



Uniwersytet Mikołaja Kopernika w Toruniu
Wydział Matematyki i Informatyki
Hibatűrő Rendszerek Kutatócsoport

Wykrywanie występowanie chorób serca z wykorzystaniem uczenia maszynowego nadzorowanego

PRACA INŻYNIERSKA

Autor
Magdalena Szulc

Promotor
Piotr Przymus

September 17, 2019

Contents

Oświadczenie plagiatowe

Treść oświadczenia

Toruń, September 17, 2019

Magdalena Szulc
student

Abstrakt

Abstrakt pl

Abstract

Abstrakt eng

**Wykrywanie występowanie chorób
serca, wykorzystanie uczenia
maszynowego nadzorowanego**

Wstęp

Sztuczna inteligencja wśród szerokiego zakresu swoich zastosowań może zostać wykorzystana do analizy bardziej lub mniej złożonych danych medycznych, w celu przewidzenia wystąpienia choroby u konkretnej osoby, bez udziału procesu myślowego od stony specjalisty.

Do tego przeznaczenia idealnie nadaje się uczenie nadzorowane (ang. *supervised learning*) tj. rodzaj uczenia maszynowego zakładający przykład, dane wejściowe będące wzorcem na podstawie którego wyszukiwane są zależności.

Zależności typujące osoby chore, zakwalifikowane na podstawie konkretnych objawów lub wyników badań.

W dzisiejszych czasach choroby sercowo-naczyniowe stanowią najczęstszą przyczynę zgonów, a liczba osób cierpiących na te dolegliwości stale rośnie. Głównymi przyczynami zachorowalności diagnozowanymi przez specjalistów są niski poziom świadomości i profilaktyki chorób serca. Objawy są tym silniejsze im gorszy jest stan chorobowy pacjenta.

Dlatego prowadzone są intensywne prace nad zwiększeniem dostępności badań, które wspomogą diagnostykę kardiologiczną na jak najwcześniejszym etapie.

Powodem szukania dokładniejszych sposobów diagnozowania są również wysokie koszty leczenia generowane przez choroby układu krwionośnego. Według analityków firmy konsultingowej KPMG ¹ w 2011 r. koszty diagnostyki i terapii chorób serca wyniosły ponad 15 miliardów polskich złotych.

Nadzieją jaką pokładana jest w machine learningu jest szybsza diagnostyka chorób ułatwiająca oraz przyspieszająca proces ich leczenia. Zastosowanie uczenia maszynowego w medycynie, pozwala również na przetwarzanie historycznych danych, w celu poszerzenia zasobów informacji które mogą zostać wykorzystane przez lekarza.

Słowa kluczowe: uczenie maszynowe, uczenie nadzorowane

¹międzynarodowa sieć firm audytorsko-doradczych ze szczególnym uwzględnieniem branży dóbr konsumpcyjnych, usług finansowych, nieruchomości i budownictwa, technologii informacyjnych, mediów i komunikacji (TMT), transportowej (TSL), produkcji przemysłowej, a także sektora publicznego

Cel i zakres pracy

Celem pracy jest porównanie wybranych algorytmów uczenia maszynowego nadzorowanego, przy założeniu że dane wejściowe są wybrakowane, a w rezultacie zbudowanie modelu który na podstawie danych medycznych wystawia diagnozę o występowaniu zaburzeń sercowo-naczyniowych lub ich braku.

Dane medyczne wyróżniają się tym, że trudno uzyskać do nich dostęp, najczęściej nie są to informacje, które się udostępnia do użytku publicznego. Z tego powodu, kluczowym krokiem jest wybór cech branych pod uwagę przy tworzeniu modelu.

Zatem odpowiedź na pytanie jak wybrakowanie danych mocno wpływa na rezultat i czy istnieją różnice między zastosowaniem wybranych algorytmów nauczania nadzorowanego wymaga przedstawienia porównania łatwości tworzenia modelu, dokładności, złożoności oraz czasu uzyskania odpowiedzi.

