

**Akademia Górniczo-Hutnicza  
im. Stanisława Staszica w Krakowie**

---

Wydział Elektrotechniki, Automatyki, Informatyki i Elektroniki



**PRACA MAGISTERSKA**

**ŁUKASZ DZIEDZIA**

**PORÓWNANIE JAKOŚCI UOGÓLNIENIA I  
EFEKTYWNOŚCI DZIAŁANIA RÓŻNYCH METOD  
KLASYFIKACJI Z SIECIAMI NEURONOWYMI  
SONN**

PROMOTOR:

dr Adrian Horzyk

Kraków 2009

## **OŚWIADCZENIE AUTORA PRACY**

OŚWIADCZAM, ŚWIADOMY ODPOWIEDZIALNOŚCI KARNEJ ZA  
POŚWIADCZENIE NIEPRAWDY, ŻE NINIEJSZĄ PRACĘ DYPLOMOWĄ  
WYKONAŁEM OSOBIŚCIE I SAMODZIELNIE, I NIE KORZYSTAŁEM ZE  
ŹRÓDEŁ INNYCH NIŻ WYMIENIONE W PRACY.

.....

PODPIS

**AGH**  
**University of Science and Technology in Krakow**

---

Faculty of Electrical Engineering, Automatics, Computer Science and  
Electronics



**MASTER OF SCIENCE THESIS**

**ŁUKASZ DZIEDZIA**

**COMPARISON OF GENERALIZATION QUALITY  
AND PROCESSING EFFICIENCY OF VARIOUS  
CLASSIFICATION METHODS WITH SONN  
NEURAL NETWORKS**

**SUPERVISOR:**

Adrian Horzyk Ph.D

Krakow 2009



# Spis treści

<b>1. Wstęp</b> .....	8
1.1. Motywacja.....	8
1.2. Cele oraz zarys pracy .....	8
1.2.1. Biblioteka algorytmów klasyfikacji.....	9
1.2.2. Metodyka tworzenia porównań, aplikacja.....	10
1.2.3. Porównanie wybranych metod klasyfikacji z algorytmem SONN.....	10
1.3. Opis rozdziałów .....	10
<b>2. Eksploracja danych i zagadnienie klasyfikacji</b> .....	11
2.1. Eksploracja danych .....	11
2.1.1. Geneza .....	11
2.1.2. Obszary zastosowań .....	12
2.2. Metody eksploracji danych .....	12
2.2.1. Reguły asocjacyjne .....	13
2.2.2. Regresja .....	13
2.2.3. Analiza skupień .....	13
2.3. Klasyfikacja.....	13
2.3.1. Matematyczne sformułowanie zagadnienia klasyfikacji .....	14
<b>3. Metody klasyfikacji</b> .....	16
3.1. Metoda k najbliższych sąsiadów .....	17
3.1.1. Opis metody.....	17
3.1.2. Szczegóły implementacyjne oraz parametry metody .....	18
3.2. Support Vector Machines .....	19
3.2.1. Opis metody.....	19
3.2.2. Soft Margin.....	22
3.2.3. Kernel Trick.....	22
3.2.4. Zastosowanie SVM dla zagadnień niebinarnych.....	22

3.2.5.	Szczegóły implementacyjne oraz parametry metody .....	24
3.3.	Radial Basis Function Networks .....	25
3.3.1.	Trening warstwy ukrytej .....	26
3.3.2.	Trening warstwy wyjściowej .....	27
3.3.3.	Szczegóły implementacyjne oraz parametry metody .....	27
3.4.	Probabilistic Neural Networks .....	27
3.4.1.	Szczegóły implementacyjne oraz parametry metody .....	29
3.5.	Self Optimizing Neural Networks .....	29
3.5.1.	Ontogeniczne Sieci Neuronowe .....	29
3.5.2.	Sieci SONN - opis .....	30
3.5.3.	Architektura .....	32
3.5.4.	Opis działania algorytmu .....	35
3.5.5.	Szczegóły implementacyjne .....	37
<b>4.</b>	<b>Metodologia porównań metod klasyfikacji .....</b>	<b>38</b>
4.1.	Proces porównywania .....	38
4.2.	Zbiory danych .....	40
4.3.	Wstępne przetwarzanie danych .....	46
4.3.1.	Normalizacja .....	46
4.3.2.	Standaryzacja .....	46
4.3.3.	Principal Component Analysis .....	47
4.4.	Metody walidacji .....	47
4.4.1.	k-fold cross-validation .....	48
4.4.2.	Leave-one-out cross-validation .....	49
4.4.3.	Podział procentowy .....	49
4.4.4.	Test na danych treningowych .....	49
4.5.	Kryteria porównania .....	50
4.5.1.	Jakość uogólnienia .....	51
4.5.2.	Efektywność działania .....	53
4.5.3.	Skalowalność .....	55
4.5.4.	Interpretowalność .....	55
4.5.5.	Inne kryteria .....	56
<b>5.</b>	<b>Część praktyczna .....</b>	<b>57</b>
5.1.	Założenia .....	57
5.2.	Wybrane aspekty inżynierii oprogramowania .....	58

5.3. Stosowane technologie.....	59
5.4. Szczegóły implementacyjne.....	60
5.4.1. Struktura projektu.....	60
5.4.2. Podstawowe interfejsy .....	61
5.5. Prezentacja aplikacji .....	62
5.5.1. Wybór treningowego zbioru danych.....	62
5.5.2. Wybór metod wstępnego przetwarzania danych .....	63
5.5.3. Wybór oraz parametryzacja algorytmów.....	64
5.5.4. Konfiguracja symulacji.....	64
5.5.5. Przeglądanie wyników.....	66
<b>6. Wyniki porównań metod klasyfikacji.....</b>	<b>68</b>
6.1. Stanowisko testowe i uwagi o procesie porównywania .....	68
6.2. Porównania dla zbioru Wine .....	69
6.3. Porównania dla zbioru Iris .....	77
6.4. Porównania dla zbioru Yeast.....	80
6.5. Porównania dla zbioru Congressional Voting Records.....	83
6.6. Porównania dla zbioru Car Evaluation .....	86
6.7. Porównania dla zbioru Ionosphere.....	89
6.8. Porównania dla zbioru Arrhythmia.....	92
6.9. Skalowalność - porównanie wybranych wskaźników.....	95
<b>7. Podsumowanie .....</b>	<b>98</b>
7.1. Wnioski .....	98
7.2. Możliwości rozwoju pracy .....	100
<b>A. Zawartość dołączonej płyty CD-ROM.....</b>	<b>101</b>

# 1. Wstęp

Fascynacja zagadnieniem klasyfikacji danych nie jest zjawiskiem nowym. Ludzkość od stuleci próbuje klasyfikować rozmaite obiekty, dzięki czemu poznaje ich naturę i wykorzystuje zdobytą wiedzę do najróżniejszych celów. Wraz z rozwojem techniki komputerowej możliwości gromadzenia, przetwarzania i analizy danych nieustannie się polepszają, co stawia nowe wyzwania dla twórców oraz użytkowników algorytmów klasyfikacji. W związku z tym konieczne jest konstruowanie nowych oraz ulepszanie istniejących algorytmów tak, aby coraz skuteczniej radzić sobie z ogromną ilością danych i pozyskiwać cenne informacje w nich zawarte. Ponadto, wiele z istniejących algorytmów w swoim działaniu lub budowie przypomina niedościgniony wzór „naturalnego komputera” – ludzki mózg, co sprawia, że badania nad tymi metodami klasyfikacji są niezwykle interesujące.

## 1.1. Motywacja

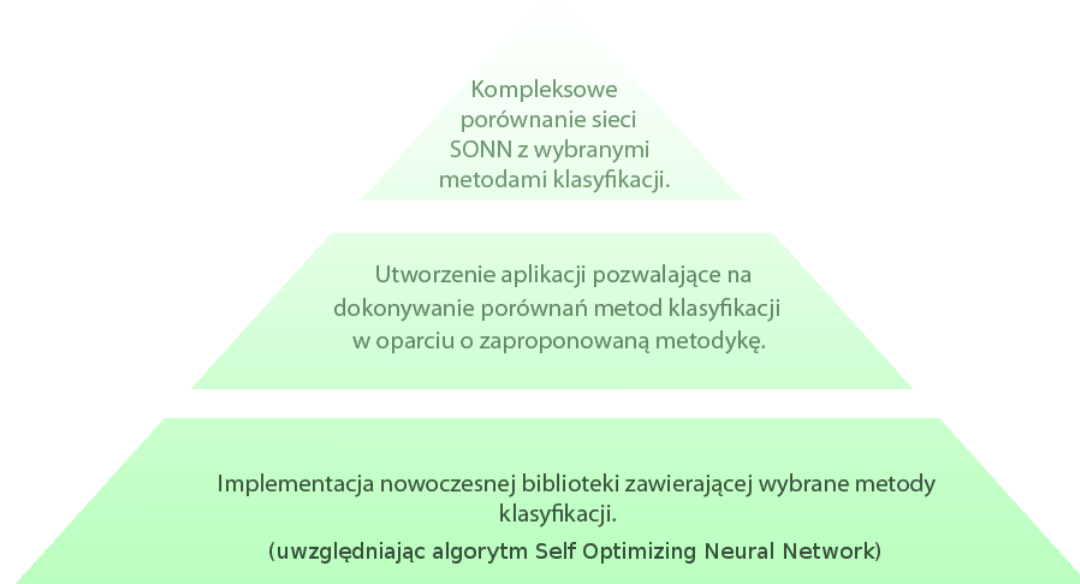
Osobistą motywacją do podjęcia tego tematu pracy było zainteresowanie metodami sztucznej inteligencji, dziedziną eksploracji danych oraz związanymi z tym zagadnieniami. Praca ta jest naturalną konsekwencją wybranego kierunku studiów oraz projektów i prac podejmowanych w trakcie edukacji. Dobór technologii oraz wybranych metod inżynierii oprogramowania, dzięki którym udało się zrealizować pracę, stanowi element podsumowania zdobytej wiedzy i doświadczeń. Realizując tą pracę intencją autora było zwiększenie oraz ugruntowanie wiadomości z podejmowanych dziedzin oraz dodanie skromnego wkładu do badań naukowych prowadzonych na Akademii Górniczo-Hutniczej w Krakowie.

## 1.2. Cele oraz zarys pracy

Tematem pracy jest dokonanie analizy porównawczej jakości uogólnienia oraz efektywności działania wybranych metod klasyfikacji z algorytmem *Self Optimizing Neural Networks* (SONN) autorstwa dr Adriana Horzyka, który od lat rozwijany jest w ramach badań naukowych na AGH w Krakowie. W niniejszej pracy analizowany jest algorytm SONN w wersji 3 (SONN-3) i ilekroć w pracy będzie mowa o algorytmie SONN to dotyczy to wersji SONN-3.



Wyniki uzyskane w pracy mają dostarczyć informacji o możliwościach i ograniczeniach algorytmu, które mogą posłużyć do rozwoju kolejnych wersji tego algorytmu.



Rysunek 1.1: Etapy pracy nad aplikacją, główne cele

Czynności niezbędne do wykonania analizy porównawczej wybranych metod klasyfikacji z sieciami SONN można opisać, jako proces składający się z trzech głównych etapów (patrz rysunek 1.1) opisanych w poniższych sekcjach.

### 1.2.1. Biblioteka algorytmów klasyfikacji

W celu uzyskania miarodajnych wyników porównań algorytmy muszą zostać zaimplementowane w tej samej technologii. Istnieje wiele porównań rozmaitych algorytmów, ale najczęściej przeprowadzane są z wykorzystaniem różnych aplikacji, które różnią się technologią wykonania, jakością, czy innymi szczegółami wpływającymi na powstawanie różnic w zakresie czasu tworzenia klasyfikatora. Ponadto zdarza się, że algorytmy w istniejących narzędziach wspierających eksplorację danych nie są dokładnie opisane, przez co trudno użytkownikom zweryfikować otrzymane rezultaty.

Istotnym elementem tej pracy jest etap implementacji spójnej biblioteki zawierającej wybrane algorytmy na wspólnej platformie technologicznej, przy jednoczesnym zadbaniu o udostępnienia informacji o parametrach i szczegółach implementacyjnych algorytmów. W ramach tego etapu konieczna jest analiza oraz implementacja wybranych algorytmów klasyfikacji danych. Jakość wykonania czynności w tym etapie jest kluczowa dla dalszych prac, stąd bardzo istotne jest zastosowanie odpowiednich narzędzi inżynierii oprogramowania opisanych w rozdziale 5. „Część praktyczna”.

### 1.2.2. Metodyka tworzenia porównań, aplikacja

Na tym etapie prac najważniejszym zadaniem jest analiza zagadnień związanych z procesem porównywania metod klasyfikacji. W oparciu o zebrane informacje konieczne jest zaproponowanie i opisanie metodyki dokonywania porównań pozwalającej na rzetelne skonfrontowanie algorytmu SONN z wybranymi metodami klasyfikacji. Na podstawie nakreślonej metodyki należy zaimplementować praktyczną, nowoczesną aplikację wspierającą proces porównania algorytmów. Aplikacja wspierająca porównywanie powinna wykorzystywać wcześniej zaimplementowaną bibliotekę algorytmów.

### 1.2.3. Porównanie wybranych metod klasyfikacji z algorytmem SONN

Ostatnim etapem realizowanym w ramach pracy jest wykorzystanie utworzonej aplikacji do wykonania porównań metod klasyfikacji, będących istotą pracy. Dokonując porównań należy odpowiednio dobrać zbiory danych, różne konfiguracje algorytmów oraz metod wstępnego przetwarzania danych. Praca powinna zostać zakończona przedstawieniem rezultatów oraz ich interpretacją, a także wskazaniem potencjalnych możliwości udoskonalenia algorytmu SONN.

## 1.3. Opis rozdziałów

Rozdział 1 „Wstęp” – wprowadzenie do tematyki pracy, przedstawienie celów oraz zakresu pracy.

Rozdział 2 „Eksploracja danych i zagadnienie klasyfikacji” – wprowadzenie teoretyczne do dziedziny eksploracji danych, a w szczególności zagadnienia klasyfikacji będącego przedmiotem pracy.

Rozdział 3 „Metody klasyfikacji” – opis metod klasyfikacji danych implementowanych oraz analizowanych w ramach pracy, z uwzględnieniem algorytmu *Self Optimizing Neural Networks*.

Rozdział 4 „Metodologia porównań algorytmów klasyfikacji” – opis zagadnień związanych z porównywaniem metod klasyfikacji. Opis procesu oraz kluczowych koncepcji wykorzystanych w pracy.

Rozdział 5 „Część praktyczna” – prezentacja aplikacji, techniczne aspekty zaimplementowanej biblioteki oraz aplikacji, opis stosowanych metod inżynierii oprogramowania.

Rozdział 6 „Wyniki porównań metod klasyfikacji” – prezentacja wyników przeprowadzonych porównań.

Rozdział 7 „Podsumowanie” – zebranie wniosków dotyczących porównań, analiza możliwości rozwoju systemu oraz algorytmu SONN.

Dodatek A – zawartość dołączonej płyty CD.

## 2. Eksploracja danych i zagadnienie klasyfikacji

Algorytmy klasyfikacji, w tym sieci SONN, stanowią istotny fragment dziedziny nazywanej *data mining*, która w związku z lawinowym przyrostem informacji staje się obszarem nauki coraz bardziej popularnym i pożądanym. W języku polskim występuje wiele tłumaczeń zwrotu *data mining*. W niniejszym rozdziale i całym opracowaniu używana będzie nazwa eksploracja danych.

W rozdziale „Eksploracja danych i zagadnienie klasyfikacji” zaprezentowana zostanie dziedzina eksploracji danych, omówione zostaną związane z nią metody oraz szczegółowo omówione zostanie zagadnienie klasyfikacji, będące przedmiotem niniejszej pracy.

### 2.1. Eksploracja danych

Eksploracja danych jest czynnością naturalną dla każdego człowieka. Ludzki mózg potrafi analizować i rozpoznawać nawet bardzo skomplikowane wzorce danych. Dla przykładu można wymienić zdolność do rozpoznawania ludzi po głosie, rysach twarzy czy innych cechach osobniczych. Zatem dlaczego powstała osobna dziedzina nauki poświęcona tak naturalnej dla człowieka czynności?

#### 2.1.1. Geneza

Ludzkość generuje oraz przechowuje coraz większe ilości danych. Dane pochodzą z różnych źródeł, do których można zaliczyć chociażby centra badawcze, instytucje rządowe, firmy, czy popularny Internet. Aktualnie szacuje się, że co każde trzy lata ilość przechowywanych danych zostaje podwojona [1]. Bardzo prawdopodobne jest również to, że z czasem będzie się skracał wraz z rozwojem technologii.

Ogromne ilości danych, których analiza wykracza poza możliwości człowieka stały się powodem do powstania dziedziny eksploracji danych, której zadaniem jest dostarczanie metod i narzędzi pomocnych w procesie analizy danych pod kątem wyszukiwania istotnych informacji w dużych zbiorach danych.

### 2.1.2. Obszary zastosowań

Eksploracja danych znajduje zastosowanie w wielu dziedzinach nauki stając się istotnym czynnikiem wpływającym na ich rozwój. Do obszarów zastosowań eksploracji danych należą:

- Medycyna – Zagadnienie eksploracji danych w medycynie jest chętnie wykorzystywane, tym bardziej, że wpływa na zwiększenie skuteczności leczenia, czy zapobiegania rozwojowi groźnych chorób.

Przykład: Zastosowanie niektórych metod pozwala na tworzenie systemów wspierających diagnostykę medyczną (jednym z przykładów zastosowania może być diagnoza płuc w celu wykrycia zmian nowotworowych).

- Biznes – jest to dziedzina, w której eksploracja danych znalazła wiele praktycznych zastosowań, a umiejętnie stosowana może doprowadzić do zwiększenia zysków oraz dynamiki rozwoju przedsiębiorstwa.

Przykład: Bardzo przydatną cechą nowoczesnych systemów wspierających zarządzanie relacjami z klientami (*Customer Relationship Management, CRM*) jest możliwość dokonania analizy skupień grup klientów, czy dokonywania analizy koszykowej w celu usprawnienia kampanii marketingowej przedsiębiorstwa.

- Technika – podobnie jak w przypadku wyżej wymienionych obszarów zastosowań metody eksploracji danych znalazły swoje zastosowanie w wielu gałęziach techniki.

Przykład: Diagnoza stanu skomplikowanej infrastruktury technicznej (np. sieci energetyczne) odbywa się za pomocą metod eksploracji danych.

- Biotechnologia – ze względu na dużą ilość danych występujących w badaniach biotechnologicznych eksploracja danych jest bardzo skutecznym narzędziem wspierającym badania w zakresie tej dziedziny.

Przykład: eksploracja danych w biotechnologii jest pomocna przy analizie zależności między różnymi strukturami DNA a prawdopodobieństwem wystąpienia raka.

- Inne – eksploracja danych jest lub może być wykorzystywana wszędzie tam, gdzie mamy do czynienia ze zbiorami danych, których rozmiary sprawiają, że niemożliwa staje się ich analiza przez człowieka.

## 2.2. Metody eksploracji danych

Rozróżnia się kilka metod występujących w wachlarzu narzędzi eksploracji danych. Najważniejsze z nich zostaną opisane poniżej.

### 2.2.1. Reguły asocjacyjne

Reguły asocjacyjne stosowane są w celu znalezienia zależności pomiędzy występowaniem wzorców w zbiorze danych. Częstym przykładem zastosowania reguł asocjacyjnych jest analiza koszykowa, dzięki której sprzedawca może wyznaczyć reguły występujące w procesie kupna jego produktów.

Przykład: Sprzedawca oferujący produkty A, B, C, D, E po zastosowaniu algorytmów wyszukujących reguły asocjacyjne na zbiorze archiwalnych transakcji może uzyskać informacje o zależnościach występujących przy kupnie produktów, np. kupując produkty B, C kupowany jest także produkt E. Takie algorytmy wykorzystywane są w wielu nowoczesnych serwisach zakupowych (np. amazon.com).

### 2.2.2. Regresja

Z regresją mamy do czynienia, gdy na podstawie informacji o próbkach ze zbioru danych (gdzie każda z próbek określona jest przez zestaw cech i związana jest z pewną, ciągłą wartością wynikową) chcemy znaleźć zależność pozwalającą na przewidzenie, jaka wartość wynikowa związana będzie z nową, nieznaną wcześniej próbką określoną przez taki sam zestaw cech.

Przykład: Regresja jest wykorzystywana do przewidywania szeregów czasowych (na przykład w analizie giełdowej).

### 2.2.3. Analiza skupień

Kolejną metodą, dzięki której wyszukać możemy cenne informacje w dużym zbiorze danych jest analiza skupień. Dzięki niej próbki danych opisane przez wiele cech możemy pogrupować zgodnie z pewnym, nieznanym wcześniej podobieństwem (wydzielić klastry).

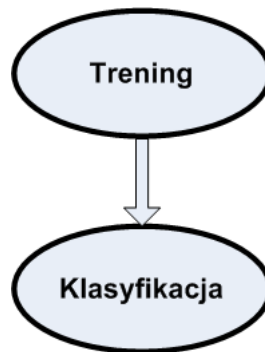
Przykład: Analiza skupień wykorzystywana jest między innymi przez przeglądarkę internetową Clusty.org [2] do semantycznego grupowania wyników wyszukiwań.

## 2.3. Klasyfikacja

Zagadnienie klasyfikacji możemy scharakteryzować jako próbę znalezienia mechanizmu, pozwalającego na przyporządkowanie próbki danych określonej przez pewien zestaw cech do jednej ze znanych klas. Taki mechanizm nosi nazwę klasyfikatora i jego działanie uzależnione jest od wybranej koncepcji, natomiast proces tworzenia klasyfikatora nazywa się najczęściej uczeniem lub treningiem. Matematyczne sformułowanie zagadnienia klasyfikacji znajduje się w podrozdziale 2.3.1, natomiast opis metod analizowanych w ramach pracy w rozdziale 3.

Warto zwrócić uwagę na podobieństwo klasyfikacji z zagadnieniem analizy skupień, dla obu podejść celem jest określenie przynależności nieznanego wzorca do pewnej grupy. Najistotniejszą różnicą pomiędzy tymi dwoma podejściami jest brak zdefiniowanego zbioru klas dla zagadnienia analizy skupień (*ang. unsupervised learning*), podczas gdy metody klasyfikacji wymagają sprecyzowania klasy dla każdego z występujących w problemie wzorców (*ang. supervised learning*).

Bardzo często w literaturze problem klasyfikacji definiuje się jako zagadnienie dwuetapowe, przedstawione schematycznie na rysunku 2.1.



Rysunek 2.1: Dwuetapowość procesu klasyfikacji

Najczęstszym sposobem wykorzystywania algorytmów klasyfikacji jest wytrenowanie najlepszego klasyfikatora dla danego problemu (zbioru danych), a następnie wielokrotne stosowanie go do klasyfikacji danych wcześniej nieznanymi, na których klasyfikator nie był uczony, licząc na jego zdolność do uogólnienia zdobytej w trakcie treningu wiedzy.

Podobnie jak wcześniej opisywane metody, również klasyfikacja danych stosowana jest w wielu dziedzinach. Potwierdzeniem tego może być zestaw zbiorów danych wykorzystywanych przy porównaniach metod klasyfikacji w tej pracy, który został opisany w rozdziale 4.2. Nie jest kompletna lista możliwych obszarów zastosowań, a raczej jej skromny wycinek, który jednak powinien obrazować wszechstronność w zakresie możliwości wykorzystania zagadnienia klasyfikacji.

### 2.3.1. Matematyczne sformułowanie zagadnienia klasyfikacji

Jeżeli:

$X$  - przestrzeń cech

$Y$  - zbiór znanych klas

Dany jest  $n$ -elementowy zbiór treningowy:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}, \text{ gdzie:}$$

$\mathbf{x}_i \in X$  - wektor cech określający  $i$ -ty element zbioru danych (wzorec)

$y_i \in Y$  - element ze zbioru dostępnych klas, mówiący do jakiej klasy  $i$ -ty obiekt należy  
To zagadnienie klasyfikacji polega na znalezieniu takiego odwzorowania (klasyfikatora):

$$k : X \longrightarrow Y,$$

które dla dowolnego nieznanego wektora  $\mathbf{x} \in X$  określi poprawnie, do jakiej klasy  $y \in Y$  dany wektor należy.

### 3. Metody klasyfikacji

Spośród wielu popularnych metod klasyfikacji wybranych zostało kilka, które zostaną porównane z metodą SONN. Podczas wybierania algorytmów klasyfikacji do zestawu zaimplementowanego oraz stosowanego w porównaniach stosowano następujące kryteria:

- Popularność metody – istotne jest, aby dokonywana analiza było miarodajna oraz porównywalna. Osiągnięcie tego jest możliwe dzięki możliwości odniesienia wyników testów metody SONN do wyników algorytmów, które są rozpoznawalne i często występują w różnego rodzaju publikacjach i testach.
- Skuteczność metody – należy pamiętać, że algorytmy klasyfikacji stanowią narzędzia, które mają skutecznie wspierać inne dziedziny. Jest to istotny argument przemawiający za koniecznością tworzenia porównań takich, jak opisane w niniejszej pracy. Takie podejście wymusza skupienie uwagi na metodach osiągających dobre rezultaty, które mają szansę zostać zastosowane w praktycznych zagadnieniach.
- Własne zainteresowanie – spośród wielu istniejących algorytmów ostatecznie do implementacji wybrane zostały te, które okazały się szczególnie interesujące. Do takich metod niewątpliwie należą SVM i RBFN. W porównaniach nie znalazła się metoda perceptronu wielowarstwowego (*Multilayer Perceptron*, MLP), która była przedmiotem innego realizowanego przeze mnie projektu [3].

Ponieważ nie sposób było zaimplementować wszystkich istniejących metod, duży nacisk położony został na implementację biblioteki tworzonej w ramach pracy tak, aby możliwe było rozszerzanie istniejącego kodu o nowe algorytmy. Szczegóły tego rozwiązania omówione zostaną w rozdziale 5.

W dalszej części tego rozdziału opisane zostaną metody zaimplementowane w ramach pracy. Na szczególną uwagę zasługują informacje o szczegółach implementacyjnych, parametrach oraz różnych wariantach metod. Gdy będzie to możliwe zilustrowane zostaną przykłady działania algorytmu.



## 3.1. Metoda k najbliższych sąsiadów

Metoda k najbliższych sąsiadów (*k Nearest Neighbours*, kNN)[4] jest jedną z najprostszych metod klasyfikacji. Dzięki swojej prostocie wysoka interpretowalność<sup>1</sup> nadaje się doskonale jako narzędzie dydaktyczne podczas poznawania zagadnienia klasyfikacji. Również dzięki prostocie algorytm ten stał się bardzo popularny w dziedzinie eksploracji danych, co przemawia za sensownością porównywania innych metod z tym algorytmem.

### 3.1.1. Opis metody

Metoda k najbliższych sąsiadów bazuje na prostej zasadzie:

„wzorzec danych należy do tej klasy, do których należy największa liczba spośród k sąsiadujących z nim wzorców”.

Trening metody k najbliższych sąsiadów sprowadza się do zapamiętania zbioru uczącego (metoda *lazy*). Jest to dość specyficzne dlatego, że większość algorytmów klasyfikacji w fazie treningu wykonuje czasochłonne obliczenia. Konsekwencją tej specyfiki jest wydłużony czas klasyfikacji w porównaniu z innymi metodami.

Klasyfikacja konkretnego wzorca polega na wyznaczeniu k próbek ze zbioru treningowego znajdujących się najbliżej (w kontekście stosowanej metryki) badanego wzorca. Do najczęściej stosowanych metryk w algorytmie *k Nearest Neighbours* należą:

- Metryka euklidesowa
- Metryka Manhattan
- Metryka Hamminga

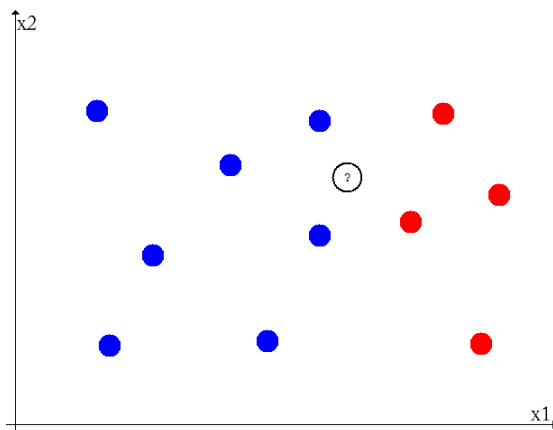
Ostatnim etapem jest określenie klasy, do jakiej należy analizowana próbka. W praktyce do najczęściej stosowanych metod wyboru klasy należą:

- *Majority Voting* (Głosowane na równych prawach) – próbka należy do tej klasy, do której należy najwięcej spośród k znalezionych sąsiadów.
- *Inverse Distance Voting* (Głosowanie z uwzględnieniem odległości) – dla każdej z klas występujących wśród znalezionych k sąsiadów obliczana jest suma odwróconych odległości od analizowanej próbki. Próbka danych zostaje zakwalifikowana do klasy, dla której obliczona suma jest największa.

