



Neuronowy model dawcy krwi

KWD

Krzysztof MICHALSKI

Rafał POKRYWKA

27.01.2020

Streszczenie

W tym projekcie zajmowaliśmy się przewidywaniem, czy pacjent odda krew w marcu 2007, na podstawie danych z Blood Transfusion Service Center Data Set. Było to zatem zadanie klasyfikacji binarnej. Zastosowaliśmy kilka modeli - regresję liniową, regresję z użyciem drzew decyzyjnych i lasów losowych oraz w pełni połączoną sieć neuronową. Modele zostały poddane analizie skuteczności i na jej podstawie najlepszy okazał się ... z dokładnością ...

Spis treści

| | | |
|----------|---|-----------|
| 1 | Wprowadzenie | 1 |
| 1.1 | Opis problemu | 1 |
| 2 | Opis metody | 3 |
| 2.1 | Wprowadzenie teoretyczne | 3 |
| 2.2 | Badania symulacyjne | 4 |
| 2.3 | Wyniki symulacji dla różnych modeli | 9 |
| 2.3.1 | Dane niezrównoważone | 9 |
| 2.3.2 | Dane zrównoważone | 11 |
| 3 | Podsumowanie | 13 |
| A | Kod programu | 14 |

Rozdział 1

Wprowadzenie

Dane do projektu uzyskaliśmy ze strony Blood Transfusion Service Center Data Set. Składały się one z dwóch plików *transfusion.data* i *transfusion.names*.

W pliku *transfusion.data* znajdowało się 748 przykładów opisanych 4 cechami:

- R (Recency - months since last donation)
- F (Frequency - total number of donation)
- M (Monetary - total blood donated in c.c.)
- T (Time - months since first donation)

oraz jedną etykietą - binary variable representing whether he/she donated blood in March 2007 (1 stand for donating blood; 0 stands for not donating blood). Zatem zadanie można zaliczyć do kategorii klasyfikacji binarnej. Klasy te były słabo zrównoważone - klasę 1 posiadało około 24% przykładów, a klasę 0 - 76%.

W pliku *transfusion.names* znajdował się opis danych - źródło, opis atrybutów oraz autorzy.

1.1 Opis problemu

Zadanie polega na klasyfikacji binarnej przykładów na podstawie 4 cech numerycznych. Celem klasyfikacji, jest określenie, czy nowy pacjent kliniki odda krew po pewnym okresie czasu. W tym celu można wykorzystać różne techniki uczenia maszynowego takie jak regresja liniowa, drzewa decyzyjne, lasy losowe czy sieci neuronowe. Ze względu na niewielką ilość przykładów oraz cech je opisujących, skuteczność takiej klasyfikacji może być niewielka

- zbiory danych o takiej wielkości są niewystarczające dla dobrego wytrenowania na przykład sieci neuronowej. Ze względu na niezrównoważenie klas, modele mogą mieć tendencję do częstszego przewidywania dominującej klasy - bo zapewnia to dużą dokładność takiego modelu. Można spróbować zrównoważyć zbiór wybierając podzbiór przykładów o takiej samej reprezentacji klasy pozytywnej i negatywnej.

Rozdział 2

Opis metody

2.1 Wprowadzenie teoretyczne

W celu rozwiązania zadania wykorzystaliśmy kilka metod uczenia maszynowego:

Regresję logistyczną - czyli metodę służącą do klasyfikacji binarnej, pozwalającą określić prawdopodobieństwo przynależności przykładu do jednej z dwóch klas. Jeżeli prawdopodobieństwo to przekracza 50% oznacza to, że model przewiduje, że dany przykład należy do tej klasy. Do nauki regresora logistycznego wykorzystywana jest:

- Logarytmiczna funkcja straty, która wykonuje nieliniową transformację wyjściowych wartości na przedział od 0 do 1, tak, żeby suma wartości obu klas wynosiła 1.
- Metoda spadku wzdłuż gradientu prostego - pozwalająca zbliżyć się do minimum globalnego powyższej funkcji straty.
- Funkcję logistyczną - pozwalającą określić wartość prawdopodobieństwa dla każdej klasy.