W pracy opisano następujące algorytmu uczenia nadzorowanego:

- lasy decyzyjne (ang. *decisions-forests*)
- metoda wektorów nośnych (ang. *support vector machines*, SVM)
- k-najbliższych sąsiadów (ang. *k-neares neighbours*, KNN)

Praktyczna część pracy napisana została w języku Python z wykorzystaniem scikit-learn, obsługującym wiele algorytmów maszynowego uczenia się w tym uczenia nadzorowanego i docelowo wybranych algorytmów przedstawionych w teoretycznej części pracy.

Biblioteka opiera się o Numerical Python, zestaw narzędzi do obliczeń na macierzach, wektorach oraz o pakiet Science python umożliwiający metody numeryczne takie jak całkowanie, różniczkowanie itp. .

Do przygotowania danych wykorzystano zestaw narzędzi Pandas, ułatwiający tworzenie struktur danych i ich analizę. W celu wizualizacji wyników w postaci wykresów zastosowano Matplotlib.

Chapter 1

Wprowadzenie teoretyczne

Uczenie maszynowe (ang. *machine learning*, ML) to dziedzina zajmująca się zestawem algorytmów, które analizując duże zbiory danych wysatwiają predykcję na temat zadanego problemu. Uczenie maszynowe zależnie od sposobu *trenowania* algorytmu wyróżnia min. uczenie nadzorowane oraz uczenie bez nadzoru. Dane oraz wynik który przewidywanie ma osiągnąć uzależniają wybór kategorii.

Uczenie maszynowe nadzorowane (ang. *supervised learning*) to klasa algorytmów uczenia maszynowego, która bazuje na poetykietowanych już danych wejściowych. Ten typ uczenia świetnie nadaje się do rozwiązywania problemów z zakresu klasyfikacji. Nadzór polega na porównaniu rezultatów działania modelu z wynikami które są zawarte w danych wejściowych (*dane oznaczone*). Algorytm po osiągnięciu żądanej efektywności jest w stanie dokonać klasyfikacji przykładu dla którego nie posiada odpowiedzi. Sprawdza się to obecnie w rekomendacji produktów oraz diagnozie chorób.

Uczenie maszynowe bez nadzoru (ang. *unsupervised learning*) to klasa algorytmów uczenia maszynowego która głównie rozwiązuje problemy grupowania. Dane dostarczane do modelu nie zawierają *oznaczeń*, zatem nauczanie polega na wyciąganiu konkluzji z poprzednio wykonanych iteracji. Na skuteczność modeli budownych w oparciu o uczenie bez nadzoru wpływ ma rozmiar dostarczonego do nauki zbioru danych, im jest on większy tym bardziej wzrasta efektywność. Takie zbiory można uzyskać rejestrując dane na bieżąco dlatego do najczęstszych zastosowań tej klasy algorytmów, można zaliczyć rozpoznawanie mowy czy obrazu.

Podział osób na kategorie cierpiące na choroby sercowo-naczyniowe oraz zdrowe, to dylemat klasyfikacyjny nadający się do rozwiązania za pomocą algorytmów uczenia maszynowego nadzorowanego i na nich skupia się dalsza część pracy.

