

UNIwersytet Zielonogórski

Wydział Informatyki, Elektrotechniki i Automatyki

Praca dyplomowa

Kierunek: Informatyka

ANALIZA PORÓWNAWCZA BIBLIOTEK  
UCZENIA MASZYNOWEGO JĘZYKA C++ NA  
POTRZEBY ZASTOSOWAŃ W BIOSTATYSTYCE

inż. Kacper Wojciechowski

Promotor:

Prof. dr hab. inż. Dariusz Uciński

Pracę akceptuję:

.....

(data i podpis promotora)

Zielona Góra, czerwiec 2023



## Streszczenie

Niniejsza praca ma na celu analizę i porównanie dostępnych w języku C++ bibliotek uczenia maszynowego, pod kątem ich zastosowania w pracy na danych biostatystycznych. W kolejnych rozdziałach czytelnik zapoznawany jest z:

- Ogólną postać problemów napotykaną w procesie implementacji rozwiązań uczenia maszynowego;
- Charakterystyką wybranego zestawu danych biostatystycznych wykorzystanych do testów omawianych bibliotek;
- Typami oraz uzyskiwanymi wynikami wybranych pod kątem danych eksperymentalnych metod uczenia maszynowego w środowisku prototypowym;
- Bibliotekami Tensorflow, Shark, Caffe i PyTorch wraz z metodami implementacji poszczególnych metod wzorcowych;
- Zbiorczym podsumowaniem funkcjonalności oferowanych przez wyżej wymienione biblioteki.

**Słowa kluczowe:** uczenie maszynowe, C++, biblioteka, sieci neuronowe, głębokie uczenie maszynowe, płytkie uczenie maszynowe.

# Spis treści

<b>1. Wstęp</b>	<b>1</b>
1.1. Wprowadzenie . . . . .	1
1.2. Cel i zakres pracy . . . . .	2
1.3. Struktura pracy . . . . .	2
<b>2. Uczenie maszynowe w ujęciu praktycznym</b>	<b>3</b>
2.1. Problemy współczesnego uczenia maszynowego . . . . .	3
2.2. Język C++ jako narzędzie do rozwiązania problemów uczenia maszy- nowego . . . . .	5
2.3. Cel powstania bibliotek . . . . .	6
<b>3. Inżynieria danych eksperymentalnych i testowe szablony modeli</b>	<b>7</b>
3.1. Omówienie danych eksperymentalnych . . . . .	7
3.2. Charakterystyka i przetwarzanie danych . . . . .	8
3.2.1. Analiza rozkładu danych . . . . .	8
3.2.2. Czyszczenie i normalizacja rozkładu danych . . . . .	9
3.3. Szablony docelowych modeli dla zadanych danych eksperymentalnych	12
3.3.1. Regresja logistyczna . . . . .	12
3.3.2. Głęboka sieć neuronowa . . . . .	14
3.3.3. Maszyna wektorów nośnych . . . . .	16
<b>4. Biblioteka TensorFlow</b>	<b>17</b>
4.1. Wprowadzenie . . . . .	17
4.2. Metody wdrożenia modeli . . . . .	17
4.2.1. API w języku C++ . . . . .	17
4.2.2. Środowisko wykonawcze ONNX . . . . .	18
4.2.3. Zestaw narzędzi TensorFlow Lite . . . . .	19
4.2.4. Za pomocą Protobuf . . . . .	19
4.3. Formaty źródeł danych . . . . .	19
4.4. Metody przetwarzania i eksploracji danych . . . . .	20
4.4.1. Operacje strukturalne . . . . .	20
4.4.2. Eksploracja danych . . . . .	21
4.4.3. Sterowanie przepływem obliczeń . . . . .	22
4.5. Modele uczenia maszynowego . . . . .	22
4.5.1. Regresja liniowa . . . . .	22
4.5.2. Regresja logistyczna . . . . .	23
4.5.3. Regresja Deminga . . . . .	23
4.6. Metody analizy modeli . . . . .	24
4.7. Dostępność dokumentacji i źródeł wiedzy . . . . .	24

<b>5. Biblioteka Shark</b>	<b>25</b>
5.1. Wprowadzenie . . . . .	25
5.2. Metody wdrożenia modeli . . . . .	25
5.3. Formaty źródeł danych . . . . .	25
5.4. Metody przetwarzania i eksploracji danych . . . . .	25
5.5. Modele uczenia maszynowego . . . . .	25
5.6. Metody analizy modeli . . . . .	25
5.7. Dostępność dokumentacji i źródeł wiedzy . . . . .	25
<b>6. Biblioteka Caffe</b>	<b>26</b>
6.1. Wprowadzenie . . . . .	26
6.2. Metody wdrożenia modeli . . . . .	26
6.3. Formaty źródeł danych . . . . .	26
6.4. Metody przetwarzania i eksploracji danych . . . . .	26
6.5. Modele uczenia maszynowego . . . . .	26
6.6. Metody analizy modeli . . . . .	26
6.7. Dostępność dokumentacji i źródeł wiedzy . . . . .	26
<b>7. Biblioteka PyTorch</b>	<b>27</b>
7.1. Wprowadzenie . . . . .	27
7.2. Metody wdrożenia modeli . . . . .	27
7.3. Formaty źródeł danych . . . . .	27
7.4. Metody przetwarzania i eksploracji danych . . . . .	27
7.5. Modele uczenia maszynowego . . . . .	27
7.6. Metody analizy modeli . . . . .	27
7.7. Dostępność dokumentacji i źródeł wiedzy . . . . .	27
<b>8. Podsumowanie</b>	<b>28</b>
8.1. Nakład pracy . . . . .	28
8.2. Funkcjonalności . . . . .	28
8.3. Dostępność źródeł wiedzy i dokumentacja . . . . .	28

# Spis rysunków

2.1. Schemat perceptronu - Simplelearn . . . . .	4
2.2. Multithreading in modern C++ - Modernes C++ . . . . .	5
3.1. Histogram rozkładu zmiennej odpowiedzi . . . . .	8
3.2. Przykłady histogramów zmiennych decyzyjnych . . . . .	9
3.3. Przykład analizy obserwacji odstających dla poszczególnych klas zmiennej odpowiedzi . . . . .	9
3.4. Porównanie rozkładu danych przed i po transformacji logarytmicznej.	10
3.5. Porównanie rozkładów danych przed i po zastosowaniu transformacji pierwiastkiem sześciennym. . . . .	11
3.6. Porównanie uzyskanych rozkładów danych przed i po odwrotnej transformacji Arrheniusa. . . . .	12
3.7. Wykres p-wartości dla całego zestawu zmiennych decyzyjnych. . . . .	13
3.8. Wykres i p-wartości istotnych zmiennych decyzyjnych . . . . .	14
3.9. Krzywa charakterystyczna odbiornika (ROC) dla modelu regresji logistycznej . . . . .	14
3.10. Schemat struktury sieci . . . . .	15
3.11. Krzywa charakterystyczna odbiornika dla zestawu testowego . . . . .	15
3.12. Krzywa charakterystyczna odbiornika dla danych walidacyjnych . . . . .	15
3.13. Krzywa charakterystyczna odbiornika dla danych uczących modelu SVM . . . . .	16
3.14. Krzywa charakterystyczna odbiornika dla danych walidacyjnych modelu SVM . . . . .	16

# Spis tabel

3.1. Lista istotnych regresorów . . . . .	13
3.2. Struktura modelu sieci neuronowej . . . . .	15
3.3. Wartości składowych X modelu dla poszczególnych zmiennych decy- zyjnych . . . . .	16

# Rozdział 1

## Wstęp

### 1.1. Wprowadzenie

We współczesnym stanie techniki coraz częściej można spotkać się z urządzeniami i programami o inteligentnych funkcjach, takich jak predykcja zjawisk na podstawie zestawu danych, rozpoznawanie obrazu, analiza mowy, czy przetwarzanie języka naturalnego. Znajdują one zastosowanie w różnych dziedzinach codziennego życia, m.in. w medycynie. W zależności od potrzeb, techniki uczenia maszynowego można wykorzystać do zastosowań medycznych, jak np. rozpoznawanie komórek rakowych na skanach rezonansem magnetycznym, podejmowanie decyzji na podstawie zbioru objawów obecnych u pacjenta, lub przewidywanie norm związków naturalnie występujących w organizmie ludzkim w zależności od okoliczności i wyników pomiarów.

Jedną z istotnych dziedzin medycyny jest biostatystyka, polegająca na wykorzystaniu analizy statystycznej do wnioskowania na podstawie zbiorów danych, takich jak rezultaty przeprowadzonych badań (np. morfologicznych, poziomu poszczególnych hormonów we krwi, itp.), informacji o nawykach żywieniowych oraz stylu życia pacjenta. Szczególnie istotną formą systemów operujących w tej dziedzinie są systemy eksperckie, wykorzystujące techniki płytkiego i głębokiego uczenia maszynowego w celu wspierania diagnozy stawianej przez wykwalifikowanych lekarzy.

U podstaw wyżej wymienionych zagadnień leży implementacja rozwiązań opartych o teorię uczenia maszynowego, oraz wszelkie związane z tym problemy. W związku z tym na przestrzeni lat powstało wiele gotowych narzędzi, takich jak biblioteki i *frameworki*, mające na celu wsparcie programistów w szybkim i prawidłowym wprowadzaniu rozwiązań sztucznej inteligencji na różne platformy docelowe oraz w różnych językach, począwszy od języka C++, przez Python, po środowiska takie jak Matlab.