Drzewa decyzyjne - czyli wszechstronny algorytm uczenia maszynowego, który można wykorzystywać do różnych zadań - zarówno klasyfikacji (binarnej i wieloklasowej) jak i do zadań regresji. Podejście tego modelu do klasyfikacji binarnej polega na dzieleniu każdej z cech danych na dwa zbiory w sposób liniowy - tak, żeby zbiory te były jak najlepiej rozdzielone. Oznacza to, że w podzbiorach przykładów po każdej ze stron tej tak zwanej granicy decyzyjnej, znadowało się jak najwięcej przykładów należących do jednej klasy - podzbiory te powinien być jak najmniej zanieczyszczone przykładami należącymi do drugiej grupy. Wysokość takiego drzewa określa ile takich liniowych podziałów cech przykładu może wykonać model.

Lasy losowe - jest to jedna z metod uczenia zespołowego. Polega ona na wytrenowaniu pewnej ilości drzew decyzyjnych, a następnie zastosowaniu jednej z technik głosowania drzew w celu określenia rezultatu. Głosowanie można przeprowadzać na przykład w sposób większościowy - w takim przypadku, każdy model określa, do której klasy według niego należy dany przykład - klasa określana przez zespół to ta, na którą głosowało najwięcej modeli.

Sieci neuronowe - czyli algorytm wielokrotnej zmiany reprezentacji cech określających przykłady, w celu znalezienia reprezentacji pozwalającej - w przypadku klasyfikacji - w jak największym stopniu rozdzielić od siebie przykłady należące do różnych klas. Sieci składają się z:

- **Neuronów** - które są elementami wykonującymi liniową transformację danych otrzymywanych na wejściu - sumę ważoną wejść wraz z tak zwanym członem obciążenia, niezależnym od wartości wejść.
- **Warstw** - które są złożone z neuronów. Neurony w ramach jednej warstwy nie są ze sobą połączone.
- **Połączeń pomiędzy warstwami** - czyli połączeniem danych wejściowych do neuronów pierwszej warstwy, a następnie wyjść neuronów jednej warstwy z wejściami neuronów kolejnej warstwy.
- **Funkcji aktywacji** - określających sposób modyfikacji reprezentacji danych przekazywanych pomiędzy warstwami. W przypadku klasyfikacji ostatnia warstwa stosuje funkcję softmax zwracającą prawdopodobieństwa przynależności danego przykładu do danej klasy.

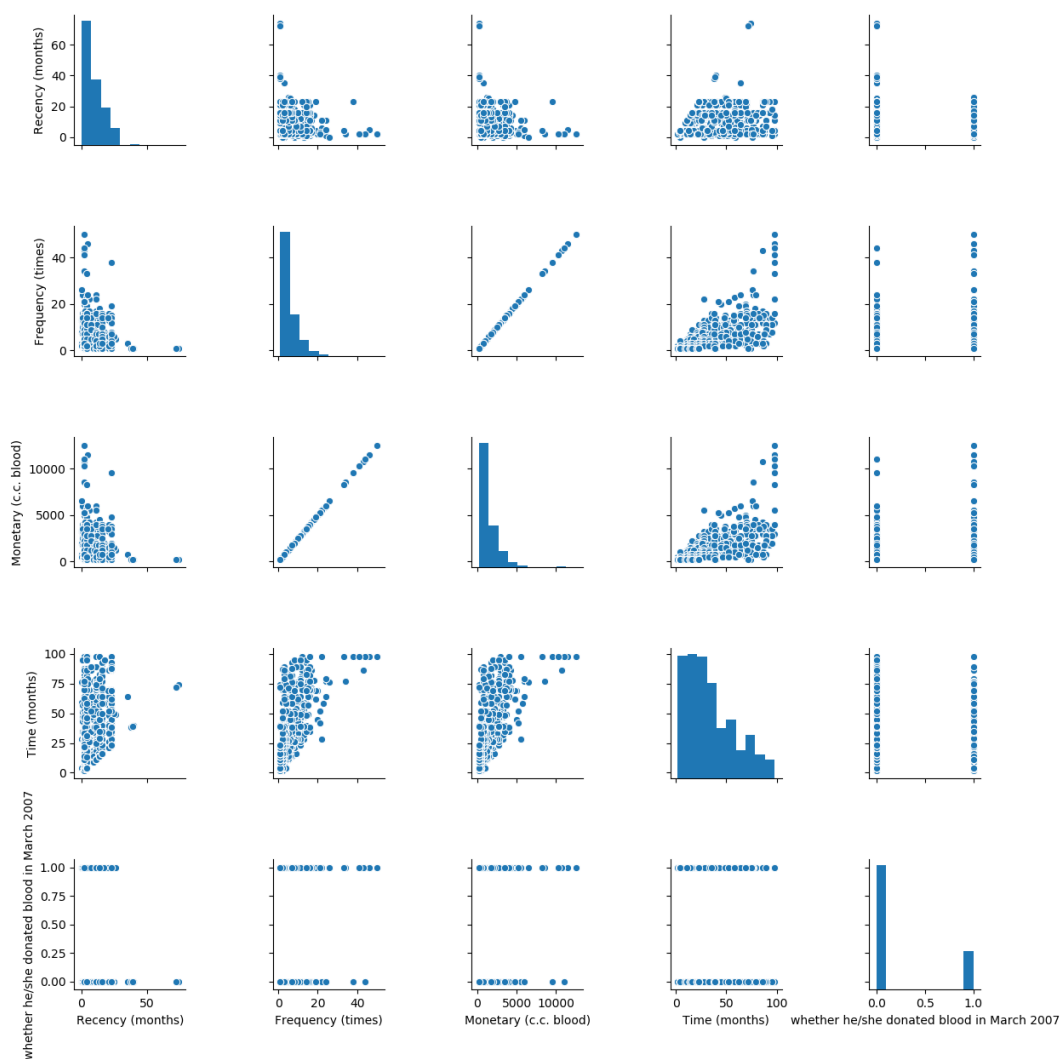
Sieci neuronowe są trenowane z użyciem metody spadku wzdłuż gradientu prostego. Każdy krok takiej nauki określa się jako epokę danej sieci.

