

Politechnika Poznańska
Wydział Informatyki i Zarządzania
Instytut Informatyki

Praca dyplomowa magisterska

**OPTYMALIZACJA KLASYFIKATORA SVM ZA POMOCĄ PROGRAMOWANIA
GENETYCZNEGO**

Tomasz Ziętkiewicz

Promotor
dr hab. Krzysztof Krawiec

Poznań, 2013

Spis treści

1	Wprowadzenie	1
1.1	Cel i zakres pracy	1
1.2	Struktura pracy	1
2	Podstawy teoretyczne	2
2.1	Uczenie maszynowe	2
2.1.1	Systemy klasyfikujące	3
	Formalizacja problemu klasyfikacji	3
	Metody oceny skuteczności klasyfikacji	3
2.1.2	Miary skuteczności klasyfikacji	4
2.2	SVM - Maszyny wektorów wspierających	5
2.2.1	Klasyfikatory liniowe	5
2.2.2	Maszyny wektorów wspierających	8
	Funkcje jądrowe	10
2.2.3	Algorytmy wykorzystujące funkcje jądrowe	11
2.2.4	Klasyfikacja więcej niż dwóch klas	11
	Jeden przeciw wszystkim	11
	Klasyfikacja parami	11
2.3	Obliczenia ewolucyjne	12
	Selekcja	13
	Krzyżowanie	14
	Mutacja	14
	Genotyp	15
2.3.1	Programowanie genetyczne	15
	Inicjalizacja populacji	16
	Krzyżowanie	16
	Mutacja	16
	Selekcja	16
2.4	Optymalizacja parametrów SVM - przegląd literatury	18
2.4.1	Miary przystosowania (fitness)	19
2.4.2	Optymalizacja parametrów	19
2.4.3	Ewolucja kerneli	19
	Genetic Kernel SVM	20
	KTree	20
	Evolutionary Kernel Machine	21
	Kernel GP	21

3	Algorytm Kernel GP	23
3.1	Opis algorytmu	23
3.1.1	Inicjalizacja populacji	23
	Generowanie funkcji	24
3.1.2	Ewaluacja kerneli	25
3.1.3	Selekcja	26
3.1.4	Krzyżowanie i mutacja	26
3.1.5	Walidacja rozwiązania	26
3.2	Implementacja	27
3.3	Złożoność obliczeniowa	27
4	Wyniki działania algorytmu na popularnych zbiorach danych	28
4.1	Metodologia pomiarów	28
4.2	Opis zbiorów danych	29
4.3	Fitness	29
4.4	Wyniki klasyfikacji zbioru walidującego	30
4.4.1	Monotoniczność funkcji trafności	35
4.5	Czas wykonania	38
4.6	Użycie pamięci	38
4.7	Podsumowanie wyników	38
5	Case study - klasyfikacja danych ADHD 200	39
5.1	Opis zbioru danych	39
5.1.1	Surowe dane	39
5.1.2	Preprocessing	40
5.2	Konstrukcja i selekcja cech	40
5.3	Wyniki klasyfikacji	40
5.3.1	Kernel GP	40
5.3.2	Porównanie z innymi algorytmami	40
6	Podsumowanie	41
	Literatura	42
	Zasoby internetowe	44

Rozdział 1

Wprowadzenie

1.1 Cel i zakres pracy

Niniejsza praca ma dwa podstawowe cele:

- Stworzenie algorytmu programowania genetycznego optymalizującego parametry klasyfikatora SVM
- Zastosowanie stworzonego algorytmu do klasyfikacji danych ze zbioru ADHD-200

Realizacja drugiego z powyższych celów służyć ma przede wszystkim sprawdzeniu efektywności stworzonego algorytmu, ale jest też wyzwaniem samym w sobie. Zbiór danych ADHD-200 nie poddaje się łatwo klasyfikacji za pomocą metod uczenia maszynowego, dlatego każda poprawa wyników klasyfikacji względem wyników dotychczas osiągniętych będzie sporym sukcesem.

1.2 Struktura pracy

Struktura pracy jest następująca: rozdział drugi przedstawia ważniejsze zagadnienia teoretyczne związane z pracą oraz zawiera przegląd literatury. W rozdziale trzecim opisano zaimplementowany algorytm Kernel GP oraz przedstawiono sposób jego implementacji. Rozdział czwarty przedstawia wyniki działania algorytmu na standardowych zbiorach danych używanych do testowania algorytmów maszynowego uczenia. W rozdziale piątym prezentowane są wyniki działania algorytmu na zbiorze ADHD-200. Rozdział szósty zawiera podsumowanie.

Rozdział 2

Podstawy teoretyczne

2.1 Uczenie maszynowe

Uczenie maszynowe (ang. *Machine Learning*) to dziedzina informatyki zajmująca się konstruowaniem *systemów uczących się* [KS03]. Podstawową cechą takich systemów jest to, że potrafią one zmieniać sposób swojego działania w miarę jak napływają do nich kolejne dane. Zmiana działania systemu może mieć różną skalę - od zmiany pojedynczych parametrów programu, przez zapamiętywanie danych wejściowych po całkowitą zmianę wykonywanego algorytmu. Niezależnie od skali każda taka zmiana powinna mieć wpływ na jego przyszłe działanie i powinna mieć na celu uzyskanie jak najwyższej *oceny* pracy systemu. Jak ujmuje to Tom Mitchell [Mit97]:

System uczy się z doświadczenia E ze względu na pewną klasę zadań T i ocenę wykonania P jeśli ocena wykonania zadań należących do klasy T rośnie wraz z doświadczeniem E .

Systemy uczące się mają wiele zastosowań, między innymi [KS03]

- Rozpoznawanie mowy ludzkiej
- Rozpoznawanie tekstu pisanego (OCR, ang. *Optical Character Recognition*)
- Diagnostyka medyczna
- Klasyfikacja tekstów, np. na potrzeby filtrowania niechcianych wiadomości
- Automatyczna identyfikacja zagrożeń na podstawie obrazu z kamer przemysłowych
- Kierowanie autonomicznymi pojazdami
- Prognozowanie pogody
- Prognozowanie zmian kursów akcji na giełdzie
- Wykrywanie podejrzanych transakcji finansowych
- Biometria - identyfikacja ludzi na podstawie cech takich jak głos, wygląd twarzy, odciski palców, sposób chodzenia
- Wspomaganie podejmowania decyzji

2.1.1 Systemy klasyfikujące

Jednym z typów systemów uczących się są *systemy klasyfikujące* (inaczej *klasyfikatory*). Operują one na zbiorach *przykładów* opisanych za pomocą pewnego zbioru *atrybutów*. *Przykłady* (zwane też *obserwacjami*) reprezentują pewne obiekty, które różnią się od siebie wartościami atrybutów. Każdy przykład jest całkowicie scharakteryzowany przez swoje wartości atrybutów, co oznacza, że dwa przykłady o identycznych wartościach atrybutów są z punktu widzenia systemu klasyfikującego nieodróżnialne. Przykładem zbioru przykładów może być np. zbiór pacjentów, zaś zbiorem atrybutów zbiór cech takich jak wiek, płeć, wzrost, wyniki testów laboratoryjnych. Wśród zbioru atrybutów wyróżnia się jeden specjalny atrybut zwany atrybutem decyzyjnym (w odróżnieniu od pozostałych - atrybutów warunkowych) zwany też klasą lub etykietą obiektu. Zazwyczaj wartość tego atrybutu nie jest znana bezpośrednio, jej zdobycie stanowi pewną wartość. W przytoczonym przypadku pacjentów takim atrybutem może być na przykład diagnoza choroby. Uczenie systemu klasyfikującego polega na dostarczeniu do systemu zbioru przykładów z przypisanymi etykietami. Zbiór taki nazywamy *zbiorem trenującym / uczącym*. Na podstawie przykładów ze zbioru uczącego system wytwarza wewnętrzną reprezentację, która następnie umożliwia przypisanie nieznanym klasyfikatorowi etykiet/klas nowym przykładom, które nie występowały w zbiorze uczącym.

Formalizacja problemu klasyfikacji

W celu uściślenia dalszych rozważań konieczne jest wprowadzenie notacji formalnej opisującej problem klasyfikacji [KS03]. Zbiór wszystkich możliwych obiektów x , których dotyczy dany problem klasyfikacji, nazywany jest dziedziną i jest oznaczany przez U . Atrybut $a_i(x) : U \rightarrow V_{ai}$ to dowolna funkcja określona na dziedzinie U z przeciwdziedziną A_i . Zbiór wszystkich atrybutów oznaczamy przez $A = \{a_1, a_2, \dots, a_n\}$.

Każdy przykład $x \in X$ można opisać jako wektor w n -wymiarowej przestrzeni atrybutów Ω , czyli $\langle a_1(x), a_2(x), \dots, a_n(x) \rangle$.

Problem klasyfikacji polega na znalezieniu odwzorowania, które każdemu $x_i \in D$ (gdzie D to zbiór danych wejściowych) przypisuje jego klasę c_i . W przypadku klasyfikacji binarnej $c_i \in \{-1, 1\}$.

Metody oceny skuteczności klasyfikacji

W celu oceny skuteczności systemu należy za jego pomocą dokonać klasyfikacji przypadków, które nie były użyte podczas jego uczenia i których etykiety są znane (choć nie dostępne klasyfikatorowi). Wyniki klasyfikacji porównuje się z właściwymi etykietami i w ten sposób szacuje skuteczność klasyfikacji. W tym celu można wydzielić ze zbioru przykładów specjalny podzbiór, zwany *zbiorem testującym*, który jest używany do testowania a w fazie uczenia klasyfikator nie ma do niego dostępu. Czasami wydziela się też *zbiór walidujący*, który jest używany w trakcie uczenia w celu optymalizacji parametrów algorytmu. Stałego podziału zbioru przykładów na zbiór trenujący, testujący i walidujący można dokonać tylko wtedy, gdy zbiory te są wystarczająco liczne. W przeciwnym przypadku może okazać się, że nie są one wystarczająco reprezentatywne i na przykład rozkład przykładów z poszczególnych klasy jest mocno skrzywiony w którymś ze zbiorów. Aby tego uniknąć można posłużyć się metodą *k-krotnej walidacji krzyżowej*. Polega ona na podzieleniu zbioru na k podzbiorów i następnie powtarzanych k -razy fazach uczenia i testowania klasyfikatora, przy czym za każdym razem k -ty podzbiór służy jako zbiór testujący/walidujący a pozostałym podzbiory jak zbiór uczący. Skuteczności klasyfikacji oblicza się wtedy jako średnią sprawność osiąganą we wszystkich k testach.

2.1.2 Miary skuteczności klasyfikacji

Do oceny jakości klasyfikacji można używać różnych miar. W przypadku klasyfikacji binarnej (czyli kiedy rozróżniamy tylko dwie klasy przykładów) większość z nich można wyrazić za pomocą stosunku kilku z czterech wartości wyrażających liczbę przypadków klasyfikowanych w określony sposób. Wartości te są odnoszą się zawsze do jednej z klas, która jest w pewien sposób wyróżniona. Na przykład w przypadku diagnozy medycznej zazwyczaj taką klasą jest grupa osób chorych na jakąś chorobę. Przypadki zaklasyfikowane jako należące do tej klasy określane są jako zaklasyfikowane *pozytywnie* (ang. *positive*) natomiast przypadki zaklasyfikowane jako do niej nienależące jako zaklasyfikowane *negatywnie* (ang. *negative*). Słowa ang. "True" oraz ang. "False" odnoszą się odpowiednio do przypadków zaklasyfikowanych prawidłowo i nieprawidłowo:

- **True Positive (TP)** - liczba przypadków **poprawnie** zaklasyfikowanych jako **należące** do wyróżnionej klasy,
- **True Negative (TN)** - liczba przypadków **poprawnie** zaklasyfikowanych jako **nienależące** do wyróżnionej klasy,
- **False Positive (FP)** - liczba przypadków **niepoprawnie** zaklasyfikowanych jako **należące** do wyróżnionej klasy (inaczej błąd pierwszego rodzaju)
- **False Negative (FN)** - liczba przypadków **niepoprawnie** zaklasyfikowanych jako **nienależące** do wyróżnionej klasy (inaczej błąd drugiego rodzaju)

Poniżej zostały opisane miary, o których będzie mowa w dalszej części pracy. Wszystkie one zawierają się w przedziale $\langle 0, 1 \rangle$.