Istotnym krokiem w przygotowywaniu oprogramowania wykorzystującego sztuczną inteligencję jest prawidłowy wybór wspomnianych wcześniej narzędzi dokonywany na etapie projektowania, tak, aby oferowały one możliwości adekwatne do wymagań funkcjonalnych. Niniejsza praca dokonuje analizy porównawczej bibliotek uczenia maszynowego dla języka C++ w kontekście zastosowań w dziedzinie biostatystyki, celem umożliwienia czytelnikowi trafnego wyboru odpowiedniego narzędzia do realizacji projektu badawczego.



## 1.2. Cel i zakres pracy

Celem pracy jest przeprowadzenie analizy i przygotowanie zestawienia bibliotek do uczenia maszynowego dla języka C++, obrazując przykłady bazujące na zestawie danych biostatystycznych.

Zakres pracy obejmował:

- Przegląd dostępnych bibliotek języka C++;
- Inżynierię i kształtowanie danych;
- Płytkie i głębokie uczenie nadzorowane;
- Kwestie wydajnościowe w dopasowywaniu i wdrażaniu modeli;
- Badania praktyczne w oparciu o zestaw danych medycznych i biologicznych.

## 1.3. Struktura pracy

Pierwszy rozdział przedstawia ogólnym zagadnieniem dotykany przez pracę, porzucając dziedziny problemu i jej zastosowań, do istoty tematu pracy. Dodatkowo omawiany jest cel i zakres realizacji pracy, oraz jej strukturę.

Kolejny rozdział wprowadza czytelnika do tematu uczenia maszynowego, oraz napotykanym w nim problemów dotyczących złożoności obliczeniowej oraz zużycia zasobów. Stanowią one podstawę do zaproponowania języka C++ jako technologii wspierającej ich rozwiązanie przy pomocy bibliotek.

Tematem rozdziału trzeciego jest przygotowanie elementów testowych do wykorzystania w późniejszej analizie porównawczej. Składa się na nie wybranie i przygotowanie do zestawu danych biostatystycznych do procesu uczenia oraz wybrane wzorcowych rozwiązań. Czytelnik przeprowadzony jest przez normalizację danych i selekcję najlepiej dopasowanych regresorów, oraz zostaje zapoznany z przykładowymi wynikami rozwiązań wzorcowych.

Dalsza część pracy składa się z bloku omówienia i analizy wybranych bibliotek uczenia maszynowego pod kątem określonych kryteriów. Pierwszą z nich, opisaną w rozdziale czwartym, jest biblioteka TensorFlow autorstwa Google. W poszczególnych sekcjach poruszone zostały tematy takie jak możliwe sposoby wdrożenia utworzonych w niej modeli, wymagana forma danych, dostępne funkcjonalności związane z danymi oraz modelami, a także dostępność źródeł informacji na temat sposobu pracy z biblioteką.

# Rozdział 2

## Uczenie maszynowe w ujęciu praktycznym

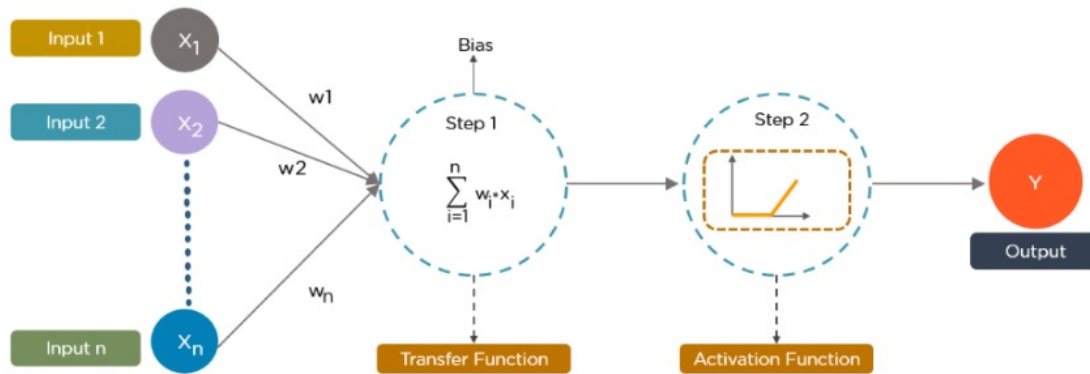
### 2.1. Problemy współczesnego uczenia maszynowego

Na uczenie maszynowe składają się zaawansowane techniki algorytmiczne i złożone struktury danych przeprowadzające obliczenia na zadanym przez użytkownika zestawie danych uczących, testujących, i danych otrzymywanych w trakcie użytkowania wytworzonego modelu.

Do podstawowych form modeli należą modele produkowane w wyniku technik takich jak regresja liniowa i nieliniowa, regresja logistyczna czy liniowa analiza dyskryminacyjna. W ich wyniku tworzone są modele w postaci wielomianów, które później wymagają stosunkowo bardzo małych nakładów mocy obliczeniowej w celu ewaluacji wyników na podstawie zadanego zestawu danych.

Bardziej zaawansowanymi metodami uczenia maszynowego są drzewa decyzyjne, stanowiące strukturę opartą o logikę drzewa. Każdy z poziomów drzewa odpowiada najlepszemu na danym etapie predyktorowi z dostępnych regresorów, powodując rozgałęzienie na poszczególne wartości lub zakresy. Proces obliczania wartości zmiennej wyjściowej odbywa się poprzez przejście przez drzewo od korzenia do jednego z końcowych liści.

Do najbardziej zaawansowanych, aczkolwiek także najbardziej wymagających obliczeniowo i pamięciowo technik uczenia maszynowego należą techniki uczenia głębokiego wykorzystujące sieci neuronowe, jak np. głębokie sieci neuronowe (ang. *Deep Neural Network*, *DNN*) i konwolucyjne sieci neuronowe (ang. *Convolutional Neural Network*, *CNN*). U podstaw tych metod leży struktura sieci neuronowej, składająca się z warstwy wejściowej, jednej lub więcej warstw ukrytych posiadających perceptrony, oraz jednej warstwy wyjściowej. Każdy węzeł z poprzedniej warstwy połączony jest z każdym węzłem w następnej warstwie, lecz perceptrony znajdujące się w tej samej warstwie są wzajemnie niezależne. Każde połączenie posiada przypisaną wagę użytą do przeliczenia wartości wchodzącej do danego perceptronu z danego sąsiada z poprzedniej warstwy. Wewnątrz perceptronu obliczana jest suma iloczynów wyjść z poprzednich perceptronów i wag odpowiadających połączeniom, a następnie dla uzyskanej sumy obliczana jest wartość funkcji aktywacyjnej, która stanowi wartość wyjściową perceptronu. Przykładowa sieć wykorzystująca pojedynczy perceptron w pojedynczej warstwie ukrytej przedstawiona została na rys. 2.1.



**Rysunek 2.1.** Schemat perceptronu - Simplelearn

Bardziej rozbudowane metody wykorzystujące sieci neuronowe, jak np. CNN, wymagają dodatkowych kroków obliczeniowych związanych z wstępnym przetworzeniem danych wejściowych, aby były one przyswajalne dla wykorzystywanej sieci.

Analizując struktury danych wymagane przez poszczególne omówione powyżej rodzaje modeli, wyróżnić można następujące problemy napotykane podczas implementacji metod uczenia maszynowego:

- Wymagania wydajnościowe – są one ściśle powiązane ze złożonością obliczeniową wykorzystanych metod, wydajnością zastosowanego języka i wydajnością zastosowanej platformy sprzętowej. Docelowym efektem jest minimalizacja czasu wymaganego na uczenie modelu (choć tutaj tolerowane są także długie czasy, szczególnie w przypadku dużych zestawów danych uczących) i czasu propagacji modelu (w przypadku czego minimalizacja czasu propagacji stanowi priorytet).
- Wymagania pamięciowe – wynikają one z wykorzystywanych platform sprzętowych i ich ograniczeń pamięciowych. Przykładem powyższego dylematu jest zastosowanie modeli uczenia maszynowego na platformach mobilnych i platformach systemów wbudowanych, gdzie obecne rozmiary pamięci RAM i pamięci masowej (szczególnie w przypadku platform wbudowanych) potrafią być wyraźnie ograniczone w stosunku do systemów komputerowych.

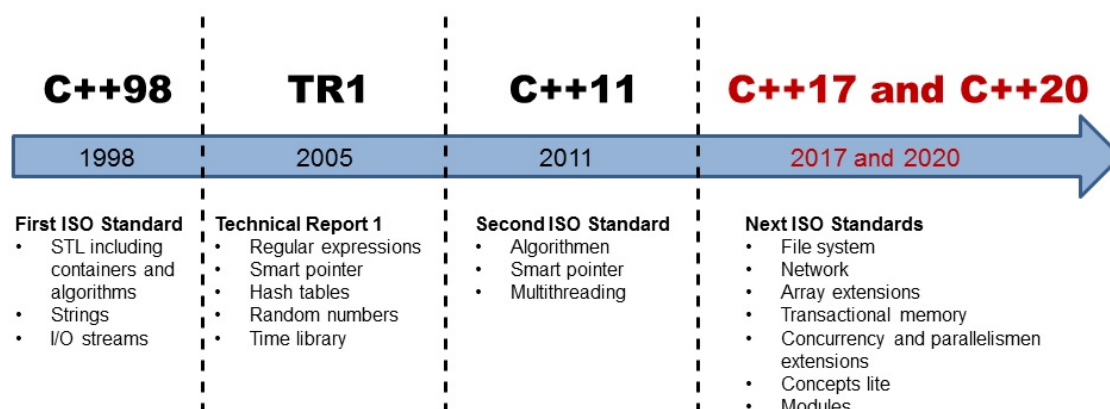
W trakcie rozwoju technologii uczenia maszynowego, postawiono stanowcze kroki w kierunku rozwiązywania powyższych problemów, aby sprostać narastającym wymaganiom związanym z coraz to nowymi i bardziej skomplikowanymi zastosowaniami sztucznej inteligencji. Dokonywano tego poprzez między innymi optymalizację algorytmów, dobór platform sprzętowych o wysokim taktowaniu, możliwym zrównolegleniu operacji, oraz wykorzystaniu wysoko wydajnych języków programowania, w szczególności języków mających możliwość wykorzystania wsparcia ze strony niskopoziomowych operacji.