---

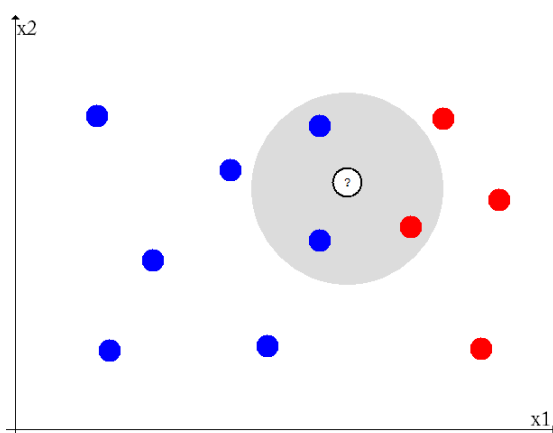
<sup>1</sup>Kryterium interpretowalności zostało zarysowane w rozdziale 4.5.4.

Działanie metody k najbliższych sąsiadów zostało zilustrowane dla przykładu dwuwymiarowego na rysunkach 3.1 oraz 3.2. Na pierwszym z nich widoczne są wzorce należące do dwóch różnych klas (kolory niebieski i czerwony) oraz klasyfikowany punkt (oznaczony znakiem zapytania), dla którego chcemy określić przynależność do jednej z klas.



Rysunek 3.1: Koncepcja działania metody kNN, problem

Na rysunku 3.2 zaznaczony został obszar zawierający  $k$  (w przykładzie  $k = 3$ ) najbliższych sąsiadów.



Rysunek 3.2: Koncepcja działania metody kNN, rozwiązanie

Stosując metodę *Majority Voting* można stwierdzić, że dla  $k = 3$  badany punkt będzie należał do klasy oznaczonej kolorem niebieskim.

### 3.1.2. Szczegóły implementacyjne oraz parametry metody

Zaimplementowany algorytm korzysta z metryki euklidesowej jako sposobu do wyznaczania odległości pomiędzy wektorami danych.

Zastosowana metoda wyboru klasy to bardzo popularna metoda *Majority Voting*.

Parametry związane z metodą k najbliższych sąsiadów przedstawione są w tabeli 3.1.

Tablica 3.1: Parametry metody kNN

Nazwa	Opis
k	Ilość wyszukiwanych najbliższych sąsiadów. Wraz ze wzrostem ilości sąsiadów zmniejsza się ryzyko zbyt zgrubnego oszacowania przynależności do klasy oraz zwiększa czas obliczeń. Należy jednak starannie dobierać wielkość tego parametru tak, aby zasięg poszukiwania sąsiedztwa nie był zbyt szeroki, co może doprowadzić do pogorszenia wyników klasyfikacji.

## 3.2. Support Vector Machines

Metoda *Support Vector Machines* (Maszyny Wektorów Wspierających, SVM)[5, 6, 7, 8, 9, 10, 11, 12] jest kolejną metodą klasyfikacji wybraną i zaimplementowaną w ramach niniejszej pracy. Jest to algorytm, który osiąga bardzo dobre wyniki w wielu testach, dlatego jest często stosowany w praktycznych zagadnieniach.

### 3.2.1. Opis metody

Istotą metody *Support Vector Machines* jest odnalezienie hiperpłaszczyzny, która podzieli przestrzeń cech na podprzestrzenie reprezentujące każdą z klas. Analizę tej metody warto rozpocząć od wyprowadzenia niezbędnych podstaw, na przykładzie najprostszego przypadku, a następnie rozszerzyć analizę o bardziej zaawansowane koncepcje, takie jak *Soft Margin*, *Kernel Trick*, czy metody osiągania wieloklasowości.

Rozważmy najprostszy przypadek. Mamy  $m$ -elementowy zbiór danych:

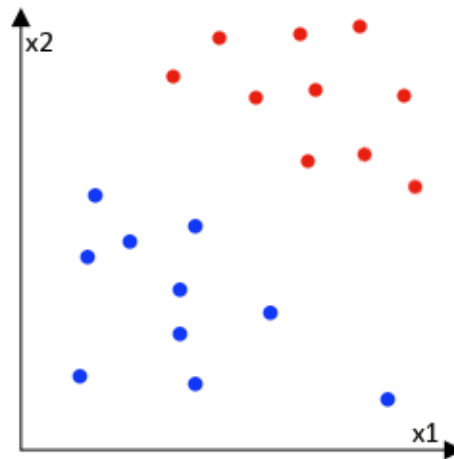
$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}, \text{ gdzie:}$$

$\mathbf{x}_i$  - wektor określający  $i$ -ty element zbioru

$y_i$  - klasa  $i$ -tego elementu zbioru

Analizowany zbiór jest dwuklasowy (binarne zagadnienie klasyfikacji) oraz dane są liniowo separowane (założenie). Ilustracją takiego przykładu dla  $n = 2$ , gdzie  $n$  jest liczbą wymiarów przestrzeni cech może być rysunek 3.3.

Dla takiego przypadku istnieje nieskończenie wiele hiperpłaszczyzn, które potrafią podzielić przestrzeń cech na podprzestrzenie reprezentujące każdą z klas. Przykłady takich rozwiązań zostały zamieszczone na rysunku 3.4 jako proste w kolorze zielonym. Punkty znajdujące się na hiperpłaszczyźnie spełniają zależność:



Rysunek 3.3: Koncepcja działania metody SVM, liniowo separowalny (dwuklasowy) zbiór danych

$$\mathbf{w}\mathbf{x}_i * b = 0, \text{ gdzie:}$$

$\mathbf{w}$  - normalna względem hiperpłaszczyzny

$\mathbf{x}$  - punkt w przestrzeni

Dzięki założeniu liniowej separowalności danych wszystkie punkty ze zbioru uczącego spełniają odpowiednio:

$$\mathbf{w}\mathbf{x}_i + b \geq 1, \text{ dla } y_i = 1 \quad (3.1)$$

$$\mathbf{w}\mathbf{x}_i + b \leq -1, \text{ dla } y_i = -1 \quad (3.2)$$

Co dla wszystkich punktów zbioru uczącego można zapisać jako:

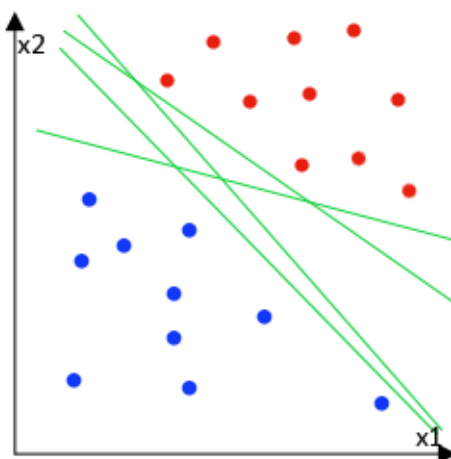
$$y_i(\mathbf{w}\mathbf{x}_i + b) - 1 = 0 \quad (3.3)$$

Aby wybrać najlepsze rozwiązanie poszukuje się klasyfikatora, który podzieli przestrzeń z największym marginesem. Zastosowanie takiego rozwiązania motywowane jest założeniem, że wraz ze wzrostem szerokości marginesu zmniejsza się ryzyko popełnienia błędu w przypadku klasyfikacji nieznanymi danymi.

Można udowodnić [7], że maksymalizacja marginesu sprowadza się do minimalizacji funkcji celu:

$$Q(\mathbf{w}) = \mathbf{w}^T \mathbf{w}, \text{ przy ograniczeniach:} \quad (3.4)$$

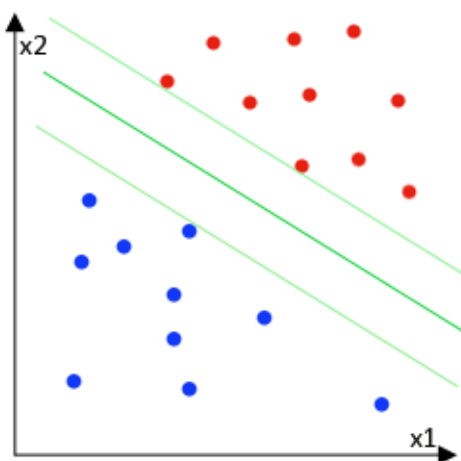
$$y_i(\mathbf{w}\mathbf{x}_i + b) - 1 = 0 \quad (3.5)$$



Rysunek 3.4: Koncepcja działania metody SVM, przykłady możliwych rozwiązań

Tak postawiony problem można rozwiązać wieloma sposobami (metoda najszybszego spadku, metoda symulowanego wyżarzania). Bardzo często stosowanym sposobem jest zastosowanie metody programowania kwadratowego (*ang. Quadratic Programming*), które może być rozwiązane za pomocą algorytmu *Sequential Minimal Optimization* (SMO) [9, 10] opisanego w dalszej części tego rozdziału.

Znalezione zgodnie z tą metodą rozwiązanie przedstawion zostało na rysunku 3.5. Opisywany do tej pory przypadek nosi nazwę *Linear SVM* (LSVM).



Rysunek 3.5: Koncepcja działania metody SVM, klasyfikator liniowy z największym marginesem

Założmy, że dane nie są liniowo separowalne, co w praktyce jest bardzo częstym zjawiskiem. W takim przypadku zastosowanie metody maksymalnego marginesu przedstawione wcześniej nie jest uzasadnione. Istnieją dwa podejścia pozwalające na rozwiązanie problemu danych nieseparowalnych liniowo.

### 3.2.2. Soft Margin

Pierwszym podejściem stosowanym w przypadku danych nieseparowalnych liniowo jest metoda *Soft Margin*. Polega ona na dodaniu czynnika kary za odległość od hiperpłaszczyzny, dla każdego błędnie sklasyfikowanego wzorca. Zmodyfikowana funkcja celu ma postać:

$$Q(\mathbf{w}) = \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^R \varepsilon_i, \text{ gdzie:} \quad (3.6)$$

$C$  - współczynnik kary

$R$  - ilość błędnie sklasyfikowanych próbek

$\varepsilon_i$  - odległość  $i$ -tej błędnie sklasyfikowanej próbki od marginesu

### 3.2.3. Kernel Trick

Inne podejście pozwalające stosować metodę SVM dla zbiorów zawierających dane nieseparowalne liniowo nazywane jest *Kernel Trick* [11]. Podejście to polega na mapowaniu nie separowalnych liniowo próbek do przestrzeni o większej liczbie wymiarów, w której możliwe jest zastosowanie klasyfikatora liniowego (dane są separowalne liniowo). W tym celu stosuje się jedną z funkcji jądrowych (*ang. kernel function*). Do najczęściej używanych metod jądrowych w algorytmie SVM należą:

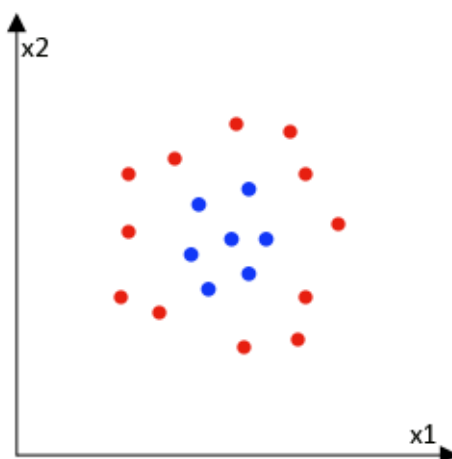
- funkcja gaussowska  $K(\mathbf{u}, \mathbf{v}) = e^{-\left(\frac{\|\mathbf{u}-\mathbf{v}\|^2}{2\sigma^2}\right)}$
- funkcja wielomianowa  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}\mathbf{v} + 1)^d$  lub  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}\mathbf{v})^d$
- funkcja sigmoidalna  $K(\mathbf{u}, \mathbf{v}) = \tanh(\eta\mathbf{u}\mathbf{v} + \rho)$

Na rysunkach 3.6 i 3.7 przedstawiono schematycznie zasadę działania *Kernel Trick* dla przypadku dwuwymiarowego. Rysunek 3.6 przedstawia zbiór danych nieseparowalnych liniowo, w którym jak w poprzednich przypadkach występują dwie klasy (zagadnienie binarne).

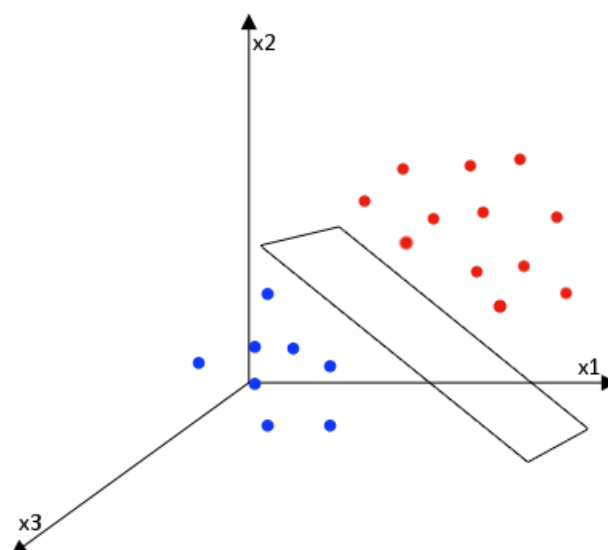
Na kolejnym rysunku 3.7 przedstawiono odpowiednio przetransformowane dane do przestrzeni o większej liczbie wymiarów oraz hiperpłaszczyznę wyznaczającą nowe obszary klasyfikacji.

### 3.2.4. Zastosowanie SVM dla zagadnień niebinarnych

Przedstawione powyżej rozważania dotyczą binarnego zagadnienia klasyfikacji. Najczęściej klasyfikowane dane należą do większej ilości klas, stąd istnieje konieczność adaptacji algorytmu SVM do rozwiązywania problemów wieloklasowych [12]. Rozwiązaniem tego problemu jest stworzenie wielu klasyfikatorów binarnych i odpowiednie ich wyspecjalizowanie. Przy założeniu, że  $k$  jest ilością klas występujących w problemie stosowane podejścia można opisać jako:



Rysunek 3.6: Koncepcja działania metody SVM, dane nieseparowalne liniowo



Rysunek 3.7: Koncepcja działania metody SVM, „Kernel Trick”

- *One Versus One* – dla każdej z istniejących par klas tworzony jest klasyfikator binarny potrafiący rozróżnić, do której z tych dwu klas analizowana próbka danych należy. Przynależność do klasy określana jest na podstawie największej ilości rozpoznań próbki jako reprezentanta danej klasy. Ilość utworzonych klasyfikatorów binarnych wynosi  $\binom{k}{2}$ .
- *One Versus All* – tworzonych jest  $k$  klasyfikatorów. Każdy z nich specjalizuje się w rozpoznawaniu przynależności do jednej z klas. Ostatecznie próbka zostaje sklasyfikowana do klasy, której klasyfikator wyznaczył największą wartość wyjściową (*winner takes all strategy*). Ilość utworzonych klasyfikatorów binarnych dla tego przypadku wynosi  $k$ .

### 3.2.5. Szczegóły implementacyjne oraz parametry metody

Algorytm SVM został zaimplementowany w oparciu o koncepcję Johna C. Platta [9, 10]. Koncepcja ta nosi nazwę *Sequential Minimal Optimization* (SMO). Zaproponowany przez Platta algorytm dotyczy rozwiązania zagadnienia programowania kwadratowego wynikającego z przedstawienia problemu SVM jako minimalizacji funkcji celu (patrz wzór 3.4). Koncepcja ta polega na zastąpieniu pierwotnego zagadnienia szeregiem mniejszych problemów programowania kwadratowego, które mogą zostać rozwiązane bez konieczności wykonywania skomplikowanych i czasochłonnych obliczeń numerycznych. Nie sposób w ramach niniejszej pracy analizować dokładnie tej ciekawej koncepcji, warto natomiast zwrócić uwagę na wyniki algorytmu SVM zaimplementowanego w wersji SMO z innymi podejściami [9, 10], co uzasadnia wybór tej metody.

Ze względu na szerokie zainteresowanie metodą SVM, w niniejszej pracy zostały zaimplementowane trzy odmiany tej metody:

1. Linear SVM – implementacja koncepcji LSVM opisanej we wcześniejszej części tego rozdziału.
2. RBF SVM – implementacja wykorzystująca funkcję gaussowską jako funkcję jądrową metody.
3. Polynomial SVM – implementacja wykorzystująca funkcję wielomianową jako funkcję jądrową metody.

Parametry stosowane w metodzie *Support Vector Machines* przedstawione zostały w tabeli 3.2.

Tablica 3.2: Parametry metody SVM

Nazwa	Opis
eps	Dla wszystkich metod. Parametr związany z dokładności obliczeń. Najczęściej jest to liczba rzędu 0.1 – 0.001.
cost	Dla wszystkich metod. Parametr wysokości kosztu za błędne sklasyfikowanie próbki. Jego wpływ wynika z zależności w opisanej funkcji celu(odnośnik).
sigma	Dla metody RBF SVM. Wartość odchylenia standardowego występująca w stosowanej funkcji Gaussa.
d	Dla metody Polynomial SVM. Parametr stopnia stosowanego wielomianu.



### 3.3. Radial Basis Function Networks

Istnieje wiele odmian sieci neuronowych, które są wykorzystywane w zagadnieniu klasyfikacji. W ramach pracy zaimplementowana została interesująca ich odmiana, a mianowicie *Radial Basis Function Networks* (Sieci Radialnych Funkcji Bazowych, RBFN) [13, 14, 15, 16]. *Radial Basis Function Networks* są sieciami bez sprzężeń zwrotnych (*ang. feedforward*). Ten rodzaj sieci składa się z trzech warstw:

- Warstwa wejściowa – rola neuronów w tej warstwie ogranicza się do przekazywania sygnałów wejściowych do warstwy ukrytej.
- Warstwa ukryta – neurony w tej warstwie wykorzystują radialną funkcję bazową do wygenerowania sygnału wyjściowego.
- Warstwa wyjściowa – sumują wygenerowane w warstwie ukrytej sygnały wyjściowe (uwzględniając wagę każdego z połączeń), na tej podstawie obliczane są sygnały wyjściowe całej sieci.

Schemat budowy sieci typu RBF został przedstawiony na rysunku 3.8. Występujące na tym rysunku symbole oznaczają kolejno:

$x_i$  –  $i$ -ty element wektora wejściowego danych

$R_j$  –  $j$ -ty neuron warstwy ukrytej, wartość wyjściowa neuronu  $R_0$  to zawsze 1

$y_m$  –  $m$ -te wyjście sieci odpowiadające klasie  $m$

$n$  – wymiar przestrzeni cech

$k$  – ilość neuronów warstwy ukrytej

$c$  – ilość neuronów warstwy wyjściowej (ilość klas występujących w problemie klasyfikacji)

Radialna funkcja jest funkcją spełniająca zależność (przyjęta norma najczęściej oznacza normę euklidesową):

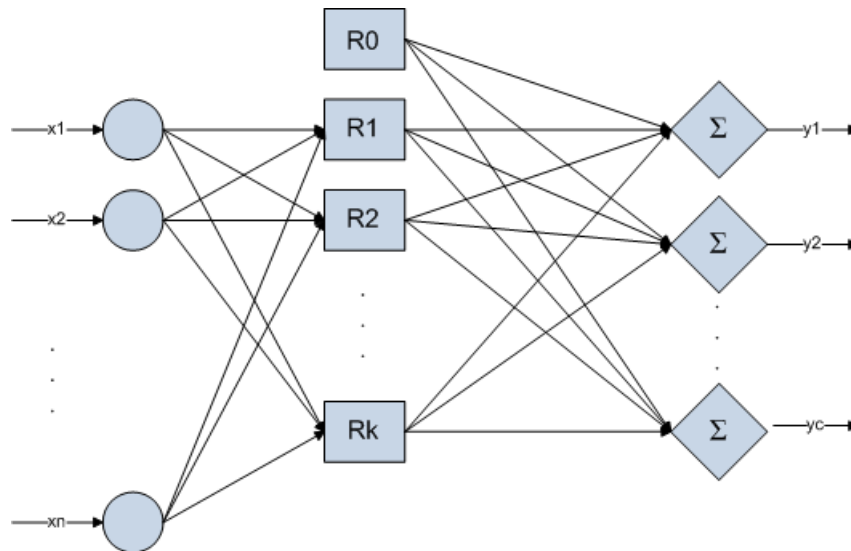
$$\phi(\mathbf{x}, \mathbf{c}) = \phi(\|\mathbf{x} - \mathbf{c}\|)$$

Najczęściej wykorzystywane funkcje radialne w sieciach RBF to:

- funkcja gaussowska  $f(\mathbf{x}, \mathbf{c}) = e^{-(\beta\|\mathbf{x}-\mathbf{c}\|^2)}$
- funkcja wielokwadratowa (*ang. multiquadratic*)  $f(\mathbf{x}, \mathbf{c}) = \sqrt{\|\mathbf{x} - \mathbf{c}\|^2 + \beta^2}$

Trening *Radial Basis Function Networks* podzielić można na dwa etapy:

- Trening warstwy ukrytej – wybór centrów dla funkcji radialnych oraz innych parametrów funkcji. Ta część treningu odbywa się bez nadzoru (*ang. unsupervised learning*).



Rysunek 3.8: Schemat architektury Radial Basis Function Network

- Trening warstwy wyjściowej – dobranie współczynników wag dla każdego z neuronów warstwy wyjściowej. Dobieranie współczynników jest nazywane nauką z nadzorem (*ang. supervised learning*).

Możliwość podziału treningu tak, aby warstwa ukryta i wyjściowa trenowane były oddzielnie jest istotne z punktu widzenia wydajności. Dzięki takiej konstrukcji sieci, w tym algorytmie nie występuje wymagający obliczeniowo proces wstecznej propagacji błędów (*ang. backpropagation*).

### 3.3.1. Trening warstwy ukrytej

Pierwszym z zadań związanych z uczeniem warstwy ukrytej sieci RBF jest znalezienie centrów funkcji radialnych. Istnieje kilka podejść do sposobu wyboru centrów neuronów warstwy ukrytej. Do najczęściej stosowanych należą:

- Wybór losowy
- Wybór za pomocą algorytmów analizy skupień
- Wybór za pomocą modeli probabilistycznych

Z zagadnieniem treningu warstwy ukrytej związany jest także dobór pozostałych parametrów wymaganych przez funkcje radialne. Dla obu wymienionych wcześniej funkcji radialnych konieczne jest określenie parametru dyspersji ( $\beta$ ), który najczęściej jest wspólny dla każdego z neuronów warstwy ukrytej, a jego wartość to średnia odległość między kolejnymi centrami.

### 3.3.2. Trening warstwy wyjściowej

W kolejnym etapie treningu sieci RBF konieczne jest określenie współczynników wag dla każdego z połączeń warstwy ukrytej z warstwą wyjściową. Wagi muszą zostać dobrane tak, aby minimalizować wartość błędu będącego różnicą między oczekiwaną wartością wyjścia sieci a wartością obliczoną przez sieć dla każdego z neuronów warstwy wyjściowej. Tak postawiony problem można sprowadzić do minimalizacji funkcji błędu:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left( \sum_{k=0}^K \mathbf{w}_k \phi_k^n - y^n \right)^2, \text{ gdzie} \quad (3.7)$$

$N$  - ilość próbek w zbiorze treningowym

$K$  - ilość neuronów warstwy ukrytej

$\mathbf{w}_k$  - waga powiązana z wartością wyjściową  $k$ -tego neuronu warstwy ukrytej

$\phi_k^n$  - wartość funkcji radialnej  $k$ -tego neuronu warstwy ukrytej dla  $n$ -tego wektora

$y^n$  - oczekiwana wartość wyjściowa dla  $n$ -tego wektora zbioru treningowego

### 3.3.3. Szczegóły implementacyjne oraz parametry metody

Zaimplementowane w ramach pracy sieci RBF pozwalają na wykorzystanie funkcji Gaussa oraz funkcji wielokwadratowej, jako funkcji radialnych występujących w neuronach warstwy ukrytej. Centra dla funkcji radialnych mogą zostać wylosowane ze zbioru treningowego lub znalezione za pomocą algorytmów analizy skupień (*K-mean clustering*).

Parametry związane z metodą RBFN przedstawione zostały w tabeli 3.3.

Tablica 3.3: Parametry metody RBFN

Nazwa	Opis
k	Ilość neuronów w warstwie ukrytej. Wielkość ta powinna zostać dostosowana do ilości wzorców w zbiorze treningowym. Przy zbyt dużej ilości neuronów warstwy ukrytej istnieje ryzyko przeuczenia sieci ( <i>ang. overfitting</i> ).

## 3.4. Probabilistic Neural Networks

Metoda *Probabilistic Neural Networks* (Probabilistyczne Sieci Neuronowe, PNN) [17, 18, 19, 20] jest inspirowana bayesowskimi metodami klasyfikacji i stanowi ich rozwinięcie. Głównym założeniem tej metody jest próba aproksymacji funkcji gęstości prawdopodobieństwa rozkładu wektorów ze zbioru treningowego.

Schemat budowy *Probabilistic Neural Networks* przedstawiony został na rysunku 3.9. Występujące na rysunku symbole oznaczają kolejno:

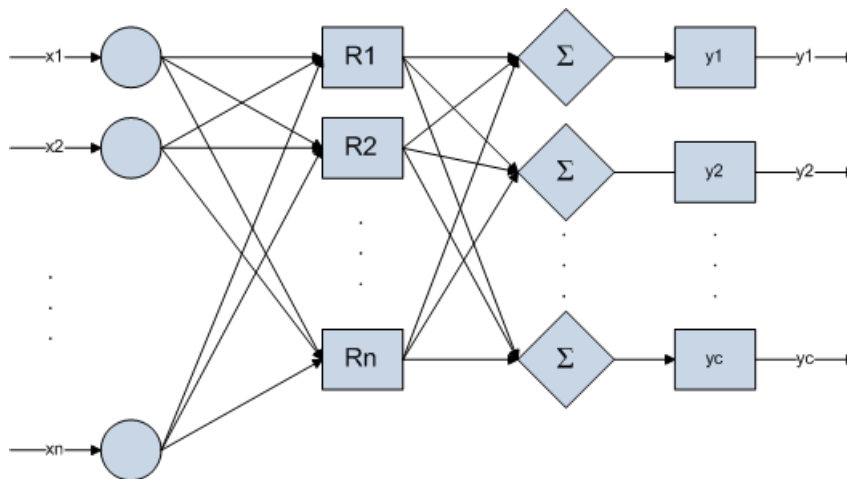
$x_i$  –  $i$ -ty element wektora wejściowego danych

$P_i$  –  $i$ -ty neuron warstwy ukrytej

$y_m$  –  $m$ -te wyjście sieci odpowiadające klasie  $m$

$n$  – wymiar przestrzeni cech, ilość neuronów w warstwie ukrytej

$c$  – ilość neuronów warstwy wyjściowej (ilość klas występujących w problemie klasyfikacji)



Rysunek 3.9: Schemat architektury Probabilistic Neural Network

Sieci PNN zbudowane są z następujących warstw:

- Warstwa wejściowa – składa się z neuronów, których liczba równa jest wymiarowi przestrzeni cech, jej zadaniem jest propagacja sygnałów wejściowych do każdego neuronu warstwy ukrytej.
- Warstwa ukryta – każdy z neuronów w tej warstwie skojarzony jest z jednym wektorem występującym w ciągu uczącym. Porównując działanie tej warstwy z warstwą ukrytą sieci RBF (patrz podrozdział 3.3.1), można powiedzieć, że każdy z neuronów odpowiada centrum zlokalizowanemu w punkcie określonym przez próbkę z ciągu uczącego, z którą jest skojarzony.
- Warstwa wyjściowa – zawiera dokładnie tyle neuronów, ile znanych jest klas w danym problemie klasyfikacji. Neurony w tej warstwie sumują wartości obliczone w warstwie ukrytej dla danej klasy. Prawdopodobieństwo przynależności do klasy reprezentowanej przez neuron warstwy wyjściowej może zostać obliczone po znormalizowaniu wartości wyjściowych do przedziału  $(0, 1]$ .

### 3.4.1. Szczegóły implementacyjne oraz parametry metody

Implementacja metody *Probabilistic Neural Networks* zastosowana w niniejszej pracy inspirowana jest opracowaniem tej metody autorstwa Nikolaya Y. Nikolaeva [17]. Zastosowana funkcja jądrowa to funkcja gaussowska.

Parametry stosowane w metodzie PNN przedstawione zostały w tabeli 3.4.

Tablica 3.4: Parametry metody PNN

Nazwa	Opis
sigma	Parametr odchylenia standardowego dla funkcji Gaussa stosowanej w neuronach warstwy ukrytej.

## 3.5. Self Optimizing Neural Networks

Spośród wielu metod klasyfikacji szczególną rolę odgrywają sieci neuronowe. Istnieje wiele rodzajów sieci stosowanych w zagadnieniu klasyfikacji, jak na przykład Perceptron Wielowarstwowy (MLP), czy omawiane w poprzednich rozdziałach Sieci Radialnych Funkcji Bazowych (RBFN) lub Sieci Probabilistyczne (PNN). Modele te charakteryzują się jednak pewnymi ograniczeniami, które często dyskwalifikują te metody w praktycznych zastosowaniach. Do takich ograniczeń należą:

- Potrzeba określenia topologii sieci – konieczność określania liczby warstw oraz liczby neuronów. Konieczność ustalenia architektury sieci a priori lub w sposób doświadczalny sprawia, że proces ten jest bardzo czasochłonny, a najczęściej prowadzi do nieoptymalnego doboru struktury, co wpływa na obniżenie zdolności generalizacyjnych modelu.
- „Przekleństwo wymiarowości” (*ang. curse of dimensionality*) – wraz ze wzrostem liczby wymiarów przestrzeni cech nieliniowo rośnie liczba neuronów (lub warstw) niezbędnych do stworzenia sieci realizującej poprawną klasyfikację danych.

Problemy te były jednym z powodów opracowania koncepcji ontogenicznych sieci neuronowych.