- **Precyzja** (ang. *precision*) - określa jaka część przypadków zaklasyfikowanych jako należące do wyróżnionej klasy rzeczywiście do niej należy. Dana jest wzorem:

$$precision = \frac{TP}{TP + FP}$$

- **Kompletność** (ang. *recall*) - określa jaka część przypadków należących do wyróżnionej klasy została prawidłowo zaklasyfikowana jako należące do niej. Dana jest wzorem:

$$recall = \frac{TP}{TP + FN}$$

- **Trafność (lub dokładność)** (ang. *Accuracy*) - stosunek liczby przypadków ze zbioru walidującego, które zostały zaklasyfikowane poprawnie do liczby wszystkich przypadków w zbiorze testującym. Może być wyrażona jako:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Miara F_1** (ang. *F_1 measure*) - miara uwzględniająca zarówno precyzję (ang. *precision*) jak i kompletność (ang. *recall*). Miara ta nie uwzględnia wartości *TN*. Jej wartość jest dana wzorem:

$$F_1 = 2 \times \frac{precision \times recall}{precision + recall}$$

- **MCC** ang. *Matthews correlation coefficient* - miara, która w przeciwieństwie do miary F_1 bierze pod uwagę wszystkie cztery wartości (*TP*, *TN*, *FP* i *FN*). Dana wzorem:

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

- *Średnie prawdopodobieństwo wyboru właściwej klasy* - niektóre klasyfikatory zamiast przypisywać każdemu z przykładów jedną z klas potrafią zwrócić dla każdego przykładu rozkład przynależności do wszystkich rozważanych klas. Jakość klasyfikacji można wtedy obliczyć jako uśrednioną po wszystkich przykładach wartość prawdopodobieństwa przypisanego klasie, do której przykład należy. Wartość taka może wahać się od wartości 0 (kiedy dla każdego przykładu do jego właściwej klasy zostało przypisane prawdopodobieństwo 0) do wartości 1 (kiedy dla każdego przykładu do jego właściwej klasy zostało przypisane prawdopodobieństwo 1).

W przypadku problemów, w których wyróżnia się $k > 2$ klas miara korzystająca z wartości TP, TN, FP i FN jest obliczana jako średnia wartość tej miary dla k problemów binarnych polegających na zaklasyfikowaniu przykładów jako należących lub nienależących do wybranej klasy.

2.2 SVM - Maszyny wektorów wspierających

Maszyna wektorów wspierających (SVM, ang. *Support Vector Machine*) to rodzaj klasyfikatora binarnego. Stanowi on rozszerzenie *klasyfikatora liniowego*, lecz w przeciwieństwie do niego jest w stanie poprawnie klasyfikować dane *nieseparowalne liniowo*. Jest to możliwe dzięki dokonywanej przez SVM transformacji danych do wyższych wymiarów za pomocą *funkcji jądrowych*.

2.2.1 Klasyfikatory liniowe

Jednym z najprostszych klasyfikatorów jest *klasyfikator liniowy*. Rozwiązuje on problem klasyfikacji binarnej poprzez znalezienie w przestrzeni atrybutów Ω hiperpłaszczyzny, która dzieli ją na dwie części odpowiadające dwóm klasom decyzyjnym: $\{-1, 1\}$.

Definicja 2.2.1 *Hiperpłaszczyzna w przestrzeni Ω to zbiór:*

$$\{x \in \Omega | \langle w, x \rangle + b = 0\}, w \in \Omega, b \in R \quad (2.1)$$

$\langle x, y \rangle$ oznacza iloczyn skalarny wektorów x i y :

$$\langle x, y \rangle = \sum_{i=1}^N [x]_i [y]_i$$

gdzie $[x]_i$ to i -ta wartość wektora x .

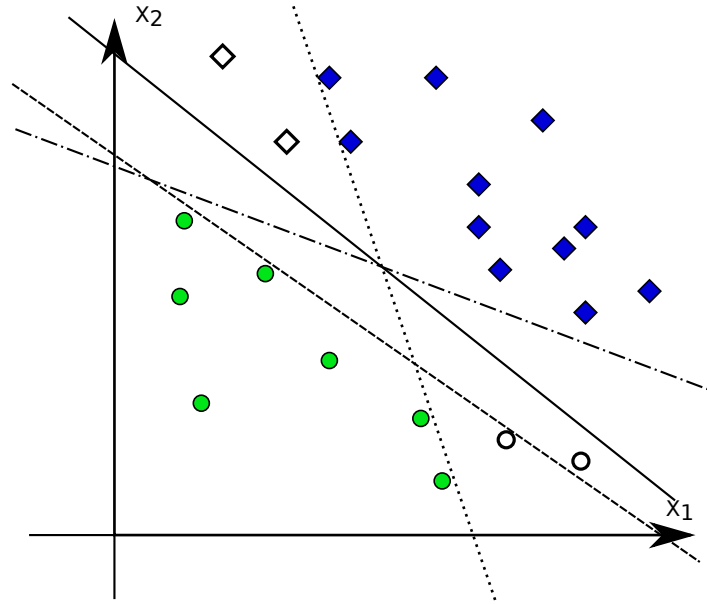
Wektor w we wzorze 2.1 to wektor wag, normalny do hiperpłaszczyzny, $\|w\|$ to norma euklidesowa tego wektora, czyli jego długość, a $b/\|w\|$ to odległość płaszczyzny od początku układu współrzędnych. Oba te parametry można razem dowolnie przeskalowywać, to jest pomnożyć w i b przez tę samą stałą zachowując tę samą hiperpłaszczyznę. Dlatego wprowadza się ograniczenie, po którego zastosowaniu otrzymujemy tak zwaną postać kanoniczną hiperpłaszczyzny:

Definicja 2.2.2 *Dla danego zbioru obserwacji $x_1, x_2, \dots, x_m \in \Omega$ wektor w i parametr b wyznaczają kanoniczną postać hiperpłaszczyzny jeśli:*

$$\min_{i=1..m} |\langle w, x_i \rangle + b| = 1 \quad (2.2)$$

Hiperpłaszczyzna dana wzorem 2.1 definiuje funkcję decyzyjną, która każdemu przypadkowi z Ω przypisuje klasę decyzyjną:

$$\begin{aligned} f_{w,b} : \Omega &\rightarrow \{\pm 1\} \\ x &\mapsto f_{w,b}(x) = \text{sgn}(\langle w, x \rangle + b) \end{aligned} \quad (2.3)$$



RYСУNEK 2.1: Hiperpłaszczyzny separujące dwa zbiory punktów w przestrzeni dwówymiarowej. Każda z nich poprawnie separuje punkty ze zbioru uczącego - zielone koła i niebieskie kwadraty. Przykłady ze zbioru testowego (puste kwadraty i kółka) są poprawnie separowane jedynie przez dwie proste (prosta narysowana linią ciągłą i prosta narysowana kropkami i kreskami).

Wynikiem uczenia klasyfikatora liniowego jest znalezienie hiperpłaszczyzny i odpowiadającej jej funkcji decyzyjnej, która przykładom ze zbioru uczącego $(x_i, y_i) \in \Omega$ przypisuje prawidłowe etykiety, czyli dla każdego (jeśli zbiór jest liniowo separowalny), lub dla jak największej liczby przypadków x_i zachodzi $f_w, b(x_i) = y_i$. Zazwyczaj kilka hiperpłaszczyzn równie dobrze rozdziela przypadki ze zbioru uczącego, mogą się jednak one różnić zdolnością do klasyfikacji zbioru testowego, co pokazano na rysunku 2.2.1. Optymalna hiperpłaszczyzna separująca to taka, która charakteryzuje się największym marginesem, czyli odległością hiperpłaszczyzny do najbliższych położonych obserwacji [SS02].

Definicja 2.2.3 Dla hiperpłaszczyzny danej wzorem $x \in \Omega | \langle w, x \rangle + b = 0$ oraz zbioru obserwacji $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ marginesem tego zbioru od hiperpłaszczyzny nazywamy minimalną odległość hiperpłaszczyzny od punktów z tego zbioru:

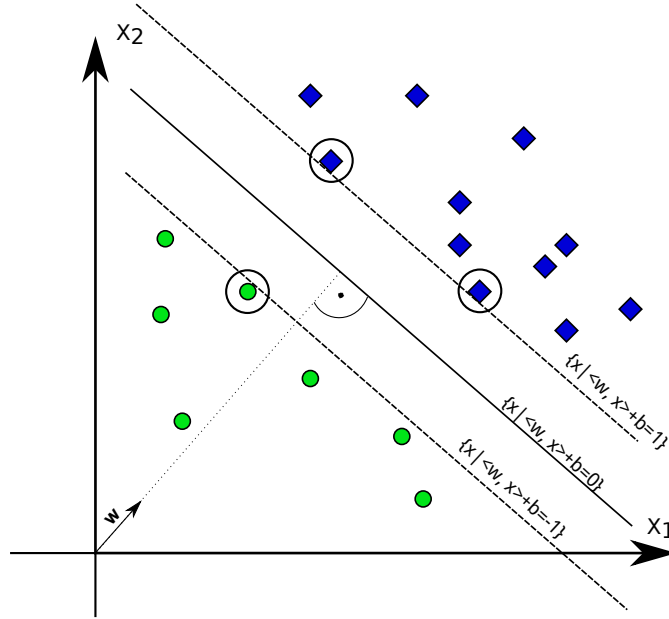
$$\rho_{w,b} := \min_{i=1..m} y_i \times (\langle w, x_i \rangle + b) / \|w\| \quad (2.4)$$

Żeby maksymalizować margines powinniśmy minimalizować $\|w\|$, zachowując warunek 2.2. Problem znalezienia optymalnej hiperpłaszczyzny separującej zbiór przykładów $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ można zatem zapisać jako problem optymalizacyjny:

$$\begin{aligned} \min_{w \in \Omega, b \in \mathbb{R}} \quad & \tau(w) = \frac{1}{2} \|w\|^2 \\ \text{p.o.} \quad & y_i (\langle w, x_i \rangle + b) \geq 1 \quad \text{dla } i = 1..m \end{aligned} \quad (2.5)$$

Ograniczenia w powyższym problemie zapewniają, że wartość funkcji decyzyjnej $f_{w,b}(x_i)$ będzie równa y_i , czyli, że hiperpłaszczyzna poprawnie odseparuje przykłady z dwóch grup. Osiągnięcie celu optymalizacji zapewnia znalezienie hiperpłaszczyzny o maksymalnym marginesie.

Powyższy problem programowania matematycznego jest podany w tak zwanej *formie prymalnej*. W praktyce rozwiązuje się wersję *dualną* problemu, która ma przyjmując następującą postać:



RYСУNEK 2.2: Hiperpłaszczyzna separująca dwa zbiory punktów w przestrzeni dwuwymiarowej wraz z marginesami. Przykłady w kółku to wektory podpierające.

Definicja 2.2.4 Wersja dualna problemu znalezienia optymalnej hiperpłaszczyzny:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^m} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \\ \text{p.o.} \quad & \alpha_i \geq 0, \text{ dla } i = 1..m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned} \quad (2.6)$$

Dla formy dualnej problemu optymalizacji funkcja decyzyjna przyjmuje postać:

$$f(x) = \text{sgn} \left(\sum_{i=1}^m y_i \alpha_i \langle x, x_i \rangle + b \right) \quad (2.7)$$

We wzorze 2.5 dla większości przykładów wartość współczynników α_i wyniesie 0. Przykłady te nie mają wpływu na wynik optymalizacji i otrzymaną hiperpłaszczyznę. Pozostałe przykłady, dla których $\alpha_i > 0$ noszą nazwę *wektorów wspierających*. Leżą one dokładnie na marginesie (na hiperpłaszczyźnie wyznaczonej przez margines, czyli w przypadku kanonicznej postaci hiperpłaszczyzny separującej w odległości $1/\|w\|$ od niej). Wektory wspierające są jedynymi przykładami ze zbioru trenującego, które należy zapamiętać w celu wyznaczenia hiperpłaszczyzny separującej. Ilość wektorów wspierających wyznacza złożoność hipotezy, którą stanowi otrzymana funkcja decyzyjna i pozwala oszacować górną granicę oczekiwanego prawdopodobieństwa błędu klasyfikacji. Górna granica prawdopodobieństwa błędu klasyfikacji klasyfikatora uczonego na zbiorze o wielkości $m-1$ równa się liczbie wektorów otrzymanych przy uczeniu na zbiorze o wielkości m podzielonej przez m .