## 2.2. Język C++ jako narzędzie do rozwiązywania problemów uczenia maszynowego

Dostępne są różne języki i środowiska wspierające uczenie maszynowe, począwszy od języków takich jak Python, C++, Java czy Matlab. Jednak spośród wymienionych kandydatów szczególnie istotnym wyborem jest język C++.

C++ to język imperatywny charakteryzujący się silnym typowaniem, łączący programowanie niskopoziomowe dla konkretnych architektur z wysokopoziomym programowaniem, w związku z czym oferuje programistom dużą kontrolę nad wykorzystaniem pamięci i możliwość optymalizacji w postaci m.in. dostosowywania wykorzystanych typów danych do wymagań funkcjonalnych tworzonej sieci, kontroli lokalizacji zmiennych (programista decyduje czy zmienna lub struktura znajdzie się na stosie czy stercie) oraz optymalizację czasów wywołań funkcji poprzez sugerowanie kompilatorowi utworzenia funkcji inline. W przeciwieństwie do języków skryptowych których kod jest interpretowany w trakcie wykonywania, takich jak Python i język środowiska Matlab, C++ jest językiem kompilowanym. Oznacza to, że program napisany w C++ przetwarzany jest z postaci tekstu do wykonawczego kodu binarnego dostosowanego do wybranej architektury procesora. Usuwa to całkowicie nadmiar złożoności obliczeniowej wykonywanego programu związanej z interpretacją poleceń i tłumaczeniem ich na język procesora danej platformy w trakcie wykonywania programu, gdyż jest to wykonywane tylko raz, na etapie kompilacji, dodatkowo pozwalając na zastosowanie przez kompilator mechanizmów optymalizacji dostępnych dla wybranej platformy.

Część mechanizmów z języka C++, wywodzących się jeszcze z języka C, pozwala na wykorzystanie wstawek kodu źródłowego w języku Assembler dla wybranego procesora, co zwiększa wydajność programu kosztem przenośności kodu. Dodatkowo niektóre platformy oferują API modułów akceleracji sprzętowej (jak np. system Android udostępniający *Neural Networks API*, *NNAPI* dla sieci neuronowych), co oferuje dodatkowe przyspieszenie czasu działania programu.



Rysunek 2.2. Multithreading in modern C++ - Modernes C++

Jedną z popularnych technik mających na celu znaczne zwiększenie wydajności modeli sztucznej inteligencji jest zrównoleglenie przetwarzania. Dostępność mechanizmów wielowątkowych dla procesorów (wprowadzonych w standardzie C++11 i

dalej rozwijanych, jak przedstawiono na rys. 2.2), oraz kompatybilność języka C++ z językiem CUDA pozwala wykonywać wiele obliczeń równolegle poprzez wykorzystanie wielu rdzeni lub oddelegowaniu części przetwarzania do karty (lub wielu kart) graficznej (gdzie ilość procesorów GPU znacząco przewyższa ilość rdzeni CPU). Dodatkowym atutem wykorzystania języka C++ przy tworzeniu modelu sztucznej inteligencji jest łatwa integracja z programami dedykowanymi do wysokiej wydajności, napisanymi w tym języku.

Wymienione wyżej mechanizmy i cechy charakterystyczne języka umożliwiają programistom znaczną optymalizację przygotowywanych rozwiązań sztucznej inteligencji, co przekłada się na bardziej efektywne zużycie pamięci, zabezpieczenie przed przeładowaniem stosu procesora, oraz krótsze czasy propagacji utworzonych modeli.

## 2.3. Cel powstania bibliotek

Implementacja mechanizmów pozwalających na tworzenie rozwiązań sztucznej inteligencji, z racji na swoją złożoność, wymagania dotyczące kompetencji twórców oraz konieczność optymalizacji jest czasochłonna i kosztowna. Tu z pomocą przychodzą biblioteki utworzone przez korporacje oraz społeczność programistów *open source*. Stanowią one gotowe zbiory mechanizmów (najczęściej pisane w sposób obiektowy, a więc ubrane w klasy posiadające określone zestawy metod), które są na bieżąco optymalizowane przez grupy programistów wykorzystujące je w prywatnych projektach lub pracy zawodowej. Oferują one możliwość wykorzystania gotowych modeli utworzonych w innych technologiach, a czasem także bezpośrednie przygotowanie modelu na podstawie odpowiednio sformatowanego i odpowiednio przystosowanego zestawu danych.

Użycie gotowych bibliotek nie tylko oszczędza kosztu i przyspiesza tworzenie pożądanego rozwiązania sztucznej inteligencji, lecz także zapewnia większą niezawodność, gdyż elementy zawarte w bibliotece są implementowane, dokładnie testowane i poprawiane przez programistów o wysokich kompetencjach, jak m.in. w przypadku biblioteki TensorFlow posiadającej wsparcie od pracowników Google.

Większość bibliotek przeznaczonych do uczenia maszynowego, nawet wykorzystywanych w językach takich jak Python, napisana jest w języku C++, oferując API dostępne dla określonych języków docelowych. Niestety nie wszystkie biblioteki napisane w ten sposób oferują dostęp do całego API w języku C++ dla wykorzystujących je programów zewnętrznych, lub bywa on utrudniony i skomplikowany, co sprawia że w powszechnej praktyce część bibliotek dedykowanych dla języka C++ operuje na modelach przygotowanych w ramach innej, lub czasem nawet tej samej biblioteki, napisanych w innym języku. Częstym przypadkiem jest tutaj wykorzystanie właśnie języka Python do utworzenia grafu modelu lub modelu w formacie ONNX (ang. *Open Neural Network Exchange*).

W ramach analizy porównawczej w niniejszej pracy, porównywane będą biblioteki oferujące zarówno tworzenie modeli w ramach języka C++, jak i wymagające wykorzystania modeli z innego źródła.

# Rozdział 3

## Inżynieria danych eksperymentalnych i testowe szablony modeli

### 3.1. Omówienie danych eksperymentalnych

W celu zestawienia funkcjonalnego bibliotek uczenia maszynowego w języku C++ i przedstawienia przykładów konieczne było wybranie danych eksperymentalnych możliwych do wykorzystania jako porównawczy punkt odniesienia. Jako w/w dane wybrano bazę dotyczącą diagnostyki raka piersi „*Wisconsin Diagnostic Breast Cancer*” z listopada 1995 roku, w której zamieszczono wyniki obrazowania określone w sposób liczbowy. Autorami zestawu są Dr. Wiliam H. Wolberg, W. Nick Street oraz Olvi L. Mangasarian z Uniwersytetu Wisconsin [1]. Baza ta jest dostępna do pobrania z repozytorium Uniwersytetu Californii [2]. Dane mają następującą strukturę:

- 1) ID - numer identyfikacyjny pacjentki;
- 2) Diagnosis [*Malignant* - *M* / *Benign* - *B*] - charakter nowotworu, **zmienna odpowiedzi**;
- 3) Dane klasyfikujące:
  - a) *Radius* - średnica guza;
  - b) *Texture* - tekstura guza;
  - c) *Perimeter* - obwód guza;
  - d) *Area* - pole guza;
  - e) *Smoothness* - gładkość, miara lokalnych różnic w promieniu guza;
  - f) *Compactness* - zwartość, wykorzystywana do oceny stadium guza;
  - g) *Concavity* - stopień wklęsłości miejsc guza;
  - h) *Concave points* - punkty wklęsłości guza;
  - i) *Symmetry* - symetria guza, pomagająca w ocenie charakteru przyrostu guza.

- j) *Fractal dimension* („coastline approximation” - 1) - wymiar fraktalny pozwalający na ilościowy opis złożoności komórek nerwowych, umożliwiającą stwierdzenie nowotworzenia się zbioru komórek.

Dla każdej ze zmiennych odpowiedzi została zebrana średnia wartość, odchylenie standardowe oraz średnia trzech największych pomiarów, gdzie każdy zestaw ustawiony jest sekwencyjnie (np. kolumna 3 - średni promień, kolumna 12 - odchylenie standardowe promienia, kolumna 22 - średnia trzech największych pomiarów promienia). Każda ze zmiennych ma charakter ciągły. Zredukowany zestaw danych, zawierający jedynie zmienne decyzyjne informujące o średnich wartościach znaleźć można jako dodatek do książki „*Biostatistics Using JMP: A Practical Guide*” autorstwa Trevora Bihla [3].