### 3.5.1. Ontogeniczne Sieci Neuronowe

Sieci ontogeniczne charakteryzują się zdolnością do automatycznej konstrukcji swej topologii (struktury). Sieci tego rodzaju możemy sklasyfikować ze względu na charakterystykę algorytmu konstrukcyjnego:

- Modele powiększające strukturę – to podejście polega na rozpoczęciu procesu uczenia od pewnej minimalnej struktury i stopniowym powiększaniu jej (dodawanie neuronów, warstw i połączeń synaptycznych) tak, aby polepszać zdolności generalizacji sieci.
- Modele zmniejszające strukturę – w tym podejściu proces uczenia rozpoczyna się ustalenia początkowej struktury, która dobrze realizuje zadanie generalizacji na określonym zbiorze (dalsze etapy uczenia nie wpływają na redukcję błędu). Po ustaleniu wstępnej struktury następuje proces redukcji sieci najczęściej poprzez usuwanie połączeń międzyneuronowych oraz neuronów, które mają nieistotny wpływ na realizowane zadanie klasyfikacji.

Analizowana w niniejszej pracy metoda *Self Optimizing Neural Networks* autorstwa dr Adriana Horzyka [21, 22, 23, 24, 25, 26] należy do grupy sieci ontogenicznych. Obok niej, jako przykłady innych podejść, można wymienić takie modele jak: RAN (*Resource Allocation Network*) [27], IncNet [28], upstart, czy algorytm korelacji kaskadowej (*ang. cascade-correlation network*).

### 3.5.2. Sieci SONN - opis

Celem autora metody SONN było zaprojektowanie uniwersalnego narzędzia klasyfikacji. Metoda ta realizuje postulat ontogeniczności, czyli automatycznego konstruowania topologii oraz doboru wag połączeń synaptycznych. Model SONN można sklasyfikować jako sieć powiększającą strukturę (główny algorytm konstrukcyjny), przy czym warto zauważyć, że algorytm ten posiada także możliwość dokonywania rewizji struktury sieci i redukowaniu zbędnych neuronów.

Algorytm *Self Optimizing Neural Networks* rozwijany jest od kilku lat w ramach badań [21, 22, 23, 24, 25, 26] prowadzonych przez dr Adriana Horzyka na Akademii Górniczo-Hutniczej w Krakowie. W trakcie badań zaprezentowano kilka wersji algorytmu, których krótka historia zamieszczona jest poniżej.

- SONN-1 (rok 2001) – wersja ta ma początki w opracowywanej w pracy doktorskiej [26] dr Adriana Horzyka metodzie automatycznej konfiguracji sieci neuronowych dla problemów rozpoznawania wzorców binarnych. Do głównych możliwości sieci należy umiejętność automatycznej konstrukcji topologii sieci. Podstawowym ograniczeniem tego modelu było wymaganie stosowania wyłącznie wzorców binarnych w zbiorze treningowym.
- SONN-2 (rok 2005) – poprawa oraz optymalizacja algorytmu konstrukcyjnego (szacowana złożoność obliczeniowa:  $O([\text{wymiar przestrzeni cech}] * [\text{ilość wektorów treningowych}]^2)$ ), zdolność do redukcji wymiaru przestrzeni cech

poprzez automatyczne wartościowanie najbardziej znaczących cech z punktu widzenia dyskryminacji wzorców.

- SONN-3 (rok 2008) – najnowsza wersja, która jest przedmiotem niniejszej pracy. Metoda SONN została wzbogacona o inteligentny sposób tworzenia warstwy wejściowej sieci (algorytm *Automatic Discriminative Lossy Binarization Conversion Algorithm*, ADLBCA), która potrafi klasyfikować dane o wartościach rzeczywistych, całkowitych oraz binarnych. Oprócz tej zmiany dokonano dalszej optymalizacji algorytmu konstrukcyjnego.
- SONN-4 – wersja, która aktualnie jest rozwijana. Planowane zmiany to umożliwienie konstruowania sieci w oparciu o dane symboliczne oraz umiejętność odnajdywania kombinacji dyskryminujących poszczególne klasy już na etapie tworzenia warstwy wejściowej.

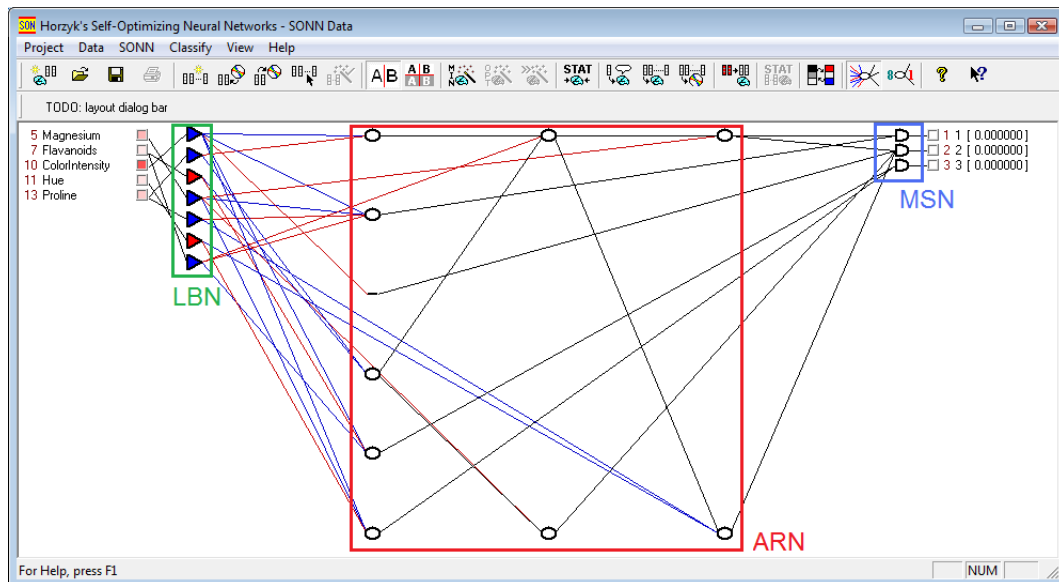
Najistotniejsze cechy *Self Optimizing Neural Networks* przedstawione zostały w poniższym zestawieniu.

- Zdolność do automatycznego kreowania struktury sieci (ontogeniczność).
- Automatyczny dobór parametrów w procesie uczenia.
- Umiejętność automatycznej konwersji rzeczywistych danych wejściowych do postaci binarnych w warstwie wejściowej sieci (algorytm ADLBCA).
- Zdolność do klasyfikacji z zerowym błędem wzorców zastosowanych podczas konstruowania sieci (jeżeli wzorce zbioru treningowego nie są sprzeczne).
- Automatyczna redukcja wymiaru przestrzeni cech – podczas konstruowania sieci cechy, które nie mają istotnego znaczenia w procesie klasyfikacji zostają pominięte, co korzystnie wpływa na zdolność uogólniania oraz czas działania algorytmu.
- Dzięki dążeniu do zachowania jak najmniejszej struktury sieci, metoda SONN posiada dobre właściwości generalizacyjne, jest niewrażliwa na przeuczenie (*ang. overfitting*).
- Algorytm SONN jest niewrażliwy na rozkład klas w zbiorze treningowym.
- Umiejętność rozpoznawania negatywów wzorców danych.

Ciekawym aspektem działania sieci SONN jest stosowanie wartości binarnych jako pobudzeń generowanych przez neurony LBN sieci, co stanowi analogię do sposobu funkcjonowania ludzkiego mózgu.

### 3.5.3. Architektura

Sieci SONN są sieciami wielowarstwowymi. Ilość warstw i neuronów w każdej warstwie jest determinowana w procesie konstrukcyjnym<sup>2</sup>. Przykładową sieć SONN utworzoną dla zbioru danych Wine przedstawia rysunek 3.10. Na tym rysunku oznaczono miejsca występowania neuronów każdego z typów występujących w sieciach *Self Optimizing Neural Networks*.



Rysunek 3.10: Przykładowa struktura sieci SONN utworzona dla zbioru Wine (program autorstwa dr Adriana Horzyka)

W sieciach SONN wyróżnić można trzy rodzaje neuronów wyspecjalizowanych w pełnieniu określonych funkcji, które występują w kolejnych warstwach sieci opisanych w poniższych sekcjach.

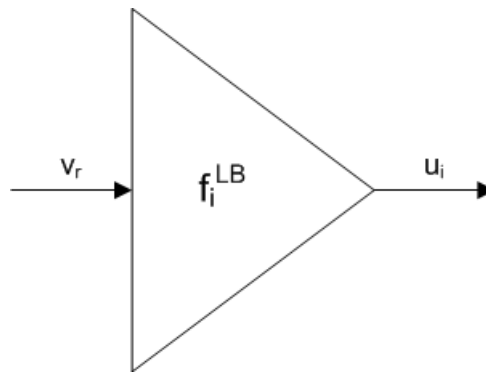
#### Warstwa wejściowa (Lossy Binarizing Subnetwork)

Warstwa wejściowa sieci SONN składa się wyłącznie z neuronów typu LBN (*Lossy Binarizing Neuron*). Ich zadaniem jest przekształcenie rzeczywistych danych wejściowych na binarne sygnały przesyłane do kolejnych warstw sieci tak, aby już w początkowym etapie tworzenia sieci wyróżnić istotne cechy. Na podstawie tych cech utworzony zostanie model sieci. Schemat neuronu LBN znajduje się na rysunku 3.11. Za tworzenie warstwy wejściowej sieci odpowiedzialny jest algorytmu binaryzacji (np. ADLBCA), którego działanie opisane zostanie w podrozdziale 3.5.4.

Każdy neuron LBN odpowiada za przekształcenie wektora danych rzeczywistych  $v_r$  na bipolarny sygnał  $u_i = f_i^L B(v_r, R_i)$ . Neuron taki związany jest z konkretną cechą  $k$  zbioru

<sup>2</sup>Proces konstrukcyjny w sieciach ontogenicznych jest odpowiednikiem procesu uczenia (treningu) dla innych algorytmów klasyfikacji.





Rysunek 3.11: Self Optimizing Neural Network, schemat LBN

uczącego oraz przedziałem wartości tej cechy, który jednoznacznie wskazuje na przynależność do jednej z klas problemu. Funkcja wyznaczająca wartość wyjścia neuronu LBN przedstawiona jest poniżej.

$$f_i^{LB}(v_r, R_i) = \begin{cases} +1 & \text{dla } v_r^k \in R_i^k, \\ -1 & \text{w przeciwnym przypadku.} \end{cases}, \text{ gdzie} \quad (3.8)$$

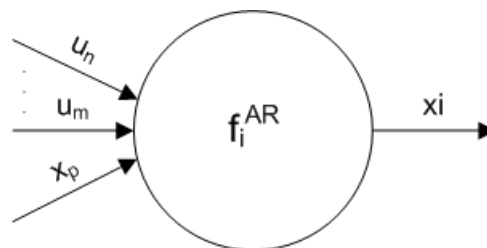
$f_i^{LB}$  - wartość wyjścia  $i$ -tego neuronu LBN

$v_r^k$  - wartość  $k$ -tej cechy wektora liczb rzeczywistych

$R_i^k$  - przedział reprezentowany przez  $i$ -ty neuron LBN dla  $k$ -tej cechy

### Warstwy ukryte (Aggregation Reinforcement Subnetworks)

Warstwy ukryte powstają w wyniku działania głównego algorytmu konstrukcyjnego sieci. W każdej z warstw ukrytych występują wyłącznie neurony typu ARN (*Aggregation Reinforcement Neuron*), których schemat zilustrowany został na rysunku 3.12. Celem warstw ukrytych jest odszukiwanie i wzmacnianie podobieństw między wzorcami występującymi w zbiorze treningowym.



Rysunek 3.12: Self Optimizing Neural Network, schemat ARN

Neurony ARN na wejściu mogą przyjmować sygnały wygenerowane przez neurony LBN w warstwie wejściowej ( $u_m, \dots, u_n$ ) lub sygnały wyjściowe z innego (pojedynczego) neuronu ARN spośród neuronów występujących w warstwach wcześniejszych ( $x_p$ ). Każde z wejść

związane jest z automatycznie obliczaną wartością wagi. Model matematyczny neuronu ARN przedstawiony został w wzorze 3.9.

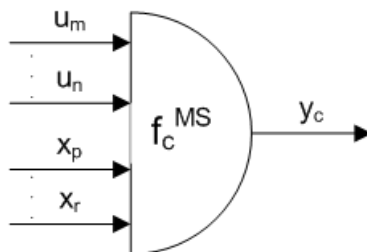
$$x_i = f_i^{AR}(u_m, \dots, u_n, x_p) = w_0^{x_p} x_p + \sum_{k \in \{i, \dots, j\}} w_k^{u_k} u_k \quad (3.9)$$

Liczba warstw oraz neuronów wzmacniających agregację zależna jest od stopnia korelacji danych treningowych. Wyróżnić można dwie reguły określające zależność pomiędzy stopniem korelacji danych a strukturą sieci SONN:

1. Wraz ze wzrostem stopnia korelacji wzorców reprezentujących różne klasy w zbiorze treningowym, struktura sieci staje się coraz bardziej skomplikowana (rośnie liczba warstw i neuronów) i odwrotnie.
2. Wraz ze zmniejszaniem się stopnia korelacji wzorców reprezentujących te same klasy w zbiorze treningowym, struktura sieci staje się coraz bardziej skomplikowana i odwrotnie.

### Warstwa wyjściowa (Maximum Selection Subnetwork)

Warstwa wyjściowa sieci SONN jest najprostszym strukturalnie oraz logicznie elementem modelu sieci. W skład tej warstwy wchodzi wyłącznie neuron typu MSN (*Maximum Selection Neuron*). Liczba neuronów zależy ściśle od liczby klas występujących dla danego zbioru danych (dla każdej z klas tworzony jest jeden neuron MSN). Zadaniem neuronów ostatniej warstwy jest wyznaczenie maksimum z otrzymywanych na wejściu sygnałów. Wejściowe sygnały mogą pochodzić bezpośrednio z warstwy wejściowej ( $u_m, \dots, u_n$ ) lub z dowolnej liczby neuronów ARN ( $x_p, \dots, x_r$ ). Warto zauważyć, że wejścia neuronów MSN nie są związane z żadną wagą. Schemat neuronu MSN przedstawiony został na rysunku 3.13, a jego matematyczny opis w wzorze 3.10.



Rysunek 3.13: Self Optimizing Neural Network, schemat MSN

$$y_c = f_m^{MS}(u_m, \dots, u_n, x_p, \dots, x_r) = \max(u_m, \dots, u_n, x_p, \dots, x_r) \quad (3.10)$$

### 3.5.4. Opis działania algorytmu

Działanie metody jest zdecydowanie bardziej wyszukane od wielu algorytmów opisywanych we wcześniejszych częściach pracy. Metody tej nie sposób opisać prostymi regułami, jest to raczej zbiór algorytmów wyspecjalizowanych w konkretnych zadaniach. Opis działania algorytmu zostanie przedstawiony zgodnie z porządkiem tworzenia kolejnych warstw sieci, począwszy od pojawieniu się na wejściu danych uczących.

#### Tworzenie warstwy wejściowej – Lossy Binarizing Subnetwork

Za tworzenie warstwy wejściowej sieci SONN odpowiedzialny jest algorytm *Automatic Discriminative Lossy Binarization Conversion Algorithm* (ADLBCA) [25]. Celem tego algorytmu jest utworzenie grupy neuronów odpowiedzialnych za przekształcenie rzeczywistych danych wejściowych do postaci binarnych (bipolarnych) sygnałów. Dokonuje się tego poprzez wyszukanie przedziałów na zbiorze danych dla każdej z cech, które najlepiej determinują przynależność do konkretnej klasy występującej w danym problemie klasyfikacji. Dąży się do minimalizacji ilości neuronów w tej warstwie poprzez wyszukiwanie jak największych przedziałów (algorytm zachłanny). Warto zauważyć, że już na tym etapie dokonuje się selekcja cech, które mają największy wpływ na określenie przynależności do klas.

Dla każdego ze znalezionych przedziałów tworzony jest neuron LBN (patrz podrozdział 3.5.3), który na swym wyjściu da sygnał  $+1$ , jeżeli podany wzorec ma wartość cechy w przedziale, który dany neuron reprezentuje lub  $-1$  w przeciwnym przypadku.

#### Tworzenie warstw ukrytych – Aggregation Reinforcement Subnetwork

Po utworzeniu neuronów LBN następuje główny algorytm konstrukcyjny modelu SONN. Omawianie tego kroku należy rozpocząć od wyjaśnienia pojęcia *Global Discrimination Coefficient* (globalny współczynnik rozróżnienia, GDC).

Globalny współczynnik rozróżnienia wyznaczany jest dla każdego neuronu LBN (każdej zbinaryzowanej cechy używanej przez sieć SONN). Celem wprowadzania tego współczynnika jest konieczność dokładnego określenia zdolności dyskryminacyjnych neuronów doprowadzających sygnały do warstw ukrytych, dzięki temu możliwe jest przeprowadzenie algorytmu konstrukcyjnego i optymalny dobór wag połączeń synaptycznych. W celu wyznaczenia współczynników GDC konieczna jest analiza całego zbioru treningowego. Skutkuje to koniecznością przebudowy sieci po zmianie zbioru uczącego, co generalnie jest cechą charakterystyczną większości klasyfikatorów.

$$d_{(k+)}^n = \begin{cases} \frac{P_k^m}{(M-1)Q^m} \sum_{h=1 \wedge h \neq m}^M \left(1 - \frac{P_k^h}{Q^h}\right) & \text{gdy } u_k^n \geq 0 \wedge u^n \in C^m, \\ 0 & \text{gdy } u_k^n < 0 \wedge u^n \in C^m. \end{cases} \quad (3.11)$$

$$d_{(k-)}^m = \begin{cases} \frac{N_k^m}{(M-1)Q^m} \sum_{h=1 \wedge h \neq m}^M \left(1 - \frac{N_k^h}{Q^h}\right) \text{ gdy } u_k^n \leq 0 \wedge u^n \in C^m, \\ 0 \text{ gdy } u_k^n > 0 \wedge u^n \in C^m. \end{cases} \quad (3.12)$$

$$P_k^m = \sum_{u_k^n \in \{u \in U \vee C^m : u_k^n > 0, n \in \{1, \dots, Q\}\}} u_k^n \quad (3.13)$$

$$N_k^m = \sum_{u_k^n \in \{u \in U \vee C^m : u_k^n < 0, n \in \{1, \dots, Q\}\}} -u_k^n \quad (3.14)$$

$$Q^m = \|\ u \in U \cap C^m : n \in \{1, \dots, Q\} \ \| \quad (3.15)$$

Symbole w powyższych wzorach oznaczają:

$M$  - liczba klas dla konkretnego problemu klasyfikacji

$Q$  - liczba wzorców w zbiorze treningowym

$C^m$  - identyfikator  $m$ -tej klasy

$u_k^n$  - wartość binarnego wyjścia  $n$ -tego neuronu związanego z  $k$ -tą cechą

Po wyznaczeniu globalnych współczynników dyskryminacji konieczne jest określenie, które z nich posłużą w procesie konstrukcyjnym tak, aby wszystkie wzorce treningowe mogły zostać sklasyfikowane poprawnie. Takie współczynniki noszą nazwę obligatoryjnych (*Obligatory*) i uwzględnienie ich jest konieczne. Rozróżnia się jeszcze współczynniki opcjonalne, których uwzględnienie może nie wpłynąć na rozmiar sieci. Algorytm wyznaczania obligatoryjnych i opcjonalnych współczynników dyskryminacji przebiega w następujących krokach:

1. Dla każdego ze współczynników, które nie zostały uznane za obligatoryjne, należy oszacować ilość nierozróżnionych wcześniej wzorców, które mogą zostać rozróżnione przez ten współczynnik. Obliczenie iloczynu wartości GDC i szacowanej ilości wzorców.
2. Wybór największego z iloczynów obliczonych w poprzednim kroku.
3. Oznaczenie wzorców, które mogą zostać rozróżnione przez współczynnik wybrany w kroku 2.
4. Jeżeli pozostały niedyskryminowane wzorce, to powrót do kroku 2. W przeciwnym wypadku - koniec.

Zastosowanie powyższego algorytmu gwarantuje zdolność sieci SONN do bezbłędnej klasyfikacji wzorców uczących. Z drugiej strony, poszukiwanie cech o największych zdolnościach określających przynależność do klas korzystnie wpływa na właściwości generalizacji w ten sposób zbudowanego modelu.

Neurony ARN w podsieci ARS są łączone z dowolną liczbą neuronów LBN oraz pojedynczym - jeżeli taki istnieje - neuronem ARN z poprzedniej warstwy. Celem tworzenia

połączeń jest wzmacnianie (*ang. reinforcement*) zdolności do określania przynależności do konkretnej klasy.

### **Tworzenie warstwy wyjściowej – Maximum Selection Subnetwork**

Liczba neuronów warstwy wyjściowej bezpośrednio zależy od ilości klas występujących w danym problemie klasyfikacji. Neurony w warstwie wyjściowej mają za zadanie wybrać maksymalną wartość uzyskaną od połączonych z nimi neuronami warstw wcześniejszych. Wartość wyjściowa każdego z neuronów MSN może być interpretowana jako prawdopodobieństwo, z jakim sieć zakwalifikowała klasyfikowany wzorzec do klasy związanej z danym neuronem MSN. Warto zwrócić uwagę, że wejścia neuronów warstwy wejściowej nie są powiązane z żadną wagą, wszystkie niezbędne obliczenia dokonują się w warstwach ukrytych i warstwie wejściowej. Najwyższa wygenerowana wartość wyjściowa spośród neuronów MSN jest traktowana jako wartość wyjściowa sieci SONN, co obrazuje równanie 3.10.

### **3.5.5. Szczegóły implementacyjne**

Implementacja algorytmu SONN, która jest elementem biblioteki algorytmów, bazuje na doświadczalnym kodzie algorytmu stworzonym przez dr Adriana Horzyka. Kod został przepisany na platformę .NET (C#) i zawiera kilka usprawnień w stosunku do wersji oryginalnej.

Metoda SONN nie jest parametryzowana.

## 4. Metodologia porównań metod klasyfikacji

W niniejszym rozdziale omówiona zostanie metodologia dokonywania porównań algorytmów klasyfikacji. Przedstawiony zostanie proces badawczy prowadzący do otrzymania wyników oraz zostaną omówione zagadnienia związane z tym procesem. Zaprezentowane zostaną wyselekcjonowane zbiory danych, dla których sprawdzone zostanie działanie klasyfikatorów. Szczególny nacisk położony zostanie na kryteria porównawcze, wskaźniki stosowane do porównań oraz mechanizmy walidacji. Zostaną także omówione sposoby wstępnego przetwarzania danych, stosowane nie tylko przy porównywaniu algorytmów, ale też w innych dziedzinach.

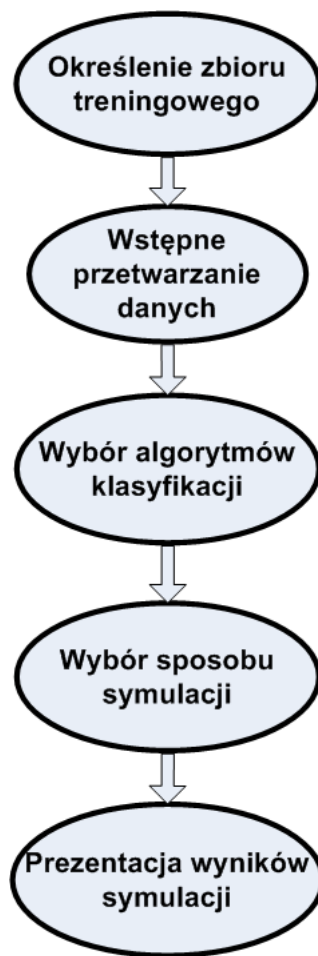
Motywacją do napisania tego rozdziału w obecnym kształcie był fakt istnienia znacznej ilości publikacji poruszających zagadnienie klasyfikacji i oferujących porównanie poszczególnych metod, przy braku opisu metodologii zastosowanej do przeprowadzenia tych porównań.

Wreszcie rozdział ten może stanowić podsumowanie istniejącego i stosowanego w praktyce zestawu koncepcji i narzędzi związanego ze stosowaniem i porównywaniem metod klasyfikacji.

### 4.1. Proces porównywania

Dokonując porównań metod klasyfikacji użytkownicy muszą wykonać szereg czynności, które wymagane są do uzyskania finalnych wyników. Charakterystyka algorytmów klasyfikacji sprawiła, że opisany poniżej proces (patrz rysunek 4.1) jest standardowy i występuje w takim lub podobnym kształcie w wielu aplikacjach związanych z eksploracją danych i metodami klasyfikacji [29, 30]. Kolejne etapy tego procesu to:

1. Określenie zbioru treningowego – celem dokonywania porównań jest wybranie najlepszego klasyfikatora, który zgodnie z przyjętymi kryteriami jakości, pozwala na realizację zagadnienia klasyfikacji dla konkretnego problemu. Przez problem rozumiemy zestaw danych (zbiór danych) podlegających klasyfikacji. Zbiory danych, dla których przeprowadzone zostaną badania opisane będą w podrozdziale 4.2.
2. Wstępne przetwarzanie danych – wybrany zbiór danych może pochodzić z różnych źródeł. Często przed przystąpieniem do treningu klasyfikatora przeprowadzane są czyn-



Rysunek 4.1: Schemat procesu porównywania metod klasyfikacji

ności pozwalające na transformację zbioru danych do innej postaci, co często pozwala na uzyskanie lepszych wyników klasyfikacji. Metody przetwarzania danych zaimplementowane w ramach pracy zostaną zaprezentowane w podrozdziale 4.3.

3. Wybór algorytmów klasyfikacji – w tym etapie wybierane są algorytmy, które mają zostać sprawdzone na wybranym i przetworzonym wcześniej zbiorze danych. Istotna jest możliwość sprawdzenia konkretnego algorytmu w różnych wariantach, jeżeli takie występują oraz przy zastosowaniu różnej konfiguracji parametrów występujących w algorytmie. Algorytmy klasyfikacji, które znalazły się w bibliotece zaimplementowanej w ramach pracy opisane zostały w rozdziale 3.
4. Wybór sposobu symulacji – z porównywaniem metod klasyfikacji związane jest zagadnienie walidacji, dzięki któremu określić można, jakie wyniki osiągane są na znanym zbiorze treningowym, a jakich można spodziewać się po wprowadzeniu danych nieznanymi podczas uczenia. Stosowane sposoby walidacji opisane zostały w podrozdziale 4.4.

5. Prezentacja wyników – ostatnim krokiem omawianego procesu jest prezentacja wyników otrzymanych po zakończeniu symulacji. Poprzez wyniki rozumiane są statystyki algorytmów oraz kryteria jakości, które opisane zostały w podrozdziale 4.5.

## 4.2. Zbiory danych

Kluczowym elementem przy dokonywaniu ogólnego porównania metod klasyfikacji jest wybór zbiorów danych. Wybierając zbiory danych do niniejszej pracy brane były pod uwagę następujące kryteria:

1. Rozpoznawalność – istnieje wiele publikacji przedstawiających porównanie metod klasyfikacji. Bardzo często takie porównania przeprowadza się na tych samych zbiorach, co czyni je rozpoznawalnymi i dzięki temu możliwa jest konfrontacja konkretnego algorytmu z innymi, opisanymi w innych pracach. Do takich często stosowanych zbiorów niewątpliwie należą Wine i Iris, które wykorzystane zostaną również w niniejszej pracy.
2. Zróżnicowanie – w celu uzyskania rzetelnych wyników koniecznym jest zróżnicowanie zestawu stosowanych zbiorów. W związku z tym wybrane zostały zbiory charakteryzujące się:
  - (a) Różną wielkością – wielkość zbioru wpływa na czas treningu i klasyfikacji. Wśród stosowanych zbiorów znajdują się zarówno zbiory małe (Iris, wielkość wolumenu = 600), jak i zbiory stosunkowo duże (Arrhythmia, wielkość wolumenu = 126108).
  - (b) Różną liczbą cech – liczba cech dla wielu algorytmów jest istotnym parametrem wpływającym na strukturę klasyfikatora, dlatego wybrane zostały zbiory opisane przez kilka (Iris, liczba cech = 4) oraz kilkaset cech (Arrhythmia, liczba cech = 279).
  - (c) Różnym rozkładem ilości reprezentantów każdej z klas – wśród stosowanych zbiorów znajdują się takie, które zawierają równomierny podział reprezentantów każdej z klas (Iris), jak i te, dla których rozkład ten jest nierównomierny (np. Yeast).
  - (d) Rodzajem danych – wśród wybranych zbiorów większość opisana jest przez dane rzeczywiste, ale znajdują się też dane w formacie binarnym bipolarnym<sup>1</sup> (Congressional Voting Record).
  - (e) Dziedziną pochodzenia – w celu zaprezentowania obszarów wykorzystywania metod klasyfikacji wybrane zostały zbiory związane z różnymi dziedzinami. Wśród

---

<sup>1</sup>W rzeczywistości są to dane trójwartościowe, w których występują także wartości nieokreślone (1 - głos na tak, -1 - głos na nie oraz 0 - wstrzymanie się od głosu)



tych dziedzin znajduje się medycyna (Arrhythmia), biochemia (Wine, Yeast), botanika (Iris), technika (Ionosphere), a nawet polityka (Congressional Voting Records).

3. Powszechność – realizowane porównania algorytmów mogą być zestawiane z podobnymi próbami. Aby było to możliwe konieczne jest stosowanie publicznie dostępnych zestawów danych. Zbiory danych stosowane w niniejszej pracy pochodzą z popularnego i ogólnodostępnego repozytorium danych [31].