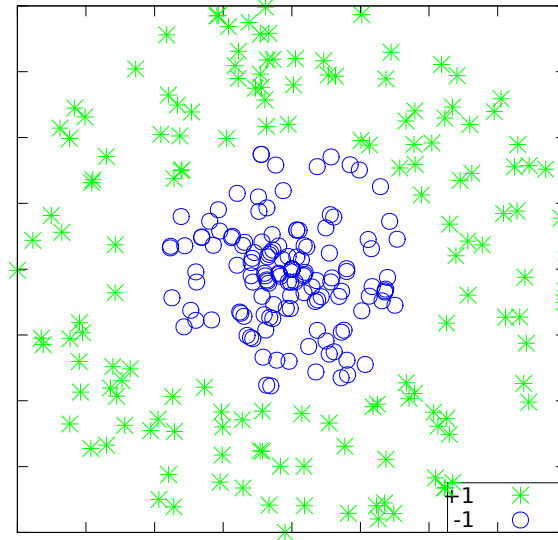
W przypadku, w którym dane nie są separowalne liniowo ze względu na szum, tzn. nie istnieje taka hiperpłaszczyzna, która dla danego zbioru przykładów spełnia ograniczenia ze wzoru 2.5, do problemu optymalizacyjnego wprowadza się tak zwane *zmienne osłabiające* (ang. *slack variables*) $\zeta_i, i = 1..m$, oraz stałą $C > 0$, która jest parametrem algorytmu. Wzór 2.5 przyjmuje wówczas postać:

$$\begin{aligned} \min_{w \in \Omega, b \in \mathbb{R}} \quad & \tau(w) = \frac{1}{2} \|w\|^2 + \frac{C}{m} \sum_{i=1}^m \zeta_i \\ \text{p.o.} \quad & y_i (\langle w, x \rangle + b) \geq 1 - \zeta_i \quad \text{dla } i = 1..m \\ & \zeta_i \geq 0 \quad \text{dla } i = 1..m \end{aligned} \quad (2.8)$$

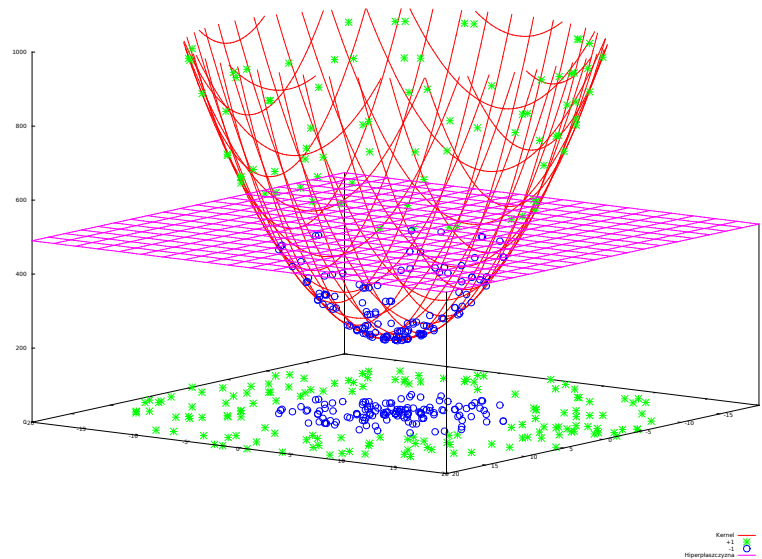
Natomiast wzór 2.6:

$$\begin{aligned} \max_{\alpha \in R^m} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \\ \text{p.o.} \quad & 0 \leq \alpha_i \leq \frac{C}{m}, \text{ dla } i = 1..m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned} \quad (2.9)$$

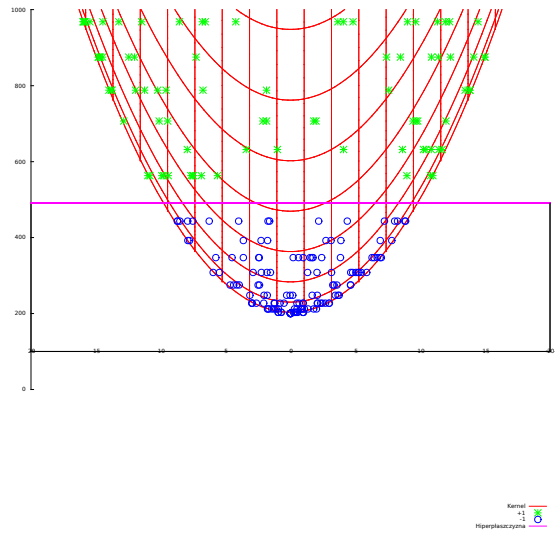
2.2.2 Maszyny wektorów wspierających



RYSUNEK 2.3: Zbiór danych opisanych w przestrzeni 2D. Klas nie da się odseparować za pomocą hiperpłaszczyzny



RYSUNEK 2.4: Zbiór z rysunku 2.2.2 po dokonaniu transformacji do przestrzeni 3-wymiarowej. Na wysokości 0 pokazano dane przed transformacją. W nowej przestrzeni klasy mogą być rozdzielone hiperpłaszczyzną.



RYСУNEK 2.5: Zbiór z rysunku 2.2.2 po dokonaniu transformacji do przestrzeni 3-wymiarowej, rzut wzdłuż osi X

Opisane powyżej klasyfikatory są w stanie klasyfikować tylko dane, które są separowalne liniowo. Jeśli dane cechuje pewna inna niż liniowa zależność między wartościami atrybutów decyzyjnych i warunkowych (na przykład taka jak przedstawiona na rysunku 2.2.2), to klasyfikator liniowy nie poradzi sobie z ich klasyfikacją. Aby rozwiązać ten problem można dokonać mapowania danych wejściowych do przestrzeni o większej liczbie wymiarów, w której te dane mogą stać się liniowo separowalne. Na rysunku 2.2.2 pokazano dane z rysunku 2.2.2 po transformacji do przestrzeni trój-wymiarowej za pomocą funkcji wielomianowej. Jak widać w nowej przestrzeni dane są separowalne przez hiperpłaszczyznę, co widać dokładnie na rysunku 2.2.2.

Mapowanie z oryginalnej przestrzeni atrybutów do nowej przestrzeni jest dokonywane przez funkcję $\Phi: \Omega^m \mapsto \Omega^{m+k}$, która dane n -wymiarowe zamienia w dane $n+k$ -wymiarowe. Na przykład mapowanie wektora $x = [x_1, x_2]$ z przestrzeni 2-wymiarowej do przestrzeni 3-wymiarowej może wyglądać następująco: $\Phi(x) = \Phi([x_1, x_2]) = [x_1^2, x_2^2, x_1 x_2]$. Aby zastosować mapowanie w klasyfikatorze wystarczy zamienić wszystkie wystąpienia wektorów wynikiem ich mapowań, czyli na przykład cel optymalizacji ze wzoru 2.6 przyjąłby postać:

$$\max_{\alpha \in R^m} W(\alpha) = \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \Phi(x_i), \Phi(x_j) \rangle \quad (2.10)$$

natomiast funkcja decyzyjna ze wzoru 2.7:

$$f(x) = \text{sgn} \left(\sum_{i=1}^m y_i \alpha_i \langle \Phi(x), \Phi(x_i) \rangle + b \right) \quad (2.11)$$

Powyższe podejście charakteryzuje się dużą złożonością obliczeniową. Zarówno obliczanie mapowania dla każdego przykładów jak i obliczenie iloczynu skalarnego przykładów w nowej przestrzeni, która może mieć bardzo dużo wymiarów, stanowi spory narzut obliczeniowy. Dokonywanie wszystkich mapowań a następnie obliczanie ich iloczynu nie jest jednak konieczne, można obie te operacje wykonać w jednym kroku, stosując tak zwany ang. *Kernel Trick* [BGV92]. Polega on na zastosowaniu

podstawienia $\langle \Phi(x), \Phi(x') \rangle \mapsto k(x, x')$ gdzie $k(x, x')$ to *funkcja jądrowa* (ang. *kernel*). Po zastosowaniu *kernel trick* otrzymujemy następujący problem optymalizacyjny:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^m} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{p.o.} \quad & \alpha_i \geq 0, \text{ dla } i = 1..m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned} \quad (2.12)$$

oraz następującą funkcję decyzyjną:

$$f(x) = \text{sgn} \left(\sum_{i=1}^m y_i \alpha_i k(x, x_i) + b \right) \quad (2.13)$$

Funkcje jądrowe

Funkcje jądrowe pozwalają obliczyć iloczyn skalarny wyników mapowania ($\Phi()$) dwóch wektorów (x, x') do wysoko wymiarowej przestrzeni atrybutów (Ω^{m+k}), bez bezpośredniego dokonywania mapowania:

$$\begin{aligned} \Phi : \Omega^m &\mapsto \Omega^{m+k} \\ \langle \Phi(x), \Phi(x') \rangle &\mapsto k(x, x') \end{aligned}$$

Najczęściej używane funkcje jądrowe to:

- Wielomianowa: $k(x, y) = \langle x, y \rangle^d$
- Gausowska (RBF): $k(x, y) = e^{-\sigma * \|x - y\|^2}$
- Sigmoidalna: $k(x, y) = \tanh(\gamma \langle x, y \rangle + \tau)$

Jako funkcję jądrową w SVM można użyć jednak dowolnej symetrycznej funkcji, która jest dodatnio określona (ang. *positive definite*), czyli dla której macierz:

$$K_{ij} := k(x_i, x_j)$$

spełnia warunek:

$$\sum_{i,j} c_i c_j K_{ij} \geq 0 \quad (2.14)$$

dla wszystkich $c_i \in \mathbb{R}$ co jest równoważne temu, że jej wszystkie wartości własne są dodatnie [SS02].

Funkcje jądrowe, które nie są dodatnio określone, a jedynie warunkowo dodatnio określone (ang. *conditionally positive definite*) również mogą być użyte w metodach jądrowych, nie gwarantują jednak, że znalezione rozwiązanie będzie unikalne. Wyniki eksperymentalne pokazują jednak, że w praktyce funkcje, które są jedynie warunkowo dodatnio określone, mogą dawać równie dobre wyniki co funkcje dodatnio określone [BTB05]. Funkcja jest warunkowo dodatnio określona, jeśli spełnia warunek ze wzoru 2.14 przy czym $\sum_{i,j} c_i = 0$.

Funkcje jądrowe cechuje własność domknięcia ze względu na pewne operacje. Oznacza to, że funkcje powstałe poprzez połączenie poprawnych funkcji jądrowych za pomocą tych operatorów również są poprawnymi funkcjami jądrowymi [STC04]. Jeżeli $k_1(x, z)$ i $k_2(x, z)$ są kernelami, to również poniższe funkcje są kernelami:

$$\begin{aligned} k(x, z) &= k_1(x, z) + k_2(x, z) \\ k(x, z) &= k_1(x, z) \times k_2(x, z) \\ k(x, z) &= a \times k_1(x, z) \\ k(x, z) &= e^{k_2(x, z)} \\ k(x, z) &= p(k_1(x, z)) \end{aligned} \quad (2.15)$$

gdzie $a \in \mathbb{R}^+$, a $p()$ to wielomian o dodatnich współczynnikach.

2.2.3 Algorytmy wykorzystujące funkcje jądrowe

SVM nie jest jedynym algorytmem wykorzystującym funkcję jądrową do dokonania mapowania danych do wysoko wymiarowych przestrzeni. Ich użycie jest zasadne wszędzie tam, gdzie zamiast działać wprost na wejściowych wektorach danych można operować na ich iloczynach skalarnych [SS02]. W związku z tym funkcje jądrowe znalazły zastosowania w wielu algorytmach, zwanych zbiorczo jak metody jądrowe (ang. *kernel methods*):

- *Kernel-NN* [YJZ02] - modyfikacja *k-NN* korzystająca z funkcji jądrowych do obliczania odległości między przykładami
- *k-PCA* [SS02] - modyfikacji algorytmu analizy głównych składowych *Principal Component Analysis* (PCA), który może służyć m.in. do konstrukcji cech opisujących przykłady, klasyfikowane potem za pomocą jakiegoś algorytmu klasyfikacji.
- *SVR* - *Support Vector Regression* - regresja za pomocą SVM
- *Kernel Fisher Discriminant* (KFD [SS02]) - zastosowanie funkcji jądrowych w liniowym dyskryminatorze Fishera, który może służyć do redukcji wymiarów opisujących przykłady i/lub do klasyfikacji.

2.2.4 Klasyfikacja więcej niż dwóch klas

Klasyfikator SVM jest ze swej natury klasyfikatorem binarnym, to znaczy potrafi separować jedynie dwie klasy przykładów. W praktyce jednak wiele problemów klasyfikacji dotyczy zbiorów, w których wyróżniono więcej niż dwie klasy. Aby zastosować SVM do rozwiązania takich problemów stosuje się zazwyczaj jedną z dwóch technik: *jeden przeciw wszystkim* (ang. One versus the Rest) oraz *klasyfikacja parami* (ang. pairwise classification).

Jeden przeciw wszystkim

Metoda *jeden przeciw wszystkim* polega na nauczaniu M klasyfikatorów binarnych, gdzie M oznacza liczbę klas. Każdy z nich separuje jedną klasę od pozostałych $M - 1$ klas, traktując te ostatnie jak jedną klasę. Klasyfikacja polega na zastosowaniu wszystkich M klasyfikatorów na testowanym przykładzie i wybraniu tego, który maksymalizuje wartość ze wzoru 2.13, przed zastosowaniem funkcji sgn , czyli:

$$\operatorname{argmax}_{j=1..M} g^j(x), \text{ gdzie } g^j(x) = \sum_{i=1}^m y_i a_i^j k(x, x_i) + b^j$$

. Jako wady tej metody wymienia się to, że pod-problemy przez nią rozwiązywane są niesymetryczne (jeśli w wyjściowym problemie dystrybucja przykładów do klas jest równomierna, to w każdym pod-problemie klasa agregująca "pozostałe" klasy jest $M - 1$ razy większa). Z problemem asymetryczności lepiej radzi sobie druga z wymienionych metod.