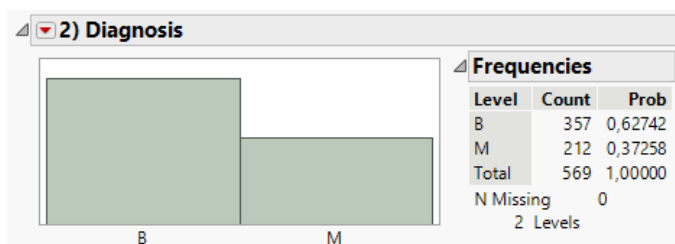
## 3.2. Charakterystyka i przetwarzanie danych

W celu przeprowadzenia procesu uczenia maszynowego, jednym z najistotniejszych kroków jakie należy podjąć jest wstępne zaznajomienie się z zestawem danych i jego analiza pod kątem rozkładu poszczególnych zmiennych oraz prawdopodobieństw. W tym celu wykorzystane zostało oprogramowanie JMP.

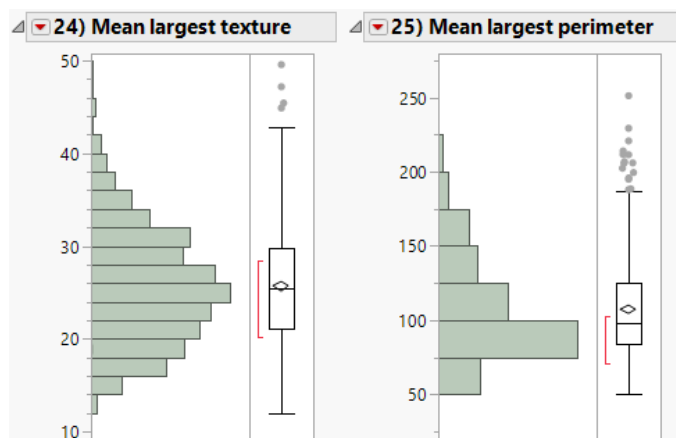
### 3.2.1. Analiza rozkładu danych

Proces analizy rozkładu rozpoczęty został od przyjrzenia się zmiennej odpowiedzi (*Diagnosis*). Rysunek 3.1 przedstawia uzyskany histogram, wraz z tabelą określającą ilość obserwacji danej klasy i współczynnik prawdopodobieństwa przynależności odpowiedzi do danej klasy. Zauważyć można, że dla użytego zestawu danych ilość zarejestrowano 357 obserwacji łagodnego raka piersi, a jego prawdopodobieństwo przynależności do klasy *Benign* wynosi  $\approx 62,7\%$ , natomiast do klasy *Malignant* przynależało 212 obserwacji z prawdopodobieństwem  $\approx 37,3\%$ .

Podczas analizy histogramów zmiennych decyzyjnych, stwierdzono że znaczna ilość ma charakter prawostronnie skośny oraz występują dla nich obserwacje odstające, o czym informuje znajdujący się po prawej stronie histogramu wykres okienkowy (ang. *box graph*), co przedstawiono na rysunku 3.2. Wyjątkiem okazała się zmienna *Mean Largest Concave Points*, która mimo lekkiej skośności, okazała się nie posiadać obserwacji odstających. Na podstawie tych informacji stwierdzono, że aby przygotować dane w odpowiedni sposób do procesu uczenia należy przeprowadzić ich czyszczenie oraz normalizację rozkładu.



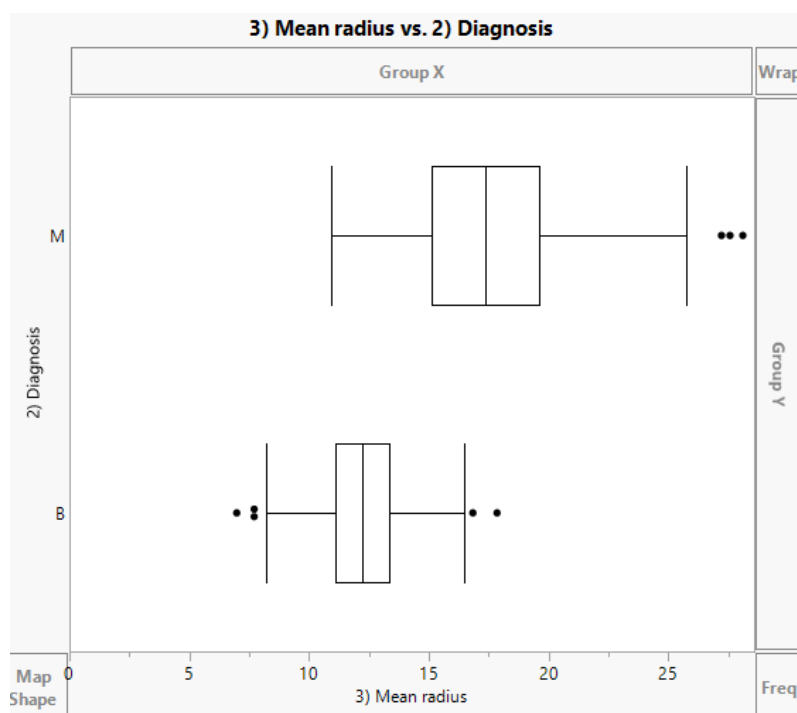
Rysunek 3.1. Histogram rozkładu zmiennej odpowiedzi



Rysunek 3.2. Przykłady histogramów zmiennych decyzyjnych

### 3.2.2. Czyszczenie i normalizacja rozkładu danych

Na pełny zestaw danych składa się 569 obserwacji. Podczas wstępnej analizy stwierdzono istnienie 13 brakujących wartości dla regresora *Std err concave points*, dla których przyjęto wartość średnią z całej kolumny. Głównym problemem okazały się obserwacje odstające oraz skośności rozkładu. Do analizy obserwacji odstających wykorzystano wykresy okienkowe, gdzie oś Y reprezentowała zmienną odpowiedzi, natomiast oś X czyszczoną zmienną decyzyjną. Przykładowy wykres został przedstawiony na rysunku 3.3. Ze względu na bardzo małą ilość obserwacji zdecydowano się rozpocząć proces przystosowywania danych do uczenia poprzez normalizację ich rozkładu, aby zminimalizować lub wyeliminować konieczność usunięcia danych odstających.

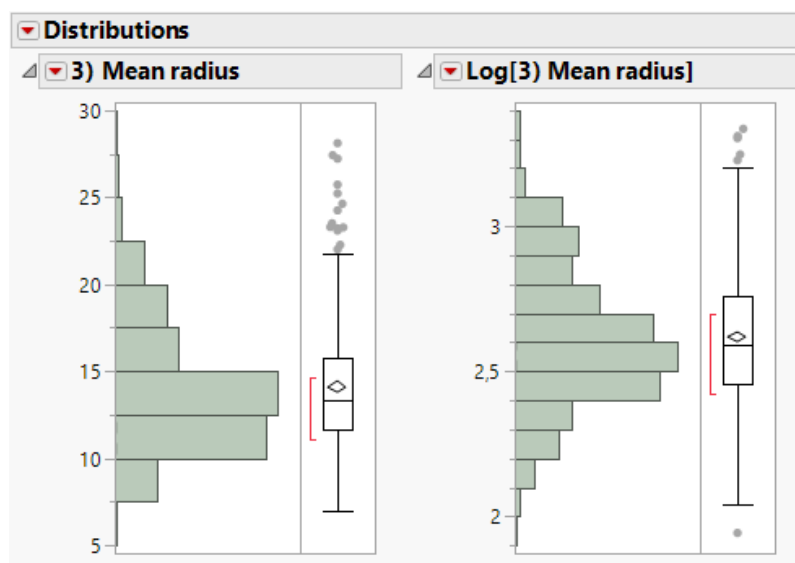


Rysunek 3.3. Przykład analizy obserwacji odstających dla poszczególnych klas zmiennej odpowiedzi



W pierwszym podejściu zdecydowano się na zastosowanie transformacji logarytmicznej dla wszystkich zmiennych decyzyjnych i porównanie charakterystyk uzyskanych rozkładów z oryginalnymi. Zmienna *Mean largest concave points* okazała się posiadać rozkład bardzo zbliżony do standardowego, w związku z czym wyłączono ją z dalszej analizy normalizacji. Przykładowe wyniki przedstawiono na rysunku 3.4. Transformacja ta okazała się skutecznym rozwiązaniem jedynie dla następujących zmiennych:

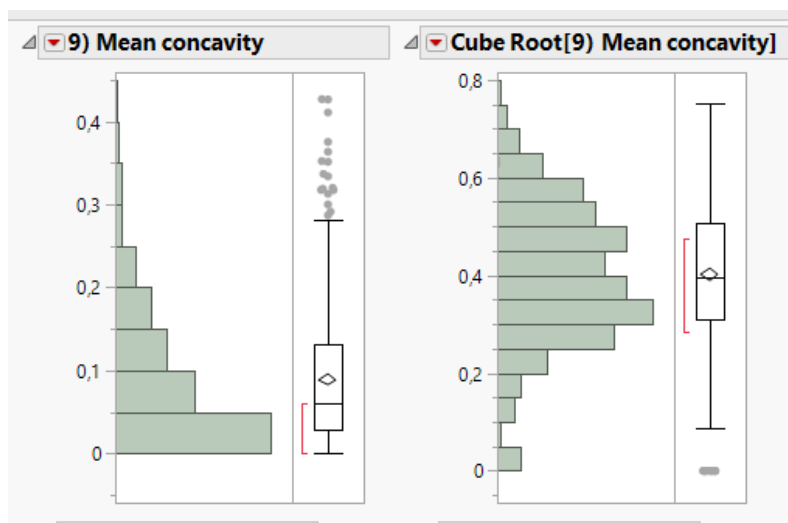
1. *Mean radius*;
2. *Mean texture*;
3. *Mean perimeter*,
4. *Mean area*;
5. *Mean smoothness*;
6. *Mean symmetry*;
7. *Std err texture*;
8. *Std err smoothness*;
9. *Std err compactness*;
10. *Std err concave points*;
11. *Mean largest texture*;
12. *Mean largest smoothness*;
13. *Mean largest compactness*.