Wybrane zbiory danych zostały zaprezentowane wraz z informacjami interesującymi z punktu widzenia klasyfikacji (liczba cech, liczba wektorów). Do opisu dodany został też wskaźnik „wielkość wolumenu” obrazujący wielkość zbioru danych (iloczyn liczby cech i liczby wektorów).

Tablica 4.1: Charakterystyka zbioru danych Iris

<b>Nazwa</b>	Iris
<b>Liczba cech</b>	4
<b>Liczba wektorów treningowych</b>	150
<b>Wielkość wolumenu</b>	600
<b>Rozkład danych wg. klas</b>	Iris-Setos = 50, Iris-Virginica = 50, Iris-Versicolor = 50
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Iris">http://archive.ics.uci.edu/ml/datasets/Iris</a>
<b>Opis</b>	Zbiór ten zawiera informacje o trzech gatunkach kwiatów z rodzaju kosaćców (irysów). Popularny zbiór danych występujący w wielu publikacjach dotyczących zagadnienia klasyfikacji. Cechuje się niewielką liczbą danych i równomiernym rozkładem klas.

Tablica 4.2: Charakterystyka zbioru danych Wine

<b>Nazwa</b>	Wine
<b>Liczba cech</b>	13
<b>Liczba wektorów treningowych</b>	178
<b>Wielkość wolumenu</b>	2314
<b>Rozkład danych wg. klas</b>	1 = 59, 2 = 71, 3 = 48
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Wine">http://archive.ics.uci.edu/ml/datasets/Wine</a>
<b>Opis</b>	Zbiór danych zawierający wybranych 13 parametrów zebranych na podstawie analizy chemicznej trzech 3 rodzajów win produkowanych w pewnym regionie Włoch. Jest to popularny, często występujący w publikacjach zbiór danych. Jest uznawany za nieskomplikowany.

Tablica 4.3: Charakterystyka zbioru danych Yeast

<b>Nazwa</b>	Yeast
<b>Liczba cech</b>	8
<b>Liczba wektorów treningowych</b>	1484
<b>Wielkość wolumenu</b>	11872
<b>Rozkład danych wg. klas</b>	CYT = 463, NUC = 429, MIT = 244, ME3 = 163, ME2 = 51, ME1 = 44, EXC = 37, VAC = 30, POX = 20, ERL = 5
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Yeast">http://archive.ics.uci.edu/ml/datasets/Yeast</a>
<b>Opis</b>	Zbiór zawiera informacje o lokalizacji protein w analizowanych komórkach drożdży. Zbiór ten został wybrany ze względu na nierównomierny rozkład klas w zbiorze danych.

Tablica 4.4: Charakterystyka zbioru danych Congressional Voting Records

<b>Nazwa</b>	Congressional Voting Records
<b>Liczba cech</b>	16
<b>Liczba wektorów treningowych</b>	435
<b>Wielkość wolumenu</b>	6960
<b>Rozkład danych wg. klas</b>	democrate = 267, republican = 168
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Congressional+Voting+Records">http://archive.ics.uci.edu/ml/datasets/Congressional+Voting+Records</a>
<b>Opis</b>	Zbiór danych zawierający wyniki 16 głosowań kongresu Stanów Zjednoczonych Ameryki Północnej z roku 1984. Zbiór też został wybrany, ponieważ składa wyłącznie z danych o wartościach: +1, 0, -1.

Tablica 4.5: Charakterystyka zbioru danych Car Evaluation

<b>Nazwa</b>	Car Evaluation
<b>Liczba cech</b>	6
<b>Liczba wektorów treningowych</b>	1728
<b>Wielkość wolumenu</b>	10368
<b>Rozkład danych wg. klas</b>	unacc = 1210, acc = 384, good = 69, vgood = 65
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Car+Evaluation">http://archive.ics.uci.edu/ml/datasets/Car+Evaluation</a>
<b>Opis</b>	Zbiór zawierający dane charakteryzujące przykładowe samochody. Klasa w tym zbiorze określa jeden z czterech stopni akceptowalności danego modelu.
<b>Uwagi</b>	Tekstowe dane opisowe zostały zamienione na postać liczbową.

Tablica 4.6: Charakterystyka zbioru danych Ionosphere

<b>Nazwa</b>	Ionosphere
<b>Liczba cech</b>	34
<b>Liczba wektorów treningowych</b>	351
<b>Wielkość wolumenu</b>	11934
<b>Rozkład danych wg. klas</b>	g = 225, b = 126
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Ionosphere">http://archive.ics.uci.edu/ml/datasets/Ionosphere</a>
<b>Opis</b>	Zbiór danych zawierający wyniki testowych pomiarów jakości działania radarów. Każdy z pomiarów może zostać sklasyfikowany do jednej z dwóch klas (good, bad). Zbiór ten jest przydatny do testowania szczególnego przypadku – binarnego zagadnienia klasyfikacji.

Tablica 4.7: Charakterystyka zbioru danych Arrhythmia

<b>Nazwa</b>	Arrhythmia
<b>Liczba cech</b>	279
<b>Liczba wektorów treningowych</b>	452
<b>Wielkość wolumenu</b>	126108
<b>Rozkład danych wg. klas</b>	01 = 245, 02 = 44, 03 = 15, 04 = 15, 05 = 13, 06 = 25, 07 = 3, 08 = 2, 09 = 9, 10 = 50, 14 = 4, 15 = 5, 16 = 22
<b>Źródło</b>	<a href="http://archive.ics.uci.edu/ml/datasets/Arrhythmia">http://archive.ics.uci.edu/ml/datasets/Arrhythmia</a>
<b>Opis</b>	Zbiór zawiera dane osobnicze pacjentów (wiek, płeć, wzrost, waga) oraz dane pobrane z aparatury medycznej. Klasami w tym zbiorze są grupy określające występowanie bądź brak arytmii w pracy serca. Zbiór został wybrany ze względu na dużą liczbę cech, z których większość może być zredukowana, co stanowi dobre sprawdzenie zdolności do redukcji zbędnych wymiarów przestrzeni cech.
<b>Uwagi</b>	Oryginalny zbiór zawiera brakujące wartości dla niektórych cech, które zostały uzupełnione wartościami średnimi dla danej cechy.

## 4.3. Wstępne przetwarzanie danych

Przetwarzanie danych jest niezmiernie istotnym zagadnieniem towarzyszącym procesowi klasyfikacji, jak również innym metodom eksploracji danych. Celem stosowania wstępnego przetwarzania danych w zagadnieniu klasyfikacji jest taka transformacja zbioru danych, której zastosowanie utworzy nowy zbiór, dla którego algorytmy klasyfikacji rozwiążą problem w krótszym czasie lub z mniejszym błędem. Spośród wielu sposobów transformacji danych omówione zostaną te, które zostały zaimplementowane w ramach pracy.

### 4.3.1. Normalizacja

Jest jedną z najczęściej stosowanych metod wstępnego przetwarzania danych. Działanie tej metody polega na takim przekształceniu wartości wybranej cechy, aby wartości po transformacji zawierały się w przedziale  $[min, max]$ , gdzie najczęściej stosowanymi wartościami są  $min = 0$  oraz  $max = 1$ . Zależność między danymi przed i po normalizacji określona jest wzorem 4.1.

$$x'_i = min + \frac{(max - min)(x_i - x_{min})}{x_{max} - x_{min}}, \text{ gdzie} \quad (4.1)$$

$x'_i$  - wartość cechy dla  $i$ -tego wektora po normalizacji

$x_i$  - wartość cechy dla  $i$ -tego wektora przed normalizacją

$x_{min}$  - minimalna wartość cechy

$x_{max}$  - maksymalna wartość cechy

$max$  - oczekiwana wartość maksymalna po normalizacji

$min$  - oczekiwana wartość minimalna po normalizacji

### 4.3.2. Standaryzacja

Kolejnym sposobem transformacji danych zaimplementowanym w ramach pracy jest standaryzacja. Standaryzacja pozwala nam na uzyskanie zbioru danych, w którym wartości dla danej cechy mają oczekiwaną wartość średnią równą zero oraz odchylenie standardowe równe jedności. W ramach pracy zaimplementowana została metoda standaryzacji  $Z$ , w której zależność między danymi przed i po przetworzeniu wyraża się wzorem 4.2.

$$x'_i = \frac{x_i - \mu}{\sigma} \quad (4.2)$$

$x'_i$  - wartość cechy dla  $i$ -tego wektora po normalizacji

$x_i$  - wartość cechy dla  $i$ -tego wektora przed normalizacją

$\mu$  - średnia wartość cechy przed standaryzacją

$\sigma$  - odchylenie standardowe cechy przed standaryzacją

### 4.3.3. Principal Component Analysis

Zastosowanie wyżej wymienionych metod wstępnego przetwarzania danych może mieć wpływ na jakość uogólnienia klasyfikatorów. Metodę *Principal Component Analysis* (PCA, Analiza głównych składowych, *Proper Orthogonal Decomposition*, czasem nazywana także dyskretną transformacją Karhunen-Loève)[32] można z kolei stosować w celu skrócenia czasu treningu i klasyfikacji algorytmu poprzez redukcję wymiaru przestrzeni cech, przy zachowaniu najistotniejszych informacji zawartych w zbiorze danych.

Algorytm *Principal Component Analysis* należy do grupy metod analizy statystycznej. Ogólne zadanie metody PCA można sformułować jako wyznaczenie obrotu układu współrzędnych tak, aby pierwsze jego współrzędne związane były ze zmiennymi, które niosą najwięcej informacji o danych w zbiorze. Dzięki tej operacji można wybrać składowe główne (zmienne niosące najwięcej informacji) i usunąć pozostałe, dzięki czemu dokonuje się redukcja wymiaru problemu i skraca czas obliczeń.

Algorytm PCA składa się z następujących kroków:

1. Transformacja wstępna. Taka transformacja polega na odjęciu średniej wartości każdej cechy od każdej pozycji dla danej cechy. Dzięki tej operacji uzyskujemy zbiór danych, którego średnia wynosi zero.
2. Wyznaczenie macierzy kowariancji dla zbioru uzyskanego po wstępnej transformacji.
3. Obliczenie wartości własnych macierzy kowariancji oraz ich posortowanie w porządku malejącym.
4. Wybór składowych głównych. Spośród wartości własnych wybierane są te, które uznane zostały za składowe główne.
5. Utworzenie macierzy składającej się z wektorów własnych macierzy kowariancji dla wybranych wartości własnych.
6. Transformacja względem składowych głównych. Mając wybrane wektory własne, odpowiadające składowym niosącym najwięcej informacji, możliwa jest transformacja oryginalnego zbioru danych do postaci, w której pierwszymi składowymi są składowe dla największych znalezionych wartości własnych macierzy kowariancji.

## 4.4. Metody walidacji

Klasyfikatory są tworzone do rozwiązywania realnych problemów. Celem osób stosujących metody klasyfikacji w praktyce jest uzyskanie jak najlepszego mechanizmu dla danego pro-

blemu, w związku z tym pojawia się konieczność walidowania algorytmów w odniesieniu do konkretnego zbioru danych.

Naturalną sytuacją jest, że badacze próbujący zbudować mechanizm klasyfikacji, dysponują zbiorem danych, który reprezentuje część możliwych przypadków występujących w obrębie badanej dziedziny. Dla ustalenia uwagi przytoczę hipotetyczną sytuację, w której grupa lekarzy współpracująca ze specjalistami z dziedziny eksploracji danych próbuje zbudować mechanizm automatyzujący diagnozę raka płuc. Lekarze dysponują zbiorem przebadanych dotychczas przypadków, a celem tworzonego mechanizmu jest jak najdokładniejsze diagnozowanie przypadków jeszcze nieznanych. Pojawia się problem „jak na podstawie istniejących danych określić, który klasyfikator będzie najlepszy”. W praktyce stosuje się kilka prostych metod, które zostaną opisane poniżej. Większość z nich wywodzi się ze statystycznej metody sprawdzianu krzyżowego (walidacji krzyżowej). Algorytmy te zostaną udostępnione w ramach aplikacji wspierającej porównywanie sieci SONN z wybranymi metodami klasyfikacji.

Tworząc podzbiory treningowe warto zwrócić uwagę na rozkład wzorców według klas w zbiorze oryginalnym i podzbiorze testowym. Najczęściej dąży się do zachowania proporcji między przedstawicielami kolejnych klas. Takie podejście nosi nazwę walidacji stratyfikowanej (*ang. stratified*) i zostało uwzględnione podczas implementacji aplikacji wspierającej proces porównania metod klasyfikacji.

Opisując metody walidacji podawane przykłady dotyczyć będą dziesięcioelementowego zbioru przedstawionego na rysunku 4.2. W przykładach kolorem żółtym oznaczone zostaną wzorce ze zbioru uczącego, natomiast próbki ze zbioru walidacyjnego oznaczone są kolorem zielonym.

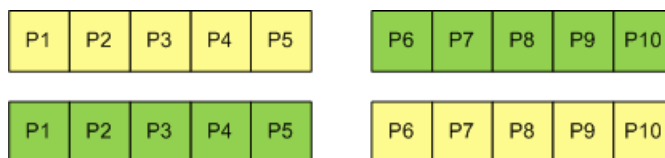
P1	P2	P3	P4	P5	P6	P7	P8	P9	P10
----	----	----	----	----	----	----	----	----	-----

Rysunek 4.2: Metody walidacji, zbiór danych

#### 4.4.1. k-fold cross-validation

Metoda *k-fold cross-validation* jest bardzo często stosowaną metodą walidacji, pozwalającą na uzyskanie informacji o możliwościach uogólnienia modelu, przy niskim nakładzie obliczeniowym. W metodzie *k-fold cross-validation* zbiór treningowy dzielony jest na  $k$  podzbiorów tak, że do walidacji przeznaczona jest każda z próbek tylko raz. Na rysunkach 4.3 oraz 4.4 przedstawiono uzyskane podziały dzięki tej metodzie dla różnych parametrów.





Rysunek 4.3: Metody walidacji, k-fold cross validation (k=2)



Rysunek 4.4: Metody walidacji, k-fold cross validation (k=5)

#### 4.4.2. Leave-one-out cross-validation

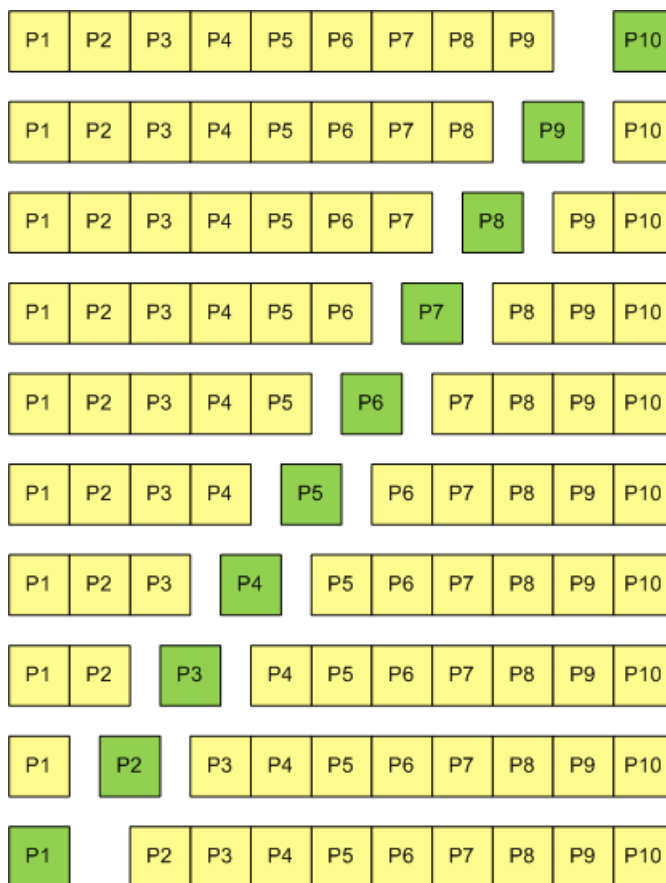
Szczególnym przypadkiem *k-fold cross-validation*, który często występuje jako osobna metoda walidacji jest *leave-one-out cross-validation* (LOOCV). Metodę tą można opisać jako  $n$ -krotną (gdzie  $n$  to ilość próbek danych w zbiorze uczącym) walidację dla każdej próbki, gdzie do uczenia modelu zostały wybrane pozostałe wzorce. Zatem jest to przypadek *k-fold cross-validation*, dla parametru  $k$  równego ilości próbek w zbiorze. Przykład podziału zbioru uczącego na podzbiory przedstawiony został na rysunku 4.5.

#### 4.4.3. Podział procentowy

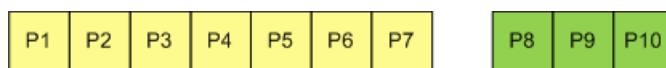
Często stosowaną metodą walidacji jest podział zbioru treningowego na dwie części: treningową i walidacyjną, przy czym część treningowa zawiera najczęściej więcej elementów (popularne proporcje to: 90 : 10, 75 : 25). Można porównać tą metodę do wykorzystania jednego z podziałów uzyskanych metodą *k-fold cross-validation*. Na poniższym przykładzie (patrz rysunek 4.6) przedstawiony został podział 70 : 30.

#### 4.4.4. Test na danych treningowych

Kolejnym testem jest sprawdzenie zachowania klasyfikatora na wszystkich dostępnych próbkach ze zbioru treningowego. Dzięki tej metodzie można stwierdzić, czy klasyfikator potrafi nauczyć się poprawnie rozpoznawać wszystkie znane przypadki treningowe.



Rysunek 4.5: Metody walidacji, LOOCV



Rysunek 4.6: Metody walidacji, Podział procentowy (70:30)

## 4.5. Kryteria porównania

Porównując metody klasyfikacji konieczne jest określenie, jakie kryteria zostaną zastosowane przy dokonywaniu porównań. Tematem pracy jest porównanie jakości uogólnienia i efektywności działania, dlatego te kryteria zostaną omówione szczegółowo. Konieczne jest analizowanie tych dwóch kryteriów jednocześnie, ponieważ można sobie wyobrazić algorytmy, które osiągają doskonałe wyniki przy ogromnym zużyciu zasobów lub zużywają minimalne ilości zasobów, ale fatalnie radzą sobie z zadaniem klasyfikacji. Opisane zostaną także inne kryteria, które mogą mieć wpływ na wybór metody klasyfikacji.

W praktyce, podczas porównań interesujący jest klasyfikator osiągający najlepsze rezultaty w kontekście stosowanej miary. Dla celów badawczych w tworzonej aplikacji oraz niniejszej pracy podawane będą także wartości średnie oraz maksymalne i minimalne wartości miar uzyskane dla danego algorytmu oraz różnych sposobów walidacji.

### 4.5.1. Jakość uogólnienia

Poprzez jakość uogólnienia (generalizacji) rozumiana jest skuteczność metody podczas rozpoznawania wzorców wybranych zgodnie z określoną metodą walidacji (patrz podrozdział 4.4). O uogólnieniu możemy mówić, gdy dana metoda nauczy się rozpoznawać klasę wzorca wykorzystując pewne podobieństwa do innych wzorców, które były obecne w procesie uczenia, przy jak najmniejszym stopniu bezpośredniego porównywania wzorców (przeuczenie, *ang. overfitting*).

Aby zobrazować tę różnicę można posłużyć się popularną analogią do nauki ze zrozumieniem, która pozwala na radzenie sobie z problemami wcześniej nieznanymi oraz nauki „na pamięć”, która może sprawdzać się w problemach znanych, ale najczęściej jest zupełnie nieskuteczna w konfrontacji z nowym zagadnieniem. Istnieje wiele miar stosowanych w procesie walidacji, których wybór zależy od obszaru zastosowania klasyfikatora. Opisywane w tej pracy miary jakości uogólnienia podzielić można na dwie kategorie:

- Miary bazujące na poprawności określenia przynależności do klasy – współczynnik jakości określany jest wyłącznie na podstawie tego, czy klasyfikator poprawnie określił przynależność do klasy. Do takich miar należą  $Q\_COR$  (współczynnik poprawności),  $Q\_ERR$  (współczynnik błędu).
- Miary bazujące na pewności<sup>2</sup> określenia przynależności do klasy – miary te analizują współczynnik prawdopodobieństwa, z jakim klasyfikator określił przynależność do konkretnej klasy. Do takich miar należą:  $Q\_SSE$  (*sum squared error*, sumaryczny błąd kwadratowy),  $Q\_MSE$  (*mean squared error*, średni błąd kwadratowy),  $Q\_RMSE$  (*root mean squared error*, pierwiastek średniego błędu kwadratowego),  $Q\_ARV$  (*average relative variance*, średnia wariancja względna).

Miary należące do grupy miar jakości uogólnienia zostaną oznaczone przedrostkiem „Q”, w celu łatwiejszego rozróżnienia podczas analizy wyników porównań.

#### Współczynnik poprawności, $Q\_COR$

$$Q\_COR = \frac{\text{ilość wzorców sklasyfikowanych poprawnie}}{\text{ilość wzorców w zbiorze danych}}$$

lub wartość procentową:

$$Q\_COR = \frac{\text{ilość wzorców sklasyfikowanych poprawnie}}{\text{ilość wzorców w zbiorze danych}} * 100\%$$

<sup>2</sup>Pewność określenia przynależności czasem bywa nazywana prawdopodobieństwem. Warto zwrócić uwagę, że to prawdopodobieństwo może nie spełniać aksjomatycznej definicji Kołmogorowa, a jedynie oznaczać potocznie rozumianą szansę tego, że badany wzorec należy do pewnej klasy.

**Współczynnik błędów, Q\_ERR**

$$Q\_ERR = \frac{\text{ilość wzorców sklasyfikowanych błędnie}}{\text{ilość wzorców w zbiorze danych}} * 100\%$$

**Sumaryczny błąd kwadratowy, Q\_SSE**

$$Q\_SSE = \sum_{i=1}^n (f(x_i) - y_i)^2$$

**Średni błąd kwadratowy, Q\_MSE**

$$Q\_MSE = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 = \frac{1}{n} Q\_SSE$$

**Pierwiastek średniego błędów kwadratowych, Q\_RMSE**

$$Q\_RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2} = \sqrt{Q\_MSE}$$

**Średnia wariancja względna, Q\_ARV**

$$Q\_ARV = \frac{\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2}{\frac{1}{n} \sum_{i=1}^n \left( \frac{1}{n} \sum_{i=1}^n y_i - y_i \right)^2}$$

**Macierz niepewności**

Macierz niepewności (*ang. confusion matrix*) jest sposobem wizualizacji wyników działania klasyfikatora (lub innych metod eksploracji danych). Dzięki tej metodzie można uzyskać przejrzystą interpretację ile wzorców zostało zakwalifikowanych do poszczególnych klas oraz do jakiej klasy powinny być zostać zaklasyfikowane.

Na rysunku 4.7 przedstawiona została macierz niepewności wygenerowana w aplikacji zrealizowanej w części praktycznej (odnośnik). Ten konkretny przykład dotyczy jednego etapu walidacji krzyżowej algorytmu kNN na zbiorze danych Wine. Kolumny macierzy reprezentują wynik klasyfikacji (określenie przynależności do klasy), natomiast wiersze macierzy reprezentują oczekiwany wynik klasyfikacji (zakładając pełną poprawność). *Confusion matrix* dostarcza wygodnej wizualizacji, która dla tego przykładu pozwala łatwo stwierdzić, że dwa wzorce, które powinny zostać skojarzone z klasą „2” zostały niepoprawnie sklasyfikowane. Przy całkowitej skuteczności algorytmu macierz niepewności jest macierzą diagonalną.

		Sklassyfikowano jako		
		1	2	3
Oczekiwana klasa	1	6	0	0
	2	1	5	1
	3	0	0	5

Rysunek 4.7: Macierz niepewności dla jednego etapu walidacji krzyżowej algorytmu kNN na zbiorze Wine

### 4.5.2. Efektywność działania

Efektywność działania to ogólna nazwa grupy kryteriów, które związane są z czasem i wykorzystaniem zasobów przez mechanizmy klasyfikacji. Kryteria stosowane w porównaniach w ramach pracy zostaną wyjaśnione w kolejnych sekcjach.

Miary należące do grupy miar efektywności działania zostaną oznaczone przedrostkiem „E\_”, w celu łatwiejszej analizy wyników porównań.

#### Czas treningu, E\_TTR

Kryterium efektywności działania najczęściej utożsamiane jest z czasem treningu klasyfikatora. Czas treningu określa jak szybko algorytm potrafi odnaleźć zależności w zbiorze danych niezbędne do przeprowadzenia klasyfikacji. Dąży się do minimalizacji tego wskaźnika (w pracy został oznaczony symbolem E\_TTR).

Podczas doświadczeń przeprowadzanych w ramach pracy zauważono, że takie zdefiniowanie kryterium efektywności działania bywa niewystarczające.

#### Czas klasyfikacji, E\_TCL

Mając na uwadze dwuetapowość zadania klasyfikacji (rysunek 2.1) konieczna jest również analiza czasu klasyfikacji (wskaźnik oznaczony symbolem E\_TCL). Okazuje się, że dla niektórych algorytmów (np. kNN, co można zaobserwować analizując wyniki porównań metod klasyfikacji) etap treningu jest mniej skomplikowany obliczeniowo od części klasyfikacyjnej. Z drugiej strony należy pamiętać o najczęstszym sposobie wykorzystywania metod klasyfikacji, które polegają na jednorazowym wytrenowaniu algorytmu (często w warunkach, gdzie czas treningu nie ma dużego znaczenia). Mając wytrenowany klasyfikator jest on wielokrotnie używany do klasyfikacji wzorców, co wskazuje na bardziej istotną wartość czasu klasyfikacji (taki przypadek będzie miał miejsce w systemach czasu rzeczywistego realizujących wiele operacji klasyfikacji dla wytrenowanego klasyfikatora).

Chcąc obiektywnie porównać metody klasyfikacji, należy brać pod uwagę zarówno czas klasyfikacji, jak i treningu. Dodatkowo wspomnianą zaletą takiego podejścia jest możliwość skonfrontowania dowolnego z czasów z realnymi potrzebami uwarunkowanymi przez obszar zastosowań.

Naturalnie najlepszym klasyfikatorem z punktu widzenia efektywności działania będzie ten, który osiągnie najmniejsze czasy treningu i klasyfikacji. Przeprowadzając wiele symulacji warto także analizować średnie czasy dla każdego problemu i klasyfikatora, co zostało również uwzględnione w implementowanej w ramach pracy aplikacji.

### Czas całkowity operacji, $E_{TOT}$

Miarą agregującą wartości czasu treningu i czasów klasyfikacji jest całkowity czas operacji ( $E_{TOT}$ ), pokazujący jak dużo czasu zajęło danej metodzie utworzenie klasyfikatora i wykonanie wszystkich przewidzianych operacji klasyfikacji. Sposób jego obliczania obrazuje wzór 4.3.

$$E_{TOT} = E_{TTR} + \sum_i E_{TCL_i} \quad (4.3)$$

### Wykorzystanie pamięci

Często pomijanym kryterium, które ma szczególne znaczenia w systemach wbudowanych (*ang. embedded*) jest ilość pamięci wykorzystywanej przez klasyfikator. Możemy rozróżnić dwie miary dla tego kryterium:

- **Złożoność pamięciowa** – zależność wymaganej pamięci RAM używanej w procesie treningu klasyfikatora od wielkości zbioru treningowego. Ze względu na zastosowanie technologii wykorzystującej automatyczne zwalnianie pamięci<sup>3</sup> (*ang. garbage collecting*) niemożliwy był pomiar tej wielkości, przez co nie znajdzie się ona w wynikach porównań.
- **Objętość pamięciowa wytrenowanego klasyfikatora** – realna wielkość zajmowanego miejsca w pamięci dla wytrenowanego klasyfikatora. Oznacza ona jak wiele miejsca w pamięci potrzeba do przechowywania obiektu utworzonego klasyfikatora. Pomiar tego wskaźnika może być przeprowadzony poprzez serializację<sup>4</sup> do formatu binarnego obiektu klasyfikatora. Taka wielkość jest udostępniana w symulacjach oraz występuje w porównaniach, jest ona oznaczona symbolem  $E_{MEM}$ .

<sup>3</sup>Więcej informacji o mechanizmie *garbage collecting* w środowisku .NET znaleźć można na stronach MSDN (<http://msdn.microsoft.com/en-us/magazine/bb985010.aspx>, <http://msdn.microsoft.com/en-us/library/0xy59wtx.aspx>)

<sup>4</sup>Więcej informacji o serializacji w środowisku .NET znaleźć można na stronach MSDN (<http://msdn.microsoft.com/en-us/library/system.runtime.serialization.aspx>)

### 4.5.3. Skalowalność

Kryterium skalowalności (*ang. scalability*) dla metod klasyfikacji określa czy konkretny algorytm może zostać zastosowany dla dowolnie dużego zbioru danych. Jeżeli algorytm można zastosować dla dużego zbioru danych i osiąga on dobre wyniki, to o takim algorytmie mówimy, że jest dobrze skalowalny lub, że dobrze się skaluje. W praktyce, ciężko określić miarę, która w sposób obiektywny potrafiła porządkować klasyfikatory pod kątem skalowalności. Z tego powodu skalowalność jest rzadko wymienianym kryterium podczas porównywania metod klasyfikacji.