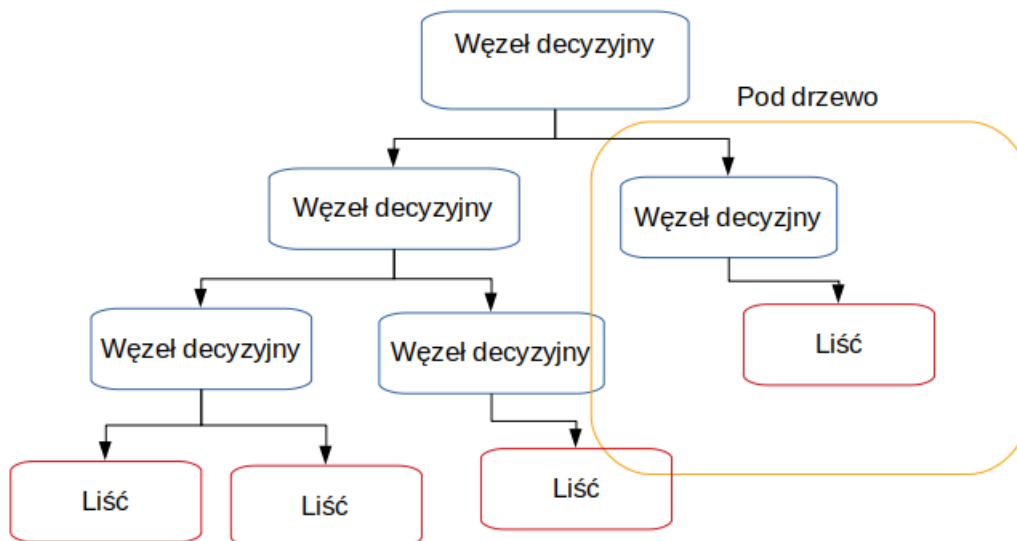


Figure 1.1: Schemat 1

1.1 Wybrane algorytmy uczenia maszynowego nadzorowanego

Drzewa decyzyjne (ang. *decisions trees*) są uznawane za najprostszyszy i najbliższy ludzkiemu zrozumieniu algorytm uczenia, który swoją nazwę zawdzięcza graficznej reprezentacji w postaci drzewa. Każdy węzeł oznacza atrybut, na podstawie którego następuje rozróżnienie. W modelu kluczowa jest kolejność cech, które występują po sobie ponieważ determinuje to otrzymany rezultat.

Prawie każdy algorytm uczenia maszynowego nadzorowanego można podzielić na dwa etapy. W pierwszym opracowywany jest wzorzec, na którym bazują późniejsza predykcja. Etap nauki dla drzewa decyzyjnego polega na typowaniu atrybutów, które stają się węzłami decyzyjnymi, dzielącymi rekordy na dwa mniejsze zestawy i tak aż nie ma możliwości dalszego podziału.

O metodologii drzew decyzyjnych oparta jest dokładniejsza forma nauczania nadzorowanego: *losowe lasy decyzyjne*.

Losowe lasy decyzyjne (ang. *random decision forests*) to technika polegająca na połączeniu wielu drzew decyzyjnych w celu uniknięcia problemu z *nadmiernym dopasowaniem* do treningowego zestawu danych na którym został przeszkolony. Utworzony szablon aby poprawnie działać na danych testowych i służących weryfikacji, nie może stać się charakterystycznym przypadkiem rozwiązującym przypadek testowy.

W tym celu dla losowych lasów decyzyjnych najpierw stosuje się **agregację bootstrap'ową**.

Z treningowego zestawu danych losuje się, co ważne z możliwymi powtórzeniami, wiersze