**Rysunek 3.4.** Porównanie rozkładu danych przed i po transformacji logarytmicznej.

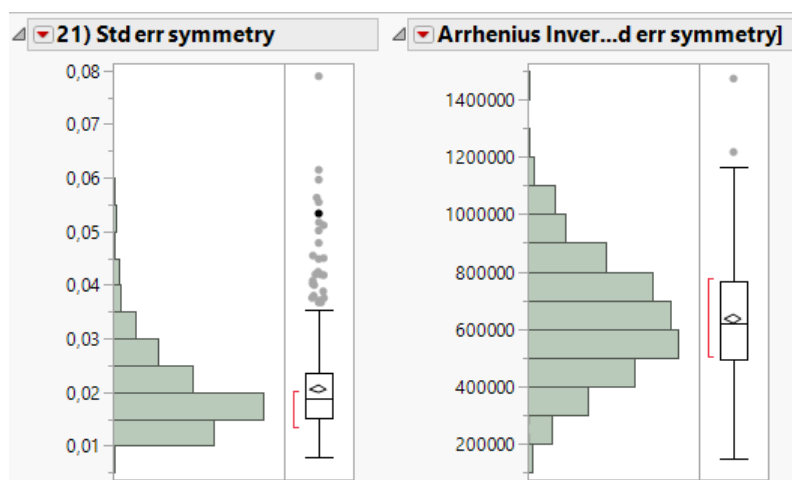
W drugim kroku podjęto próbę wykorzystania transformacji pierwiastkiem sześciennym dla pozostałych zmiennych decyzyjnych, ze względu na jej skuteczność dla danych o rozkładzie prawoskośnym. Rysunek 3.5. przedstawia porównanie rozkładu zmiennej *Mean concavity* przed i po transformacji pierwiastkiem sześciennym. Pomysłynie znormalizowano rozkład następujących zmiennych:

1. *Mean compactness*;
2. *Mean concavity*;
3. *Mean concave points*;
4. *Std err concavity*;
5. *Mean largest radius*;
6. *Mean largest perimeter*;
7. *Mean largest concavity*;
8. *Mean largest symmetry*.



**Rysunek 3.5.** Porównanie rozkładów danych przed i po zastosowaniu transformacji pierwiastkiem sześciennym.

Ostatecznym krokiem okazało się zastosowanie odwrotnej transformacji Arrheniusa. Niestety część z uzyskanych zmodyfikowanych zmiennych decyzyjnych zachowała częściowy skośny rozkład, jednak inne przetestowane transformacje, jak m.in. pierwiastek kwadratowy, potęga kwadratowa, logarytm  $x+1$ , logarytm dziesiętny, funkcja potęgowa, funkcja wykładnicza, przyniosły rezultaty porównywalne lub gorsze od uzyskanego w wyniku w/w odwrotnej transformacji Arrheniusa. Rysunek 3.6 przedstawia porównanie uzyskanych rozkładów.



**Rysunek 3.6.** Porównanie uzyskanych rozkładów danych przed i po odwrotnej transformacji Arrheniusa.

Ze względu na bardzo małą ilość obserwacji, zdecydowano się na zachowanie wszystkich obserwacji odstających, aby zapobiec utracie informacji i zmianie uzyskanych w procesie normalizacji rozkładów.

### 3.3. Szablony docelowych modeli dla zadanych danych eksperymentalnych

Ze względu na dychotomiczny charakter zmiennej odpowiedzi, wybrany został przedstawiony poniżej zestaw metod dla których wykonano i przedstawiono testy praktyczne. Szablony struktury rozwiązań, takie jak np. wybór zmiennych uczestniczących w procesie uczenia, lub struktura sieci neuronowej zostały ustalone w sposób empiryczny z wykorzystaniem programu do uczenia maszynowego JMP.

#### 3.3.1. Regresja logistyczna

Badanie zależności w modelu regresji logistycznej odbyło się z wykorzystaniem wykresu wpływu zmiennej decyzyjnej na zmienną odpowiedzi opartego o p-wartość. Jako próg pozwalający na odrzucenie hipotezy zerowej (hipotezy o braku wpływu zmiennej na odpowiedź) przyjęto 0.05 jednostek. Rysunek 3.7 przedstawia w/w wykres wraz z p-wartościami dla poszczególnych zmiennych. Zauważyć można, że dla części zmiennych nie została wyznaczona p-wartość – oznacza to, że część zmiennych jest ze sobą skorelowanych.

Pierwszym krokiem w wybraniu istotnych zmiennych było usunięcie zmiennych skorelowanych, drugim natomiast stopniowe usuwanie zmiennych o p-wartości powyżej określonego progu. Rysunek 3.8 przedstawia listę wraz z wykresem kolumnowym istotnych regresorów. Ich lista, wraz z odpowiadającymi im p-wartościami została umieszczona w tabeli 3.1.

Nazwa zmiennej	p-wartość
<i>Log mean largest texture</i>	0,00000
<i>Log mean largest compactness</i>	0,00000
<i>Cube root mean largest symmetry</i>	0,00001
<i>Arrhenius inverse std err symmetry</i>	0,00005
<i>Arrhenius inverse std err radius</i>	0,00018
<i>Cube root mean concave points</i>	0,00056
<i>Cube root mean largest concavity</i>	0,00069
<i>Log std err texture</i>	0,00252
<i>Cube root mean largest perimeter</i>	0,00526
<i>Log mean smoothness</i>	0,04867
<i>Log mean radius</i>	0,04884

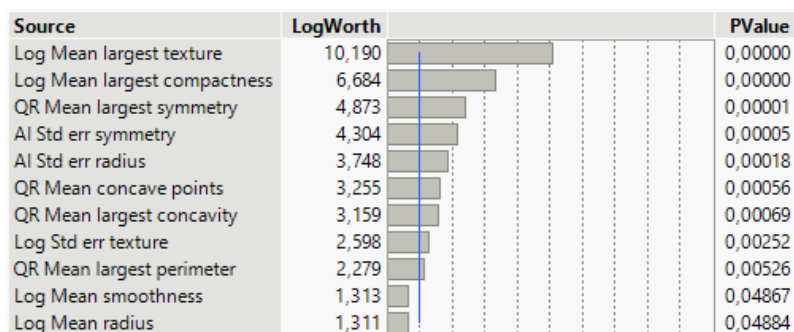
Tabela 3.1. Lista istotnych regresorów

Source	LogWorth		PValue
AI Mean largest area	951,928		0,00000
QR Mean concavity	610,420		0,00000
AI Std err area	476,593		0,00000
Log Std err smoothness	358,978		0,00000
AI Mean fractal dimation	218,578		0,00000
AI Mean largest fractal dimation	.		0,00000
QR Mean largest symmetry	.		.
Mean largest concave points	.		.
QR Mean largest concavity	.		.
Log Mean largest compactness	.		.
Log Mean largest smoothness	.		.
QR Mean largest perimeter	.		.
Log Mean largest texture	.		.
QR Mean largest radius	.		.
AI Std err fractal dimation	.		0,00000
AI Std err symmetry	.		0,00000
Log Std err concave points	.		0,00000
QR Std err concavity	.		.
Log Std err compactness	.		0,00000
AI Std Err perimeter	.		0,00000
Log Std err texture	.		0,00000
AI Std err radius	.		0,00000
Log Mean symmetry	.		0,00000
QR Mean concave points	.		.
QR Mean compactness	.		.
Log Mean smoothness	.		.
Log Mean area	.		.
Log Mean perimeter	.		.
Log Mean texture	.		0,00000
Log Mean radius	.		0,00000

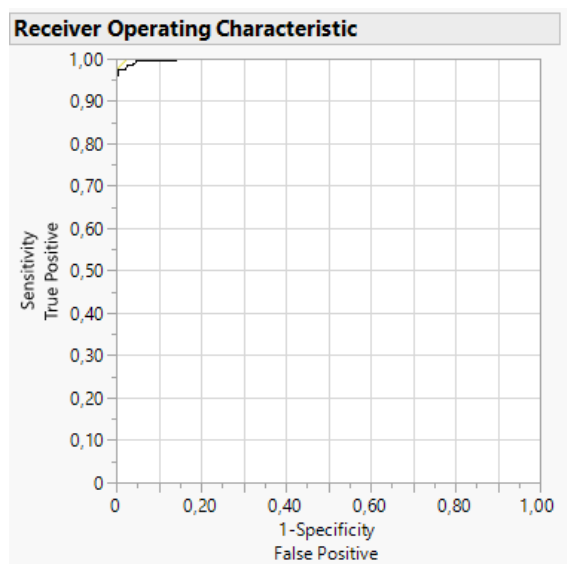
Rysunek 3.7. Wykres p-wartości dla całego zestawu zmiennych decyzyjnych.

Dla wybranego zestawu zmiennych model osiągnął dokładność na poziomie  $R^2 = 0.9401$ . Zgodnie z macierzą pomyłek, 207 obserwacji typu *Malignant* oraz 335 obserwacji *Benign* zostało zaklasyfikowanych poprawnie. Oznacza to, że model uży-

skalał tylko 2 wyniki typu *false-positive* (prawdopodobieństwo 0,6%) i 5 wyników typu *false-negative* (prawdopodobieństwo 2,4%) dla danych treningowych. Ze względu na mały zestaw obserwacji, ryzyko przeuczenia jest znikome, w związku z czym nie wytypowano zestawu danych walidacyjnych. Rysunek 3.9 przedstawia krzywą charakterystyczną odbiornika dla modelu.