Klasyfikacja parami

Metoda *klasyfikacji parami* polega na stworzeniu wszystkich możliwych kombinacji klasyfikatorów binarnych: dla każdej z M klas tworzy się $M - 1$ klasyfikatorów oddzielających ją od poszczególnych klas, co daje łącznie $(M - 1) * M/2$ klasyfikatorów. Każdy z nich uczony jest tylko na przykładach należących do 2 wybranych klas. Klasyfikacja przykładu odbywa się poprzez głosowanie: każdy

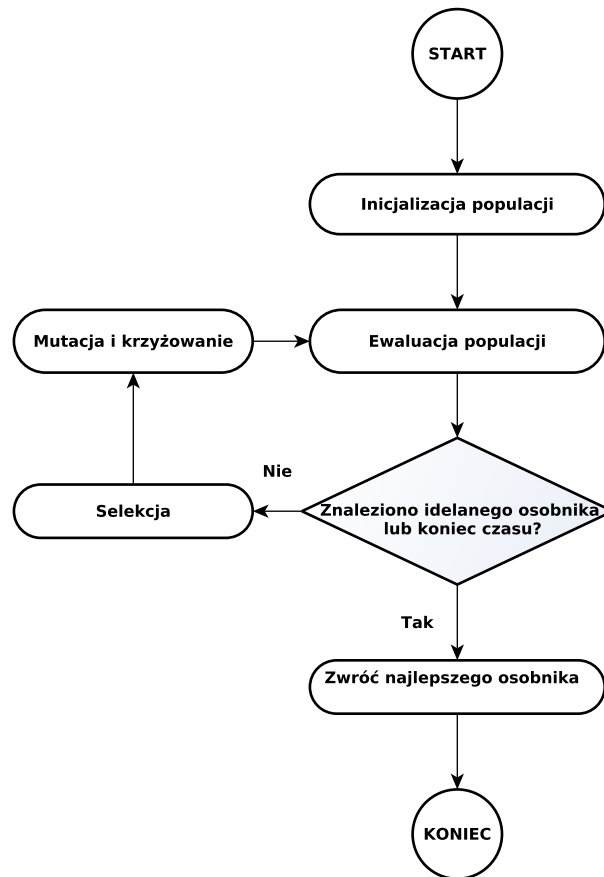
z $(M - 1) * M/2$ klasyfikatorów oddaje głos na klasę, do której przypisał dany przykład. Przykładowi przypisuje się klasę, która zdobyła największą liczbę głosów. Choć w tej metodzie liczba klasyfikatorów, które trzeba nauczyć a potem użyć do klasyfikacji jest znacznie większa niż w metodzie *jeden przeciw wszystkim*, to jednak uczenie odbywa się na mniejszej liczbie przykładów a klasyfikacja z użyciem mniejszej liczby wektorów wspierających, co może przyspieszyć sprawia, że etapy te przebiegają szybciej niż w metodzie *jeden przeciw wszystkim*. Ponieważ jednak liczba klasyfikatorów rośnie potęgowo względem liczby klas M , to dla dużych M *klasyfikacja parami* może okazać się wolniejsza niż *jeden przeciw wszystkim*.

2.3 Obliczenia ewolucyjne

Obliczenia ewolucyjne (ang. *Evolutionary computation (EC)*) to grupa inspirowanych biologicznie technik obliczeniowych. Algorytmy należące do tej grupy to *Algorytmy Ewolucyjne* (ang. *Evolutionary Algorithms (EA)*). Algorytmy te stanowią rodzaj metaheurystyk, czyli uniwersalnych algorytmów służących rozwiązywaniu różnych problemów optymalizacji, często wykorzystujących w trakcie działania bardziej specyficzne algorytmy. Algorytmy ewolucyjne wyróżnia wśród innych metaheurystyk między innymi przynależność do grupy algorytmów populacyjnych, które charakteryzują się równoległym przeszukiwaniem wielu rozwiązań, między którymi zachodzą interakcje wpływające na ocenę poszczególnych rozwiązań jak i kierunek ich zmian. Idea obliczeń ewolucyjnych czerpie swoje inspiracje z biologicznych teorii ewolucji, z nich również zapożycza terminologię. Poniżej przedstawiono znaczenie terminów używanych w kontekście obliczeń ewolucyjnych [Luk09]:

- *osobnik* (ang. individual - jedno z potencjalnych rozwiązań problemu
- *genotyp* (ang. genotype) - dane opisujące osobnika, podlegające procesom mutacji i krzyżowania
- *relacja rodzic-dziecko* ang. parent-parent - dziecko to osobnik powstały poprzez mutację rodzica lub w wyniku krzyżowania dwóch rodziców
- *mutacja* (ang. mutation) - zmiana genotypu pojedynczego osobnika w wyniku której powstaje osobnik o nowym genotypie
- *krzyżowanie* (ang. crossover) - stworzenie jednego lub więcej nowych osobników, których genom jest kombinacją genomów rodziców
- *populacja* - zbiór osobników item *pokolenie* - kompletny cykl operacji wykonanych na pokoleniu, również pokolenie będące jego wynikiem
- *przystosowanie / wartość funkcji przystosowania* (ang. fitness) - wartość określająca jakość osobnika ze względu na cel optymalizacji przyjęty w rozwiązywanym problemie
- *selekcja* (ang. selection) - wybór podzbioru osobników z populacji na podstawie ich przystosowania
- *powielanie* - utworzenie nowego osobnika poprzez skopiowanie istniejącego osobnika bez zmian
- *reprodukcja* (ang. breeding) - proces tworzenia nowego pokolenia (nowej generacji) osobników z osobników już istniejących poprzez procesy mutacji, krzyżowania i selekcji.

Przebieg typowego algorytmu ewolucyjnego przedstawiono na rysunku 2.6, poniżej przedstawiono jego opis:



RYSUNEK 2.6: Diagram przepływu algorytmu ewolucyjnego GP.

1. Utwórz początkową populację osobników korzystając z losowego mechanizmu lub użyj istniejącej populacji
2. Oblicz wartość *funkcji dopasowania* każdego osobnika z populacji (dokonaj ewaluacji populacji).
3. Jeśli znaleziono idealnego osobnika (wartość fitness wyniosła 1) lub skończył się czas, zwróć najlepszego znalezione dotąd osobnika
4. Dokonaj selekcji osobników
5. Utwórz nową populację poprzez mutację, krzyżowanie lub powielanie wybranych w poprzednim kroku osobników
6. Wróć do punktu 2

Selekcja

Selekcja polega na wybraniu podzbioru P osobników z populacji. Jej celem jest promowanie dobrych rozwiązań, które powstały w wyniku generowania nowych osobników, mutacji oraz krzyżowania. Dzięki dokonywaniu selekcji geny osobników lepiej przystosowanych do rozwiązania postawionego przed algorytmem problemu są "przekazywane dalej", biorąc udział w tworzeniu nowych osobników. "Dobroć rozwiązań jest oceniana poprzez funkcję przystosowania, która przyjmuje wartości z przedziału $\langle 0, 1 \rangle$, gdzie 0 oznacza najgorsze rozwiązanie a 1 rozwiązanie idealne. Najprostszym, naiwnym", sposobem przeprowadzenia selekcji jest wybór z niej określonej liczby osobników o najwyższej wartości fitness. Rozwiązanie to jednak nie jest optymalne - czasami pożądane fragmenty genotypu są

"ukryte" w osobnikach, które nie osiągają najwyższych wartości fitness, dlatego warto dopuścić również niektóre ze słabszych osobników do etapu krzyżowania. Głównym parametrem wpływającym na proces selekcji jest *napór selekcyjny*. Im większy napór selekcyjny tym trudniej osobnikom przejść selekcję co prowadzi do ujednolicania się osobników w populacji, ponieważ do następnego pokolenia przechodzi tylko wąska grupa najlepszych osobników, zmniejszając tym samym pulę genów pokolenia. Zmniejszanie się różnorodności jest niepożądanym procesem, ponieważ prowadzi do stagnacji procesu ewolucyjnego - kiedy wszystkie osobniki mają ten sam genotyp jedyne zmiany w populacji możliwe są poprzez mutację. Stagnacja procesu ewolucyjnego oznacza, że populacja z pokolenia na pokolenie się nie zmienia, ponieważ przez etap selekcji przechodzą zawsze te same, podobne do siebie osobniki. Jest to równoważne z utknięciem w optimum lokalnym (które może, choć nie musi być optimum globalnym rozwiązywanego problemu). Dlatego, żeby nie dopuścić do zbyt szybkiego ujednolicenia populacji stosuje się różne, bardziej wyrafinowane niż opisany powyżej naiwny, algorytmy selekcji:

- *ruletka* - metoda ruletki, zwana też wyborem losowym powtórzeniami ang. *Fitness-Proportionate Selection* polega na losowaniu osobników z prawdopodobieństwem wylosowania proporcjonalnym do ich wartości fitness. Metaforycznie można ją przedstawić w ten sposób, że każdemu osobnikowi przypisuje się pole na kole od ruletki o wielkości proporcjonalnej do jego wartości fitness i kręcąc ruletką losuje ustaloną liczbę osobników. Ta metoda selekcji dopuszcza wybranie jednego osobnika kilka razy a także umożliwia przypuszczenie osobników o małych wartościach fitness, co jest korzystne dla różnorodności populacji. Wadą tej metody jest ukryte w niej założenie o ilorazowym charakterze skali wartości fitness - osobnik A o k -krotnie większym fitness od osobnika B będzie średnio wybierany k razy częściej niż osobnik B, choć w ogólności nie musi być tak, że k -krotnie większy fitness oznacza k -krotnie lepszego lub pożądanego z punktu widzenia puli genów populacji osobnika. Z tym problemem radzi sobie następna opisana metoda.
- *turniej* - metoda turniejowa (ang. *tournament selection*) - jest najbardziej popularną i jednocześnie najprostszą z metod selekcji [Luk09]. Polega na losowaniu ze zwracaniem t osobników i wybraniu tego, który ma największą wartość fitness. Parametr t to tak zwany *rozmiar turnieju*, zwiększając go zwiększamy presję selekcyjną. Dla $t = 1$ selekcja sprowadza się do losowego wyboru osobników, dla $t = P$ otrzymamy populację składającą się z P kopii osobnika o najwyższym fitness.

Krzyżowanie

Osobniki wybrane w procesie selekcji z ustalonym jako parametr procesu ewolucyjnego prawdopodobieństwem ulegają krzyżowaniu, czyli wymieszaniu ich genotypów, w czego wyniku powstaje nowy osobnik o przypuszczalnie nieobecnym wcześniej genotypie.

Mutacja

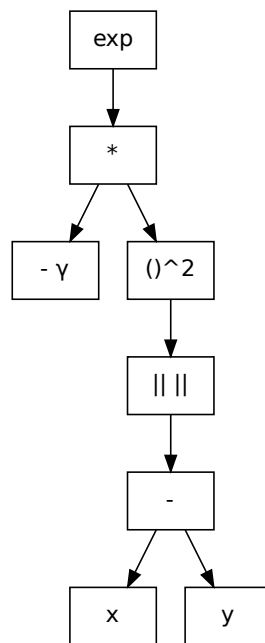
Mutacja polega na losowej, zazwyczaj nieznacznej, zmianie genotypu osobnika. Wpływa korzystnie na różnorodność populacji wprowadzając fragmenty genomu, które mogły wcześniej nie być obecne w puli genów populacji. Mutacja podobnie jak krzyżowanie jest stosowana z pewnym prawdopodobieństwem. Im to prawdopodobieństwo większe a zmiany dokonywane przez mutację bardziej rozległe, tym bardziej eksploratywny staje się algorytm ewolucyjny - trudniej mu "utknąć" w lokalnym optimum.

Genotyp

Poszczególne algorytmy ewolucyjne różnią się od siebie przede wszystkim ze względu na sposób reprezentacji osobników i powiązane z nim metody mutacji i krzyżowania. Selekcja osobników także może przebiegać na różne sposoby. Opisanie wszystkich możliwych parametrów i metod używanych w algorytmach ewolucyjnych nie jest celem tej pracy, ich opis można znaleźć na przykład w [Luk09]. Poniżej opisano jeden z rodzajów algorytmów ewolucyjnych jaką jest programowanie genetyczne, ze względu na jej wykorzystanie w niniejszej pracy.

2.3.1 Programowanie genetyczne

Programowanie genetyczne (GP, ang. *Genetic Programming*) to rodzaj algorytmów ewolucyjnych, które wyróżniają się specyficzną reprezentacją - ewolucji podlegają wyrażenia reprezentowane za pomocą drzew, na przykład wyrażenia matematyczne, logiczne lub kod programu komputerowego. Najczęściej w definicji programowania genetycznego podaje się właśnie ten ostatni rodzaj osobników [Luk09] [PLM08], lecz ewolucja innych wyrażeń reprezentowanym za pomocą drzew przebiega w algorytmie programowania genetycznego w identyczny sposób.



RYSUNEK 2.7: Przykładowe drzewo, które może być wynikiem działania algorytmu GP.
Drzewo przedstawia wyrażenie $e^{-\gamma * ||x-y||^2}$

Drzewo generowane przez algorytm GP (np. takie jak pokazane na rysunku 2.7) składa się z *terminali* (liście) oraz *funkcji* (węzły wewnętrzne oraz korzeń). Terminale to wartości wejściowe generowanego wyrażenia - stałe i zmienne - na rysunku są to zmienne x i y oraz stała λ . Funkcje to operacje wykonywane wartościach pochodzących z terminali lub innych funkcji, w przypadku drzewa z rysunku 2.7 są to: funkcja wykładnicza $exp()$, norma euklidesowa $|||$ oraz funkcja potęgowa (korzeń). Każdy węzeł wewnętrzny posiada luk(i) wejściowe, którymi funkcja "pobiera" wartości swoich argumentów oraz luk wyjściowy, którym przekazuje wartość do kolejnych funkcji (lub w przypadku korzenia zwraca wartość obliczoną przez całe wyrażenie).