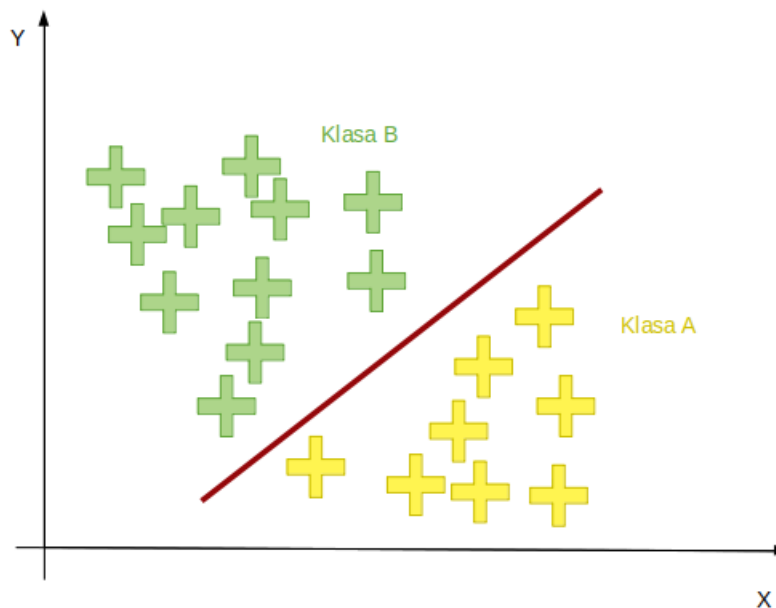


Figure 1.2: Schemat 2 ¹

danych dla których trenowany będzie model. Jako rezultat brana jest większość lub średnia wartości uzyskanych wyników dla poszczególnych drzew decyzyjnych. Dodatkowo dla drzew decyzyjnych w lasach losowych, atrybuty odpowiadające za kategoryzację są wybierane z wylosowanego podzbioru.

Wśród zalet lasów losowych należy wyróżnić iż potrafią one trafnie wykalkulować brakujące wartości cech. Idealnie znajdują zastosowanie dla realnych danych, których zasadniczym problemem jest ich niekompletność.

Dane medyczne posiadają szeroką wariację zmiennych z dużym prawdopodobieństwem wybrakowania, zastosowanie do nich lasów decyzyjnych ma potencjał na pozytywne rezultaty.

Metoda wektorów nośnych (ang. *support vector machines*, skr. **SVM**) to algorytm uczenia maszynowego nadzorowanego, który każdy parametr z dostępnych cech dla danych wejściowych, traktuje jako punkt w przestrzeni. Na podstawie ułożenia punktów dzieli się je na 2 klasy. Graficznie jest to reprezentowane przez prostą dla której odległość między najbliższymi dwoma punktami dla wektorów jest możliwie największa. Taka prosta nazywana jest *prostą marginalną* i powstaje ona poprzez generowanie i selekcję tych prostych które rzetelnie szufladkują klasy danych.

Techinka ta gwarantuje precyzyjniejsze rezultaty niż drzewa decyzyjne, niestety dla dużych zbiorów danych czas trwania szkolenia znacznie się wydłuża oraz istnieją przypadki dla których podział jedną prostą jest niewykonalny, taki przypadek reprezentuje rozkład na schemacie nr. 2.

¹Schemat wzorowany na http://res.cloudinary.com/dyd911kmh/image/upload/f_auto,q_auto:best/v1526288453/index3_s

²Schemat wzorowany na http://res.cloudinary.com/dyd911kmh/image/upload/f_auto,q_auto:best/v1526288453/index3_s