Rysunek 3.8. Wykres i p-wartości istotnych zmiennych decyzyjnych



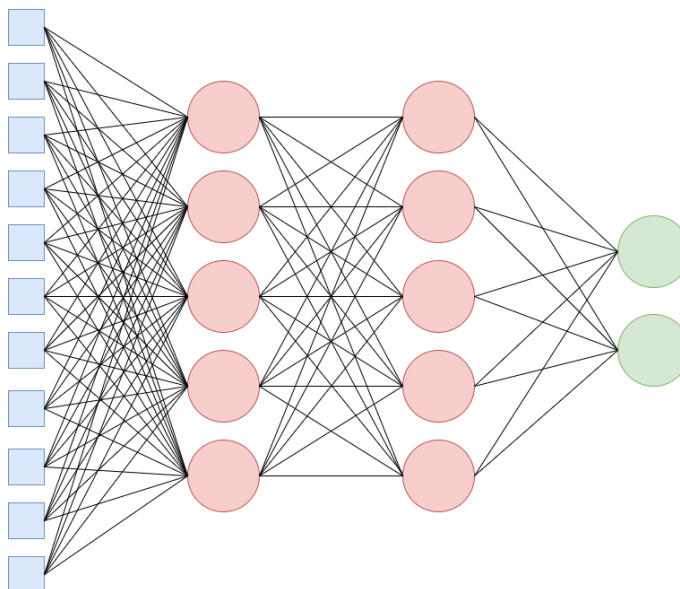
Rysunek 3.9. Krzywa charakterystyczna odbiornika (ROC) dla modelu regresji logistycznej

### 3.3.2. Głęboka sieć neuronowa

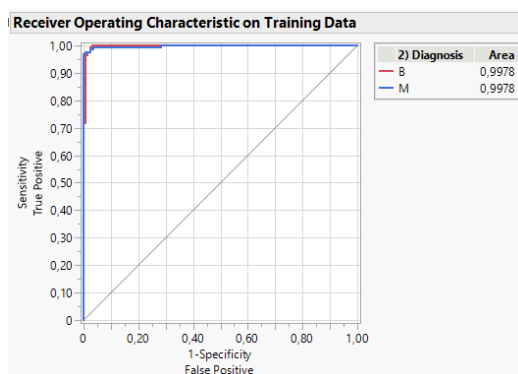
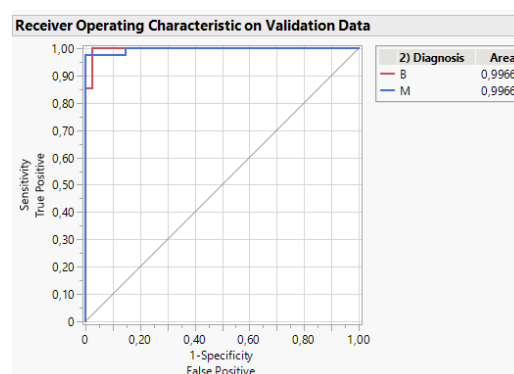
Do przygotowania sieci neuronowej wykorzystano zestaw zmiennych zawartych w tabeli 3.1. Dane zostały losowo podzielone na dane uczące i walidacyjne w proporcji 80% do 20%. W wyniku prób i błędów, optymalny model uzyskano przy strukturze przedstawionej w tabeli 3.2. Graficzny schemat struktury został także przedstawiony na rysunku 3.10.

Środowisko JMP nie udostępnia informacji o funkcji aktywacji warstwy wyjściowej, w związku z czym w tabeli 3.2 została ona pominięta. Dla ziarna o wartości 1234 uzyskano model którego statystyka  $R^2$  dla danych treningowych wyniosła 0.966268, natomiast dla danych testowych 0.9924547. Trafność dla losowo wybranego zestawu

Typ warstwy	ilość neuronów	aktywacja
ukryta	5	tangens hiperboliczny
ukryta	5	tangens hiperboliczny
wyjściowa	2	—

**Tabela 3.2.** Struktura modelu sieci neuronowej**Rysunek 3.10.** Schemat struktury sieci

testowego wyniosła 100%, natomiast dla danych uczących napotkano 5 przypadków *false-negative* (prawdopodobieństwo 3%) oraz 1 przypadek *false-positive* (prawdopodobieństwo 0,4%). Rysunki 3.11 oraz 3.12 przedstawiają krzywe charakterystyczne odbiornika dla zestawu testowego i walidacyjnego.

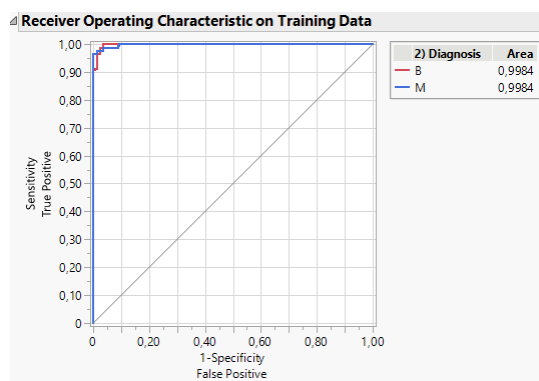
**Rysunek 3.11.** Krzywa charakterystyczna odbiornika dla zestawu testowego**Rysunek 3.12.** Krzywa charakterystyczna odbiornika dla danych walidacyjnych

### 3.3.3. Maszyna wektorów nośnych

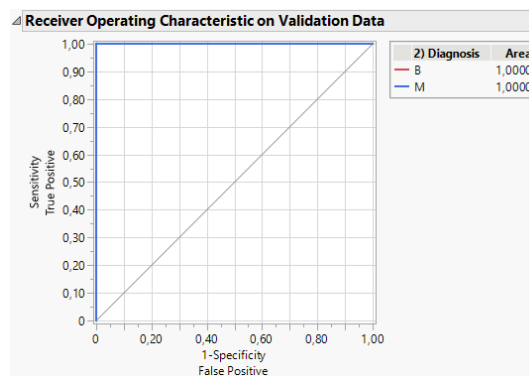
Do predykcji diagnozy wykorzystano ten sam zestaw regresorów, zawartych w tabeli 3.1. Ponownie w celu walidacji użyto metody wybrania losowego zestawu walidacyjnego spośród dostarczonych danych, w proporcji 80% obserwacji uczących i 20% testowych, z użyciem wartości 1234 dla ziarna generatora liczb pseudolosowych. Jako funkcję jądra maszyny wektorów nośnych (ang. Support Vector Machine, SVM) wybrano *Radial Basis Function*, która jest domyślnym wyborem dla SVM w środowisku JMP.

Zmienna decyzyjna	wartość X
<i>Log mean largest texture</i>	3,217
<i>Log mean largest compactness</i>	-1,5504
<i>Cube root mean largest symmetry</i>	0,65891
<i>Arrhenius inverse std err symmetry</i>	635100
<i>Arrhenius inverse std err radius</i>	38170
<i>Cube root mean concave points</i>	0,33665
<i>Cube root mean largest concavity</i>	0,5951
<i>Log std err texture</i>	0,1049
<i>Cube root mean largest perimeter</i>	4,7045
<i>Log mean smoothness</i>	-2,3502
<i>Log mean radius</i>	2,6191

**Tabela 3.3.** Wartości składowych X modelu dla poszczególnych zmiennych decyzyjnych



**Rysunek 3.13.** Krzywa charakterystyczna odbiornika dla danych uczących modelu SVM



**Rysunek 3.14.** Krzywa charakterystyczna odbiornika dla danych walidacyjnych modelu SVM

Utworzony w ten sposób model posiada generalizowaną statystykę  $R^2$  na poziomie 0.97161 dla zestawu walidacyjnego, i uzyskał wskaźnik błędnej klasyfikacji wynoszący 0% dla danych testowych, oraz 1,3% dla danych uczących. Tabela 3.3 przedstawia wartości X dla poszczególnych regresorów. Rysunki 3.13 oraz 3.14 przedstawiają krzywe charakterystyczne odbiornika dla uzyskanego modelu.

# Rozdział 4

## Biblioteka TensorFlow

### 4.1. Wprowadzenie

Tensorflow to darmowa biblioteka do uczenia maszynowego typu *open-source* stworzona przez Google. Oferuje ona możliwość przeprowadzenia płytkiego oraz głębokiego uczenia maszynowego, z wykorzystaniem własnych lub gotowych zestawów danych. Pozwala ona także na użycie pretrenowanych modeli jako rozwiązanie docelowe lub część składowa własnego modelu. Stanowi ona jedną z najpopularniejszych bibliotek ML, w szczególności wśród osób dopiero rozpoczynające swoje doświadczenia z tą gałęzią informatyki.

Jak podano na stronie domowej biblioteki Tensorflow[4], deklaruje ona bycie odpowiednią zarówno do przeprowadzania badań, jak i do wykorzystania w zastosowaniach przemysłowych. W momencie pisania niniejszej pracy, najnowszą dostępną wersją jest Tensorflow 2.11.