Inicjalizacja populacji

Wybranie sposobu inicjalizacji populacji sprowadza się do określenia metody losowego generowania osobników. Drzewa są budowane przez łączenie ze sobą funkcji wybieranych ze zbioru funkcji oraz terminali. Każda funkcja ma określoną liczbę argumentów oraz opcjonalnie typ danych wejściowych i wyjściowych. W najprostszym przypadku wszystkie węzły drzewa operują na tych samych typach danych, wtedy nie narzuca się ograniczeń na to które węzły mogą się ze sobą łączyć. Jeśli jednak zbiór funkcji i terminali jest heterogeniczny ze względu na wartości typów danych, konieczne jest wprowadzenie ograniczeń. Odmiana GP, w której każda funkcja ma zdefiniowane typy danych, na których operuje nazywany jest *Strongly Typed GP* STGP. Kolejnym ograniczeniem najczęściej narzucanym na generowane drzewa jest ich głębokość, czyli odległość od najdalszego liścia do korzenia. Ze względu na to ograniczenie wyróżnia się 3 podstawowe metody generowania osobników:

- *full* - generuje drzewa o maksymalnej dozwolonej głębokości, w których każdy liść jest oddalony od drzewa o maksymalną odległość. Genrowanie polega na budowaniu drzewa zaczynając od korzenia poprzez losowanie kolejnych funkcji ze zbioru funkcji i dołączenie ich do funkcji już istniejących w drzewie, tak długo aż wszystkie gałęzie osiągną oczekiwaną głębokość. Na końcu dla każdej gałęzi losuje się kończący ją terminal.
- *grow* - w tej metodzie węzły budujące drzewo są losowane z obu zbiorów: funkcji i terminali. Dzięki temu, w przypadku wylosowania terminala dana gałąź kończy swój wzrost i otrzymujemy z gałęziami o potencjalnie zróżnicowanej długości.
- *ramped half-and-half* - w metodzie tej połowę populacji inicjalizuje się za pomocą metody *grow* a drugą połowę za pomocą algorytmu *full*, dodatkowo zmieniając maksymalną możliwą głębokość. Metoda ta pozwala otrzymać populację najbardziej urozmaiconą pod względem kształtu i wielkości drzew.

Krzyżowanie

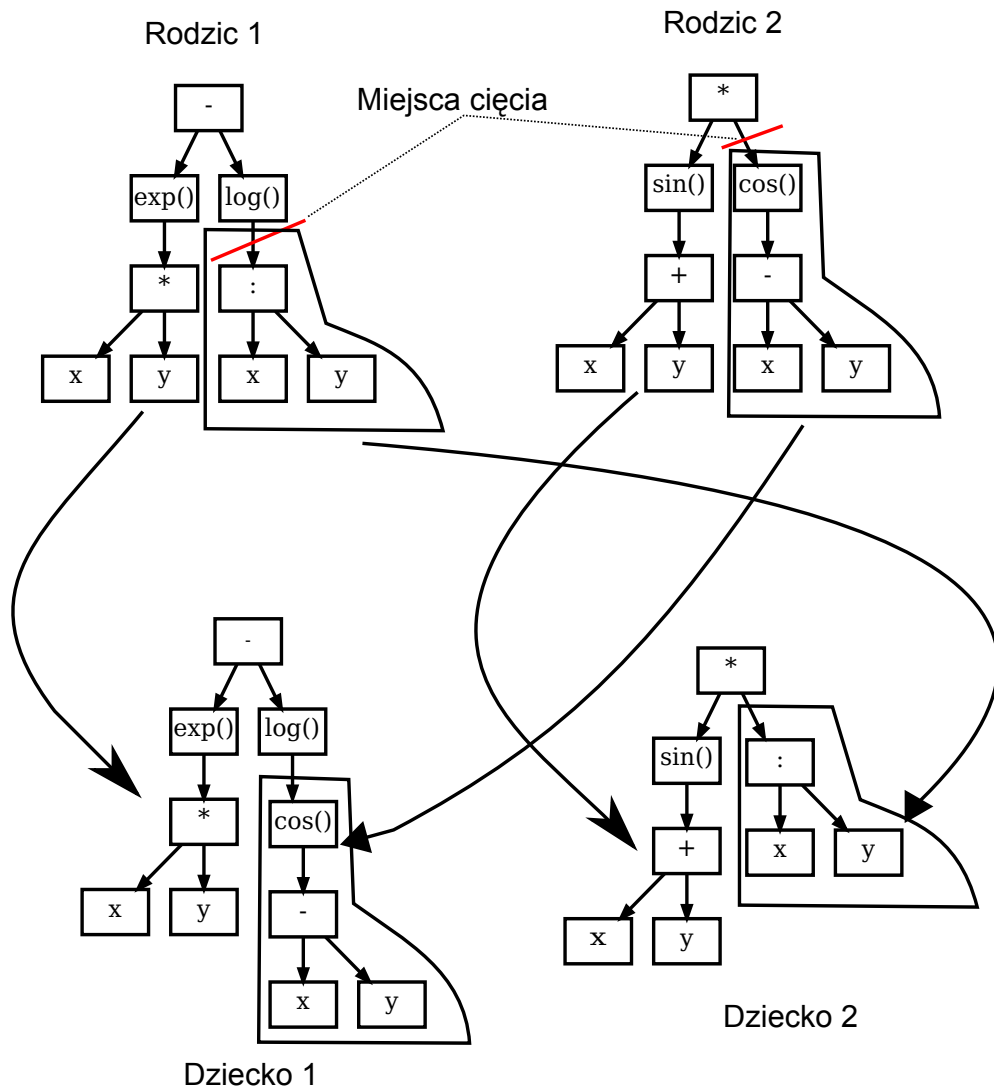
Najpowszechniejsza metoda krzyżowania w programowaniu genetycznym to *krzyżowanie poddrzew* ang. subtree crossover. Polega ona na zamianie poddrzew dwóch rodziców w czego wyniku powstaje dwóch potomków. Poddrzewa otrzymuje się poprzez wybór punktu cięcia, w którym poddrzewo jest odcinane od rodzica. Przykładowy proces krzyżowania zilustrowano na rysunku 2.3.1.

Mutacja

Mutacja drzewa polega na zamianie jego losowo wybranego fragmentu (poddrzewa) przez losowo wygenerowane drzewo odpowiedniej długości. Pod tym względem mutacja przypomina krzyżowanie, z tą różnicą, że odcięty fragment mutowanego drzewa jest zastępowany przez losowo wygenerowane drzewo a nie poddrzewo pochodzące od drugiego rodzica. Przykładową mutację drzewa pokazano na rysunku 2.3.1.

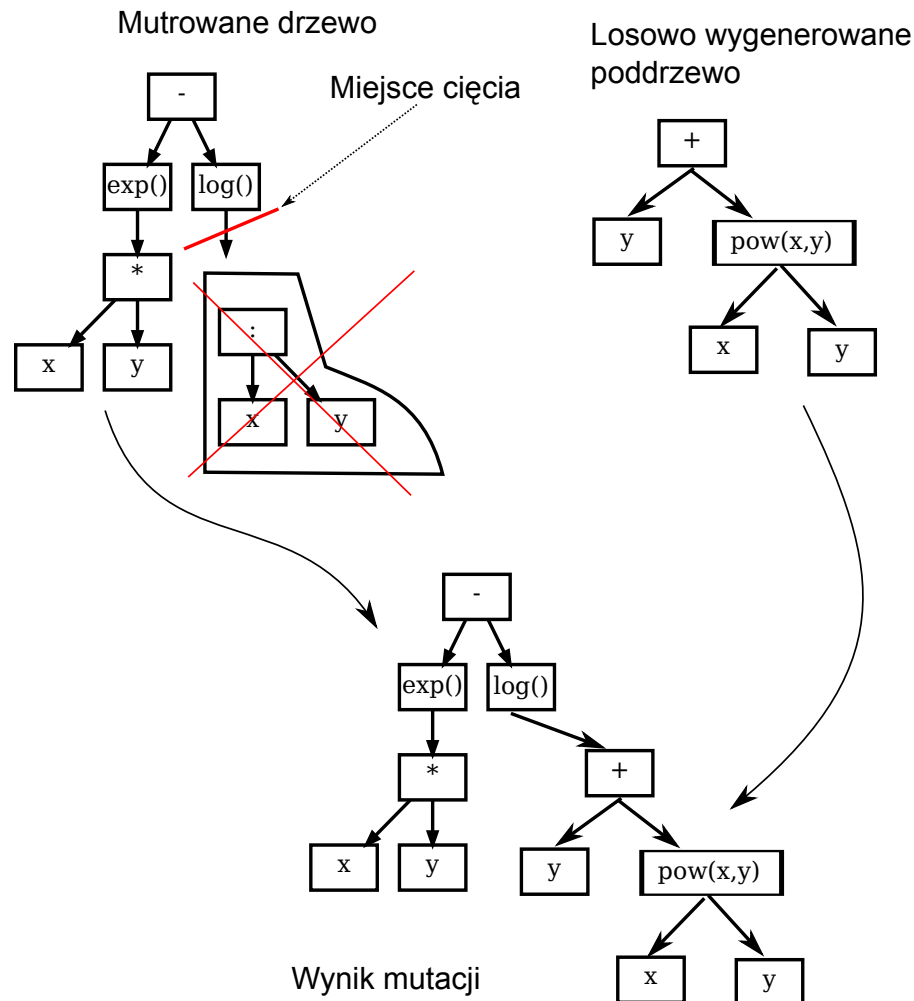
Selekcja

Selekcja w GP może być dokonywana jedną z metod używanych w innych algorytmach ewolucyjnych, jednak istnieją też metody specyficzne dla algorytmów GP. Jednym z problemów programowania genetycznego jest to, że drzewa powstałe w wyniku procesu ewolucyjnego mogą być bardzo duże, co nie jest pożądaną cechą - większe drzewo dłużej oblicza zwracaną wartość, zajmuje więcej miejsc w pamięci. Dlatego wielkość drzew należy ograniczać, jeśli wzrost drzewa nie prowadzi do zwiększenia wartości funkcji dopasowania. Wielkość generowanych drzew jest regulowana przez dwa mechanizmy.



RYSUNEK 2.8: Krzyżowanie drzew. U góry pokazano drzewa rodziców, punkty cięcia zaznaczono czerwonymi kreskami. Na dole widać dwóch potomków będących wynikiem krzyżowania.

Pierwszy to proste ograniczenie na maksymalną głębokość drzewa. Drugi mechanizm, o angielskiej nazwie *parsimony pressure*, promuje mniejsze drzewa podczas selekcji. W tym celu stosowany jest algorytm *selekcji turniejowej leksykograficznej z koszykami* (ang. Bucket Lexicographic Tournament Selection). Algorytm ten sortuje populację według przystosowania osobników, następnie grupuje je w N "koszyki". Następnie selekcja przebiega według zasad selekcji turniejowej, z tym, że porównuje się nie przystosowanie osobników, ale koszyk, do którego są przypisane. W przypadku gdy w turnieju porównywane są dwa osobniki z tego samego koszyka wygrywa ten, który jest mniejszy.



RYSUNEK 2.9: Przykładowa mutacja drzewa. U góry po lewej pokazano drzewo poddawane mutacji, u góry po prawej losowo wygenerowane poddrzewo, na dole wyniki mutacji. Odcinana część mutowanego drzewa (na rysunku przekreślona) zostaje odrzucona (nie jest już nigdzie wykorzystywana).

2.4 Optymalizacja parametrów SVM - przegląd literatury

Na wyniki osiągane przez algorytm SVM duży wpływ ma dobór stosowanej funkcji jądrowej wraz z jej parametrami (σ dla jądra RBF, d dla jądra wielomianowego, oraz γ i τ dla jądra sigmoidalnego) oraz parametrów samego algorytmu (na przykład parametru C ze wzoru 2.9) [HCL03].

Problem wyboru funkcji jądrowej i parametrów stosowanych przez klasyfikator SVM można rozwiązać na kilka sposobów. Najczęściej stosuje się po prostu jedną z popularnych funkcji - RBF, sigmoidalną lub wielomianową, kierując się doświadczeniem, "dobrymi praktykami" lub ręcznie sprawdzając wyniki osiągane z użyciem tych funkcji i różnych wartości ich parametrów. Najczęściej stosowaną funkcją jądrową wydaje się być RBF [HCL03] [HM05].

Możliwa jest również automatyzacja poszukiwania optymalnych parametrów. Proces ten nosi nazwę *optymalizacji hiperparametrów* (ang. Hyperparameter optimization). Jest to określenie zbioru metod służących optymalizacji parametrów jakiegoś algorytmu.

2.4.1 Miary przystosowania (fitness)

W przypadku wszystkich metod optymalizacji parametrów SVM konieczne jest wybranie jakiejś funkcji celu, która będzie optymalizowana. Funkcja ta informuje o tym, czy zmiany parametrów idą w dobrą stronę i pozwala oceniać i porównywać różne ich kombinacje. Oprócz miar uniwersalnych dla wszystkich problemów klasyfikacji (patrz rozdział 2.1.2), możliwe jest stosowanie miar specyficznych dla algorytmu SVM, takich jak: [HM05] [RS04]

- liczba wektorów wspierających (minimalizacja)
- promień najmniejszej sfery R , o środku w początku układu współrzędnych, obejmującej wszystkie przykłady uczące (minimalizacja): $R = \max_{1 \leq i \leq m} (K(x_i, x_i))$
- suma współczynników α : $\sum_{i=1}^N \alpha_i$ Dąży się do minimalizacji tej miary, co prowadzi do maksymalizacji marginesu hiperpłaszczyzny [HM05].
- suma *zmiennych osłabiających* (ang. *slack variables*) ζ : $\sum_{i=1}^N \zeta_i$, której minimalizacja prowadzi do mniejszej liczby niepoprawnie zaklasyfikowanych przykładów ze zbioru uczącego

oraz miar będących kombinacjami powyższych, takich jak: [RS04]

- $R^2 \sum_{i=1}^m \alpha_i + \sum_{i=1}^m \zeta_i$
- $(R^2 + 1/C)(\|\omega\|^2 + 2C \sum_{i=1}^m \zeta_i)$
- $(R^2 + 1/C) \sum_{i=1}^m \alpha_i$

Miary te pomagają zapobiec w przeuczeniu algorytmu, to znaczy zapewnić, że skuteczność klasyfikacji zbioru testującego będzie podobna do tej osiągananej na zbiorze uczącym, ponieważ nie mierzą skuteczności wprost.