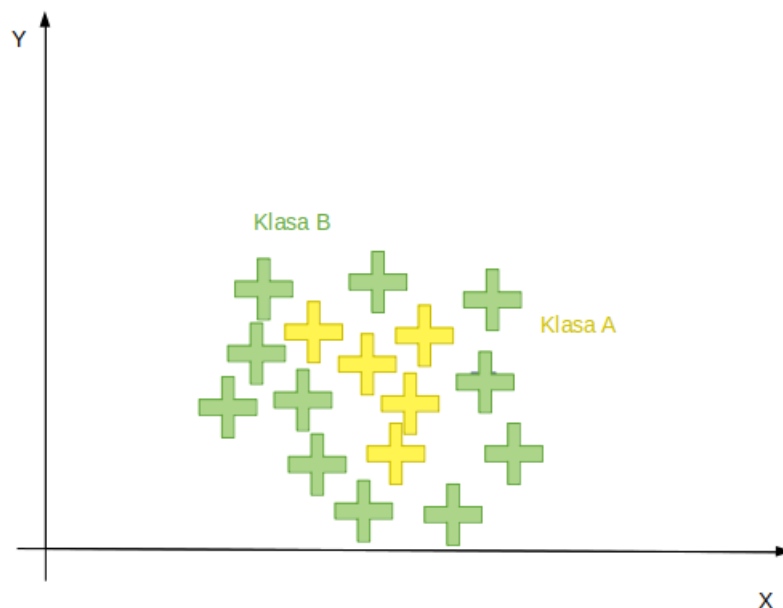


Figure 1.3: Schemat 2 ²

Zbór danych wykorzystany w pracy nie jest aż tak kolosalny by zaszkodzić wydajności, a małym kosztem można uzyskać celność rozwiązania zadanego problemu: wykrywania występowania chorób sercowo-naczyniowych. Istnieje jednak ryzyko uzyskania rozkładu wartości który wyklucza graficzną fragmentację zestawu danych na dwie części za pomocą prostej.

K najbliższych sąsiadów (ang. *k nearest neighbours*, skr. **KNN**) to algorytm uczenia maszynowego nadzorowanego operujący swoje estymacje dla konkretnego przypadku danych na wartościach jego K najbliższych sąsiadów (punktów) liczonych min. dla przestrzeni Euklidesowej, miasto (in. Manhattan) oraz Mińkowskiego.

Atrybut który nastraja proces uczenia się modelu i ma na niego największy wpływ określany jest jako hiperparametr. Dla KNN jest to liczba sąsiadów, im większa ilość jednostek mających wpływ, tym wierniejsze będą wyniki. Potęguje się wtedy niestety złożoność czasowa algorytmu, znacząco już większa od przedstawionych powyżej innych algorytmów.

W celu przewidzenia wartości dla nowych danych, należy odnaleźć K najbliższych punktów wyliczając odległości, a następnie przypisać odpowiedź implikowaną przez większość sąsiadów. Dla wartości K równej jeden, metoda ta nazywana jest algorytmem najbliższego sąsiada.

Dla lekarza wartością dodatnią jest wykrycie zależności które decydują o uznaniu lub zaprzeczeniu występowania choroby. Zastosowanie algorytmu KNN może nie tylko zakwalifikować osoby chorujące na serce, ale również ułatwić swoją graficzną reprezentacją wpływ cech na ostateczny osąd próbki.

³Schemat wzorowany na http://res.cloudinary.com/dyd911kmh/image/upload/f_auto,q_auto:best/v1526288453/index3_s

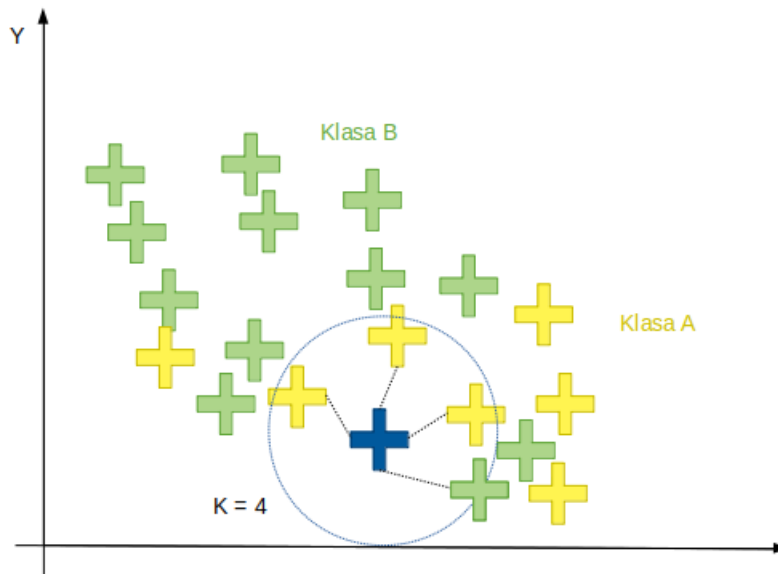


Figure 1.4: Schemat 3³

1.2 Zestawienie efektywności działania algorytmów

Konfrontacja technik uczenia maszynowego zależnie od zestawu danych będzie dawała odmienne wyniki ze względu na ich predyspozycje do zajmowania się odpowiednimi zbiorami danych.