W rzeczywistości skalowalność ściśle związana jest z kryteriami efektywności działania (rozdział 4.5.2) oraz jakości uogólnienia (rozdział 4.5.1). Ponieważ celem niniejszej pracy jest wykonanie porównań uwzględniających zbiory o różnej wielkości, w rozdziale 6.9 przedstawione zostaną zależności czasu operacji oraz potrzebnej pamięci, w funkcji stopnia trudności zbioru danych. Należy zwrócić uwagę, że taka analiza może okazać się niemiarodajna, ponieważ nie istnieje ogólny sposób wyznaczania stopnia trudności zbioru danych. Ilość elementów w zbiorze danych nie jest jedynym czynnikiem wpływającym na efektywność działania oraz jakość uogólnienia. Bardzo często kluczową rolę odgrywają zależności pomiędzy elementami zbioru danych, stopień ich skorelowania, ilość oraz rozkład klas, a także wiele innych czynników. Należy o tym pamiętać analizując takie zestawienie.

### 4.5.4. Interpretowalność

Interpretowalność jest kryterium, które bardziej związane jest z podejściem do sposobu klasyfikacji, niż z konkretnym klasyfikatorem. To kryterium porządkuje klasyfikatory według możliwości zrozumienia sposobu klasyfikacji na podstawie struktury zbudowanego klasyfikatora. Interpretowalność jest kryterium subiektywnym i niemierzalnym, natomiast istnieją przypadki, kiedy jest istotne w procesie podejmowania decyzji o wyborze klasyfikatora. Interpretowalność związana jest z zaufaniem do klasyfikatora (*ang. trust*).

Wyobraźmy sobie sytuację, w której konieczne jest zweryfikowanie zbudowanego modelu klasyfikacji przez grupę specjalistów z konkretnej dziedziny. Prawdopodobnie okaże się, że algorytm drzew decyzyjnych, w których łatwo można prześledzić proces podejmowania decyzji o przynależności do klasy oraz sporządzić graficzną wizualizację tego procesu, okaże się łatwiejszym do zinterpretowania od metody *Radial Basis Function Networks*, w której struktura sieci niesie wiele mniej informacji o procesie decyzyjnym.

### 4.5.5. Inne kryteria

Oprócz wymienionych kryteriów porównania istnieć mogą inne. Bardzo często w praktyce spotyka się porównania algorytmów w kontekście konkretnej dziedziny (np. medycyna, genetyka, mechanika).



## 5. Część praktyczna

W tym rozdziale przedstawione zostaną istotne szczegóły związane z realizacją praktycznej części pracy. Omówione zostaną założenia, które uwarunkowały ostateczny kształt stworzonych rozwiązań. W dalszej części zaprezentowane będą wybrane metody inżynierii oprogramowania, dzięki którym udało się zakończyć projekt sukcesem oraz wymienione zostaną kluczowe technologie użyte podczas realizacji części praktycznej. W podrozdziale 5.4. szczegółowo omówione zostaną istotne aspekty implementacyjne, które mogą posłużyć podczas dalszego rozwoju aplikacji. W ostatniej części rozdziału przedstawione i omówione zostaną najważniejsze elementy aplikacji wspierającej porównania.

### 5.1. Założenia

Nadrzędnym celem realizacji części praktycznej jest wykonanie narzędzi pozwalających na dokonanie porównań rozmaitych metod klasyfikacji z sieciami SONN. Z uwagi na ogrom istniejących algorytmów nie sposób było zaimplementować wszystkich z nich, dlatego duży nacisk położony został na możliwość rozszerzenia biblioteki o nowe algorytmy. Kolejnym istotnym aspektem, który został uwzględniony podczas implementacji jest całkowite oddzielenie kodu samych algorytmów od aplikacji wspierającej porównania. Dzięki temu można łatwo użyć zestawu algorytmów w dowolnej aplikacji. Mając na uwadze konieczność podziału całego projektu na dwa moduły (bibliotek algorytmów, aplikacja wspierająca porównania), można określić dokładniej wymagania w stosunku do każdego z nich.

Biblioteka algorytmów:

- Spójność – wszystkie algorytmy powinny implementować wspólny interfejs (*IAAlgorithm*), który będzie reprezentował algorytm klasyfikacji. Dzięki temu łatwiejsze będzie rozszerzanie kodu o nowe algorytmy. Ponadto, taka konwencja sprzyja poprawie jakości oraz czytelności kodu.
- Jakość – biblioteka algorytmów może być wykorzystana w innych projektach (w szczególności w aplikacji wspierającej porównania). Z tego powodu konieczne jest

napisanie zestawu testów, które oprócz wymienionego powodu powinny być przydatne podczas dalszego rozwoju biblioteki (np. SONN-4).

- Import danych z różnych źródeł – obok zestawu algorytmów konieczne jest stworzenie mechanizmów pozwalających na importowanie z rozmaitych źródeł. W rzeczywistości dane mogą pochodzić z dowolnego źródła i występować w dowolnym formacie. Konieczne jest zaimplementowanie podstawowych mechanizmów (.csv, ADO.NET), przy jednoczesnym zachowaniu możliwości rozszerzenia tej funkcjonalności (interfejs *IDataProvider*).

Aplikacja wspierająca porównania:

- Dostępność – jak wielokrotnie wspomiano w pracy istnieje wiele porównań algorytmów klasyfikacji. Za cel postawiono udostępnienie realizowanej pracy tak, aby w przyszłości osoby realizujące podobne zadania mogły wykorzystać stworzoną aplikację lub skonfrontować swoje wyniki z wynikami, które można w niej uzyskać. Mając to na uwadze warto wykorzystać zalety popularnego medium jakim jest Internet i dlatego realizowana aplikacja została wykonana jako aplikacja internetowa.
- Prostota – istnieje kilka dojrzałych narzędzi wspomagających eksplorację danych [29, 30], które zawierają wiele algorytmów, metod przetwarzania danych, mechanizmów walidacji oraz innych opcji konfiguracyjnych. Bardzo często konsekwencją tego jest nadmierny stopień skomplikowania czynności tworzenia klasyfikatora i porównania jego wyników z wynikami innych klasyfikatorów. Dlatego celem stawianym przed aplikacją realizowaną w ramach pracy jest zapewnienie prostoty oraz intuicyjności wykonywanych czynności podczas pracy z aplikacją.
- Funkcjonalność – proces porównywania metod klasyfikacji (opisany w rozdziale 4.) składa się z kilku niezbędnych etapów, z których każdy powinien zostać zaimplementowany w aplikacji.

## 5.2. Wybrane aspekty inżynierii oprogramowania

Realizując projekt informatyczny, który ma spełniać pewne wymogi odnośnie jakości, konieczne jest stosowanie sprawdzonych metod oraz praktyk, które związane są z dziedziną inżynierii oprogramowania. Poniżej wymienionych zostało kilka aspektów związanych z inżynierią oprogramowania, które odegrały szczególną rolę w realizowanym projekcie. Wskazane zostały także narzędzia, które zostały wykorzystane podczas prac.

- Zarządzanie projektem – realizując rozmaite projekty (nie tylko informatyczne) warto dzielić pracę na mniejsze etapy, kontrolować postęp w ich realizacji, śledzić błędy i wykonywać inne czynności w ramach zarządzania zadaniami, dzięki którym całość może zakończyć się z oczekiwanym rezultatem. Narzędzia które zostało wykorzystane, aby wspierać zarządzanie projektem to pakiet usługi dostępnych w ramach portalu internetowego [assembla.com](http://assembla.com).
- Wersjonowanie kodu – chcąc zabezpieczyć się przed utratą kodu aplikacji (na skutek awarii stacji roboczej lub wprowadzeniu trudnego do wykrycia błędu) warto korzystać z narzędzi umożliwiających wersjonowanie kodu. Podczas prac nad projektem stosowane było rozwiązanie Subversion (SVN).
- Testy aplikacji – jednym ze sposobów zwiększenia jakości tworzonego kodu jest utworzenie specjalnych testów sprawdzających jego poprawność. Istnieje wiele rodzajów testów aplikacji. Te, które odegrały szczególną rolę podczas prac nad projektem, to testy jednostkowe (pisane przy wsparciu zestawu bibliotek NUnit). Dzięki testom możliwe było sprawdzenie poprawności działania algorytmów, co okazało się szczególnie istotne podczas implementacji algorytmu SONN.
- Ciągła integracja (*ang. continuous integration*) – automatyzacja procesu budowania i testowania aplikacji pomaga wychwycić błąd niezwłocznie po wprowadzeniu go do repozytorium kodu. W ramach projektu stosowano popularne narzędzie Cruise Control, które umożliwiło ciągłą integrację projektu.
- Wzorce projektowe – stosowanie sprawdzonych wzorców projektowych oraz dobrych praktyk programistycznych jest niezbędną składową projektów, gdzie określone są pewne wymagania co do jakości. Kilka z zastosowanych rozwiązań zostanie przedstawionych w dalszej części tego rozdziału.

### 5.3. Stosowane technologie

Do najważniejszych technologii wykorzystanych w ramach projektu należą:

- .NET Framework 3.5 – platforma programistyczna. Jako język programowania został wybrany nowoczesny C# 3.0 wraz z mechanizmem LINQ. Do realizacji warstwy prezentacji zastosowano ASP.NET WebForms. W oparciu o ADO.NET zrealizowano jeden z providerów danych wykorzystywanych w projekcie.
- SQLServer 2005 –silnik bazy danych służącej do przechowywania stosowanych zbiorów danych.

- jQuery, AjaxControlToolkit – biblioteki wykorzystane przy programowaniu internetowego interfejsu użytkownika.
- ZedGraph – biblioteka pozwalająca na tworzenie wykresów.
- NUnit – framework wspierający tworzenie testów jednostkowych aplikacji.
- Subversion – narzędzie wspierające wersjonowanie kodu.

## 5.4. Szczegóły implementacyjne

W tym podrozdziale omówione zostaną aspekty techniczne zrealizowanego rozwiązania, dzięki którym możliwa jest rozbudowa projektu oraz integracja z innymi aplikacjami. Przedstawiona zostanie struktura projektu oraz podstawowe interfejsy.

### 5.4.1. Struktura projektu

Zrealizowane w ramach pracy rozwiązanie składa się z 9 modułów enkapsulujących podstawowe funkcje systemu.

W ramach biblioteki z algorytmami wyróżnić można następujące moduły:

- *Algorithms* - zbiór klas reprezentujących algorytmy klasyfikacji oraz inne algorytmy wykorzystywane podczas treningu klasyfikatorów (np. algorytm analizy skupień *k-mean*). W tym module znajdują się podstawowe interfejsy, które zostały opisane w dalszej części pracy.
- *Algorithms.UnitTests* - zestaw testów jednostkowych dla klas modułu *Algorithms*.
- *Data* - zbiór klas związanych z mechanizmami dostarczania (*Data.Providers*) i przetwarzania (*Data.Preprocessors*) danych używanych przez algorytmy klasyfikacji.
- *Data.UnitTests* - zestaw testów jednostkowych dla klas modułu *Data*.

Kolejne trzy moduły zawierają logikę oraz elementy warstwy prezentacji aplikacji wspierającej porównania metod klasyfikacji:

- *Web* - główny moduł aplikacji internetowej. Zawiera logikę stron, elementy interfejsu graficznego oraz ekrany systemu pomocy.
- *Web.Controls* - moduł zawierający kod niestandardowych kontrol ASP.NET (*ang. control*) używanych w aplikacji.

- *DAL* - moduł *DAL* (*Data Access Layer*) realizuje warstwę dostępu do danych. Jako dane rozumiane są zestawy algorytmów, metod przetwarzania oraz ich parametry, które używane są w aplikacji. Logika zawarta w tej części zezwala na rozszerzanie aplikacji wspierającej porównania o nowe algorytmy i metody przetwarzania danych.

Oprócz wymienionych modułów wyróżnić można także moduły realizujące funkcje pomocnicze:

- *Utils* - zbiór często wykorzystywanych w aplikacji funkcji pomocniczych (np. operacje matematyczne).
- *Utils.UnitTests* - zestaw testów jednostkowych dla klas modułu *Utils*.

### 5.4.2. Podstawowe interfejsy

Interfejs *IAlgorithm* jest kluczowym interfejsem określonym w systemie. Dzięki niemu możliwe jest korzystanie ze skomplikowanych klas algorytmów w spójny i intuicyjny sposób, nawiązujący do dwuetapowego zadania klasyfikacji (patrz rysunek 2.1).

```
public interface IAlgorithm
{
    void Train(ITrainingData trainingData);
    IClassificationResults Classify(double[] testData);
}
```

Każdy z algorytmów klasyfikacji musi realizować dwie operacje, które określone są w ramach interfejsu *IAlgorithm*:

- *Train* – trening klasyfikatora w oparciu o znany zbiór treningowy. Ta operacja powinna wywołać logikę charakterystyczną dla danego algorytmu klasyfikacji oraz przygotować tak ten obiekt, aby możliwe było wielokrotne wywołanie metody pozwalającej na sklasyfikowanie nieznanego wektora danych.
- *Classify* – operacja klasyfikacji. Dla podanego wektora liczb rzeczywistych metoda ta powinna zwrócić obiekt klasy implementującej interfejs *IClassificationResults*, który zawiera informacje o przebiegu klasyfikacji podanego wektora wartości cech.

```
public interface IAlgorithmInfoProvider
{
    Dictionary<string, string> SpecificInformation { get; }
}
```

Interfejs *IAlgorithmInfoProvider* wymaga od klas, które go implementują definicji jednej właściwości (*ang. property*), dzięki której możliwe jest uzyskanie słownika zawierającego szczegółowe informacje o obiekcie algorytmu.

Wprowadzenie powyższych interfejsów spełnia założenia wzorca projektowego fasada (*ang. facade*)[33, 34], dzięki któremu możliwe jest ukrycie skomplikowanych szczegółów biblioteki algorytmów, pozostawiając proste interfejsy, które mogą zostać użyte w zależności od wymaganej sytuacji.

Poza opisanymi interfejsami występującymi w module *Algorithms* na uwagę zasługują także inne, których znajomość może okazać się istotna podczas rozszerzania lub wykorzystywania opisywanego kodu:

- *Algorithms.IClassificationResults* - dzięki wprowadzeniu tego interfejsu możliwe jest spójne wykorzystywanie wyników klasyfikacji zwracanych przez algorytmy.
- *Data.IDataProvider* - interfejs określający sygnatury metod, które realizują importowanie danych z rozmaitych źródeł.
- *Data.IDataPreprocessor* - podobnie jak algorytmy klasyfikacji, które muszą implementować interfejs *IAlgorithm* każda z metod wstępnego przetwarzania danych musi realizować metody wymagane przez interfejs *IDataPreprocessor*.

Naturalnie, oprócz powyższych interfejsów, biblioteka zawiera inne składowe, głównie klasy realizujące często skomplikowane algorytmy (np. SONN). Nie sposób ich jednak tutaj wszystkich opisać, dlatego przedstawione zostały tylko podstawowe interfejsy, dzięki którym możliwe jest korzystanie z biblioteki w rozmaitych aplikacjach.

## 5.5. Prezentacja aplikacji

Zgodnie z założeniami, aplikacja wspierająca porównania ma być narzędziem umożliwiającym praktyczną realizację procesu przedstawionego w rozdziale 4. Użytkownik pragnący porównać wybrane metody klasyfikacji musi wykonać pięć kroków zaprezentowanych i opisanych poniżej.

Warto zwrócić uwagę na system pomocy, który dostępny jest podczas pracy z aplikacją. Można w nim znaleźć wiele cennych informacji dotyczących kolejnych stron witryny, czy zapoznać się ze specyfiką dostępnych algorytmów klasyfikacji.

### 5.5.1. Wybór treningowego zbioru danych

Pierwszym etapem pracy z aplikacją jest określenie zbioru danych, dla którego przeprowadzone zostaną porównania. Użytkownik ma następujące możliwości wprowadzenia danych:

- Import z pliku .csv - popularna metoda przechowywania zbiorów danych. Wiele programów (np. arkusze kalkulacyjne) zawierają możliwość eksportu do formatu .csv (*ang.*

*comma separated value*), dzięki czemu można przeprowadzać klasyfikację dla danych pochodzących z takich źródeł.

- Import ze źródeł danych zgodnych z ADO.NET - możliwe jest importowanie danych ze źródeł danych wspieranych przez ADO.NET[35]. Aby było to możliwe konieczne jest podanie następujących parametrów (szczegóły w systemie pomocy aplikacji):
  1. Provider - nazwa dostawcy danych (np. *System.Data.SqlClient*).
  2. ConnectionString - łańcuch połączenia (*ang. connection string*) pozwalający na połączenie ze źródłem danych.
  3. Polecenie SELECT - kwerenda zgodna ze źródłem danych, dzięki której możliwe jest pobranie danych.
- Wybór predefiniowanego zbioru - wszystkie zbiory używane podczas porównań (patrz rozdział 4.2) mogą zostać wczytane bez konieczności podawania pliku .csv lub określania specyfikacji źródła danych ADO.NET.

Po prawidłowym wczytaniu danych użytkownik zostanie powiadomiony stosownym komunikatem. Istnieje także możliwość wyświetlenia zaimportowanego zbioru danych. W tym celu należy odznaczyć opcję „Wyświetl zbiór treningowy”.

Ekran aplikacji umożliwiający wykonanie powyższych czynności przedstawiony został na rysunku 5.1.

### 5.5.2. Wybór metod wstępnego przetwarzania danych

Po wczytaniu danych istnieje możliwość wybrania oraz określenia parametrów metod wstępnego przetwarzania danych (patrz rysunek 5.2). Dostępne możliwości to:

- Normalizacja (opisano w rozdziale 4.3.1).
- Standaryzacja (opisano w rozdziale 4.3.2).
- Principal Component Analysis (opisano w rozdziale 4.3.3).
- Mix - zmiana porządku wektorów treningowych w zbiorze.

Podobnie jak w poprzednim kroku, możliwy jest podgląd wstępnie przetworzonych danych.

The screenshot shows a software application window with a tabbed interface. The active tab is 'Zbiór Treningowy'. Below the tabs, there are three main sections for data import:

- Import z pliku:** Includes a 'Pomoc' link, a text input field, and 'Wybierz' and 'Wczytaj dane' buttons.
- Import z bazy danych:** Includes a 'Pomoc' link, a 'Provider:' dropdown menu (set to 'System.Data.SqlClient'), a 'ConnectionString:' text input field (with placeholder 'tu wpisz ConnectionString ADO.NET'), a 'SELECT:' text input field (with placeholder 'Tu wpisz kwerendę pobierającą dane'), and a 'Wczytaj dane' button.
- Dostępne zbiory danych:** Includes a 'Pomoc' link and a table listing available datasets.

At the bottom of the 'Dostępne zbiory danych' section, there is a checkbox labeled 'Wyświetl zbiór treningowy' and a label 'Aktualny zbiór treningowy:'.

	Nazwa zbioru	Informacje
Pobierz	Wine	Informacje
Pobierz	Iris	Informacje
Pobierz	Car Evaluation	Informacje
Pobierz	Ionosphere	Informacje
Pobierz	Yeast	Informacje
Pobierz	Congressional Voting	Informacje
Pobierz	Arrhythmia	Informacje

Rysunek 5.1: Aplikacja - wybór treningowego zbioru danych

### 5.5.3. Wybór oraz parametryzacja algorytmów

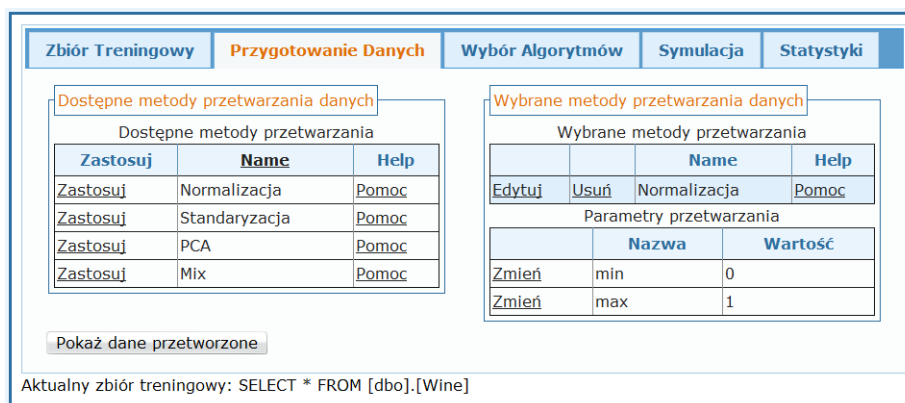
Kolejnym krokiem omawianego procesu jest wybór oraz odpowiednia parametryzacja algorytmów, które mają zostać ze sobą porównane. Dostępne są wszystkie algorytmy szczegółowo omówione w rozdziale 3. Możliwy jest wybór wielu instancji metody konkretnego typu i określenie różnej listy parametrów dla każdej z nich, dzięki temu możliwe jest porównywanie algorytmów tego samego rodzaju. Zrzut ekranu przedstawiający omawiany widok został pokazany na rysunku 5.3.

### 5.5.4. Konfiguracja symulacji

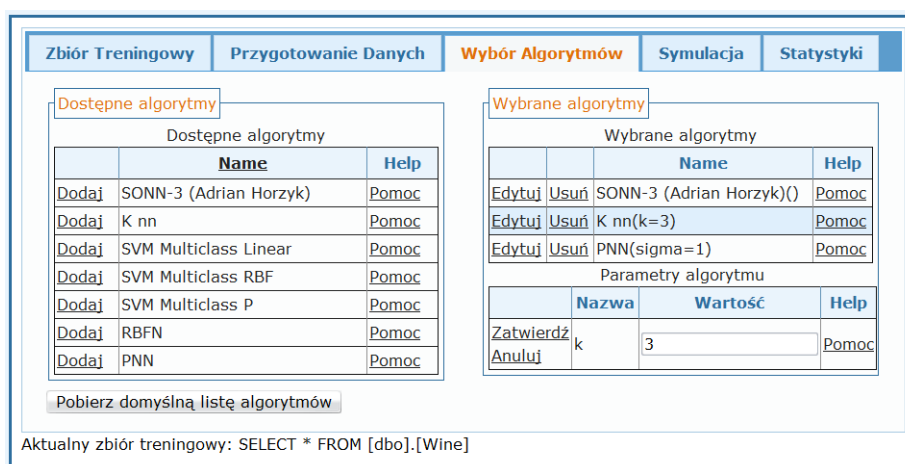
Ostatnim krokiem konfiguracyjnym przed rozpoczęciem symulacji jest wybór metod walidacji. Podobnie jak w przypadku algorytmów, wszystkie dostępne metody zostały omówione we wcześniejszej części pracy (patrz rozdział 4.4. Niektóre z nich (np. podział procentowy, *k-fold cross-validation*) mogą zostać odpowiednio sparametryzowane.

Oprócz wyboru metody walidacji możliwe jest określenie poziomu szczegółowości statystyk, które są przygotowywane w procesie symulacji. Jest to bardzo istotne podczas porównywania wielu algorytmów lub podczas korzystania z wielu metod walidacji. Ograniczenie szczegółowości statystyk przyspiesza działanie aplikacji w trakcie przeglądania wyników.





Rysunek 5.2: Aplikacja - wybór metod przetwarzania danych



Rysunek 5.3: Aplikacja - wybór algorytmów

Kolejne poziomy szczegółowości statystyk oznaczają:

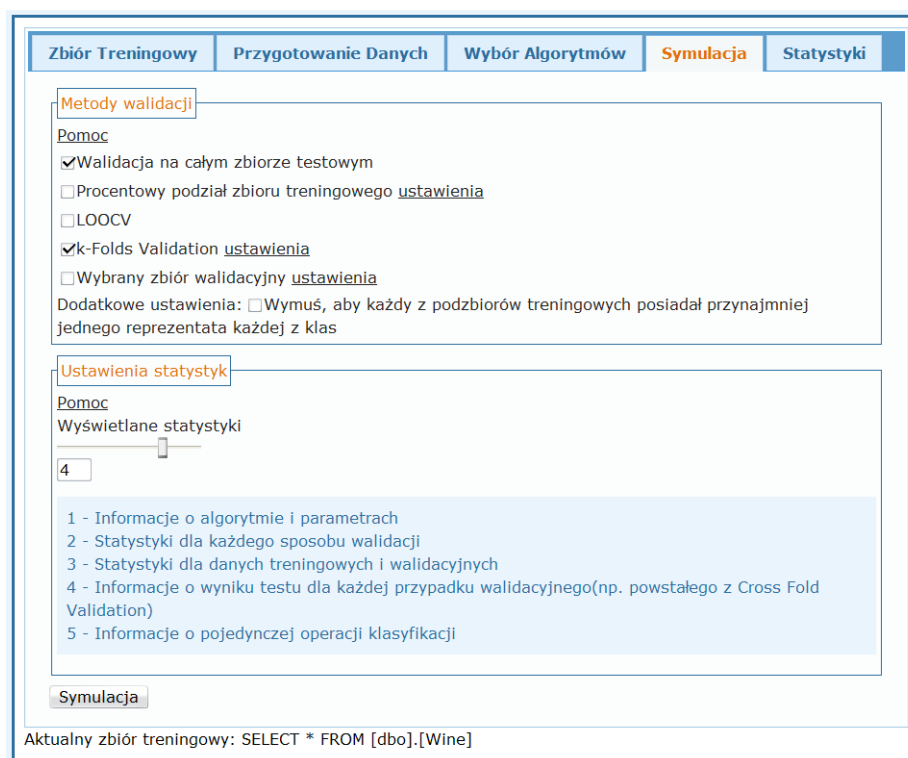
1. Poziom algorytmu - zawiera statystyki przekrojowe dla wszystkich stosowanych metod walidacji.
2. Poziom metody walidacji - dla każdej z wybranych metod walidacji dostępne są statystyki agregujące walidację na zbiorze treningowym i testowym.
3. Poziom walidacji na zbiorze treningowym/testowym - dla każdej metody walidacji przeprowadzana jest walidacja na podzbiorach treningowym i testowym. Poziom 3. pozwala na analizę statystyk dla walidacji przeprowadzonej na wymienionych podzbiorach.
4. Poziom konkretnego klasyfikatora - na tym poziomie (zależnie od metody walidacji) dostępnych jest jeden lub więcej wpisów, związanych z konkretną operacją tworzenia klasyfikatora. Jest to najbardziej istotny poziom statystyk przedstawiający takie dane jak:

informacje o klasyfikatorze, macierz niepewności, wskaźniki jakości i efektywności działania dla klasyfikatora i inne.

5. Poziom operacji klasyfikacji - najbardziej atomowy poziom statystyk zawierający informacje o pojedynczej operacji klasyfikacji wykonanej w procesie walidacji (zaleca się pominięcie tego poziomu szczegółowości w celu przyspieszenia działania aplikacji).

Po zakończeniu konfiguracji należy wcisnąć przycisk „Symulacja”, który rozpocznie proces treningu i klasyfikacji, dla każdego z wybranych algorytmów. Proces ten może trwać bardzo długo - zależnie od wybranego zbioru danych, zestawu algorytmów i metod walidacji. Po zakończeniu symulacji aplikacja przedstawi wyniki przeprowadzonych działań.

Omawiane możliwości konfiguracji dostępne są na ekranie aplikacji przedstawionym na rysunku 5.4.



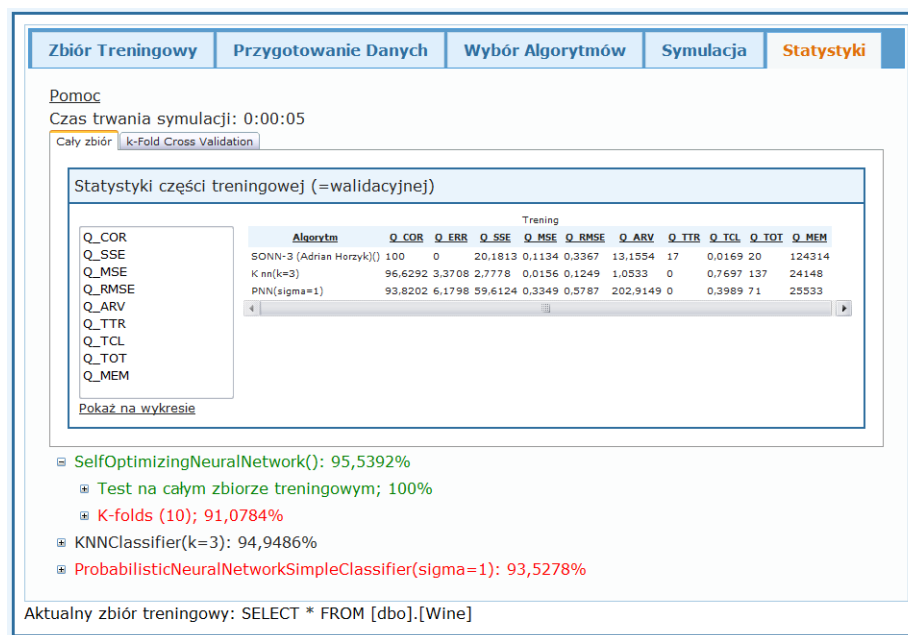
Rysunek 5.4: Aplikacja - konfiguracja symulacji

### 5.5.5. Przeglądanie wyników

Po wykonaniu symulacji aplikacja przedstawi wyniki przeprowadzonych działań. Wygląd ekranu zawierającego wyniki przedstawia rysunek 5.5. Na tym ekranie użytkownik może przeanalizować szczegółowo cały przeprowadzony przez program proces. Oprócz statystyk omówionych w poprzednim podrozdziale dostępne są także zbiorcze statystyki jakości uogólnienia i efektywności działania, zgrupowane według każdej z wybranych metod walidacji.