2.4.2 Optymalizacja parametrów

W najprostszym przypadku do dokonania hiperoptymalizacji można posłużyć się jedną z prostych metaheurystyk takich jak *Random Search* [E] lub *Grid Search* [HCL03]. *Grid Search* polega na zachłanym przeszukiwaniu przestrzeni parametrów ograniczonej przez podane zakresy dozwolonych wartości każdego parametru. Przeszukiwanie takie odbywa się z określoną rozdzielczością, czyli dla każdego parametru jest określona wielkość kroku δ , o którą zmieniany jest parametr. Ponieważ przeszukiwanie takie jest nieefektywne, zaproponowano również użycie jego ulepszonej wersji [Sta03], polegające na stopniowym zwiększaniu rozdzielczości przeszukiwania i jednoczesnym zawężaniu przeszukiwanej przestrzeni do tego podobszaru, w którym dla aktualnej rozdzielczości otrzymano najlepsze rezultaty.

Do optymalizacji parametrów SVM, używa się również algorytmów ewolucyjnych. Na przykład używano *Strategii ewolucyjnych* (ang. *Evolutionary Strategy*) w celu doboru optymalnych wartości C i σ dla jądra RBF [RS04].

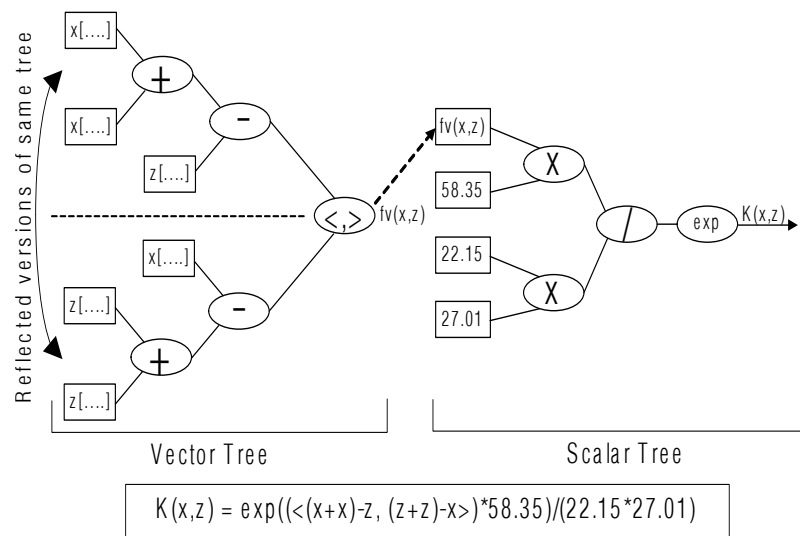
2.4.3 Ewolucja kerneli

Niemniej ważny od doboru parametrów jest dobór funkcji jądrowej użytej przez SVM do dokonania mapowania danych do wysoko wymiarowej przestrzeni. Dlatego próbuje się wykorzystywać różnego rodzaju algorytmy ewolucyjne w celu znalezienia funkcji, która daje najlepsze wyniki klasyfikacji.

Genetic Kernel SVM

W roku 2005 Howley i Madden zaproponowali użycie programowania genetycznego do projektowania funkcji jądrowych na potrzeby SVM [HM05]. Swoją metodę nazwali ang. *Genetic Kernel Support Vector Machine* (GK SVM). Algorytm GP w ich podejściu buduje funkcje ze zbioru operacji dodawania, odejmowania i mnożenia, w dwóch wersjach: skalarnej i wektorowej. Terminale to wektory opisujące przykłady uczące. Tak stworzone funkcje są służą potem do stworzenia funkcji jądrowej według wzoru: $K(x, z) = \langle treeEval(x, z), treeEval(x, z) \rangle$, gdzie $treeEval$ to wyrażenie zwrócone przez algorytm GP. Takie podejście zapewnia symetryczność funkcji jądrowej, choć nie zapewnia, że będzie ona dodatnio określona (nie musi zatem być poprawną funkcją jądrową w świetle teorii Mercera), a zatem nie zapewnia, że zostanie znalezione globalne optimum problemu optymalizacji rozwiązywanego wewnątrz SVM. Mimo to, takie funkcje mogą dawać dobre wyniki klasyfikacji SVM [BHB02]. Jako miarę przystosowania (fitness) Howley i Madden użyli błędu klasyfikacji i dodatkowo funkcji $fitness = \sum_{i=1}^m \alpha_i * R^2$ jako tak zwanego ang. *tiebreaker*, czyli funkcji służącej do uporządkowania osobników osiągających taki sam fitness. Autorzy przeprowadzili eksperymenty na sześciu zbiorach danych i porównali trafność klasyfikacji z użyciem funkcji wielomianowej, sigmoidalnej, RBF oraz stworzonych przez algorytm GP. Dla większości zbiorów wyewoluowane funkcje dawały wyniki lepsze lub porównywalne z najlepszymi wynikami otrzymanymi z pomocą jednej ze standardowych funkcji jądrowych z różnymi wartościami parametrów.

KTree



RYSunek 2.10: Przykładowe drzewo wygenerowane przez algorytm KTree (rysunek pochodzi z pracy [HM06]). Po lewej stronie widać część wektorową, po prawej część skalarną. Funkcja reprezentowana przez to drzewo to $e^{((x+x)-z, (z+z)-x) * 58.35 / (22.15 * 27.01)}$

Rok później ci sami autorzy zaprezentowali ulepszoną wersję swojego podejścia, tym razem zwaną *KTree* [HM06]. W podejściu tym drzewo reprezentujące generowany kernel składa się z dwóch części: wektorowej i skalarnej. Część wektorowa znajduje się na dole drzewa, zawiera jego liście, będące wektorami cech oraz funkcje operujące na wektorach, takie jak dodawanie i odejmowanie wektorów. Podobnie jak w poprzedniej wersji algorytmu buduje się drzewo symetryczne do wygenerowanego i oblicza ich iloczyn skalarny. Jego wynik, wraz z terminalami zawierającymi losowe stałe, stanowią liście drugiej części drzewa, która zawiera operacje wektorowe. Dzięki zastosowaniu takiej reprezentacji możliwe jest modelowanie za jej pomocą standardowych funkcji jądrowych, takich jak RBF. Przykład

takiego drzewa pokazano na rysunku 2.4.3. W pracy przebadano 3 różne funkcje fitness: dwie oparte na błędzie klasyfikacji zbioru uczącego, jako tiebreak używające kolejno wielkości drzewa oraz sumy współczynników α : $\sum_{i=1..m} \alpha_i$ oraz trzecią: opartą na 3-krotnej cross walidacji na zbiorze uczącym z wielkością drzewa jako tiebreaker-em. Dodatkowo autorzy testowali użycie "filtra Mercera" (ang. *Mercer filter*). Jeśli dla podzbioru danych uczących macierz $K = (k(x, z)_{i,j=1}^N)$ ma ujemne wartości własne, to filtr przypisuje funkcji jądrowej k najgorszą możliwość fitness tym samym eliminując funkcje nie spełniające warunków narzuconych przez teorię Mercera. W przeprowadzonych przez autorów eksperymentach metoda KTree osiągała na większości z 9 przetestowanych zbiorów danych wyniki lepsze lub porównywalne do SVM ze standardowymi kernelami. Najlepszą funkcją fitness okazała się potrójna walidacja krzyżowa. Użycie filtra Mercera okazało się nie wpływać korzystnie na fitness ani osiąganą trafność klasyfikacji, jednak autorzy ostrzegają przed użyciem algorytmu bez filtra Mercera ze względu na groźbę przeuczenia.

Evolutionary Kernel Machine

Oryginalne podejście do problemu przedstawia metoda *Evolutionary Kernel Machine* (EKM), przedstawiona w 2006 roku w [GSST06]. Funkcje jądrowe generowane są za pomocą GP, przy czym zbiór terminali stanowią funkcje operujące na wektorach a zwracające wartości skalarne (np. iloczyn skalarny, norma euklidesowa maksimum/minimum z wartości), natomiast wewnętrzne węzły pochodzą ze zbioru funkcji operujących na wartościach skalnych i zwracających wartości skalarne. Co nietypowe w tym podejściu, to funkcja fitness, która pochodzi z miary oceny marginesu zaproponowanej w [GBNT04] a używanej oryginalnie do oceny marginesu w algorytmie 1-NN:

$$F(K) = \frac{1}{m} \sum_{i=1}^m \delta_K(x'_1, y'_i) - l$$

, gdzie l to wielkość zbioru prototypów, a m to wielkość zbioru walidującego.

Wyrażenie $\delta_K(x, y)$ występujące w powyższym wzorze można przedstawić jako $\delta_K(x, y) = n(x, y) - p(x, y)$. Aby obliczyć $n(x, y)$ oraz $p(x, y)$ należy posortować przykłady ze zbioru prototypów $\{(x_1, y_1), \dots, (x_m, y_m)\}$ względem podobieństwa do przykładu (x, y) . $p(x, y)$ oznacza najniższą na posortowanej liście przykładów pozycję przykładu należącego do klasy y , czyli tej samej klasy co przykład (x, y) , $n(x, y)$ oznacza zaś najniższą pozycję przykłady należącego do klasy innej niż y . Miara podobieństwa, według której sortowane są przykłady pochodzi z użycia wygenerowanej funkcji jądrowej i wzorowana jest na algorytmie *Kernel-NN* [YJZ02]:

$$d_K(x, x')^2 = K(x, x) + K(x', x') - 2K(x, x')$$

Podsumowując: autorzy aby ocenić funkcje jądrowe na użytek algorytmu SVM, używają ich do skonstruowania klasyfikatora *Kernel-NN* i oceniają margines w tym algorytmie przy pomocy zbioru walidującego. Dodatkowo, oprócz ewolucji kerneli, autorzy stosują kooperatywną koewolucję zbioru prototypów ze zbiorem kerneli oraz koewolucję typu żywiciel-pasożyt zbioru walidującego ze zbiorem kerneli. Wyniki eksperymentów zaprezentowane w artykule pokazują porównanie błędów klasyfikacji otrzymanych za pomocą algorytmów $k-NN$, tradycyjnego SVM z jądrem RBF i metody EKM. W przypadku 2 z 6 przebadanych zbiorów danych algorytm EKM okazał się istotnie statystycznie lepszy od pozostałych dwóch algorytmów.

Kernel GP

W 2007 roku ukazał się artykuł autorstwa Seana Luke-a (autora m.in. wykorzystanej w tej pracy biblioteki ECJ [Sea10]) i Keith Sullivan opisujący algorytm *Kernel GP* (KGP) [SL07]. Algorytm ten, podobnie jak opisane powyżej algorytmy, wykorzystuje programowanie genetyczne do wygenerowania

funkcji jądrowych, jednak jako jedyny gwarantuje, że powstałe w ten sposób kernele będą spełniać warunki poprawności definiowane przez teorię Merera. Jest to możliwe dzięki własności zamknięcia zbioru kerneli ze względu na określone operacje. Własność ta umożliwia konstruowanie poprawnych funkcji jądrowych poprzez łączenie funkcji już istniejących (o których wiadomo, że są poprawne) za pomocą określonych operacji. W algorytmie *Kernel GP* zbiór terminali to zbiór trzech standardowych funkcji jądrowych: wielomianowej, sigmoidalnej i gaussowskiej. Zbiór wewnętrznych węzłów drzewa stanowią operacje dodawania, mnożenia przez stałą, mnożenia oraz funkcja wykładnicza. Ponieważ funkcje jądrowe jako argumenty przyjmują wektory a zwracają wartości skalarne, operacje łączące kernele zaś przyjmują i zwracają wartości skalarne, autorzy użyli (opisanego w części 2.3.1) silnie typowanego programowania genetycznego, żeby zapewnić zgodność typów danych przekazywanych w drzewie. Jako funkcji fitness autorzy używają wyników 10-krotnej walidacji krzyżowej na zbiorze trenującym. W eksperymentach autorzy porównali trafność klasyfikacji wydzielonego zbioru testującego przez algorytm *Kernel GP* oraz *Grid Search* (opis w części 2.4.2) użyty do SVM z jądrem gaussowskim. Dla połowy z przebadanych wzorów lepsze wyniki osiągnął *Kernel GP*, dla drugiej *Grid Search*.

Rozdział 3

Algorytm Kernel GP

3.1 Opis algorytmu

Jedną z trudności, która wiąże się z używaniem klasyfikatora SVM jest dobór odpowiedniej do zbioru danych *funkcji jądrowej*. Wymaga to doświadczenia lub przebiega na zasadzie prób i błędów. Ponadto zbiór powszechnie używanych funkcji jest ubogi - zazwyczaj ogranicza się do trzech podstawowych funkcji. Oprócz wyboru funkcji konieczne jest również ustawienie odpowiednich wartości ich parametrów.