### 4.2. Metody wdrożenia modeli

Biblioteka umożliwia następujące metody wdrożenia modeli dla języka C++ [5][6][7]:

- API TensorFlow w języku C++
- Środowisko wykonawcze ONNX
- Zestaw narzędzi TensorFlow Lite
- Za pomocą Protobuf

#### 4.2.1. API w języku C++

Z racji bycia napisanym w języku C++, biblioteka TensorFlow oferuje wsparcie do tworzenia, trenowania, i wdrażania modeli uczenia maszynowego w tym języku, poprzez udostępnienie użytkownikom interfejsu programowania aplikacji (ang. *Application Programming Interface, API*). Udostępnia ono operacje pozwalające manipulacje tensorami, wykonywanie operacji matematycznych na tensorach i ich elementach, budowę modelu w postaci grafu obliczeniowego, oraz jego uruchomienia



wykorzystując różne metody obliczeniowe, np. Adagrad oraz metodę spadku gradientowego [8]. Zastosowanie tego API pozwala na utworzenie i pracę z modelem bezpośrednio w aplikacji pisanej w języku C++.

**Listing 4.1.** Przykład utworzenia grafu obliczeniowego w TensorFlow realizującego konkatencję ciągów znakowych [9]

---

```

1  #include <tensorflow/cc/client/client_session.h>
2  #include <tensorflow/cc/ops/standard_ops.h>
3  #include <tensorflow/core/framework/tensor.h>
4
5  int main(int argc, char **argv) {
6      using namespace tensorflow;
7      using namespace tensorflow::ops;
8
9      // create a root scope
10     auto scope = Scope::NewRootScope();
11
12     // define various constants/inputs on which we
13     // will perform an operation
14     auto hello = Const(scope, std::string("hello"));
15     auto space = Const(scope, std::string(" "));
16     auto world = Const(scope, std::string("world!"));
17
18     // StringJoin operation
19     auto joinOp = StringJoin(scope, {hello, space, world});
20
21     // create a session that takes our
22     // scope as the root scope
23     ClientSession session(scope);
24
25     // Run
26     std::vector<Tensor> outputs;
27     TF_CHECK_OK(session.Run({joinOp}, &outputs));
28
29     // See our output using DebugString that tells
30     // more information about the tensor
31     std::cout << "DebugString->" << outputs[0].DebugString() << std::endl;
32
33     // we can also get the underlying data by calling flat
34     std::cout << "UnderlyingScalarValue->" << outputs[0].flat<std::string>()
35               << std::endl;
36
37     return 0;
38 }
```

---

### 4.2.2. Środowisko wykonawcze ONNX

*Open Neural Networks Exchange*, (*ONNX*) to środowisko mające na celu standaryzację reprezentacji modeli uczenia maszynowego, autorstwa firmy Microsoft [6]. Pozwala ono na uniezależnienie modelu na etapie wdrożenia od biblioteki i języka w którym został stworzony, pozwalając na ich użycie przy pomocy pojedynczego mechanizmu. Wykorzystanie tej metody wdrażania modeli polega na przygotowaniu ich za pomocą API języka okraz eksport do struktury akceptowanej przez środowisko wykonawcze ONNX.

### 4.2.3. Zestaw narzędzi TensorFlow Lite

Stanowi zbiór narzędzi (ang. *toolbox*) umożliwiających uruchomienie przygotowanych modeli na platformach mobilnych oraz urządzeniach brzegowych (takich jak moduły IoT, lub ogólnie rozumiane systemy mikroprocesorowe). Pozwala on dodatkowo m.in. na wykorzystanie akceleracji sprzętowej oraz optymalizacji modelu. Przygotowywanie modeli do uruchomienia za pomocą tej metody na wybranej platformie składa się m.in. z przygotowania modelu z użyciem biblioteki TensorFlow Lite Model Maker, wykorzystania jednego z gotowych, pretrenowanych modeli, lub konwersji modelu TensorFlow z formatu protobuf [7].

### 4.2.4. Za pomocą Protobuf

Utworzony w bibliotece TensorFlow model jest zapisywany do formatu wykorzystywanego przez mechanizm Protocol Buffers autorstwa Google, zwanego potocznie *Protobuf*. Pozwala on na przechowywanie grafu obliczeniowego modelu w notacji niezależnej od języków programowania, z której następnie narzędzie Protobuf umożliwia wygenerowanie odpowiednich klas i struktur danych dla wybranego języka, jak np. C++ czy Python [10]. Pliki Protobuf mogą mieć formę tekstową, która ułatwia edycję i analizę człowiekowi, lub formę binarną, pozwalającą na znaczne zaoszczędzenie miejsca w przypadku przechowywania danych numerycznych, jak np. wagi połączeń sieci.

**Listing 4.2.** Konstruktor klasy *Model* odczytujący strukturę modelu z pliku Protobuf [11]

```
1  #include <string>
2
3  #include "third_party/tensorflow/core/framework/graph.proto.h"
4  #include "third_party/tensorflow/core/framework/tensor.h"
5  #include "third_party/tensorflow/core/lib/io/path.h"
6  #include "third_party/tensorflow/core/platform/env.h"
7  #include "third_party/tensorflow/core/platform/init_main.h"
8  #include "third_party/tensorflow/core/platform/logging.h"
9  #include "third_party/tensorflow/core/platform/types.h"
10 #include "third_party/tensorflow/core/public/session.h"
11
12 Model(const std::string& graph_def_filename) {
13     tensorflow::GraphDef graph_def;
14     TF_CHECK_OK(tensorflow::ReadBinaryProto(tensorflow::Env::Default(),
15                                             graph_def_filename, &graph_def));
16     session_.reset(tensorflow::NewSession(tensorflow::SessionOptions()));
17     TF_CHECK_OK(session_>Create(graph_def));
18 }
```

## 4.3. Formaty źródeł danych

W języku C++ biblioteka TensorFlow wymaga podania danych w postaci obiektów klasy szablonowej *std::vector* z biblioteki STL języka [11], lub za pomocą własnych klas, np. *FeedDict*, jednak nie określa w jaki sposób muszą te dane być przechowywane. W związku z tym, możliwe jest wykorzystanie dowolnego formatu przechowywania danych do uczenia, pod warunkiem zapewnienia mechanizmu ich parsowania.

Typ przechowywany przez obiekty `std::vector` zależy jest od ilości regresorów, i może przyjmować postać typu prostego lub złożonego. Dane powinny być rozgraniczone przed procesem uczenia na dwa wektory, z których jeden przechowuje dane uczące, natomiast drugi dane walidacyjne.

**Listing 4.3.** Przykładowa funkcja realizująca krok uczenia [11]

```
1 void RunTrainStep(const std::vector<float>& input_batch,  
2                 const std::vector<float>& target_batch) {  
3     TF_CHECK_OK(session_>Run({"input", MakeTensor(input_batch)},  
4                             {"target", MakeTensor(target_batch)}},  
5                             {}, {"train"}, nullptr));  
6 }
```

## 4.4. Metody przetwarzania i eksploracji danych

API biblioteki TensorFlow dla języka C++ oferuje wiele operacji możliwych do wykonania na danych wykorzystywanych do uczenia. Operacje te można podzielić na manipulujące strukturą zapisanych danych (pozwalające na przetworzenie tensorów lub optymalizację pracy z danymi), oraz manipulujące wartościami regresorów (jak np. operacje matematyczne). Poniżej wymieniono wybrane, szczególnie istotne funkcjonalności przetwarzania i eksploracji danych, wyszczególnione w dokumentacji [8]:

### 4.4.1. Operacje strukturalne

Są to operacje mające na celu modyfikację struktury w jakich przechowywane są dane. Biblioteka oferuje szereg funkcjonalności mających na celu reformatowanie danych w celu optymalizacji wykonywania na nich obliczeń[8]. Należą do nich:

- **bitcast** - zmiana typu tensora bez kopiowania danych;
- **dekwantyzacja** - zamiana na typ zmiennoprzecinkowy z formatu stałoprzecinkowego lub całkowitoliczbowej notacji wykładniczej;
- **udawana kwantyzacja** - wykorzystywana do przygotowania danych aby móc z nich obliczyć gradient (dostępne jest kilka rodzajów udawanej kwantyzacji);
- **padding** - otoczenie wartości tensora nadmiarowymi wartościami w celu dostosowania wymiarów (możliwe uzupełnienie np. zerami, lub wartościami loszonymi);
- **stacking** - łączy N tensorów R-wymiarowych w jeden tensor R+1-wymiarowy;
- **dynamic stitching i partitioning** - łączenie danych z wielu tensorów w jeden lub rozdzielanie z jednego na wiele;
- **serializacja** - konwersja tensora na serializowany format dla Protobuf.

### 4.4.2. Eksploracja danych

W dokumentacji API TensorFlow dla C++ możliwe jest znalezienie wielu operacji pozwalających manipulować wartościami danych w celu wyliczenia pewnych cech, lub dokonania obliczeń potrzebnych w trakcie procesu uczenia[8]. Do najistotniejszych z nich należą:

- wyliczanie odcisku palca danych (ang. *fingerprint*);
- weryfikacja czy dane zawierają wartości NaN lub inf;
- wyodrębnienie fragmentów obrazów do warstwy „depth”;
- obliczenie gradientu z danych po udawanej kwantyzacji;
- obliczanie rzadkiego (ang. *sparse*) gradientu dla danego akumulatora;
- obliczenie odwrotnej permutacji tensora;
- normalizacja za pomocą kwantyzacji;
- wybór kandydatów do ewaluacji funkcji wyjściowych (ang. *candidate sampling*);
- dostosowywanie nasycenia, kontrastu i kolorystyki obrazów;
- dekodowanie obrazów z różnych formatów (np. gif, bmp, jpg, bmp);
- dekodowanie danych z różnych źródeł (np. csv, json);
- generowanie ramek (ang. *bounding box*);
- zmiana rozmiaru obrazów;
- rekodowanie obrazu z kolorystyki RGB na HSV;
- zachłanna selekcja ramek;
- różnorodne operacje matematyczne;
- pooling i konwoolucję;
- normalizację z pomocą gradientów;
- różnorodne funkcje aktywacji i straty;
- operacje na ciągach znakowych;
- algorytmy treningowe, jak np. Adagrad, Adadelata, centered RMSProp.