Ciekawą możliwością oferowaną na tym etapie pracy z aplikacją jest możliwość zapisania wybranego klasyfikatora do postaci binarnej (serializacja). Dzięki temu osoba starająca się wybrać najlepszy algorytm do praktycznych zastosowań może, po przeprowadzeniu odpowiednich porównań, wybrać i zapisać najlepszy z klasyfikatorów dla danego problemu. Tak zapisany obiekt algorytmu można odtworzyć w dowolnej aplikacji wykorzystującej biblioteki projektu wywołując statyczną metodę klasy *Algorithms.AlgorithmHelper*:

```
public static IAlgorithm Deserialize(Stream stream)
```



Rysunek 5.5: Aplikacja - wyniki

## 6. Wyniki porównań metod klasyfikacji

W niniejszym rozdziale przedstawione zostaną wyniki porównań wybranych algorytmów wykonane za pomocą zaimplementowanej aplikacji opisanej w poprzednim rozdziale. W kolejnych podrozdziałach zaprezentowane zostaną wyniki algorytmów dla każdego z wybranych zbiorów danych. Dodatkowo dla zbioru Wine porównane zostaną wyniki algorytmów po zastosowaniu popularnych metod wstępnego przetwarzania danych.

Jako metody walidacji wybrane zostały:

- test na całym zbiorze treningowym - wskaźniki jakości tak jak w rozdziale 4.5.
- dziesięcioelementowa walidacja krzyżowa (*10-fold cross-validation*) - wskaźniki opisane w rozdziale 4.5 dotyczą wartości średniej (*avg*) dla wszystkich podzbiorów powstałych w wyniku testu krzyżowego oraz wartości minimalnej (*min*) i maksymalnej (*max*), które zostały osiągnięte dla konkretnych podzbiorów. Przedstawione zostaną wyniki osiągnięte zarówno dla danych treningowych, jak i podzbioru walidacyjnego.

W podrozdziale 6.9 przeanalizowane zostaną także wskaźniki, które mogą zobrazować zachowanie konkretnych algorytmów w zależności od wielkości i stopnia skomplikowania zbioru danych.

### 6.1. Stanowisko testowe i uwagi o procesie porównywania

Wszystkie prezentowane porównania wykonane zostały na komputerze o parametrach przedstawionych w tabeli 6.1. Szczegóły techniczne aplikacji wykorzystanej do porównań opisane zostały w rozdziale 5.

W tabelach przedstawiających wskaźniki jakości uogólnienia i efektywności działania oprócz wartości liczbowych można znaleźć następujące informacje:

- NaN (*ang. not a number*) - wartość, która powstała w wyniku nieprawidłowej operacji numerycznej (np. dzielenie przez wartość bliską zero).

Tablica 6.1: Parametry techniczne komputera, na którym zrealizowano porównania

Parametr	Wartość
Procesor	Intel Core 2 Duo CPU T8300 2.40GHz
Pamięć RAM	4096 MB
System operacyjny	Windows Vista 32-bit
.NET Framework	3.5

- nieakceptowalny czas treningu - dla kilku algorytmów czas tworzenia klasyfikatora przekroczył dopuszczalną wartość (wielkość zależna od zbioru danych, najczęściej kilka godzin).
- nie udało się przeprowadzić treningu - brak możliwości przeprowadzenia treningu z powodu konieczności wykonania niedozwolonej operacji matematycznej (np. próba odwrócenia macierzy osobliwej).

Z powodu dużej ilości prezentowanych wskaźników jakości konieczne było zastosowanie skrótowych nazw algorytmów. Legenda wyjaśniająca zastosowane oznaczenia została zamieszczona w tabeli 6.2.

## 6.2. Porównania dla zbioru Wine

Jednym z najbardziej popularnych zbiorów danych, na których testowane są metody klasyfikacji, jest zbiór Wine. Dla tego zbioru przeanalizowany zostanie wpływ zastosowania różnych metod wstępnego przetwarzania danych oraz przeprowadzone zostaną porównania części treningowej i walidacyjnej powstałych w wyniku zastosowania dziesięcioelementowego testu krzyżowego.

Tabele 6.4, 6.5, 6.6 oraz 6.7 przedstawiają wyniki wybranych algorytmów klasyfikacji dla zbioru Wine poddanego zaimplementowanym metodom wstępnego przetwarzania danych. Na podstawie takiego porównania można zauważyć, że:

- Dla danych znormalizowanych (Tablica 6.5 większość algorytmów uzyskała najlepszą skuteczność (Q\_COR) oraz potrafiła utworzyć klasyfikator i przeprowadzić proces klasyfikacji w stosunkowo krótkim czasie (E\_TOT). W dalszych porównaniach wykonanych w tej pracy będzie to jedyna stosowana metoda przetwarzania danych.
- Zastosowanie algorytmu PCA do redukcji wymiaru przestrzeni cech dla wielu algorytmów istotnie skraca czas obliczeń przy jednoczesnym obniżeniu zdolności klasyfikacyjnych.

Tablica 6.2: Legenda dla algorytmów i parametrów występujących w wynikach porównań

Symbol	Algorytm	Opis parametrów
SONN-3	Self Optimizing Neural Networks	
KNN(k)	k Nearest Neighbours	(patrz tabela 3.1)
SVML(c;e)	Support Vector Machines	klasyfikator liniowy, c - cost, e - eps, (patrz tabela 3.2)
SVMR(c;s;e)	Support Vector Machines	klasyfikator nieliniowy, metoda jądrowa - f. gaussowska, s - sigma, pozostałe parametry jak dla SVML (patrz tabela 3.2)
SVMP(c;d;e)	Support Vector Machines	klasyfikator nieliniowy, metoda jądrowa: f. wielomianowa, (patrz tabela 3.2)
RBFN(k;f;c)	Radial Basis Function Networks	k - ilość neuronów warstwy ukrytej, f - funkcja radialna (1 - f. gaussowska, 2 - f. wielokwadratowa), c - metoda wyboru centrów (1 - wybór losowy, 2 - algorytm analizy skupień K-means), (patrz tabela 3.3)
PNN(sigma)	Probabilistic Neural Networks	(patrz tabela 3.4)

- Dla sieci SONN zastosowanie algorytmu PCA poskutkowało nieznacznym zwiększeniem czasu konstrukcji, znaczącym zwiększeniem używanej pamięci (Q\_MEM) oraz struktury sieci (porównanie struktury sieci SONN przy zastosowaniu normalizacji i algorytmu PCA znajduje się w tabeli 6.3). Można przypuszczać, że algorytm dyskryminacja ADL-BCA potrafi sprawniej odszukać informacje pozwalające na utworzenie klasyfikatora od algorytmu PCA.
- Sieci SONN, niezależnie od zastosowanej metody przetwarzania, potrafią bezbłędnie rozpoznawać wzorce zastosowane w procesie konstrukcji sieci.

Tablica 6.3: Struktura sieci SONN dla wybranych metod przetwarzania danych, Wine

Element struktury	Normalizacja	PCA
Liczba neuronów LBN	7	41
Liczba neuronów ARN	6	8
Liczba neuronów MSN	3	3
Liczba warstw sieci	5	5
Liczba połączeń międzyneuronowych	38	200

Tablica 6.4: Wine, brak wstępnego przetwarzania, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	20,18	0,11	0,34	13,16	55	0,06	65	124314
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	1	0,56	100	24148
KNN(k=3)	67,98	32,02	16,67	0,09	0,31	1,60	0	0,57	101	24148
KNN(k=5)	65,17	34,83	23,88	0,13	0,37	1,94	0	0,72	129	24148
KNN(k=10)	71,35	28,65	27,96	0,16	0,40	2,19	0	1,02	181	24148
KNN(k=15)	61,80	38,20	31,68	0,18	0,42	2,27	0	1,30	231	24148
KNN(k=30)	65,73	34,27	32,42	0,18	0,43	2,68	0	2,38	423	24148
SVML(c=1;e=0,1)	33,15	66,85	178,00	1,00	1,00	∞	126	0,03	132	38246
SVML(c=10;e=0,1)	33,15	66,85	178,00	1,00	1,00	∞	123	0,00	123	38246
SVML(c=50;e=0,1)	33,15	66,85	178,00	1,00	1,00	∞	120	0,01	121	38246
SVML(c=100;e=0,1)	33,15	66,85	178,00	1,00	1,00	∞	121	0,01	122	38246
SVML(c=20;e=0,1)	33,15	66,85	178,00	1,00	1,00	∞	142	0,01	144	38246
SVML(c=20;e=0,01)	94,38	5,62	16,52	0,09	0,30	4,49	6481	0,01	6482	38246
SVML(c=20;e=0,001)	nieakceptowalny czas treningu									
SVMR(c=10;s=0,5;e=0,1)	100,00	0,00	0,24	0,00	0,04	1,82	507	0,33	565	37980
SVMR(c=10;s=1;e=0,1)	100,00	0,00	0,34	0,00	0,04	2,55	584	0,31	640	37980
SVMR(c=10;s=5;e=0,1)	100,00	0,00	0,20	0,00	0,03	2,73	609	0,30	663	37980
SVMR(c=10;s=10;e=0,1)	100,00	0,00	0,75	0,00	0,06	1,90	1912	0,25	1956	37980
SVMP(c=10;d=3;e=0,01)	33,15	66,85	178,00	1,00	1,00	∞	2317	0,06	2327	38326
SVMP(c=10;d=4;e=0,01)	33,15	66,85	178,00	1,00	1,00	∞	2322	0,06	2332	38326
SVMP(c=10;d=6;e=0,01)	33,15	66,85	178,00	1,00	1,00	∞	2325	0,06	2335	38326
RBFN(k=5;f=2;c=1)	68,54	31,46	36,24	0,20	0,45	4,22	15	0,01	17	23540
RBFN(k=10;f=2;c=1)	68,54	31,46	42,68	0,24	0,49	3,09	1	0,02	4	23685
RBFN(k=15;f=2;c=1)	24,16	75,84	100,50	0,56	0,75	7,17	2	0,02	6	23830
RBFN(k=20;f=2;c=1)	42,70	57,30	74,36	0,42	0,65	5,06	3	0,02	7	23975
RBFN(k=5;f=1;c=2)	84,83	15,17	22,23	0,12	0,35	2,91	4	0,01	6	24126
RBFN(k=10;f=1;c=2)	98,88	1,12	8,52	0,05	0,22	3,14	2	0,02	5	24841
RBFN(k=15;f=1;c=2)	99,44	0,56	8,20	0,05	0,21	2,77	4	0,02	8	25556
RBFN(k=20;f=1;c=2)	100,00	0,00	6,35	0,04	0,19	2,82	5	0,03	10	26271
RBFN(k=5;f=2;c=2)	69,10	30,90	36,80	0,21	0,45	4,20	2	0,00	2	23562
RBFN(k=10;f=2;c=2)	68,54	31,46	36,29	0,20	0,45	4,21	1	0,01	3	23707
RBFN(k=15;f=2;c=2)	70,79	29,21	37,84	0,21	0,46	4,07	2	0,02	5	23852
RBFN(k=20;f=2;c=2)	69,10	30,90	45,71	0,26	0,51	3,10	4	0,02	8	23997
PNN(sigma=0,1)	100,00	0,00	0,00	0,00	0,00	NaN	0	0,49	87	25533
PNN(sigma=1)	100,00	0,00	0,00	0,00	0,00	1,01	0	0,45	80	25533
PNN(sigma=10)	90,45	9,55	14,62	0,08	0,29	1,84	0	0,43	77	25533
PNN(sigma=100)	73,60	26,40	35,06	0,20	0,44	3,31	0	0,44	79	25533

Tablica 6.5: Wine, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	20,18	0,11	0,34	13,16	63	0,05	72	124314
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	1	0,53	95	24148
KNN(k=3)	96,63	3,37	2,78	0,02	0,12	1,05	0	0,60	107	24148
KNN(k=5)	93,26	6,74	3,60	0,02	0,14	1,11	0	0,71	127	24148
KNN(k=10)	92,13	7,87	4,56	0,03	0,16	1,20	0	1,11	198	24148
KNN(k=15)	88,20	11,80	4,94	0,03	0,17	1,27	0	1,38	245	24148
KNN(k=30)	81,46	18,54	6,34	0,04	0,19	1,49	0	2,42	430	24148
SVML(c=1;e=0,1)	99,44	0,56	9,23	0,05	0,23	3,75	31	0,03	36	38246
SVML(c=10;e=0,1)	99,44	0,56	9,88	0,06	0,24	4,23	67	0,01	68	38246
SVML(c=50;e=0,1)	100,00	0,00	9,87	0,06	0,24	4,30	102	0,01	103	38246
SVML(c=100;e=0,1)	100,00	0,00	10,10	0,06	0,24	4,38	78	0,01	79	38246
SVML(c=20;e=0,1)	100,00	0,00	9,44	0,05	0,23	4,21	94	0,01	96	38246
SVML(c=20;e=0,01)	100,00	0,00	9,69	0,05	0,23	4,57	225	0,01	226	38246
SVML(c=20;e=0,001)	100,00	0,00	9,56	0,05	0,23	5,02	290	0,01	291	38246
SVMR(c=10;s=0,5;e=0,1)	100,00	0,00	1,94	0,01	0,10	3,68	344	0,11	364	37980
SVMR(c=10;s=1;e=0,1)	100,00	0,00	6,40	0,04	0,19	4,90	680	0,08	694	37980
SVMR(c=10;s=5;e=0,1)	97,75	2,25	9,75	0,05	0,23	3,49	240	0,11	259	37980
SVMR(c=10;s=10;e=0,1)	98,88	1,12	10,71	0,06	0,25	3,59	472	0,16	500	37980
SVMP(c=10;d=3;e=0,01)	100,00	0,00	8,26	0,05	0,22	4,62	1145	0,06	1156	38326
SVMP(c=10;d=4;e=0,01)	100,00	0,00	8,57	0,05	0,22	4,49	1526	0,06	1537	38326
SVMP(c=10;d=6;e=0,01)	98,31	1,69	16,37	0,09	0,30	4,08	1394	0,10	1411	38326
RBFN(k=5;f=2;c=1)	95,51	4,49	14,86	0,08	0,29	3,55	15	0,01	17	23540
RBFN(k=10;f=2;c=1)	98,88	1,12	8,68	0,05	0,22	3,35	1	0,01	3	23685
RBFN(k=15;f=2;c=1)	100,00	0,00	6,78	0,04	0,20	3,18	2	0,01	4	23830
RBFN(k=20;f=2;c=1)	100,00	0,00	5,58	0,03	0,18	3,26	4	0,02	8	23975
RBFN(k=5;f=1;c=2)	98,88	1,12	8,37	0,05	0,22	3,11	6	0,00	6	24126
RBFN(k=10;f=1;c=2)	99,44	0,56	7,96	0,04	0,21	3,13	6	0,02	9	24841
RBFN(k=15;f=1;c=2)	99,44	0,56	6,53	0,04	0,19	3,01	13	0,02	16	25556
RBFN(k=20;f=1;c=2)	99,44	0,56	5,56	0,03	0,18	3,07	15	0,02	19	26271
RBFN(k=5;f=2;c=2)	98,88	1,12	11,45	0,06	0,25	4,08	1	0,01	2	23562
RBFN(k=10;f=2;c=2)	99,44	0,56	9,59	0,05	0,23	3,59	2	0,01	4	23707
RBFN(k=15;f=2;c=2)	100,00	0,00	6,86	0,04	0,20	3,24	2	0,02	5	23852
RBFN(k=20;f=2;c=2)	100,00	0,00	5,36	0,03	0,17	3,10	2	0,03	7	23997
PNN(sigma=0,1)	100,00	0,00	0,00	0,00	0,00	1,03	0	0,49	87	25533
PNN(sigma=1)	93,82	6,18	59,61	0,33	0,58	202,91	0	0,54	96	25533
PNN(sigma=10)	93,26	6,74	78,90	0,44	0,67	3,00E+06	0	0,51	90	25533
PNN(sigma=100)	93,26	6,74	79,11	0,44	0,67	3,01E+10	0	0,48	85	25533



Tablica 6.6: Wine, dane zestandaryzowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	20,18	0,11	0,34	13,16	60	0,05	69	124314
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	1	0,60	107	24148
KNN(k=3)	75,84	24,16	10,44	0,06	0,24	1,38	0	0,57	101	24148
KNN(k=5)	76,40	23,60	14,04	0,08	0,28	1,55	0	0,72	128	24148
KNN(k=10)	71,35	28,65	17,98	0,10	0,32	1,89	0	1,02	181	24148
KNN(k=15)	67,98	32,02	21,52	0,12	0,35	2,30	0	1,29	229	24148
KNN(k=30)	56,74	43,26	28,50	0,16	0,40	3,91	0	2,34	416	24148
SVML(c=1;e=0,1)	92,70	7,30	20,55	0,12	0,34	5,64	272	0,03	278	38246
SVML(c=10;e=0,1)	96,07	3,93	14,86	0,08	0,29	4,61	993	0,01	994	38246
SVML(c=50;e=0,1)	95,51	4,49	17,63	0,10	0,31	3,70	2686	0,01	2688	38246
SVML(c=100;e=0,1)	94,94	5,06	15,05	0,08	0,29	4,29	1235	0,01	1236	38246
SVML(c=20;e=0,1)	97,75	2,25	14,13	0,08	0,28	3,88	1643	0,02	1646	38246
SVML(c=20;e=0,01)	98,88	1,12	12,27	0,07	0,26	4,74	1383	0,01	1385	38246
SVML(c=20;e=0,001)	99,44	0,56	11,00	0,06	0,25	4,46	3768	0,01	3770	38246
SVMR(c=10;s=0,5;e=0,1)	100,00	0,00	0,36	0,00	0,05	2,71	488	0,31	543	37980
SVMR(c=10;s=1;e=0,1)	100,00	0,00	0,33	0,00	0,04	3,27	364	0,30	418	37980
SVMR(c=10;s=5;e=0,1)	100,00	0,00	3,37	0,02	0,14	3,04	1905	0,13	1928	37980
SVMR(c=10;s=10;e=0,1)	98,31	1,69	9,04	0,05	0,23	3,10	1981	0,12	2002	37980
SVMP(c=10;d=3;e=0,01)	33,15	66,85	NaN	NaN	NaN	NaN	2824	0,06	2834	38326
SVMP(c=10;d=4;e=0,01)	33,15	66,85	NaN	NaN	NaN	NaN	2781	0,06	2791	38326
SVMP(c=10;d=6;e=0,01)	33,15	66,85	NaN	NaN	NaN	NaN	2854	0,06	2864	38326
RBFN(k=5;f=2;c=1)	84,83	15,17	28,98	0,16	0,40	3,39	15	0,02	18	23540
RBFN(k=10;f=2;c=1)	96,63	3,37	11,97	0,07	0,26	3,23	1	0,02	4	23685
RBFN(k=15;f=2;c=1)	98,88	1,12	9,82	0,06	0,23	3,51	2	0,01	4	23830
RBFN(k=20;f=2;c=1)	99,44	0,56	9,83	0,06	0,24	3,50	3	0,02	7	23975
RBFN(k=5;f=1;c=2)	90,45	9,55	19,88	0,11	0,33	3,56	4	0,01	5	24126
RBFN(k=10;f=1;c=2)	97,19	2,81	13,88	0,08	0,28	3,32	2	0,00	2	24841
RBFN(k=15;f=1;c=2)	98,31	1,69	10,39	0,06	0,24	3,10	4	0,02	7	25556
RBFN(k=20;f=1;c=2)	98,31	1,69	8,77	0,05	0,22	2,94	12	0,02	16	26271
RBFN(k=5;f=2;c=2)	81,46	18,54	33,69	0,19	0,44	3,97	2	0,01	4	23562
RBFN(k=10;f=2;c=2)	94,94	5,06	14,22	0,08	0,28	2,99	1	0,01	3	23707
RBFN(k=15;f=2;c=2)	98,31	1,69	9,83	0,06	0,24	3,31	1	0,02	5	23852
RBFN(k=20;f=2;c=2)	98,88	1,12	9,63	0,05	0,23	3,47	2	0,03	7	23997
PNN(sigma=0,1)	100,00	0,00	0,00	0,00	0,00	NaN	0	0,49	87	25533
PNN(sigma=1)	100,00	0,00	0,03	0,00	0,01	1,07	0	0,45	80	25533
PNN(sigma=10)	63,48	36,52	58,65	0,33	0,57	42,27	0	0,44	78	25533
PNN(sigma=100)	58,99	41,01	78,84	0,44	0,67	3,71E+05	0	0,48	85	25533

Tablica 6.7: Wine, zastosowanie algorytmu PCA i normalizacji danych, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	19,69	0,11	0,33	6,41	70	0,06	81	343926
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	0	0,49	88	7060
KNN(k=3)	65,17	34,83	19,33	0,11	0,33	1,77	0	0,53	94	7060
KNN(k=5)	64,04	35,96	27,12	0,15	0,39	2,07	0	0,67	120	7060
KNN(k=10)	67,42	32,58	30,50	0,17	0,41	2,37	0	1,01	180	7060
KNN(k=15)	65,17	34,83	32,15	0,18	0,43	2,40	0	1,26	225	7060
KNN(k=30)	65,17	34,83	32,39	0,18	0,43	2,69	0	2,39	425	7060
SVML(c=1;e=0,1)	69,10	30,90	51,65	0,29	0,54	2,84	27	0,03	33	20870
SVML(c=10;e=0,1)	68,54	31,46	44,33	0,25	0,50	4,54	26	0,01	28	20870
SVML(c=50;e=0,1)	68,54	31,46	44,20	0,25	0,50	5,55	21	0,01	23	20870
SVML(c=100;e=0,1)	68,54	31,46	43,13	0,24	0,49	6,29	32	0,01	34	20870
SVML(c=20;e=0,1)	68,54	31,46	44,57	0,25	0,50	4,47	23	0,01	24	20870
SVML(c=20;e=0,01)	68,54	31,46	44,29	0,25	0,50	4,69	35	0,01	36	20870
SVML(c=20;e=0,001)	68,54	31,46	44,36	0,25	0,50	4,58	51	0,01	53	20870
SVMR(c=10;s=0,5;e=0,1)	71,35	28,65	35,91	0,20	0,45	3,30	528	0,11	547	20892
SVMR(c=10;s=1;e=0,1)	68,54	31,46	42,98	0,24	0,49	3,78	678	0,10	696	20892
SVMR(c=10;s=5;e=0,1)	67,42	32,58	51,88	0,29	0,54	2,34	389	0,12	411	20892
SVMR(c=10;s=10;e=0,1)	67,98	32,02	53,43	0,30	0,55	2,36	721	0,12	743	20892
SVMP(c=10;d=3;e=0,01)	65,73	34,27	54,95	0,31	0,56	2,22	527	0,11	547	21238
SVMP(c=10;d=4;e=0,01)	65,73	34,27	55,48	0,31	0,56	2,00	704	0,12	725	21238
SVMP(c=10;d=6;e=0,01)	65,73	34,27	49,66	0,28	0,53	1,96	686	0,13	710	21238
RBFN(k=5;f=2;c=1)	71,91	28,09	33,26	0,19	0,43	2,56	8	0,01	9	6452
RBFN(k=10;f=2;c=1)	73,60	26,40	31,29	0,18	0,42	2,52	1	0,01	2	6597
RBFN(k=15;f=2;c=1)	72,47	27,53	33,49	0,19	0,43	2,21	2	0,01	4	6742
RBFN(k=20;f=2;c=1)	72,47	27,53	32,23	0,18	0,43	2,34	2	0,02	5	6887
RBFN(k=5;f=1;c=2)	71,35	28,65	33,18	0,19	0,43	2,37	4	0,01	5	6558
RBFN(k=10;f=1;c=2)	70,22	29,78	33,08	0,19	0,43	2,36	4	0,01	5	6793
RBFN(k=15;f=1;c=2)	71,35	28,65	33,39	0,19	0,43	2,41	3	0,01	5	7028
RBFN(k=20;f=1;c=2)	72,47	27,53	33,04	0,19	0,43	2,39	5	0,01	7	7263
RBFN(k=5;f=2;c=2)	70,22	29,78	33,16	0,19	0,43	2,32	1	0,01	2	6474
RBFN(k=10;f=2;c=2)	72,47	27,53	31,89	0,18	0,42	2,27	1	0,01	3	6619
RBFN(k=15;f=2;c=2)	72,47	27,53	33,29	0,19	0,43	2,96	1	0,01	3	6764
RBFN(k=20;f=2;c=2)	71,35	28,65	33,26	0,19	0,43	2,29	2	0,01	3	6909
PNN(sigma=0,1)	73,60	26,40	37,57	0,21	0,46	4,14	0	0,37	66	8258
PNN(sigma=1)	73,03	26,97	76,46	0,43	0,66	5,03E+03	0	0,35	63	8258
PNN(sigma=10)	73,03	26,97	79,08	0,44	0,67	5,29E+07	0	0,36	64	8258
PNN(sigma=100)	73,03	26,97	79,11	0,44	0,67	5,29E+11	0	0,36	64	8258





## 6.3. Porównania dla zbioru Iris

Zbiór Iris jest uznawany za jeden z łatwiejszych dla algorytmów klasyfikacji. Mimo tego, analiza wyników uzyskanych dzięki wybranym metodom walidacji (tabele 6.10, 6.11 oraz 6.12) może dostarczyć ciekawych wniosków.

Warto zauważyć, że popularna metoda SVM z funkcją Gaussa jako funkcją radialną, która niewątpliwie jest bardzo popularnym narzędziem klasyfikacji, uzyskuje zadowalające wyniki tylko wtedy, gdy odpowiednio dobrane są parametry. W praktyce, jest to związane z wydłużonym czasem poszukiwań najlepszego klasyfikatora dla danego problemu, ponieważ konieczne jest sprawdzenie kilku różnych konfiguracji parametrów. Bardzo często wiąże się to z doświadczeniem i wiedzą ekspercką o konkretnej dziedzinie i własnościach stosowanej metody klasyfikacji.

Tablica 6.10: Iris, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	19,28	0,13	0,36	6,28	53	0,13	73	130505
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	2	0,76	116	11446
KNN(k=3)	95,33	4,67	2,67	0,02	0,13	1,06	0	0,48	72	11446
KNN(k=5)	93,33	6,67	3,80	0,03	0,16	1,09	0	0,60	90	11446
KNN(k=10)	95,33	4,67	2,90	0,02	0,14	1,15	0	0,87	130	11446
KNN(k=15)	90,67	9,33	3,56	0,02	0,15	1,21	0	1,16	174	11446
KNN(k=30)	86,00	14,00	5,81	0,04	0,20	1,47	0	2,15	323	11446
SVML(c=1;e=0,1)	86,67	13,33	23,71	0,16	0,40	8,02	28	0,12	46	21919
SVML(c=10;e=0,1)	94,00	6,00	20,79	0,14	0,37	7,94	23	0,01	24	21919
SVML(c=50;e=0,1)	94,67	5,33	20,58	0,14	0,37	8,25	111	0,00	111	21919
SVML(c=100;e=0,1)	94,67	5,33	21,01	0,14	0,37	8,88	140	0,01	141	21919
SVML(c=20;e=0,1)	94,67	5,33	20,87	0,14	0,37	9,82	27	0,01	29	21919
SVML(c=20;e=0,01)	94,67	5,33	20,71	0,14	0,37	9,33	55	0,00	55	21919
SVML(c=20;e=0,001)	94,67	5,33	20,70	0,14	0,37	8,80	123	0,00	123	21919
SVMR(c=10;s=0,5;e=0,1)	97,33	2,67	6,29	0,04	0,20	3,50	191	0,09	204	21869
SVMR(c=10;s=1;e=0,1)	97,33	2,67	11,82	0,08	0,28	4,59	166	0,09	179	21869
SVMR(c=10;s=5;e=0,1)	82,00	18,00	24,91	0,17	0,41	7,58	328	0,11	345	21869
SVMR(c=10;s=10;e=0,1)	66,67	33,33	48,10	0,32	0,57	2,62	371	0,12	389	21869
SVMP(c=10;d=3;e=0,01)	95,33	4,67	22,11	0,15	0,38	11,43	2059	0,10	2074	22215
SVMP(c=10;d=4;e=0,01)	97,33	2,67	23,20	0,15	0,39	11,86	2543	0,09	2557	22215
SVMP(c=10;d=6;e=0,01)	96,67	3,33	25,85	0,17	0,42	13,31	6674	0,10	6689	22215
RBFN(k=5;f=2;c=1)	96,00	4,00	9,07	0,06	0,25	3,22	15	0,01	17	11062
RBFN(k=10;f=2;c=1)	98,00	2,00	4,90	0,03	0,18	2,29	1	0,01	3	11207
RBFN(k=15;f=2;c=1)	98,00	2,00	5,48	0,04	0,19	2,46	1	0,01	3	11352
RBFN(k=20;f=2;c=1)	98,67	1,33	4,62	0,03	0,18	2,26	2	0,01	4	11497
RBFN(k=5;f=1;c=2)	95,33	4,67	8,68	0,06	0,24	3,26	5	0,01	6	11288
RBFN(k=10;f=1;c=2)	97,33	2,67	6,95	0,05	0,22	2,65	3	0,01	5	11643
RBFN(k=15;f=1;c=2)	98,00	2,00	6,39	0,04	0,21	2,55	3	0,01	5	11998
RBFN(k=20;f=1;c=2)	98,00	2,00	4,57	0,03	0,17	2,11	5	0,01	7	12353
RBFN(k=5;f=2;c=2)	94,67	5,33	9,66	0,06	0,25	3,14	1	0,01	2	11084
RBFN(k=10;f=2;c=2)	98,67	1,33	6,34	0,04	0,21	2,57	1	0,01	3	11229
RBFN(k=15;f=2;c=2)	98,67	1,33	6,10	0,04	0,20	2,56	1	0,01	3	11374
RBFN(k=20;f=2;c=2)	98,67	1,33	4,87	0,03	0,18	2,45	2	0,01	3	11519
PNN(sigma=0,1)	97,33	2,67	2,93	0,02	0,14	1,30	0	0,39	59	12816
PNN(sigma=1)	91,33	8,67	53,12	0,35	0,60	375,87	0	0,38	57	12816
PNN(sigma=10)	91,33	8,67	66,52	0,44	0,67	4,76E+06	0	0,39	58	12816
PNN(sigma=100)	91,33	8,67	66,67	0,44	0,67	4,78E+10	0	0,40	60	12816



Tablica 6.12: Iris, dane znormalizowane, 10-fold cross-validation, walidacja

Table with 28 columns: Algoritm, Q\_COR(avg), Q\_COR(min), Q\_COR(max), Q\_ERR(avg), Q\_ERR(min), Q\_ERR(max), Q\_SSE(avg), Q\_SSE(min), Q\_SSE(max), Q\_MSE(avg), Q\_MSE(min), Q\_MSE(max), Q\_RMSE(avg), Q\_RMSE(min), Q\_RMSE(max), Q\_ARV(avg), Q\_ARV(min), Q\_ARV(max), E\_TTR(avg), E\_TTR(min), E\_TTR(max), F\_TCL(avg), F\_TCL(min), F\_TCL(max), E\_TOT(avg), E\_TOT(min), E\_TOT(max), F\_MEM(avg), F\_MEM(min), F\_MEM(max)

## 6.4. Porównania dla zbioru Yeast

Zbiór Yeast ze względu na nierównomierny rozkład wzorców według klas i stosunkowo dużą ilość danych może stanowić problem dla wielu algorytmów klasyfikacji. Już podczas walidacji dla całego zbioru danych (patrz tabela 6.13) można zauważyć, że większość metod osiąga wyniki dużo gorsze niż dla innych zbiorów danych. Metoda SONN dla wspomnianej metody walidacji osiągnęła stuprocentową skuteczność, jednak analiza wskaźników takich, jak Q\_MSE pozwala sądzić, że jest to związane ze specyficzną cechą tego algorytmu, dzięki której osiąga on stuprocentową skuteczność dla wzorców niesprzecznych. Zostaje to zweryfikowane podczas testu krzyżowego (tabele 6.14, 6.15), dla którego okazuje się, że metoda SONN mimo całkowitej skuteczności dla danych treningowych osiąga najgorsze rezultaty dla danych wcześniej nieznanymi.