Celem algorytmu Kernel GP jest odnalezienie optymalnej dla danego problemu funkcji jądrowej wraz z jej parametrami. Dzięki opisanej w poprzednim rozdziale własności domknięcia zbioru kerneli ze względu na pewne operacje arytmetyczne możliwe jest tworzenie nieograniczonej ilości dowolnie złożonych funkcji na podstawie kilku podstawowych kerneli. Opisywany algorytm przeszukuje przestrzeń takich funkcji za pomocą *programowania genetycznego*. Szukana jest taka funkcja, przy której użyciu klasyfikator SVM osiągnie największą *trafność (accuracy)* klasyfikacji.

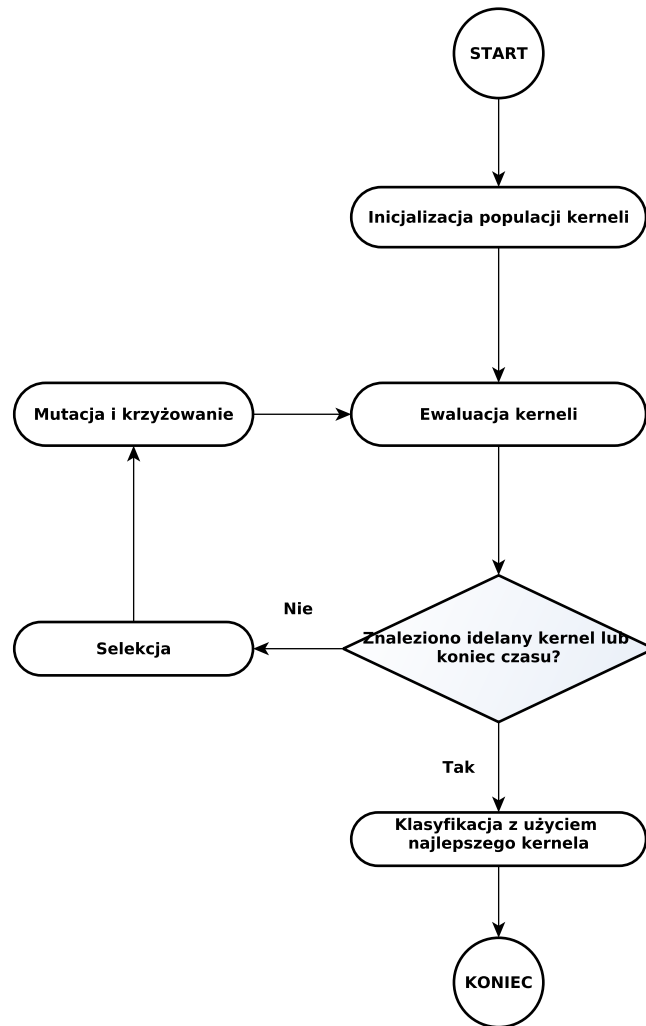
Przebieg algorytmu jest typowy dla algorytmów genetycznych:

1. Utwórz początkową populację kerneli
2. Oblicz wartość *funkcji dopasowania* każdego z kerneli: trafność klasyfikacji SVM z użyciem tego kernela
3. Jeśli znaleziono idealny kernel (wartość funkcji fitness wyniosła 1) lub skończył się czas, użyj tego kernela do klasyfikacji zbioru walidującego, zwróć wyniki klasyfikacji i zakończ algorytm.
4. Dokonaj selekcji najlepszych funkcji z populacji
5. Utwórz nową populację poprzez mutację i krzyżowanie wybranych w poprzednim kroku funkcji
6. Wróć do punktu 2

Algorytm pokazano również na diagramie przepływu na rycinie 3.1. Poszczególne kroki algorytmu zostaną opisane poniżej.

3.1.1 Inicjalizacja populacji

Podczas inicjalizacji początkowo pusta populacja jest wypełniana przez generowane w sposób losowy drzewa reprezentujące funkcje. Generowane drzewa muszą być poprawne, czyli spełniać narzucone ograniczenia na głębokość drzewa, liczbę węzłów, typ wartości zwracanych przez drzewo. Wielkość populacji jest jednym z parametrów algorytmu. Zbyt mała populacja powoduje losowe zawężenie



RYSUNEK 3.1: Diagram przepływu algorytmu Kernel GP.

przeszukiwanej przestrzeni i zmniejsza prawdopodobieństwo znalezienia optymalnej funkcji. Z drugiej strony zbyt duża wielkość populacji upodabnia algorytm genetyczny do pełnego przeszukiwania, co oczywiście zwiększa szanse znalezienia optymalnego kernela, ale wydłuża czas działania algorytmu.

Generowanie funkcji

Generowanie drzew reprezentujących funkcje jądrowe polega na łączeniu ze sobą funkcji elementarnych zgodnie z przypisanymi im ograniczeniami. Funkcje elementarne wraz z ograniczeniami zdefiniowane w algorytmie:

- Funkcje łączące - jako argument przyjmują wynik dwóch lub jednej funkcji jądrowej i ewentualnie stałą ERC. Zwracają wartość rzeczywistą. Dzięki właściwości domknięcia zbioru kerneli ze względu na operacje wykonywane przez te funkcje funkcja powstała przez połączenie dwóch kerneli funkcją łączącą jest również poprawnym kernelem [STC04].

- Dodawanie: $k(x, y) = k_1(x, y) + k_2(x, y)$
- Mnożenie: $k(x, y) = k_1(x, y) * k_2(x, y)$
- Mnożenie przez stałą: $k(x, y) = a * k_1(x, y)$

- Funkcja wykładnicza: $k(x, y) = e^{k_1(x, y)}$

Gdzie a to stała rzeczywista generowana jako stała ERC.

- Podstawowe funkcje jądrowe - jako argument przyjmują odpowiednią do funkcji liczbę stałych ERC. Zwracają wartość rzeczywistą.

- Liniowa: $k(x, y) = \langle x, y \rangle$
- Wielomianowa: $k(x, y) = \langle x, y \rangle^d$
- Gausowska *RBF): $e^{-\sigma * \|x - y\|^2}$
- Sigmoidalna: $k(x, y) = \tanh(\gamma \langle x, y \rangle + \tau)$
- Logarytmiczna: $k(x, y) = -\log(\|x - y\|^d + 1)$
- Potęgowa: $k(x, y) = (\alpha x^T z + c)^d$
- Cauchego: $k(x, y) = \frac{1}{1 + \frac{\|x - y\|^2}{\sigma}}$
- Wykładnicza: $k(x, y) = \exp\left(-\frac{\|x - y\|}{2\sigma^2}\right)$

Gdzie σ , γ , τ oraz d to wartości stałe generowane jako stałe ERC. a $\langle x, y \rangle$ to iloczyn skalarny wektorów x i y .

- Stałe ERC (ang. *Ephemeral Random Constant*) liczby rzeczywiste lub całkowite, które służą jako parametry innych funkcji. Są one liśćmi w drzewie, nie przyjmują żadnych argumentów. Mogą losowo zmieniać swoją wartość podczas mutacji.

- γ : liczba rzeczywista z zakresu $\langle 0.1, 2.0 \rangle$
- τ : liczba rzeczywista z zakresu $\langle 0.1, 1.0 \rangle$
- d : liczba całkowita z zakresu $\langle 1.0, 10.0 \rangle$
- α : liczba rzeczywista z zakresu $\langle -10.0, 10.0 \rangle$

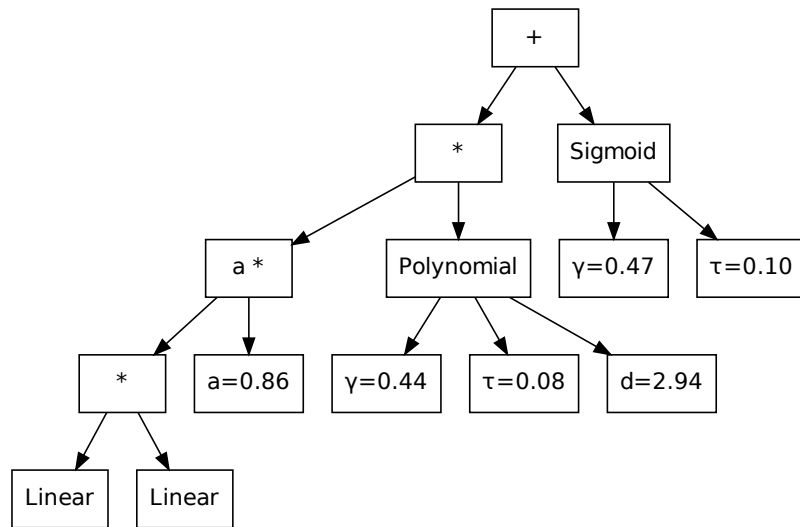
Przykładowe drzewo wygenerowane przez algorytm pokazana na ryc.??.

Wektory cech będące najważniejszymi argumentami funkcji jądrowych nie są wyodrębnione jako osobne funkcje budujące drzewo.

3.1.2 Ewaluacja kerneli

Każda wygenerowana przez algorytm GP funkcja zostaje poddana ocenie, w wyniku której zostaje jej przypisana wartość funkcji przystosowania (ang. *fitness*). W tym celu funkcja ta jest wykorzystywana przez algorytm SVM jako funkcja jądrowa a jakość wyników klasyfikacji stanowi ocenę funkcji jądrowej. Ewaluacja funkcji jądrowej może odbywać się na jeden z dwóch sposobów. Jeśli w zbiorze danych oprócz zbioru uczącego wydzielono zbiory testowe i walidujący, to sprawdzany kernel jest używany do klasyfikacji danych ze zbioru testującego. Ocena jakości klasyfikacji zostaje przeliczona na wartość *funkcji przystosowania* ewaluowanej funkcji jądrowej. Jeśli w zbiorze danych wydzielono tylko dwa podzbiory: uczący i walidujący, to zdolność klasyfikacji przez kernel jest oceniana za pomocą *walidacji krzyżowej* (ang. *cross-validation*). Walidacja krzyżowa pozwala użyć więcej danych podczas fazy uczenia, jednak wiąże się ze znacznym wzrostem złożoności obliczeniowej - zamiast jednej klasyfikacji musimy przeprowadzić k procesów uczenia i k klasyfikacji.

Do oceny jakości wyników klasyfikacji używana jest jedna z miar opisanych w części 2.1.2. Ponieważ wszystkie te miary należą do zakresu $\langle 0, 1 \rangle$, to mogą być bezpośrednio użyte jako wartość fitness ewaluowanego kernela.



RYSUNEK 3.2: Przykładowe drzewo generowane przez algorytm.

3.1.3 Selekcja

Jednym z problemów programowania genetycznego jest to, że drzewa powstałe w wyniku procesu ewolucyjnego mogą być bardzo duże, co nie jest pożądaną cechą - większe drzewo dłużej oblicza zwracaną wartość, zajmuje więcej miejsc w pamięci. Dlatego wielkość drzew należy ograniczać, jeśli wzrost drzewa nie prowadzi do zwiększenia wartości funkcji dopasowania. Wielkość generowanych drzew jest regulowana przez dwa mechanizmy. Pierwszy to proste ograniczenie na maksymalną głębokość drzewa. Wartość tę ustawiono na 6 - drzewa o większej głębokości nie zostaną w ogóle wygenerowane podczas inicjalizacji populacji czy podczas krzyżowania i mutacji. Drugi mechanizm, o angielskiej nazwie *parsimony pressure*, promuje mniejsze drzewa podczas selekcji. W tym celu stosowany jest algorytm selekcji turniejowej leksykograficznej z koszykami (ang. Bucket Lexicographic Tournament Selection). Algorytm ten sortuje populację według przystosowania osobników, następnie grupuje je w N "koszyki". Następnie selekcja przebiega według zasad selekcji turniejowej, z tym, że porównuje się nie przystosowanie osobników, ale koszyk, do którego są przypisane. W przypadku gdy w turnieju porównywane są dwa osobniki z tego samego koszyka wygrywa ten, który jest mniejszy.

3.1.4 Krzyżowanie i mutacja

Krzyżowanie polega na odcięciu dwóch losowych poddrzew z dwóch różnych osobników i zamianie ich miejscami. Wygenerowane w ten sposób drzewo musi spełniać narzucone na drzewo ograniczenia dotyczące typów i wielkości. Mutacja drzew polega na zamianie losowo wybranego poddrzewa przez losowo wygenerowane drzewo. Dodatkowo mutowane są również węzły ERC. Ich mutacja polega na dodaniu losowej wartości o rozkładzie normalnym do wartości przechowywanej w węźle. Wartość ta może być ujemna lub dodatnia.

3.1.5 Walidacja rozwiązania

Walidacja polega na użyciu najlepszego znalezionej kernela do klasyfikacji przykładów ze zbioru walidującego, które nie były używane podczas uczenia klasyfikatora SVM ani podczas ewaluacji ker-

neli. Najpierw algorytm SVM jest uczony na połączonych zbiorach trenującym i testującym, przy pomocy tej funkcji jądrowej. Następnie dokonywana jest klasyfikacja zbioru walidującego. Otrzymane w wyniku tej klasyfikacji miary jakości klasyfikacji (opisane w części 2.1.2 są miarą oceny całego algorytmu.