2.2 Badania symulacyjne

Na początku programu, w funkcji `load_and_analyze_data`, wczytaliśmy zbiór danych i wyświetliliśmy informacje na jego temat. Żeby zapoznać się z rozkładem danych, sprawdziliśmy czy i jakie są korelacje pomiędzy poszczególnymi kolumnami, oraz jakie są minimalne, średnie i maksymalne wartości.

Analiza zbioru danych w wersji niezrównoważonej:

| | R | F | M | T | D |
|---|-----------|-----------|-----------|----------|-----------|
| R | 1 | -0.182745 | -0.182745 | 0.160618 | -0.279869 |
| F | -0.182745 | 1 | 1 | 0.63494 | 0.218633 |
| M | -0.182745 | 1 | 1 | 0.63494 | 0.218633 |

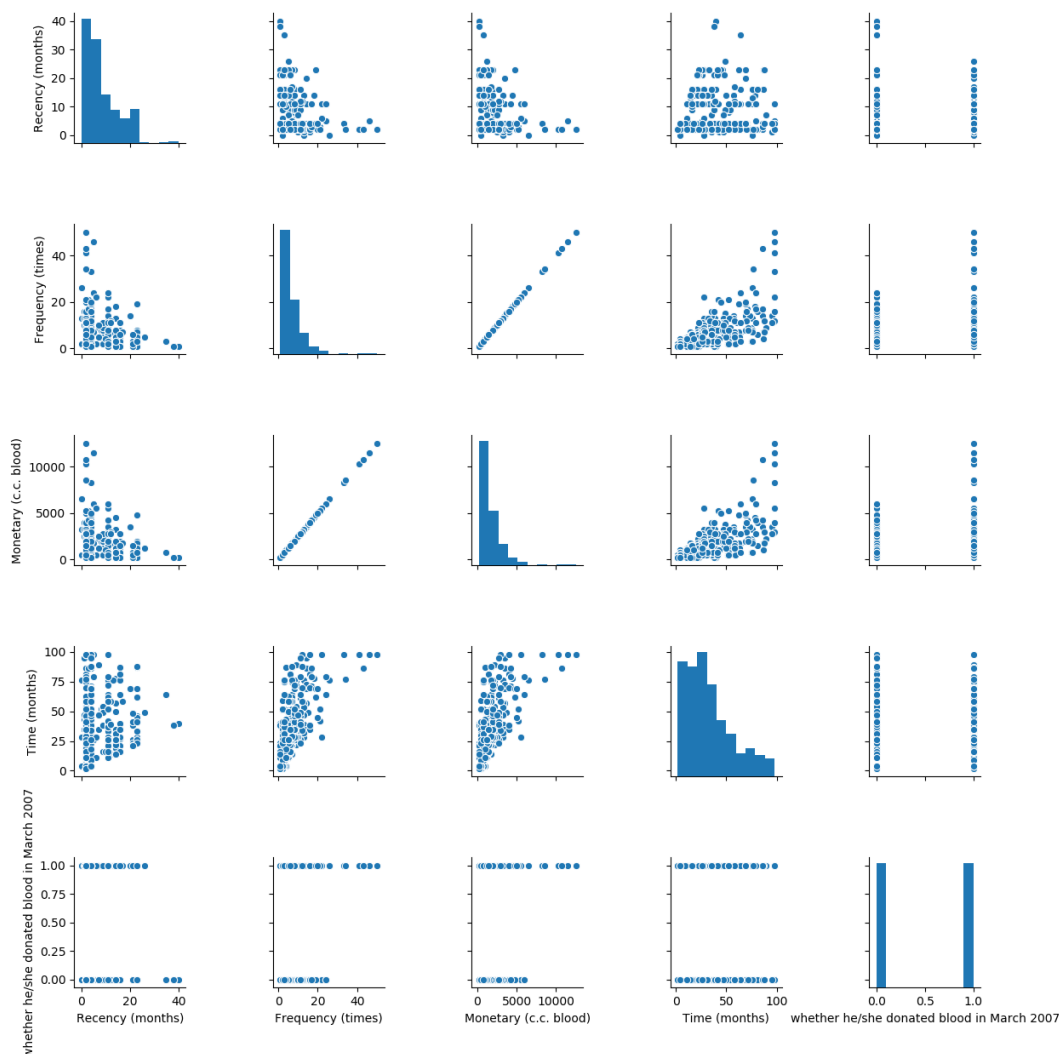


Ponieważ średnia wartość kolumny D (0 lub 1) jest równa 0.237968, to znaczy że tylko 24% wierszy ma jedynkę, czyli ilość przykładów dla klas jest nie zrównoważona. Stworzyliśmy więc dodatkowy zbiór danych, który składa się z równej ilości przykładów dla każdej z klas.

Analiza zbioru danych w wersji zrównoważonej:

| | R | F | M | T | D |
|---|-----------|-----------|-----------|------------|------------|
| R | 1 | -0.218071 | -0.218071 | 0.0949278 | -0.348375 |
| F | -0.218071 | 1 | 1 | 0.685435 | 0.216987 |
| M | -0.218071 | 1 | 1 | 0.685435 | 0.216987 |
| T | 0.0949278 | 0.685435 | 0.685435 | 1 | -0.0152334 |
| D | -0.348375 | 0.216987 | 0.216987 | -0.0152334 | 1 |

| | R | F | M | T | D |
|-------|---------|---------|---------|---------|----------|
| count | 356 | 356 | 356 | 356 | 356 |
| mean | 7.95787 | 6.34551 | 1586.38 | 33.0815 | 0.5 |
| std | 7.19436 | 6.70222 | 1675.55 | 23.8207 | 0.500704 |
| min | 0 | 1 | 250 | 2 | 0 |
| 25% | 2 | 2 | 500 | 16 | 0 |
| 50% | 4 | 5 | 1250 | 28 | 0.5 |
| 75% | 13.25 | 8 | 2000 | 46 | 1 |
| max | 40 | 50 | 12500 | 98 | 1 |



Ponieważ okazało się, że kolumny F i M są ze sobą idealnie skorelowane (F to liczba oddań krwi, a M to suma objętości oddanej krwi), to możemy pozbyć się jednej z nich, a przewidywania programu nie ulegną zmianie. Następnie dane należało przeskalować, aby znajdowały się w jednym przedziale wartości, a ich średnia wynosiła zero. Dzięki temu każda cecha będzie traktowana w równy sposób przez algorytmy uczące.