### 4.4.3. Sterowanie przepływem obliczeń

Oprócz wymienionych wcześniej poleceń optymalizujących lub przetwarzających dane, biblioteka TensorFlow posiada szereg operacji zarządzających przepływem danych w grafie obliczeniowym reprezentującym model. Należą do nich:

- operacje na zakumulowanych gradientach;
- ustawianie barier;
- operacje na mapach, kolejkach i tablicach tensorów.

## 4.5. Modele uczenia maszynowego

Biblioteka TensorFlow z racji dostarczania operacji na tensorach i realizacji przetwarzania za pomocą grafu obliczeniowego, wspiera szereg modeli zarówno płytkiego jak i głębokiego uczenia maszynowego. Należą do nich [12]:

- regresja liniowa;
- regresja logistyczna;
- regresja Deminga;
- regresja grzbietowa i lasso;
- maszyna wektorów nośnych;
- jednokierunkowa sieć neuronowa (ang. *feedforward neural network*);
- splotowa sieć neuronowa;
- rekurencyjna sieć neuronowa;
- autoenkodery;
- wzmocnione uczenie maszynowe;

Powyższe modele mogą zostać wykorzystane do rozwiązywania wielu problemów, w tym także o charakterystyce liniowej czy klasyfikacji. Poniższe fragmenty omawiają sposób implementacji poszczególnych modeli, oraz przedstawiają praktyczne przykłady dla wybranych z nich.

### 4.5.1. Regresja liniowa

Regresja liniowa przy wykorzystaniu API TensorFlow wykonywana jest poprzez metodę macierzy odwrotnej. Problem określony jest wzorem:

$$Ax = b \tag{4.1}$$

gdzie  $A$  odpowiada za macierz wartości regresorów,  $x$  za wektor współczynników poszczególnych zmiennych, a  $b$  stanowi wektor wartości oczekiwanych. Rozwiązanie w postaci ogólnej określić można wzorem:

$$x = (A^T A)^{-1} A^T b \quad (4.2)$$

Operacje te wykonać można wykorzystując transpozycję oraz iloczyn tensorów, dostępny w sekcji matematycznej API TensorFlow. Pozostałe rodzaje regresji przedstawione w dalszych sekcjach, korzystają z wyżej opisanego mechanizmu.

#### 4.5.2. Regresja logistyczna

Implementacja regresji logistycznej wykorzystuje proces regresji liniowej, jednak dokonuje na wyjściu normalizacji przewidywanego wyniku z użyciem funkcji sigmoidalnej, otrzymując prawdopodobieństwo przynależności. Obserwacja jest uznawana za należącą do sprawdzanej klasy jeżeli znormalizowana wartość przekroczy pewną zadaną wartość progową.

Implementację rozpoczęto od przygotowania mechanizmu parsowania pliku zawierającego dane eksperymentalne do postaci macierzy regresorów i wektora diagnozy. Wykorzystano w tym celu możliwość parsowania pliku CSV z dostępnego API biblioteki.

#### 4.5.3. Regresja Deminga

Polega ona na minimalizowaniu całkowitego dystansu pomiędzy punktami pomiarowymi a uzyskaną prostą, zamiast minimalizacji dystansu pionowego. W celu jej uzyskania należy wykorzystać mechanizm regresji liniowej, modyfikując funkcję straty, w celu minimalizacji następującego wzoru:

$$d = \frac{y_0 - (ax_0 + b)}{\sqrt{m^2 + 1}} \quad (4.3)$$

gdzie  $(x_0, y_0)$  oznaczają punkt pomiarowy, a prosta uzyskiwana w wyniku regresji dana jest wzorem  $y = ax + b$ . Powyższy wzór ma zastosowanie dla pomiarów w przestrzeni dwuwymiarowej (jeden regresor i jedna wartość odpowiedzi). W przypadku większej liczby regresorów należy zastosować metodę operatą o iloczynnym skalarny wektorów  $N+1$  wymiarowych, gdzie  $N$  oznacza ilość regresorów, za pomocą 3 punktów (punktu dla którego liczymy odległość, czyli obserwacji [punkt P], oraz dwóch punktów leżących na prostej uzyskanej w wyniku regresji [punkty A i B]). Metoda ta przedstawiona jest następującymi wzorami [13]:

$$\vec{pa} = P - A \quad (4.4)$$

$$\vec{ba} = B - A \quad (4.5)$$

$$t = (\vec{pa} \cdot \vec{ba}) / (\vec{ba} \cdot \vec{ba}) \quad (4.6)$$

$$d = |\vec{pa} - t * \vec{ba}| \quad (4.7)$$

Wzór 4.7 przedstawia wynikową formułę obliczania odległości, która musi być minimalizowana w ramach funkcji straty.

## 4.6. Metody analizy modeli

## 4.7. Dostępność dokumentacji i źródeł wiedzy

Niestety biblioteka TensorFlow posiada bardzo ograniczone zasoby umożliwiające pisanie z jej użyciem programów w języku C++. Dokumentacja okazuje się bardzo ograniczona i niekompletna. Dodatkowo niewielka liczba osób korzysta z tej biblioteki w języku C++, przez co ciężko znaleźć artykuły lub publikacje które prezentowałyby przykłady kodu źródłowego lub omawiały sposób pracy z udostępnionym API. Wiele dostępnych w API poleceń jak np. Placeholder, nie znalazło się w oficjalnej liście dostępnych funkcji i klas.

# Rozdział 5

## Biblioteka Shark

5.1. Wprowadzenie

5.2. Metody wdrożenia modeli

5.3. Formaty źródeł danych

5.4. Metody przetwarzania i eksploracji danych

5.5. Modele uczenia maszynowego

5.6. Metody analizy modeli

5.7. Dostępność dokumentacji i źródeł wiedzy



# Rozdział 6

## Biblioteka Caffe

6.1. Wprowadzenie

6.2. Metody wdrożenia modeli

6.3. Formaty źródeł danych

6.4. Metody przetwarzania i eksploracji danych

6.5. Modele uczenia maszynowego

6.6. Metody analizy modeli

6.7. Dostępność dokumentacji i źródeł wiedzy

# Rozdział 7

## Biblioteka PyTorch

7.1. Wprowadzenie

7.2. Metody wdrożenia modeli

7.3. Formaty źródeł danych

7.4. Metody przetwarzania i eksploracji danych

7.5. Modele uczenia maszynowego

7.6. Metody analizy modeli

7.7. Dostępność dokumentacji i źródeł wiedzy

# Rozdział 8

## Podsumowanie

8.1. Nakład pracy

8.2. Funkcjonalności

8.3. Dostępność źródeł wiedzy i dokumentacja

# Bibliografia

- [1] Olvi L. Mangasarian Dr William H. Wolberg W. Nick Street. *Wisconsin Diagnostic Breast Cancer (WDBC)*. 1995. URL: [https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(diagnostic\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(diagnostic)).
- [2] Dheeru Dua i Casey Graff. *UCI Machine Learning Repository*. 2017. URL: <http://archive.ics.uci.edu/ml>.
- [3] Trevor Bihl. *Biostatistics Using JMP: A Practical Guide*. Cary, NC: SAS Institute Inc., 2017.
- [4] Google Brain. *TensorFlow*. 2023. URL: <https://www.tensorflow.org/>.
- [5] Vitaly Bezhachev. *How to deploy Machine Learning models with TensorFlow. Part 1 - make your model ready for serving*. 2017.
- [6] Nour Islam Mokhtari. *How to Deploy TensorFlow Models in C++ in 3 different ways*. 2022.
- [7] Google Brain. *TensorFlow Lite*. 2023. URL: <https://www.tensorflow.org/lite>.
- [8] Google. *TensorFlow C++ API Reference*. 2023. URL: [https://www.tensorflow.org/api\\_docs/cc](https://www.tensorflow.org/api_docs/cc).
- [9] ksachdeva. *TensorFlow C++ Examples*. 2018. URL: [https://github.com/ksachdeva/tensorflow-cc-examples/blob/master/examples/1\\_Introduction/src/hello-world.cc](https://github.com/ksachdeva/tensorflow-cc-examples/blob/master/examples/1_Introduction/src/hello-world.cc).
- [10] Google. *A Tool Developer's Guide to TensorFlow Model Files*. 2023. URL: [https://chromium.googlesource.com/external/github.com/tensorflow/tensorflow/+r0.10/tensorflow/g3doc/how\\_tos/tool\\_developers/index.md](https://chromium.googlesource.com/external/github.com/tensorflow/tensorflow/+r0.10/tensorflow/g3doc/how_tos/tool_developers/index.md).
- [11] asimshankar. *Training TensorFlow Models in C++*. 2022. URL: <https://gist.github.com/asimshankar/5c96acd1280507940bad9083370fe8dc>.
- [12] Nick McClure. *TensorFlow Machine Learning Cookbook*. Packt Publishing, Luty 2017.
- [13] comingstorm. *Distance from point to n-dimentional line*. 2013. URL: <https://softwareengineering.stackexchange.com/questions/168572/distance-from-point-to-n-dimensional-line>.