Tablica 6.13: Yeast, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	1107,51	0,75	0,86	10910,34	20222	0,72	21288	24020508
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	2	4,11	6100	142760
KNN(k=3)	49,80	50,20	262,11	0,18	0,42	2,48	0	4,74	7037	142760
KNN(k=5)	47,64	52,36	350,00	0,24	0,49	3,29	0	6,26	9294	142760
KNN(k=10)	44,81	55,19	422,82	0,28	0,53	4,46	0	8,97	13305	142760
KNN(k=15)	46,09	53,91	454,67	0,31	0,55	5,09	0	11,33	16821	142760
KNN(k=30)	43,53	56,47	493,53	0,33	0,58	6,09	0	20,12	29856	142760
SVML(c=1;e=0,1)	55,59	44,41	1040,11	0,70	0,84	223,69	1180	0,02	1217	480911
SVML(c=10;e=0,1)	55,12	44,88	1077,01	0,73	0,85	627,19	3197	0,02	3223	480911
SVML(c=50;e=0,1)	52,63	47,37	1095,52	0,74	0,86	1027,74	13081	0,02	13105	480911
SVML(c=100;e=0,1)	49,53	50,47	1079,48	0,73	0,85	661,44	34452	0,01	34474	480911
SVML(c=20;e=0,1)	51,75	48,25	1080,97	0,73	0,85	734,39	4735	0,01	4754	480911
SVML(c=20;e=0,01)	50,74	49,26	1096,78	0,74	0,86	1026,62	11068	0,01	11089	480911
SVML(c=20;e=0,001)	46,36	53,64	1093,89	0,74	0,86	948,73	57674	0,02	57699	480911
SVMR(c=10;s=0,5;e=0,1)	64,15	35,85	1051,24	0,71	0,84	393,73	2374605	2,65	2378533	480268
SVMR(c=10;s=1;e=0,1)	60,58	39,42	1045,15	0,70	0,84	283,32	353933	2,74	357995	480268
SVMR(c=10;s=5;e=0,1)	54,65	45,35	991,39	0,67	0,82	80,56	62177	3,50	67378	480268
SVMR(c=10;s=10;e=0,1)	47,17	52,83	988,98	0,67	0,82	69,12	122436	3,13	127077	480268
SVMP(c=10;d=3;e=0,01)	nieakceptowalny czas treningu									
SVMP(c=10;d=4;e=0,01)	nieakceptowalny czas treningu									
SVMP(c=10;d=6;e=0,01)	nieakceptowalny czas treningu									
RBFN(k=5;f=2;c=1)	53,50	46,50	767,59	0,52	0,72	28,68	90	0,02	120	132180
RBFN(k=10;f=2;c=1)	58,29	41,71	738,31	0,50	0,71	27,93	21	0,04	78	132605
RBFN(k=15;f=2;c=1)	59,23	40,77	732,21	0,49	0,70	28,80	27	0,05	105	133030
RBFN(k=20;f=2;c=1)	59,97	40,03	670,90	0,45	0,67	17,99	39	0,06	135	133455
RBFN(k=5;f=1;c=2)	48,92	51,08	796,72	0,54	0,73	30,58	88	0,02	122	132566
RBFN(k=10;f=1;c=2)	57,68	42,32	749,79	0,51	0,71	30,13	115	0,04	172	133361
RBFN(k=15;f=1;c=2)	58,76	41,24	731,23	0,49	0,70	27,26	150	0,06	236	134156
RBFN(k=20;f=1;c=2)	59,37	40,63	716,25	0,48	0,69	25,25	437	0,08	549	134951
RBFN(k=5;f=2;c=2)	52,36	47,64	789,68	0,53	0,73	32,12	11	0,03	56	132202
RBFN(k=10;f=2;c=2)	58,36	41,64	737,16	0,50	0,70	29,36	18	0,04	77	132627
RBFN(k=15;f=2;c=2)	60,04	39,96	721,82	0,49	0,70	26,56	29	0,05	101	133052
RBFN(k=20;f=2;c=2)	60,98	39,02	687,06	0,46	0,68	22,16	37	0,07	135	133477
PNN(sigma=0,1)	63,27	36,73	674,22	0,45	0,67	9,53	0	10,15	15066	138947
PNN(sigma=1)	50,27	49,73	1172,04	0,79	0,89	2,43E+04	0	10,05	14919	138947
PNN(sigma=10)	48,25	51,75	1201,71	0,81	0,90	3,24E+08	0	10,21	15150	138947
PNN(sigma=100)	48,18	51,82	1202,04	0,81	0,90	3,25E+12	0	10,34	15351	138947





Tablica 6.15: Yeast, dane znormalizowane, 10-fold cross-validation, walidacja

Table with columns for Algorithm and various performance metrics such as Q\_COR, Q\_ERR, Q\_SSE, Q\_MSE, Q\_MRF, Q\_ARV, F\_TTR, F\_TCL, F\_TOT, F\_MEM, and E\_MEM. The table lists numerous models including SONN, KNN, SVM, and RBFN with their respective parameter settings and performance values.

## 6.5. Porównania dla zbioru Congressional Voting Records

Congressional Voting Records jest jednym z ciekawszych spośród analizowanych w niniejszej pracy zbiorów. Ponieważ zbiór ten składa się wyłącznie z wartości +1, -1 oraz 0 można przypuszczać, że sieci SONN, które pierwotnie były algorytmami działającymi wyłącznie dla zbiorów złożonych z cech o wartościach binarnych, osiągną dla tego zbioru bardzo dobre wyniki. Analiza wartości w tabelach 6.16, 6.17 oraz 6.18 zdaje się potwierdzać takie przypuszczenia. Szczególnie warto zwrócić uwagę na wysokie wartości wskaźników Q\_COR oraz niskie wartości wskaźnika Q\_MSE.

Można także zauważyć, że dla tego zbioru metoda radialnych sieci funkcji bazowych w pięciu z dwunastu przypadków nie była w stanie przeprowadzić procesu treningu.

Tablica 6.16: Congressional Voting Records, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	0,00	0,00	0,00	NaN	112	0,03	127	921824
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	0	1,40	611	72407
KNN(k=3)	92,64	7,36	9,89	0,02	0,15	1,11	0	1,49	648	72407
KNN(k=5)	90,80	9,20	13,72	0,03	0,18	1,13	0	1,84	800	72407
KNN(k=10)	90,80	9,20	17,10	0,04	0,20	1,17	0	2,81	1223	72407
KNN(k=15)	89,66	10,34	18,84	0,04	0,21	1,19	0	3,38	1470	72407
KNN(k=30)	89,20	10,80	22,53	0,05	0,23	1,22	0	5,93	2581	72407
SVML(c=1;e=0,1)	97,01	2,99	13,00	0,03	0,17	1,03	227	0,00	229	86876
SVML(c=10;e=0,1)	97,01	2,99	13,00	0,03	0,17	1,03	2188	0,00	2189	86876
SVML(c=50;e=0,1)	97,01	2,99	13,00	0,03	0,17	1,03	13899	0,01	13903	86876
SVML(c=100;e=0,1)	96,78	3,22	14,00	0,03	0,18	1,03	18006	0,00	18008	86876
SVML(c=20;e=0,1)	97,47	2,53	11,00	0,03	0,16	1,03	4235	0,01	4238	86876
SVML(c=20;e=0,01)	97,01	2,99	13,00	0,03	0,17	1,03	8486	0,01	8489	86876
SVML(c=20;e=0,001)	97,47	2,53	11,00	0,03	0,16	1,03	9790	0,00	9792	86876
SVMR(c=10;s=0,5;e=0,1)	100,00	0,00	0,00	0,00	0,00	NaN	4143	0,34	4292	86673
SVMR(c=10;s=1;e=0,1)	100,00	0,00	0,00	0,00	0,00	NaN	4210	0,20	4295	86673
SVMR(c=10;s=5;e=0,1)	96,78	3,22	14,00	0,03	0,18	1,03	2375	0,14	2437	86673
SVMR(c=10;s=10;e=0,1)	95,86	4,14	18,00	0,04	0,20	1,04	530	0,16	601	86673
SVMP(c=10;d=3;e=0,01)	100,00	0,00	0,00	0,00	0,00	NaN	89185	0,11	89232	86987
SVMP(c=10;d=4;e=0,01)	66,67	33,33	145,00	0,33	0,58	1,50	28428	0,09	28465	86987
SVMP(c=10;d=6;e=0,01)	38,62	61,38	435,00	1,00	1,00	∞	23846	0,09	23885	86987
RBFN(k=5;f=2;c=1)	90,34	9,66	42,00	0,10	0,31	1,11	4	0,01	8	69675
RBFN(k=10;f=2;c=1)	94,48	5,52	24,00	0,06	0,23	1,06	3	0,01	6	69780
RBFN(k=15;f=2;c=1)	95,86	4,14	18,00	0,04	0,20	1,04	4	0,01	10	69885
RBFN(k=20;f=2;c=1)	nie udało się przeprowadzić treningu									
RBFN(k=5;f=1;c=2)	92,41	7,59	33,00	0,08	0,28	1,08	16	0,00	18	70381
RBFN(k=10;f=1;c=2)	95,63	4,37	19,00	0,04	0,21	1,05	36	0,01	39	71176
RBFN(k=15;f=1;c=2)	96,78	3,22	14,00	0,03	0,18	1,03	26	0,01	32	71971
RBFN(k=20;f=1;c=2)	96,32	3,68	16,00	0,04	0,19	1,04	38	0,02	46	72766
RBFN(k=5;f=2;c=2)	90,57	9,43	41,00	0,09	0,31	1,10	2	0,00	4	69697
RBFN(k=10;f=2;c=2)	92,18	7,82	34,00	0,08	0,28	1,08	3	0,01	7	69802
RBFN(k=15;f=2;c=2)	94,94	5,06	22,00	0,05	0,22	1,05	5	0,01	11	69907
RBFN(k=20;f=2;c=2)	95,86	4,14	18,00	0,04	0,20	1,04	7	0,02	14	70012
PNN(sigma=0,1)	100,00	0,00	0,00	0,00	0,00	1,01	0	1,09	476	72661
PNN(sigma=1)	90,57	9,43	30,49	0,07	0,26	1,74	0	1,52	663	72661
PNN(sigma=10)	89,43	10,57	106,26	0,24	0,49	1,44E+04	0	1,12	489	72661
PNN(sigma=100)	89,43	10,57	108,72	0,25	0,50	1,47E+08	0	1,10	480	72661





## 6.6. Porównania dla zbioru Car Evaluation

Tabele 6.19, 6.20 oraz 6.21 przedstawiają wyniki z przeprowadzonych porównań dla zbioru Car Evaluation. Ciekawą cechą tego zbioru jest to, że domyślnie składał się z cech o wartościach symbolicznych, które zostały zamienione na wartości liczbowe. Ponieważ dane symboliczne często występują w zbiorach danych warto zastanowić się nad wpływem ich transformacji na wartości liczbowe dla badanych metod klasyfikacji.

Tablica 6.19: Car Evaluation, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	399,35	0,23	0,48	25,60	1331	0,21	1697	2091165
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	2	5,44	9404	141155
KNN(k=3)	89,18	10,82	44,11	0,03	0,16	1,20	1	7,23	12500	141155
KNN(k=5)	87,33	12,67	56,96	0,03	0,18	1,34	0	8,40	14520	141155
KNN(k=10)	84,49	15,51	81,60	0,05	0,22	1,52	0	12,17	21028	141155
KNN(k=15)	81,94	18,06	93,81	0,05	0,23	1,70	0	15,66	27061	141155
KNN(k=30)	74,77	25,23	129,02	0,07	0,27	1,99	1	27,48	47487	141155
SVML(c=1;e=0,1)	81,13	18,87	517,97	0,30	0,55	21,36	2190	0,02	2229	279310
SVML(c=10;e=0,1)	81,83	18,17	604,03	0,35	0,59	48,21	23432	0,01	23451	279310
SVML(c=50;e=0,1)	80,96	19,04	594,79	0,34	0,59	47,90	287783	0,01	287796	279310
SVML(c=100;e=0,1)	82,12	17,88	607,78	0,35	0,59	52,52	937107	0,01	937119	279310
SVML(c=20;e=0,1)	82,52	17,48	601,87	0,35	0,59	48,99	61595	0,01	61612	279310
SVML(c=20;e=0,01)	82,47	17,53	612,45	0,35	0,60	51,03	30796	0,01	30814	279310
SVML(c=20;e=0,001)	82,18	17,82	609,04	0,35	0,59	48,37	480439	0,01	480461	279310
SVMR(c=10;s=0,5;e=0,1)	99,88	0,12	178,02	0,10	0,32	7,71	1670218	1,17	1672233	279157
SVMR(c=10;s=1;e=0,1)	98,03	1,97	258,10	0,15	0,39	8,57	1561742	1,21	1563829	279157
SVMR(c=10;s=5;e=0,1)	81,60	18,40	419,49	0,24	0,49	12,06	104226	1,59	106967	279157
SVMR(c=10;s=10;e=0,1)	75,06	24,94	341,93	0,20	0,44	3,15	52179	1,68	55074	279157
SVMP(c=10;d=3;e=0,01)	nieakceptowalny czas treningu									
SVMP(c=10;d=4;e=0,01)	nieakceptowalny czas treningu									
SVMP(c=10;d=6;e=0,01)	nieakceptowalny czas treningu									
RBFN(k=5;f=2;c=1)	81,54	18,46	285,48	0,17	0,41	2,73	82	0,01	102	128215
RBFN(k=10;f=2;c=1)	81,02	18,98	272,40	0,16	0,40	2,91	14	0,02	40	128400
RBFN(k=15;f=2;c=1)	82,70	17,30	278,89	0,16	0,40	3,59	21	0,02	61	128585
RBFN(k=20;f=2;c=1)	87,15	12,85	233,55	0,14	0,37	3,21	27	0,03	79	128770
RBFN(k=5;f=1;c=2)	72,74	27,26	383,05	0,22	0,47	2,75	105	0,01	130	128521
RBFN(k=10;f=1;c=2)	81,94	18,06	264,10	0,15	0,39	2,46	160	0,02	194	128996
RBFN(k=15;f=1;c=2)	83,62	16,38	237,65	0,14	0,37	2,63	249	0,03	295	129471
RBFN(k=20;f=1;c=2)	85,94	14,06	236,30	0,14	0,37	3,31	328	0,03	383	129946
RBFN(k=5;f=2;c=2)	80,38	19,62	287,12	0,17	0,41	2,74	7	0,01	21	128237
RBFN(k=10;f=2;c=2)	81,37	18,63	270,87	0,16	0,40	2,66	18	0,02	46	128422
RBFN(k=15;f=2;c=2)	84,14	15,86	246,38	0,14	0,38	3,07	24	0,02	63	128607
RBFN(k=20;f=2;c=2)	84,90	15,10	245,40	0,14	0,38	3,25	31	0,02	72	128792
PNN(sigma=0,1)	100,00	0,00	1,64	0,00	0,03	1,07	0	7,09	12259	136235
PNN(sigma=1)	47,51	52,49	909,64	0,53	0,73	364,31	0	7,13	12315	136235
PNN(sigma=10)	41,38	58,62	971,47	0,56	0,75	3,79E+06	0	7,05	12188	136235
PNN(sigma=100)	41,32	58,68	971,99	0,56	0,75	3,79E+10	0	7,10	12265	136235







Warto zauważyć, że dla zbioru Car Evaluation metoda SONN osiągnęła korzystne rezultaty w zakresie poprawności klasyfikacji nieznanymi wzorców. Może to świadczyć o wyjątkowej skuteczności algorytmu binaryzacji ADLBCA dla zbiorów, w których wartości symboliczne zostały zamienione na wartości liczbowe.

Poza skutecznością klasyfikacji sieci SONN warto również zauważyć, że ta metoda dla zbioru Car Evaluation potrafiła w stosunkowo krótkim czasie przeprowadzić proces tworzenia klasyfikatora (warto porównać z algorytmami SVML, SVMR oraz SVMP).

## 6.7. Porównania dla zbioru Ionosphere

Tablica 6.22: Ionosphere, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	94,30	5,70	8,50	0,02	0,16	1,09	138	0,05	155	400230
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	1	1,55	545	106298
KNN(k=3)	86,32	13,68	17,33	0,05	0,22	1,19	0	1,94	682	106298
KNN(k=5)	83,19	16,81	25,48	0,07	0,27	1,25	0	2,37	832	106298
KNN(k=10)	79,49	20,51	34,32	0,10	0,31	1,34	0	3,87	1357	106298
KNN(k=15)	75,50	24,50	39,44	0,11	0,34	1,41	0	5,11	1794	106298
KNN(k=30)	70,66	29,34	50,26	0,14	0,38	1,59	0	9,42	3305	106298
SVML(c=1;e=0,1)	92,02	7,98	28,00	0,08	0,28	1,09	523	0,06	545	121648
SVML(c=10;e=0,1)	93,45	6,55	23,00	0,07	0,26	1,07	6784	0,01	6788	121648
SVML(c=50;e=0,1)	92,88	7,12	25,00	0,07	0,27	1,08	25616	0,01	25618	121648
SVML(c=100;e=0,1)	92,88	7,12	25,00	0,07	0,27	1,08	54561	0,01	54564	121648
SVML(c=20;e=0,1)	92,59	7,41	26,00	0,07	0,27	1,08	12465	0,01	12468	121648
SVML(c=20;e=0,01)	93,16	6,84	24,00	0,07	0,26	1,07	12211	0,01	12213	121648
SVML(c=20;e=0,001)	93,73	6,27	22,00	0,06	0,25	1,07	25017	0,01	25020	121648
SVMR(c=10;s=0,5;e=0,1)	100,00	0,00	0,00	0,00	0,00	NaN	3797	0,43	3949	121157
SVMR(c=10;s=1;e=0,1)	99,43	0,57	2,00	0,01	0,08	1,01	12864	0,25	12952	121157
SVMR(c=10;s=5;e=0,1)	94,59	5,41	19,00	0,05	0,23	1,06	6378	0,31	6488	121157
SVMR(c=10;s=10;e=0,1)	90,03	9,97	35,00	0,10	0,32	1,11	1538	0,44	1692	121157
SVMP(c=10;d=3;e=0,01)	97,72	2,28	8,00	0,02	0,15	1,02	2181079	0,19	2181146	121471
SVMP(c=10;d=4;e=0,01)	79,77	20,23	71,00	0,20	0,45	1,25	15232	0,09	15265	121471
SVMP(c=10;d=6;e=0,01)	64,10	35,90	351,00	1,00	1,00	∞	15980	0,08	16008	121471
RBFN(k=5;f=2;c=1)	84,90	15,10	53,00	0,15	0,39	1,18	20	0,01	25	104238
RBFN(k=10;f=2;c=1)	85,47	14,53	51,00	0,15	0,38	1,17	5	0,02	12	104343
RBFN(k=15;f=2;c=1)	91,17	8,83	31,00	0,09	0,30	1,10	8	0,03	18	104448
RBFN(k=20;f=2;c=1)	91,45	8,55	30,00	0,09	0,29	1,09	11	0,04	24	104553
RBFN(k=5;f=1;c=2)	91,45	8,55	30,00	0,09	0,29	1,09	25	0,01	29	105664
RBFN(k=10;f=1;c=2)	92,59	7,41	26,00	0,07	0,27	1,08	43	0,02	50	107179
RBFN(k=15;f=1;c=2)	90,31	9,69	34,00	0,10	0,31	1,11	57	0,03	68	108694
RBFN(k=20;f=1;c=2)	92,59	7,41	26,00	0,07	0,27	1,08	71	0,03	83	110209
RBFN(k=5;f=2;c=2)	76,07	23,93	84,00	0,24	0,49	1,31	3	0,01	7	104260
RBFN(k=10;f=2;c=2)	90,60	9,40	33,00	0,09	0,31	1,10	5	0,02	13	104365
RBFN(k=15;f=2;c=2)	91,74	8,26	29,00	0,08	0,29	1,09	7	0,03	17	104470
RBFN(k=20;f=2;c=2)	91,45	8,55	30,00	0,09	0,29	1,09	10	0,04	23	104575
PNN(sigma=0,1)	100,00	0,00	0,10	0,00	0,02	1,01	0	1,51	530	107038
PNN(sigma=1)	78,35	21,65	54,54	0,16	0,39	6,36	0	1,36	477	107038
PNN(sigma=10)	65,24	34,76	87,26	0,25	0,50	3,83E+04	0	1,36	479	107038
PNN(sigma=100)	65,24	34,76	87,75	0,25	0,50	3,82E+08	0	1,60	561	107038

Zbiór Ionosphere został uwzględniony w porównaniach ze względu na obecność wzorców sprzecznych. Dla takiego przypadku sieci SONN nie potrafią skutecznie sklasyfikować wszystkich wzorców treningowych, co można zauważyć analizując tabele 6.22, 6.23 oraz 6.24. Mimo tego, algorytm SONN dla danych walidacyjnych testu krzyżowego uzyskał rezultaty porównywalne z innymi metodami klasyfikacji.

Tablica 6.23: Ionosphere, dane znormalizowane, 10-fold cross-validation, trening

Table with 27 columns: Algoritm, Q\_COR(avg), Q\_COR(min), Q\_COR(max), Q\_ERR(avg), Q\_ERR(min), Q\_ERR(max), Q\_SSE(avg), Q\_SSE(min), Q\_SSE(max), Q\_MSE(avg), Q\_MSE(min), Q\_MSE(max), Q\_RMSE(avg), Q\_RMSE(min), Q\_RMSE(max), Q\_ARV(avg), Q\_ARV(min), Q\_ARV(max), F\_TOT(avg), F\_TOT(min), F\_TOT(max), F\_TCL(avg), F\_TCL(min), F\_TCL(max), F\_MEM(avg), F\_MEM(min), F\_MEM(max). Rows include algorithms like SONN-3, KNN(k=1), KNN(k=3), KNN(k=5), KNN(k=10), KNN(k=15), KNN(k=30), SVMML(c=1;e=0,1), SVMML(c=10;e=0,1), SVMML(c=50;e=0,1), SVMML(c=100;e=0,1), SVMML(c=20;e=0,1), SVMML(c=20;e=0,01), SVMML(c=10;e=0,5;e=0,1), SVMML(c=10;e=1;e=0,1), SVMML(c=10;e=2;e=1), SVMML(c=10;e=3;e=0,1), SVMML(c=10;e=4;e=0,1), SVMML(c=10;e=5;e=0,1), SVMML(c=10;e=6;e=0,1), SVMML(c=10;e=7;e=0,1), SVMML(c=10;e=8;e=0,1), SVMML(c=10;e=9;e=0,1), SVMML(c=10;e=10;e=0,1), SVMML(c=10;e=11;e=0,1), SVMML(c=10;e=12;e=0,1), SVMML(c=10;e=13;e=0,1), SVMML(c=10;e=14;e=0,1), SVMML(c=10;e=15;e=0,1), SVMML(c=10;e=16;e=0,1), SVMML(c=10;e=17;e=0,1), SVMML(c=10;e=18;e=0,1), SVMML(c=10;e=19;e=0,1), SVMML(c=10;e=20;e=0,1), SVMML(c=10;e=21;e=0,1), SVMML(c=10;e=22;e=0,1), SVMML(c=10;e=23;e=0,1), SVMML(c=10;e=24;e=0,1), SVMML(c=10;e=25;e=0,1), SVMML(c=10;e=26;e=0,1), SVMML(c=10;e=27;e=0,1), SVMML(c=10;e=28;e=0,1), SVMML(c=10;e=29;e=0,1), SVMML(c=10;e=30;e=0,1), SVMML(c=10;e=31;e=0,1), SVMML(c=10;e=32;e=0,1), SVMML(c=10;e=33;e=0,1), SVMML(c=10;e=34;e=0,1), SVMML(c=10;e=35;e=0,1), SVMML(c=10;e=36;e=0,1), SVMML(c=10;e=37;e=0,1), SVMML(c=10;e=38;e=0,1), SVMML(c=10;e=39;e=0,1), SVMML(c=10;e=40;e=0,1), SVMML(c=10;e=41;e=0,1), SVMML(c=10;e=42;e=0,1), SVMML(c=10;e=43;e=0,1), SVMML(c=10;e=44;e=0,1), SVMML(c=10;e=45;e=0,1), SVMML(c=10;e=46;e=0,1), SVMML(c=10;e=47;e=0,1), SVMML(c=10;e=48;e=0,1), SVMML(c=10;e=49;e=0,1), SVMML(c=10;e=50;e=0,1).



Analogicznie jak dla wcześniej analizowanego zbioru Yeast, metoda SONN wykazuje wyjątkową skuteczność dla wzorców, które zostały wykorzystane w procesie treningu oraz najgorsze rezultaty dla wzorców nieznanymi. Warto zauważyć, że wspólnymi cechami zbiorów Yeast oraz Arrhythmia są stosunkowo duża ilość klas problemu klasyfikacji oraz wielkość zbioru.

Z powodu dużej ilości klas dla zbioru Arrhythmia można zauważyć, że osiem z czternastu odmian instancji algorytmu SVM nie utworzyło klasyfikatora w dopuszczalnym czasie.