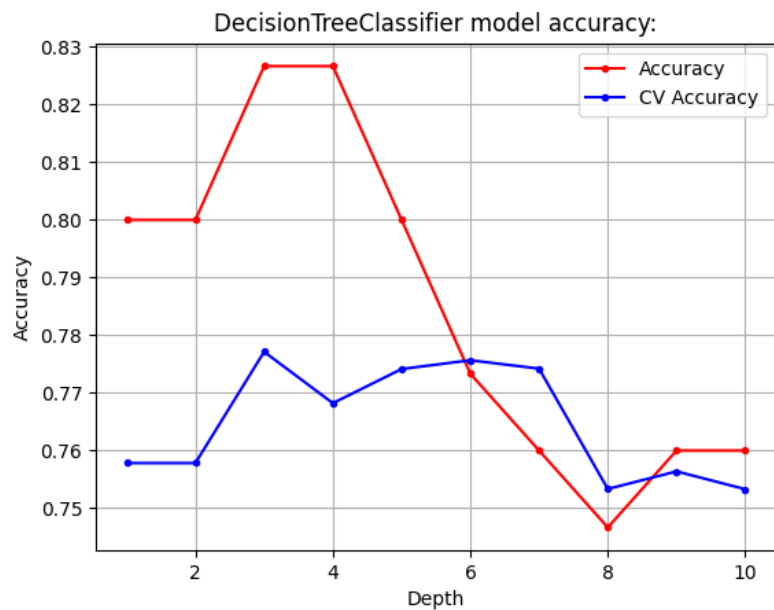
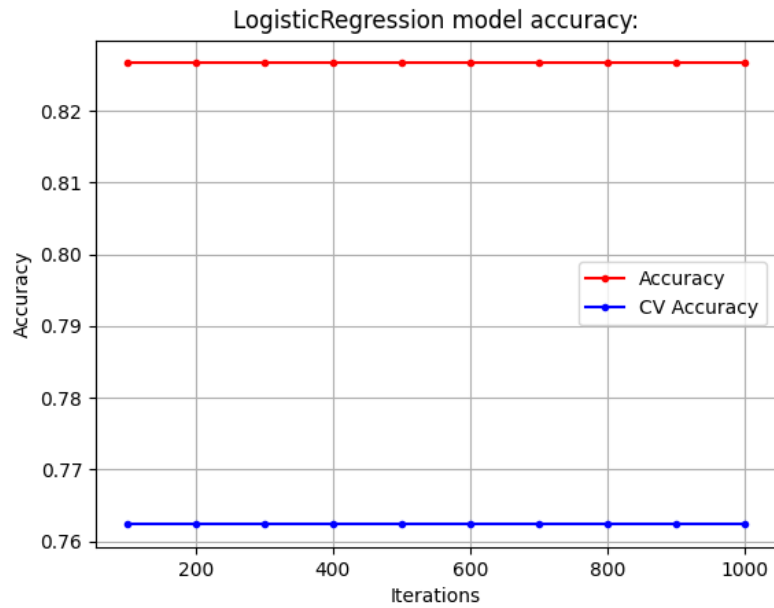
Gdy dane są już odpowiednio przygotowane, można je podzielić na dwa zbiory: trenujący i testujący. W naszym programie 90% to będą dane trenujące, a 10% testujące.