Potencjał algorytmów dla niewielkiego kompletu danych zawierającego wartości wybrakowane zostanie omówiony w późniejszych rozdziałach pracy.

Zczynając od drzew decyzyjnych, można od razu stwierdzić ich niski potencjał. Istnieje zbyt duże prawdopodobieństwo dopasowania się do modelu treningowego, gdyż wspomniany zbiór danych wejściowych nie jest wystarczająco liczny. Dlatego w dalszej części pracy omówione zostaną lasy decyzyjne.

Większej dokładności można się spodziewać po metodzie wektorów nośnych, ale jego złożoność czasowa oraz pamięciowa mogą zaniżyć jego ogólną klasyfikację.

Wskaźniki wydajności

Określenie stopnia, w jakim skonstruowany model z powodzeniem realizuje wyznaczone zadanie należy do wskaźnika wydajności. Przykładem nieprawidłowego wyboru może być próba przewidzenia wystąpienia rzadkiej choroby u pacjenta i określenie głównym miernikiem *dokładność*. W takim scenariuszu klasyfikacja wszystkich pacjentów jako zdrowych, daje niewiele odbiegającą od perfekcji dokładność, a jednocześnie błędnie osądzać każde wystąpienie choroby.

K-krotna walidacja krzyżowa (ang. *Fold Cross-Validation*) to metodyka weryfikacji poprawności modeli nauczania maszynowego. Opiera się ona na wyporze wartości swojego

hiperparamtru jakim jest K , które może przyjąć dowolną wartość mniejszą lub równą od rozmiaru danych.

Po wyborze hiperparamtru następuje segmentacja danych na K jednakowej wielkości zestawów. Wykonywanych jest k iteracji, w każdej z nich na $k-1$ kolekcjach model jest trenowany, a na pozostałej jednej weryfikowany. Procedura efektywnie pomaga ocenić poprawność działania modelu i zastosowanego algorytmu.

1.3 Model Danych

1.3.1 Omówienie danych

Budowa modelu zależna od danych ### Obróbka danych ### Budowa modelu
Metody optymalizacji #### Implementacja algorytmu 1: Losowe lasy decyzyjne
Implementacja algorytmu 2: Metoda wektorów nośnych #### Implementacja
algorytmu 3: K najbliższych sąsiadów

1.4 Wnioski i walidacja rozwiązania

1.4.1 Algorytm 1: Wyniki wniosków: Losowe lasy decyzyjne

1.4.2 Algorytm 2: Wyniki wniosków: Metoda wektorów nośnych

1.4.3 Algorytm 3 : Wyniki wniosków: K najbliższych sąsiadów

1.4.4 Porównanie algorytmów : złożoność czasowa , dokładność , złożoność implementacyjna , wpływ danych wykorzystywanych w modelu

1.4.5 Podsumowanie i opisanie wpływu danych na model

1.5 Bibliografia

1.6 Spis ilustracji

1.7 Spis tabel

Chapter 2

Podsumowanie

Treść podsumowania

Acknowledgments

Content of acknowledgments

List of Tables

List of Figures

Appendix A

Függelék

A.1 Válasz az „Élet, a világmindenség, meg minden” kérdésére

A Pitagorasz-tételből levezetve

$$c^2 = a^2 + b^2 = 42.$$

A Faraday-indukciós törvényből levezetve

$$\text{rot} E = -\frac{dB}{dt} \quad \longrightarrow \quad U_i = \oint_{\mathbf{L}} \mathbf{E} d\mathbf{l} = -\frac{d}{dt} \int_A \mathbf{B} d\mathbf{a} = 42.$$