## 6.8. Porównania dla zbioru Arrhythmia

Tablica 6.25: Arrhythmia, dane znormalizowane, test na całym zbiorze treningowym

Algorytm	Q_COR	Q_ERR	Q_SSE	Q_MSE	Q_RMSE	Q_ARV	E_TTR	E_TCL	E_TOT	E_MEM
SONN-3	100,00	0,00	341,52	0,76	0,87	5201,95	3730	0,31	3870	2877487
KNN(k=1)	100,00	0,00	0,00	0,00	0,00	NaN	2	3,26	1477	1022801
KNN(k=3)	51,33	48,67	82,67	0,18	0,43	2,20	0	3,32	1500	1022801
KNN(k=5)	49,78	50,22	112,64	0,25	0,50	2,61	0	3,67	1660	1022801
KNN(k=10)	48,45	51,55	139,75	0,31	0,56	2,94	0	4,37	1975	1022801
KNN(k=15)	45,13	54,87	152,59	0,34	0,58	3,13	0	4,97	2247	1022801
KNN(k=30)	44,91	55,09	168,43	0,37	0,61	3,42	0	7,22	3262	1022801
SVML(c=1;e=0,1)	85,84	14,16	280,82	0,62	0,79	159,73	37023	0,07	37056	1193794
SVML(c=10;e=0,1)	94,47	5,53	291,21	0,64	0,80	180,57	577838	0,04	577858	1193794
SVML(c=50;e=0,1)	97,79	2,21	305,82	0,68	0,82	204,68	3835668	0,04	3835685	1193794
SVML(c=100;e=0,1)	97,79	2,21	300,50	0,66	0,82	203,98	3121304	0,04	3121322	1193794
SVML(c=20;e=0,1)	95,58	4,42	297,01	0,66	0,81	150,96	1643353	0,04	1643373	1193794
SVML(c=20;e=0,01)	95,80	4,20	285,71	0,63	0,80	139,82	598151	0,05	598173	1193794
SVML(c=20;e=0,001)	96,46	3,54	291,17	0,64	0,80	164,99	1580671	0,05	1580694	1193794
SVMR(c=10;s=0,5;e=0,1)	100,00	0,00	57,70	0,13	0,36	29,72	824972	16,02	832212	1164754
SVMR(c=10;s=1;e=0,1)	100,00	0,00	114,87	0,25	0,50	13,35	289301	11,19	294360	1164754
SVMR(c=10;s=5;e=0,1)	82,96	17,04	260,05	0,58	0,76	88,33	165593	5,91	168266	1164754
SVMR(c=10;s=10;e=0,1)	73,89	26,11	258,78	0,57	0,76	52,20	22848	5,42	25297	1164754
SVMP(c=10;d=3;e=0,01)	nieakceptowalny czas treningu									
SVMP(c=10;d=4;e=0,01)	nieakceptowalny czas treningu									
SVMP(c=10;d=6;e=0,01)	nieakceptowalny czas treningu									
RBFN(k=5;f=2;c=1)	58,63	41,37	229,19	0,51	0,71	11,70	75	0,32	218	1020681
RBFN(k=10;f=2;c=1)	61,73	38,27	235,60	0,52	0,72	18,13	25	0,67	327	1021226
RBFN(k=15;f=2;c=1)	67,26	32,74	226,20	0,50	0,71	18,10	39	0,96	474	1021771
RBFN(k=20;f=2;c=1)	67,92	32,08	225,54	0,50	0,71	19,19	53	1,29	634	1022316
RBFN(k=5;f=1;c=2)	57,08	42,92	229,11	0,51	0,71	11,53	228	0,29	361	1031907
RBFN(k=10;f=1;c=2)	63,05	36,95	236,10	0,52	0,72	17,78	267	0,59	532	1043662
RBFN(k=15;f=1;c=2)	63,72	36,28	232,80	0,52	0,72	19,39	715	0,90	1122	1055417
RBFN(k=20;f=1;c=2)	65,93	34,07	230,11	0,51	0,71	21,06	620	1,23	1176	1067172
RBFN(k=5;f=2;c=2)	59,29	40,71	231,96	0,51	0,72	13,15	13	0,32	157	1020703
RBFN(k=10;f=2;c=2)	60,62	39,38	237,97	0,53	0,73	17,28	25	0,67	329	1021248
RBFN(k=15;f=2;c=2)	67,04	32,96	224,57	0,50	0,70	17,54	41	0,98	484	1021793
RBFN(k=20;f=2;c=2)	66,37	33,63	227,64	0,50	0,71	19,53	54	1,31	646	1022338
PNN(sigma=0,1)	100,00	0,00	0,00	0,00	0,00	1,00	0	8,32	3761	1025513
PNN(sigma=1)	74,34	25,66	275,27	0,61	0,78	14,60	0	8,16	3687	1025513
PNN(sigma=10)	54,42	45,58	384,19	0,85	0,92	4,17E+06	0	8,39	3790	1025513
PNN(sigma=100)	53,98	46,02	385,13	0,85	0,92	4,24E+10	0	8,11	3664	1025513





## 6.9. Skalowalność - porównanie wybranych wskaźników

W tym podrozdziale przeanalizowane zostaną wartości wskaźników związanych ze skalowalnością algorytmów. Na wykresach 6.1 oraz 6.2 przedstawione są wielkości wskaźników E\_TOT oraz E\_MEM. Spośród testowanych algorytmów wybrano po jednym, konkretnie sparametryzowanym, który osiągnął najlepsze wartości Q\_COR dla większości zbiorów danych. Wykaz algorytmów przedstawianych na poniższych wykresach znajduje się w tabeli 6.28. Omawiana analiza dotyczy wartości uzyskanych dla walidacji na całym zbiorze testowym.

Tablica 6.28: Algorytmy zastosowane podczas analizy skalowalności

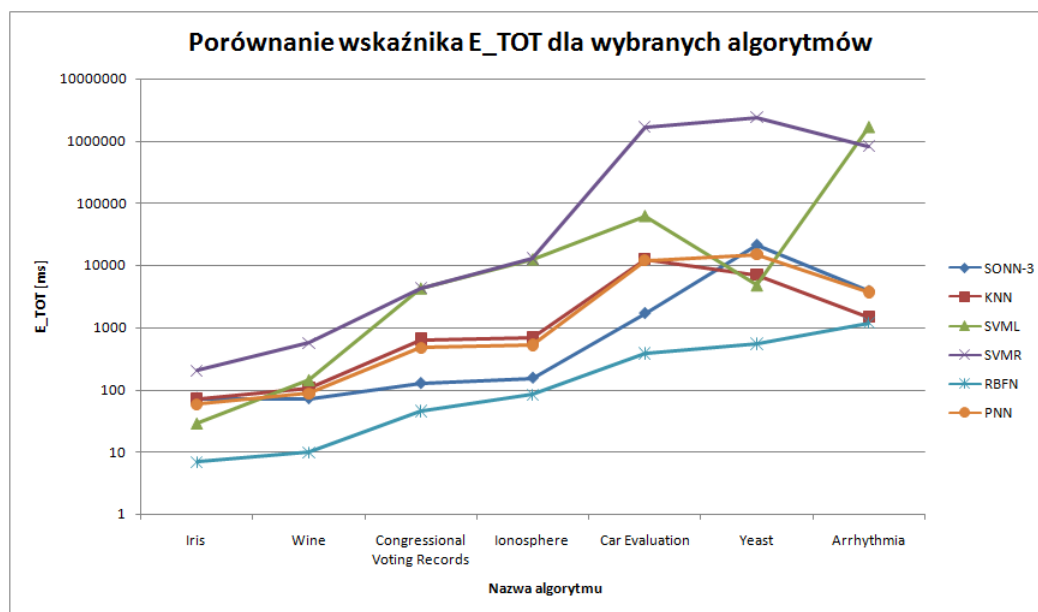
Oznaczenie	Algorytm	Parametry
SONN-3	Self Optimizing Neural Networks	
KNN	K nearest neighbours	$k = 3$
SVML	Support Vector Machines	klasyfikator liniowy, $cost = 20$ , $eps = 0,01$
SVMR	Support Vector Machines	klasyfikator nieliniowy, $cost = 10$ , $sigma = 0,5$ , $eps = 0,1$
RBFN	Radial Basis Function Networks	funkcja radialna: wielokwadratowa, metoda wyboru centrów: K-mean clustering, $k = 20$

Na poziomych osiach obu wykresów widnieją algorytmy posortowane według indeksu, który obrazuje stopień trudności zbioru danych. Jest on obliczany jako:

$$\text{Stopień trudności} = \text{ilość cech} * \text{ilość wektorów treningowych} * \text{liczba klas}$$

Warto zauważyć, że tak obliczany indeks nie uwzględnia wszystkich czynników wpływających na rzeczywisty stopień trudności zbioru. Co więcej, stopień trudności mógłby być inaczej określany dla różnych algorytmów, ponieważ każdy ma inną specyfikę i wykorzystuje inne własności danych zbioru. Niemniej jednak, tak zdefiniowany indeks pozwolił na zadowalające oszacowanie stopnia trudności zbiorów danych dla analizowanych algorytmów i wskaźników.

Na osiach pionowych przedstawiono wielkości dla każdego ze wskaźników. Warto zauważyć, że wartości na osi pionowej układają się w skali logarytmicznej.



Rysunek 6.1: Porównanie wskaźnika E\_TOT dla wybranych algorytmów

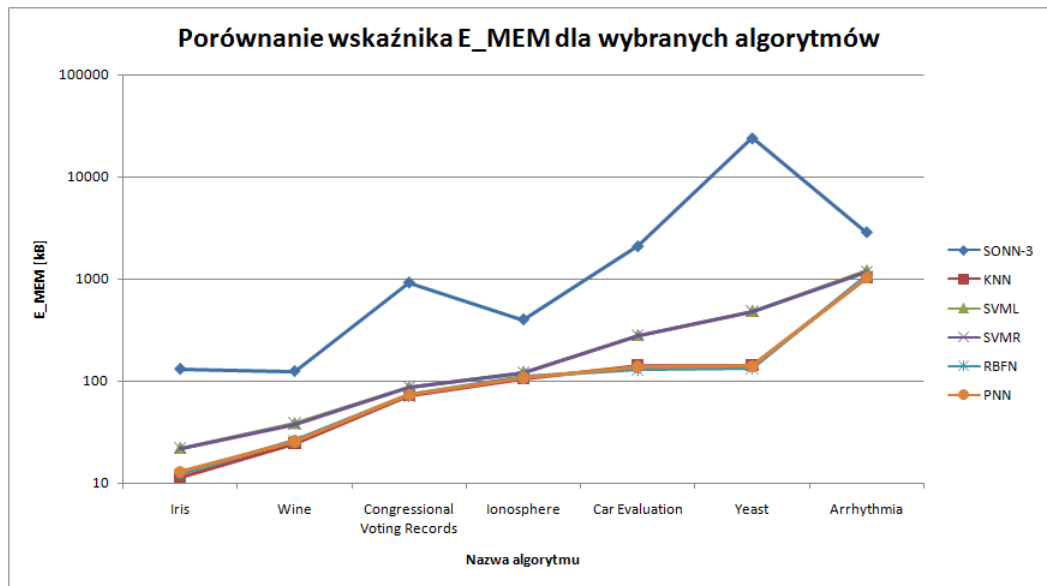
Na wykresie 6.1 przedstawiono wartości całkowitego czasu operacji E\_TOT. Istotne zależności, które mogą zostać zaobserwowane na przedstawionym wykresie, to:

- Metodą, która w najkrótszym czasie potrafi utworzyć klasyfikator okazała się metoda sieci radialnych funkcji bazowych (RBFN).
- Sieci SONN potrafią w stosunkowo krótkim czasie utworzyć klasyfikator w porównaniu z innymi metodami.
- Popularny algorytm SVM (w dwóch odmianach SVML i SVMR) wymagał więcej czasu w stosunku do pozostałych metod dla każdego ze zbiorów.
- Zaproponowany sposób obliczania współczynnika stopnia trudności wydaje się uwzględniać wystarczającą ilość informacji o zbiorze danych, dzięki czemu poprawnie udało się uporządkować zbiory według stopnia trudności.

Kolejnym istotnym z punktu widzenia skalowalności wskaźnikiem jest ilość pamięci potrzebnej do utworzenia obiektu klasyfikatora. Wartości wskaźnika E\_MEM znajdują się na wykresie 6.2. Analizując ten wykres można zauważyć, że:

- Algorytm SONN okazał się najbardziej wymagającym z punktu widzenia wymaganej pamięci.
- Zbiór Yeast jest szczególnie trudny dla sieci SONN (wskazują na to wysokie wartości wskaźników E\_MEM i E\_TOT).





Rysunek 6.2: Porównanie wskaźnika E\_MEM dla wybranych algorytmów

- Wymagana pamięć przez algorytmy SVM zależy od ilości klas występujących dla danego problemu klasyfikacji.

## 7. Podsumowanie

W ramach pracy udało się przeprowadzić porównania wybranych metod klasyfikacji z sieciami SONN. Porównania 37 algorytmów dla 7 zbiorów danych, przy użyciu najpopularniejszych metod walidacji stanowią cenny zbiór informacji o sieciach SONN oraz analizowanych algorytmach. Dzięki tym informacjom udało się zauważyć kilka interesujących wniosków, które wymienione zostaną poniżej.

Oprócz samych porównań dużą wartością pracy jest aplikacja wspomagająca porównywanie metod klasyfikacji, dzięki której możliwa jest wygodna analiza algorytmów klasyfikacji dla dowolnych zbiorów danych, najbardziej popularnych metod walidacji i wstępnego przetwarzania danych. Obok aplikacji warto wymienić także bibliotekę algorytmów, która może z łatwością zostać zastosowana w dowolnym projekcie, co umożliwi stosowanie zaimplementowanych algorytmów w praktycznych zagadnieniach.

### 7.1. Wnioski

Dzięki uzyskanemu zestawowi porównań udało się wyciągnąć następujące wnioski:

- Sieci SONN są obiecującym i ciekawym narzędziem klasyfikacji. Dla większości zbiorów danych algorytm SONN w stosunkowo krótkim czasie potrafił skonstruować strukturę sieci, która na ogół potrafiła realizować zadanie klasyfikacji z porównywalnym lub niewiele gorszym wynikiem niż popularne metody klasyfikacji.
- Ciekawą własnością sieci SONN jest osiąganie stuprocentowej skuteczności dla walidacji na całym zbiorze danych treningowych, który nie posiada wzorców sprzecznych. Jest to unikatowa własność wśród badanych metod klasyfikacji (można było ją zaobserwować również dla algorytmu kNN, dla  $k = 1$ , który jednak stanowi raczej ciekawy, szczególny przypadek tej metody, niż stosowane narzędzie klasyfikacji).
- W porównaniu objętości pamięciowej, jaka potrzebna jest do utworzenia obiektu klasyfikatora sieci SONN wykazywały najgorsze właściwości spośród badanych algorytmów. Wynika to z pewnych założeń optymalizacyjnych algorytmu SONN, dzięki którym potrafią one w krótkim czasie utworzyć klasyfikator (zmniejszenie złożoności obliczeniowej

kosztem zwiększenia zużycia pamięci). Mimo możliwości aktualnie produkowanych komputerów i nieustannej poprawie wydajności podzespołów wymagania pamięciowe sieci SONN są na tyle wysokie, że powinny stać się przedmiotem dokładnej analizy i poprawy, jeżeli algorytm ten ma być stosowany do większych niż analizowane w pracy zbiory danych.

- Sieci SONN jako jedyne spośród analizowanych algorytmów wykazują odporność na znaczne różnice w wielkościach cech (np. zbiór Wine). Pozwala to na stosowanie tego algorytmu bez konieczności wstępnego przetwarzania danych (np. normalizacja).
- Realizując porównania metod dla różnych zbiorów danych, zauważono, że niektóre metody klasyfikacji przy nieumiejętnym dobraniu parametrów nie potrafią utworzyć klasyfikatora z powodu błędu (RBFN) lub zbyt długiego czasu treningu (SVM). Brak parametrów dla algorytmu SONN jest jedną z największych zalet tego algorytmu.
- Podczas porównań niektórych metod klasyfikacji (np. Knn, PNN) warto analizować oprócz czasu treningu także czas klasyfikacji, gdyż dla tych algorytmów główne obliczenia wykonywane są na etapie klasyfikacji.
- Popularny algorytm SVM, który jest powszechnie stosowany w praktycznych zastosowaniach osiąga najlepsze rezultaty w zakresie czasu treningu dla zbiorów danych o małej liczbie klas. Jest to związane z koniecznością tworzenia wielu klasyfikatorów binarnych dla problemów wieloklasowych, co rzadko opisywane jest w publikacjach i pracach badawczych.
- Najtrudniejszym zbiorem danych dla sieci SONN (z punktu widzenia wskaźników E\_MEM i E\_TOT oraz Q\_COR) okazał się zbiór Yeast. Cechą, która wyróżnia ten zbiór spośród innych, dla których przeprowadzono porównania, jest bardzo nierównomierny rozkład wzorców dla klas problemu. Podobną własność można zaobserwować dla zbioru Arrhythmia, dla którego sieci SONN również nie osiągnęły zadowalającej skuteczności.
- Dla zbioru Car Evaluation, który składa się z wartości symbolicznych (zamienionych na kolejne liczby naturalne) metoda SONN uzyskała stosunkowo dobre wskaźniki jakości uogólnienia. Powinno to zostać wzięte pod uwagę podczas tworzenia wersji SONN-4, która według planów autora metody, ma umożliwiać przeprowadzanie klasyfikacji dla zbiorów zawierających wartości symboliczne.

## 7.2. Możliwości rozwoju pracy

Celem pracy było porównanie sieci SONN z popularnymi metodami klasyfikacji. Dzięki przeprowadzeniu odpowiednich porównań udało się zaobserwować następujące aspekty działania algorytmu, które powinny stać się przedmiotem dokładniejszej analizy i próby poprawy:

- Duże wymagania pamięciowe metody - jest to najbardziej wyraźna wada sieci SONN w porównaniu z innymi metodami klasyfikacji. Ponad dziesięciokrotnie większe wymagania pamięciowe w stosunku do innych metod mogą uniemożliwić stosowanie metody SONN dla większych niż analizowane zbiorów danych.
- Niska skuteczność dla zbiorów Yeast oraz Arrhythmia - w pierwszej kolejności warto zbadać wpływ rozkładu wzorców według klas na zdolności generalizacyjne sieci SONN.
- Skuteczność algorytmu na ogół jest nieco gorsza od innych metod klasyfikacji - możliwe, że algorytm SONN starając się utworzyć strukturę, która ze stuprocentową skutecznością rozpoznaje wzorce treningowe traci zdolność do uogólnienia (przeuczenie).

Zrealizowana w ramach pracy aplikacja wspomagająca porównania okazała się bardzo pomocnym i funkcjonalnym narzędziem. W ramach jej rozwoju na pewno warto uwzględnić:

- Funkcjonalność wizualizacji utworzonych klasyfikatorów - aktualnie podawane są tylko informacje o strukturze klasyfikatora w formie tabelarycznej.
- Zwiększenie interaktywności z użytkownikiem - rozbudowa systemu pomocy oraz dodanie kilku udogodnień takich, jak szacowanie czasu przebiegu symulacji, czy pokazywanie pozostałego czasu symulacji.

Naturalnie, można również wzbogacić pracę o kolejne porównania uwzględniające nowe zbiory danych, inne lub inaczej sparametryzowane algorytmy, metody wstępnego przetwarzania danych, czy metody walidacji. Dzięki funkcjonalnej i ogólnodostępnej aplikacji, która została zaimplementowana i wykorzystana w ramach pracy możliwe jest dokonywanie nowych porównań metod klasyfikacji, dzięki którym - zgodnie z intencją autora - możliwy będzie rozwój metod klasyfikacji, a w szczególności algorytmu *Self Optimizing Neural Networks*.

## A. Zawartość dołączonej płyty CD-ROM

Na dołączonej płycie CD-ROM znajdują się następujące pliki oraz foldery:

- plik LukaszDziedzia2009.pdf - dokument pracy magisterskiej zapisany w formacie .pdf.
- folder LukaszDziedzia2009 - źródła, które pozwalają na edycję oraz utworzenie dokumentu:
  - latex - pliki  $\LaTeX$ , które stanowią treść niniejszego dokumentu.
  - symulacje - wyniki symulacji opisanych w pracy.
- Projekt - kod źródłowy zrealizowany jako część praktyczna pracy wraz z wszystkimi referencjami niezbędnymi do uruchomienia.

# Spis rysunków

1.1	Etapy pracy nad aplikacją, główne cele . . . . .	9
2.1	Dwuetapowość procesu klasyfikacji . . . . .	14
3.1	Koncepcja działania metody kNN, problem . . . . .	18
3.2	Koncepcja działania metody kNN, rozwiązanie . . . . .	18
3.3	Koncepcja działania metody SVM, liniowo separowalny (dwuklasowy) zbiór danych . . . . .	20
3.4	Koncepcja działania metody SVM, przykłady możliwych rozwiązań . . . . .	21
3.5	Koncepcja działania metody SVM, klasyfikator liniowy z największym marginesem . . . . .	21
3.6	Koncepcja działania metody SVM, dane nieseparowalne liniowo . . . . .	23
3.7	Koncepcja działania metody SVM, „Kernel Trick” . . . . .	23
3.8	Schemat architektury Radial Basis Function Network . . . . .	26
3.9	Schemat architektury Probabilistic Neural Network . . . . .	28
3.10	Przykładowa struktura sieci SONN utworzona dla zbioru Wine (program autorstwa dr Adriana Horzyka) . . . . .	32
3.11	Self Optimizing Neural Network, schemat LBN . . . . .	33
3.12	Self Optimizing Neural Network, schemat ARN . . . . .	33
3.13	Self Optimizing Neural Network, schemat MSN . . . . .	34
4.1	Schemat procesu porównywania metod klasyfikacji . . . . .	39
4.2	Metody walidacji, zbiorów danych . . . . .	48
4.3	Metody walidacji, k-fold cross validation (k=2) . . . . .	49
4.4	Metody walidacji, k-fold cross validation (k=5) . . . . .	49
4.5	Metody walidacji, LOOCV . . . . .	50
4.6	Metody walidacji, Podział procentowy (70:30) . . . . .	50
4.7	Macierz niepewności dla jednego etapu walidacji krzyżowej algorytmu kNN na zbiorze Wine . . . . .	53

---

5.1	Aplikacja - wybór treningowego zbioru danych . . . . .	64
5.2	Aplikacja - wybór metod przetwarzania danych . . . . .	65
5.3	Aplikacja - wybór algorytmów . . . . .	65
5.4	Aplikacja - konfiguracja symulacji . . . . .	66
5.5	Aplikacja - wyniki . . . . .	67
6.1	Porównanie wskaźnika E_TOT dla wybranych algorytmów . . . . .	96
6.2	Porównanie wskaźnika E_MEM dla wybranych algorytmów . . . . .	97

# Spis tablic

3.1	Parametry metody kNN . . . . .	19
3.2	Parametry metody SVM . . . . .	24
3.3	Parametry metody RBFN . . . . .	27
3.4	Parametry metody PNN . . . . .	29
4.1	Charakterystyka zbioru danych Iris . . . . .	41
4.2	Charakterystyka zbioru danych Wine . . . . .	42
4.3	Charakterystyka zbioru danych Yeast . . . . .	43
4.4	Charakterystyka zbioru danych Congressional Voting Records . . . . .	43
4.5	Charakterystyka zbioru danych Car Evaluation . . . . .	44
4.6	Charakterystyka zbioru danych Ionosphere . . . . .	44
4.7	Charakterystyka zbioru danych Arrhythmia . . . . .	45
6.1	Parametry techniczne komputera, na którym zrealizowano porównania . . . . .	69
6.2	Legenda dla algorytmów i parametrów występujących w wynikach porównań . . . . .	70
6.3	Struktura sieci SONN dla wybranych metod przetwarzania danych, Wine . . . . .	71
6.4	Wine, brak wstępnego przetwarzania, test na całym zbiorze treningowym . . . . .	71
6.5	Wine, dane znormalizowane, test na całym zbiorze treningowym . . . . .	72
6.6	Wine, dane zestandaryzowane, test na całym zbiorze treningowym . . . . .	73
6.7	Wine, zastosowanie algorytmu PCA i normalizacji danych, test na całym zbiorze treningowym . . . . .	74
6.8	Wine, dane znormalizowane, 10-fold cross-validation, trening . . . . .	75
6.9	Wine, dane znormalizowane, 10-fold cross-validation, walidacja . . . . .	76
6.10	Iris, dane znormalizowane, test na całym zbiorze treningowym . . . . .	77
6.11	Iris, dane znormalizowane, 10-fold cross-validation, trening . . . . .	78
6.12	Iris, dane znormalizowane, 10-fold cross-validation, walidacja . . . . .	79
6.13	Yeast, dane znormalizowane, test na całym zbiorze treningowym . . . . .	80
6.14	Yeast, dane znormalizowane, 10-fold cross-validation, trening . . . . .	81



---

6.15	Yeast, dane znormalizowane, 10-fold cross-validation, walidacja . . . . .	82
6.16	Congressional Voting Records, dane znormalizowane, test na całym zbiorze treningowym . . . . .	83
6.17	Congressional Voting Records, dane znormalizowane, 10-fold cross-validation, trening . . . . .	84
6.18	Congressional Voting Records, dane znormalizowane, 10-fold cross-validation, walidacja . . . . .	85
6.19	Car Evaluation, dane znormalizowane, test na całym zbiorze treningowym . . .	86
6.20	Car Evaluation, dane znormalizowane, 10-fold cross-validation, trening . . . .	87
6.21	Car Evaluation, dane znormalizowane, 10-fold cross-validation, walidacja . . .	88
6.22	Ionosphere, dane znormalizowane, test na całym zbiorze treningowym . . . . .	89
6.23	Ionosphere, dane znormalizowane, 10-fold cross-validation, trening . . . . .	90
6.24	Ionosphere, dane znormalizowane, 10-fold cross-validation, walidacja . . . . .	91
6.25	Arrhythmia, dane znormalizowane, test na całym zbiorze treningowym . . . . .	92
6.26	Arrhythmia, dane znormalizowane, 10-fold cross-validation, trening . . . . .	93
6.27	Arrhythmia, dane znormalizowane, 10-fold cross-validation, walidacja . . . . .	94
6.28	Algorytmy zastosowane podczas analizy skalowalności . . . . .	95

## Bibliografia

- [1] praca zbiorowa, "How much information? 2003," 2003. [http://www2.sims.berkeley.edu/research/projects/how-much-info-2003/printable\\_report.pdf](http://www2.sims.berkeley.edu/research/projects/how-much-info-2003/printable_report.pdf).
- [2] "Przeglądarka internetowa clusty." <http://Clusty.org/about>.
- [3] Łukasz Dziedzia, "Symulacja procesu uczenia sztucznej sieci neuronowej," 2006. (projekt zrealizowany w ramach przedmiotu „Programowanie systemów czasu rzeczywistego”).
- [4] D. Bridge, "Topics in Artificial Intelligence." <http://www.cs.ucc.ie/~dgb/courses/tai/notes/handout4.pdf>.
- [5] C. J. Bugres, *A Tutorial on Support Vector Machines for Pattern Recognition*, p. 121–167. Data Mining and Knowledge Discovery, Kluwer Academic Publishers, 1998.
- [6] "<http://www.support-vector.net/icml-tutorial.pdf>."
- [7] A. W. Moore, "Support vector machines," 2001. <http://www.autonlab.org/tutorials/svm15.pdf>.
- [8] J. S.-T. Nello Cristianini, "An introduction to support vector machines," 2000.
- [9] J. C. Platt, "Sequential minimal optimization: A fast algorithm for training support vector machines," 1998. <http://home.mit.bme.hu/~gtakacs/download/platt1998.pdf>.
- [10] J. C. Platt, "Fast training of support vector machines using sequential minimal optimization," 2000.
- [11] "<http://omega.albany.edu:8008/machine-learning-dir/notes-dir/ker1/ker1-1.html>."
- [12] K.-B. Duan and S. S. Keerthi, "Which is the best multiclass SVM method? an empirical study," 2005.

- [13] U. Halici, "Artificial neural networks, chapter IX."
- [14] P. Schrater, "Radial basis function networks." [http://gandalf.psych.umn.edu/users/schrater/schrater\\_lab/courses/PattRecog07/RegressionIII.pdf](http://gandalf.psych.umn.edu/users/schrater/schrater_lab/courses/PattRecog07/RegressionIII.pdf).
- [15] M. Krętowska, "Sztuczne sieci neuronowe, ocena jakości sieci neuronowej, sieci RBF."
- [16] S. Osowski, *Sieci Neuronowe w ujęciu algorytmicznym*. Warszawa: PWNT, 1996.
- [17] N. Y. Nikolaev, "Probabilistic neural networks."
- [18] F. Gorunescu, "Benchmarking probabilistic neural network algorithms." <http://inf.ucv.ro/~aidc/proceedings/2006/1%20fgorunescu.pdf>.
- [19] D. F. Specht, "Probabilistic neural networks for classification, mapping or associative memory." <http://www.inf.ufrgs.br/~engel/Common/CMP121/PNN.pdf>.
- [20] "Internetowy podręcznik statystyki." [http://www.sixsigma.pl/textbook/stathome\\_stat.html?http%3A%2F%2Fwww.sixsigma.pl%2Ftextbook%2Fstneunet.html](http://www.sixsigma.pl/textbook/stathome_stat.html?http%3A%2F%2Fwww.sixsigma.pl%2Ftextbook%2Fstneunet.html).
- [21] A. Horzyk, "Self-optimizing neural network 3."
- [22] A. Horzyk, "Self optimizing neural networks sonn-3 for classification tasks."
- [23] A. Horzyk, "Self-optimizing neural networks sonn-2."
- [24] A. Horzyk and R. Tadeusiewicz, "Comparison of plasticity of self-optimizing neural networks and natural neural networks."
- [25] A. Horzyk, "Automatic discriminative lossy binary conversion of redundant real training data inputs for simplifying an input data space and data representation."
- [26] A. Horzyk, "Nowe metody uczenia sieci neuronowych bez sprzężeń zwrotnych," 1999. (praca doktorska napisana pod kierunkiem prof. zw. dr hab. inż. Ryszarda Tadeusiewicza).
- [27] J. Platt, "A resource-allocating network for function interpolation." *Neural Computation*, vol. 3, no. 2, pp. 213-225, (Summer 1991).
- [28] N. Jankowski, "Ontogeniczne sieci neuronowe w zastosowaniu do klasyfikacji danych medycznych," 1999. (praca doktorska pod kierunkiem prof. Włodzisława Duchy).
- [29] "Weka." <http://www.cs.waikato.ac.nz/ml/weka/>.
- [30] "Rapid Miner." <http://rapid-i.com/content/blogcategory/38/69/>.

- 
- [31] "UC Irvine Machine Learning Repository." <http://archive.ics.uci.edu/ml/index.html>.
- [32] L. I. Smith, "A tutorial on principal components analysis."
- [33] M. M. Robert C. Martin, "Agile, Programowanie zwinne," 2008.
- [34] J. Kerievsky, "Refaktoryzacja do wzorców projektowych," 2005.
- [35] "ADO.NET." <http://msdn.microsoft.com/en-us/library/e80y5yhx%28VS.80%29.aspx>.
- [36] "Microsoft Developer Network." <http://msdn.microsoft.com/>.
- [37] A. Kusiak, K. Kernstine, J. Kern, K. McLaughlin, and T. Tseng, "Data mining: Medical and engineering case studies," 2000. <http://www.icaen.uiowa.edu/~ankusiak/Res-in-Prog/IIE-0.pdf>.
- [38] T. Fawcett, "An introduction to ROC analysis," 2005.