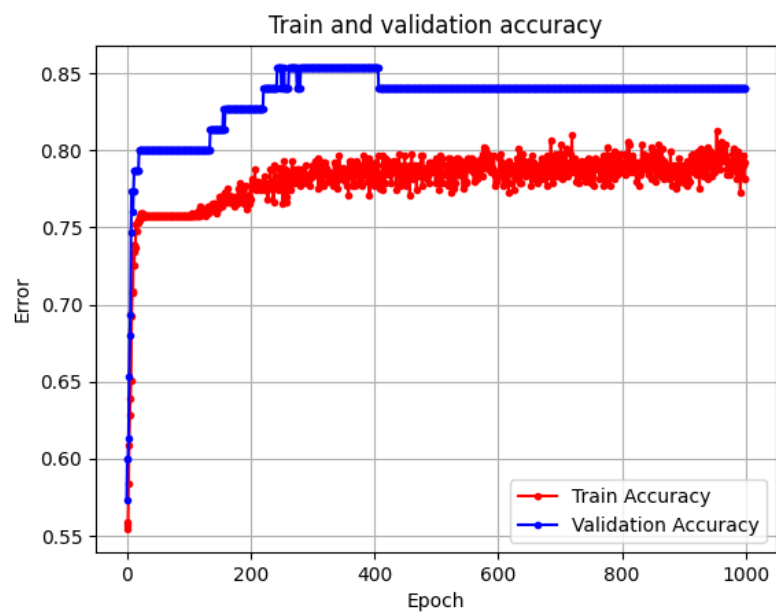
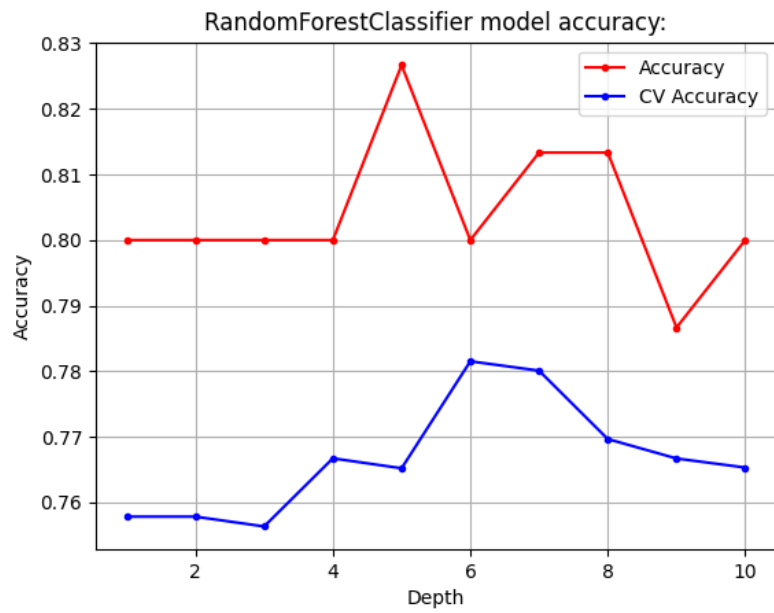
W celu porównania różnych klasyfikatorów, sprawdzamy w pętli wszystkie podane modele (LogisticRegression, DecisionTreeClassifier, RandomForestClassifier). Po wyuczeniu każdego z modeli, wypisujemy