.2	Walidacja krzyżowa	3

1 Wstęp

Celem ćwiczenia było zaimplementowanie naiwnego klasyfikatora Bayes’a. Następnie, napisany klasyfikator należało przetestować dla zbioru danych Breast Cancer Wisconsin Dataset (https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html).

2 Zasada działania naiwnego klasyfikatora Bayes’a

2.1 Sposób klasyfikacji

Podstawa naiwnego klasyfikatora Bayes’a jest wzór na prawdopodobieństwo warunkowe Bayes’a:

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

Wzór ten jest wykorzystywany w naiwnym klasyfikatorze Bayes’a do przewidywania klasy próbki pod warunkiem jej cech (B odpowiada klasie próbki, a A

- wektorze jej cech). Tak więc, prawdopodobieństwo, że próbka o parametrach $x_1, x_2, x_3, \dots, x_n$ należy do klasy X , dane jest wzorem

$$P(X|x_1, x_2, x_3, \dots, x_n) = \frac{P(x_1, x_2, x_3, \dots, x_n|X)P(X)}{P(x_1, x_2, x_3, \dots, x_n)}$$

Słowo "naiwny" w nazwie tego klasyfikatora pochodzi od faktu, że zakłada on, że zmienne (cechy) $x_1, x_2, x_3, \dots, x_n$ są od siebie niezależne. Wtedy, wzór upraszcza się do postaci.

$$P(X|x_1, x_2, x_3, \dots, x_n) = \frac{P(x_1|X)P(x_2|X)\dots P(x_n|X)P(X)}{P(x_1)P(x_2)\dots P(x_n)}$$

Klasyfikacja następuje po porównaniu prawdopodobieństw należenia próbki do danej klasy. Przewidywana klasa jest ta, której prawdopodobieństwo dla tej próbki jest największe. Zatem, jako że mianownik jest taki sam dla każdego prawdopodobieństwa dla jednej próbki, klasyfikacja próbki odbywa się na podstawie wzoru:

$$\hat{X} = \max_X (P(x_1|X)P(x_2|X)\dots P(x_n|X)P(X))$$

2.2 Trening modelu i dane ciągłe

Proces trenowania modelu polega na znalezieniu odpowiednich wartości prawdopodobieństw na podstawie danych treninowych.

O ile znalezienie prawdopodobieństw danych dyskretnych jest bardzo intuicyjne (prawdopodobieństwo klasyczne), to jeśli chodzi o ciągłe, w moim klasyfikatorze zastosowałam dla nich rozkład normalny. Na podstawie danych treningowych, dla każdej cechy z próbki program oblicza średnią i odchylenie standardowe w każdej z możliwych klas i na ich podstawie określa prawdopodobieństwa $P(x_i|X)$.

3 Hiperparametry

Model nie ma żadnych hiperparametrów.

4 Badanie jakości działania klasyfikatora

4.1 Ostateczny wynik

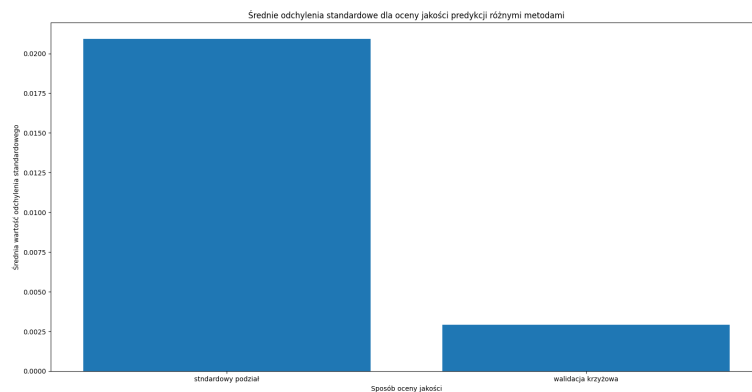
Poprawność przewidywań mojego klasyfikatora zbadałam dzieląc zbiór danych na zbiór treningowy, walidacyjny i testowy. Dla zbioru testowego, model osiągnął skuteczność około 96.49%.

0.9649122807017544

Rysunek 1: Skuteczność klasyfikatora na zbiorze testowym

4.2 Walidacja krzyżowa

Zbadałam różnice w ocenie klasyfikatora za pomocą k-krotnej walidacji krzyżowej w porównaniu ze zwykłym, stałym podziałem danych. Dla obydwu tych metod zbadałam odchylenie standardowe wyników. Uruchamiałam walidację dla stałego podziału z różnymi rozłożeniami danych w podziale, a dla walidacji krzyżowej - z inną krotnością. Rezultat widać na wykresie poniżej.



Rysunek 2: Odchylenie standardowe walidacji krzyżowej i standardowego podziału

Jak widać, dla sztywnego podziału odchylenie standardowe pomiarów okazało się dużo większe, co pokazuje większą niepewność wyników otrzymywanym właśnie tą metodą. Jakość zmierzona na zbiorze walidacyjnym może więc zależeć od tego, jak dane zostaną podzielone.

W k-krotnej walidacji krzyżowej, zbiór danych jest dzielony na k fragmentów. Modele trenuje się na k-1 fragmentach a jakość się mierzy na pozostałym. Ostateczny wynik jest uśredniony z k iteracji, gdzie kolejne z k fragmentów stają się nowym zbiorem walidacyjnym. Dzięki temu, każda próbka w zbiorze będzie oceniana, co prowadzi do większej stabilności wyników.