na wyjście podstawowe informacje o nim (np. dokładność, macierz konfuzji, ...) oraz wykres

zależności dokładności od różnych parametrów danego modelu.

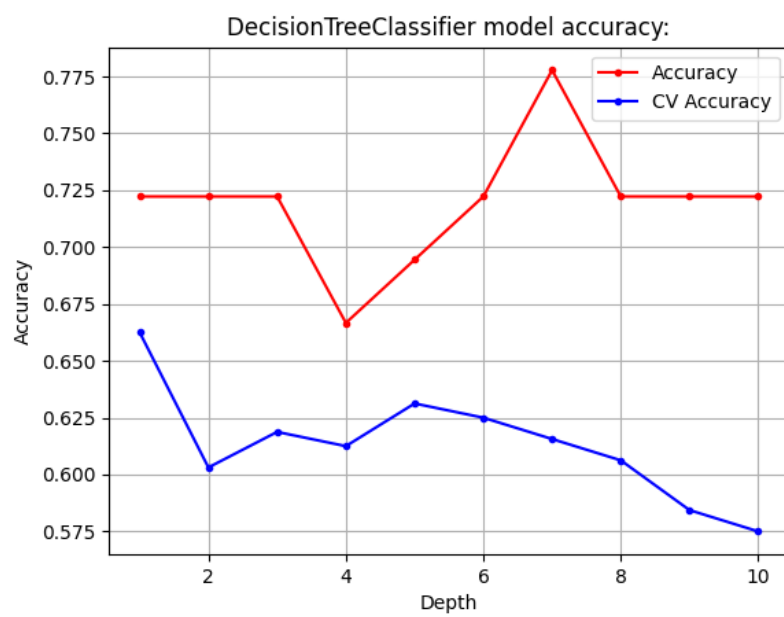
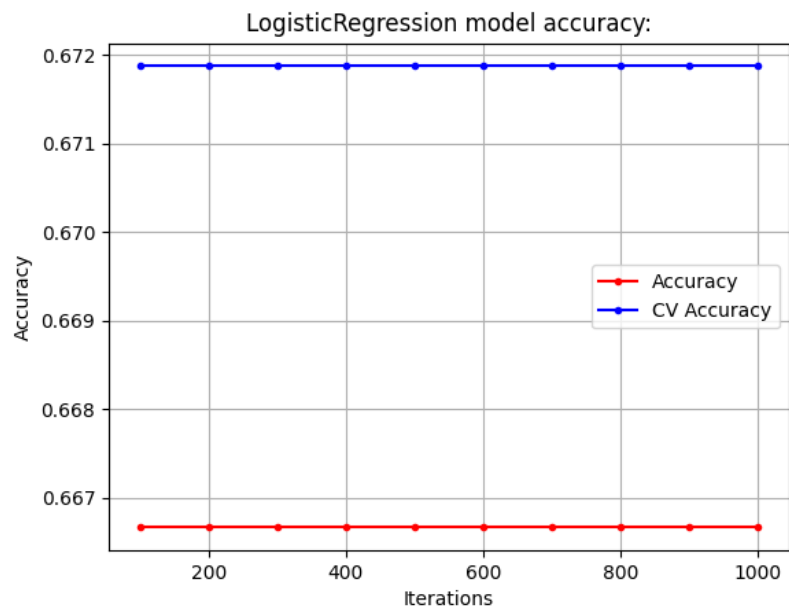
2.3 Wyniki symulacji dla różnych modeli

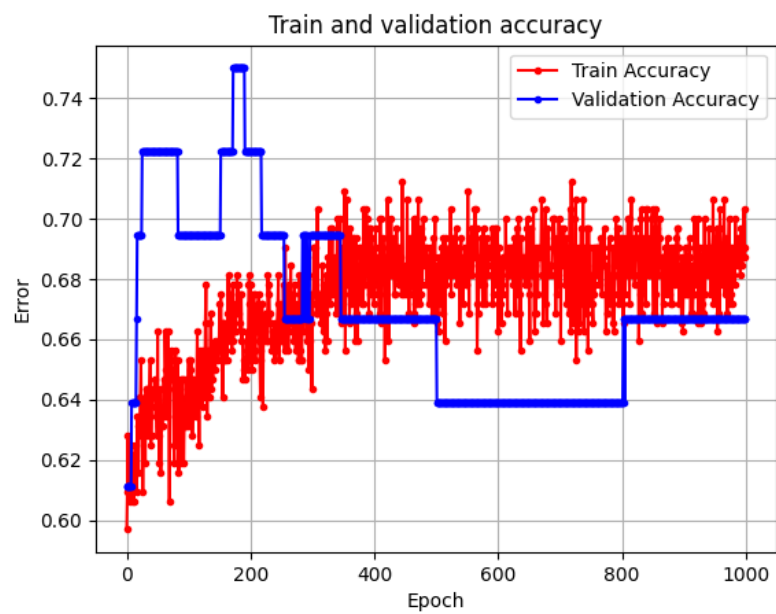
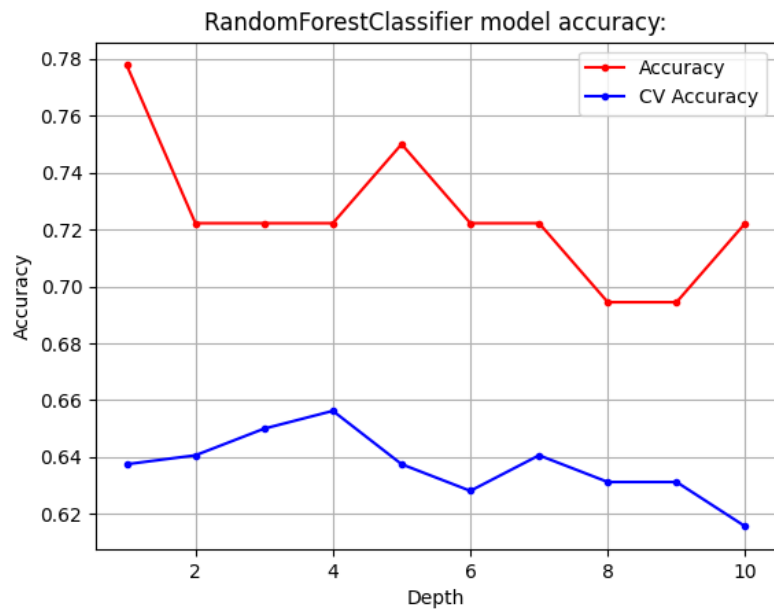
2.3.1 Dane niezrównoważone





2.3.2 Dane zrównoważone





Rozdział 3

Podsumowanie

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Dodatek A

Kod programu

```
import numpy as np

a = np.array([1,2])
b = np.array([1,2])
c = a + b
print(c) # printing results
```