3.2 Implementacja

Algorytm został napisany w języku Java z użyciem bibliotek *ECJ (Evolutionary Computing in Java)* [Sea10] oraz *LibSVM* [CL11]. Pierwsza z nich dostarcza mechanizmy *obliczeń ewolucyjnych* w tym *programowania genetycznego*. LibSVM to klasyfikator SVM napisany oryginalnie w języku C z dostępną implementacją w Javie. Mechanizmy ECJ stanowią trzon algorytmu zapewniając tworzenie populacji funkcji, ich selekcję, mutację oraz krzyżowanie. LibSVM został użyty na etapie ewaluacji wygenerowanych przez ECJ funkcji.

3.3 Złożoność obliczeniowa

Rozdział 4

Wyniki działania algorytmu na popularnych zbiorach danych

4.1 Metodologia pomiarów

Żeby oszacować trafność klasyfikacji osiąganą przez skonstruowany system konieczne było podzielenie zbioru danych na zbiór uczący i *walidujący*, a w przypadku algorytmu Kernel-GP również wydzielenie ze zbioru uczącego podzbioru *testującego*, używanego do obliczania miary przystosowania (fitness) podczas przebiegu algorytmu genetycznego. Ponieważ sposób podziału zbioru danych ma wpływ na osiąganą trafność klasyfikacji, dokonywano 5 takich podziałów a następnie wyciągano średnią oraz odchylenie standardowe z wyników otrzymanych dla tych podziałów. Ta procedura dotyczyła zarówno testowania algorytmu *Kernel-GP* jak i porównawczych testów klasyfikatora SVM z biblioteki *LibSVM*. Dla obu algorytmów stosowano te same podziały danych, przy czym w przypadku klasyfikatora *LibSVM* ze zbioru uczącego nie wydzielano zbioru testującego.

Algorytm genetyczny jest w swej naturze stochastyczny, korzysta więc z funkcji generujących liczby pseudolosowe. Aby zapewnić powtarzalność wyników i umożliwić ich porównanie ziarno generatora liczb pseudolosowych ustawiono na stałą wartość.

Aby ocenić skuteczność algorytmu genetycznego w poszukiwaniu optymalnych funkcji jądrowych oraz oszacować optymalną wielkość populacji, liczbę ewaluowanych pokoleń oraz najlepszą funkcję przystosowania przeprowadzono szereg eksperymentów obliczeniowych, w których uruchamiano algorytm dla coraz to większych wartości tych parametrów. Dla każdego przebiegu algorytmu obliczano i zapisywano kilka miar trafność klasyfikacji zbioru *walidującego* (miary te zostały opisane w części 2.1.2).

Analizując tak zebrane dane można przeanalizować na ile poszukiwanie funkcji jądrowej przez algorytm genetyczny było podobne do losowego przeszukiwania a na ile było ono zbieżne. W pierwszym przypadku na wyniki osiągane przez algorytm powinna mieć wpływ przede wszystkim wielkość populacji, w drugim również liczba populacji przez które poszukiwano rozwiązania. W szczególności ciekawym przypadkiem jest ten, gdy liczba populacji wynosi 1, czyli cały algorytm ogranicza się do wygenerowania populacji losowych osobników i wybrania jednego z nich - w tym przypadku algorytm genetyczny sprowadza się do losowego poszukiwania rozwiązania. Porównując różnicę w trafności osiąganą w trakcie jednego pokolenia i coraz większej ich liczby można ocenić czy proces ewolucyjny przebiega poprawnie.

4.2 Opis zbiorów danych

Do oceny pracy algorytmu użyto standardowych zbiorów danych służących do testowania systemów maszynowego uczenia się, dostępnych na stronie biblioteki *LIBSVM* [CL11] oraz w repozytorium UCI [D]. Zbiory zostały opisane w tabelce 4.1. Użyte nazwy zbiorów są zgodne z tymi ze strony libsvm [C].

TABLICA 4.1: Zbiory danych użyte do testowania systemu.

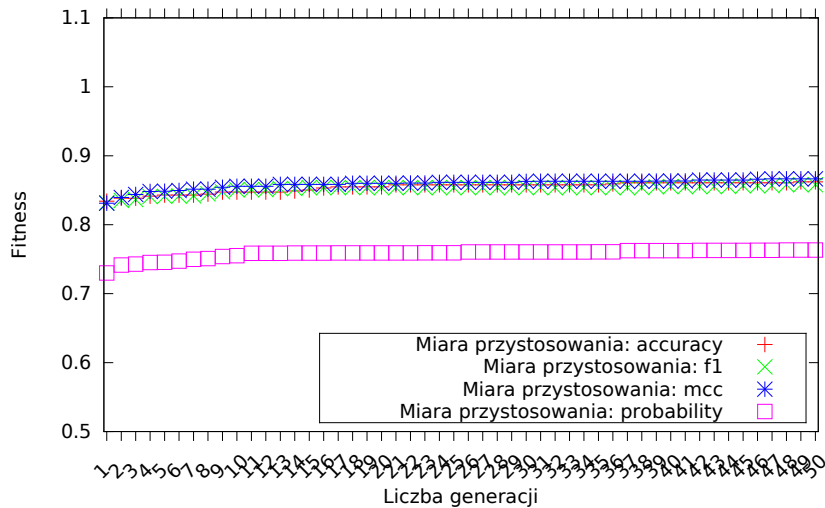
Nazwa zbioru	Liczba klas	Liczba atrybutów	Wielkość zbioru	Wielkość zbioru uczącego	Wielkość zbioru testującego	Wielkość zbioru walidującego
Iris	3	4	150	68	33	49
Letter	26	16	15000	9000	4400	6600
DNA	3	180	2000	1435	700	1051
Vowel	11	10	528	447	217	326
Breast cancer	2	10	683	343	170	170
Heart	2	13	270	136	67	67

TABLICA 4.2: Zbiory danych użyte do testowania systemu.

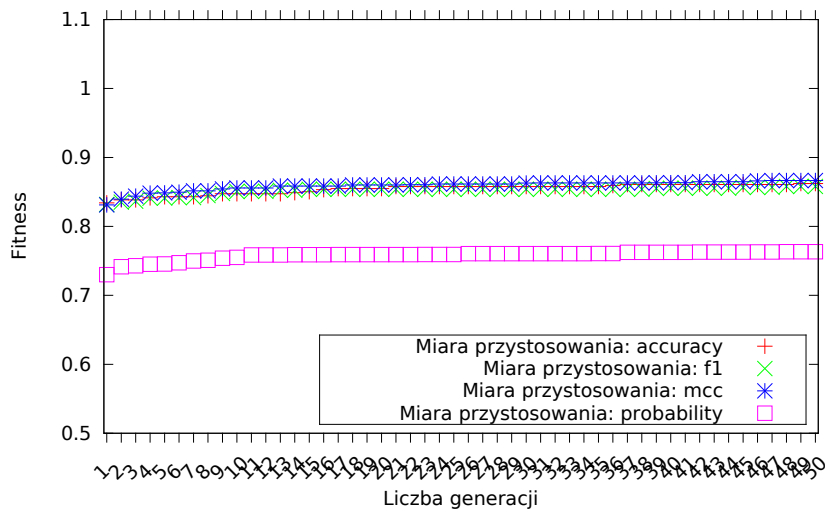
Nazwa zbioru	Liczba atrybutów ciągłych	Liczba atrybutów nominalnych	Liczba klas	Proporcje klas
Breast Cancer	10	0	2	239/444
Heart	7	6	2	150/120
DNA	0	180	3	464/485/1051
Vowel	10	0	11	Każda klasa 48 razy

4.3 Fitness

Aby ocenić dobór parametrów procesu ewolucyjnego zapisywano wartości fitness osiągane przez najlepszego osobnika w każdym pokoleniu. Wartości te zostały przedstawione na wykresach przedstawiających wartość fitness najlepszego osobnika w funkcji czasu trwania algorytmu (ilości dotychczas wygenerowanych populacji), czyli tak zwanych ang. *Fitness Graphs*. Na wykresach 4.1 - 4.2 widać jak zmienia się wartość funkcji fitness wraz z kolejnymi pokoleniami dla różnych miar jakości klasyfikacji użytych jako funkcja fitness (miary te zostały opisane w części 2.1.2 a ich użycie w części 3.1.2).



RYSUNEK 4.1: Najlepsza wartość funkcji przystosowania dla kolejnych pokoleń dla zbioru *heart*.

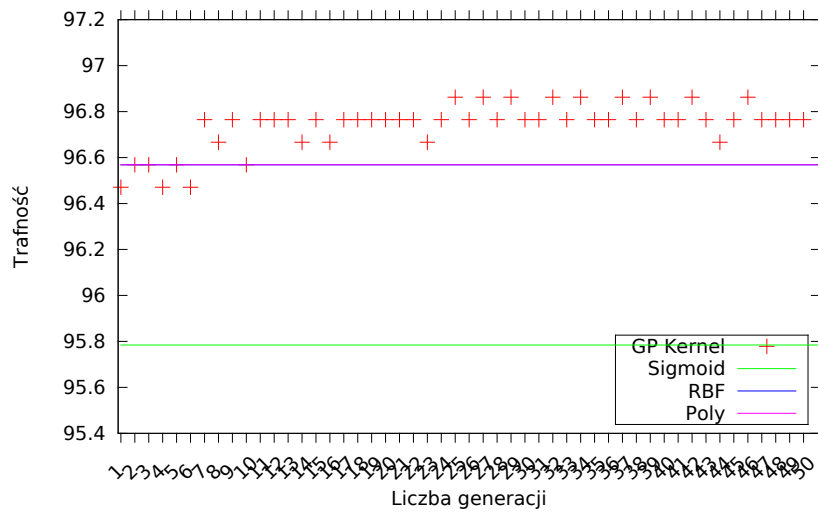


RYSUNEK 4.2: Najlepsza wartość funkcji przystosowania dla kolejnych pokoleń dla zbioru *breast*.

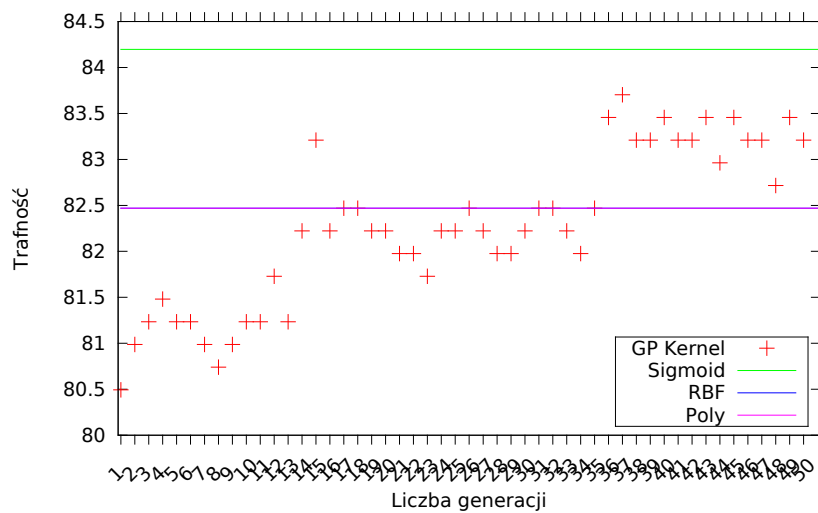
4.4 Wyniki klasyfikacji zbioru walidującego

Wyniki klasyfikacji zostały ocenione za pomocą miar opisanych w części 2.1.2. Dla każdej z tych miar przedstawiono jej wartości dla przebiegów algorytmu, w których jako funkcja fitness była wybrana właśnie ta miara. Dodatkowo na wykresach ukazano wartości danych miar uzyskane w wyniku klasyfikacji za pomocą trzech standardowych funkcji jądrowych (*Wielomianowej*, *Sigmoidalnej* i *RBF*).

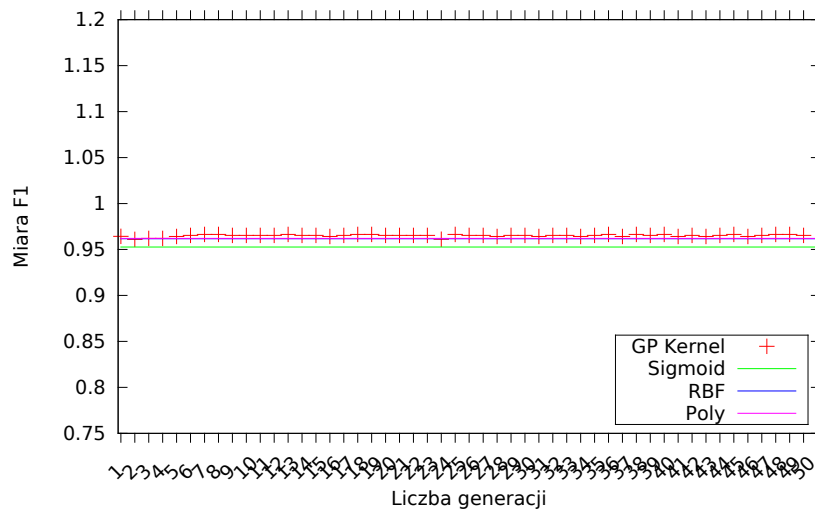
Jak widać na rysunkach 4.3-4.10 sprawność algorytmu zależy zarówno od zbioru danych jak i od wybranej miary jakości.



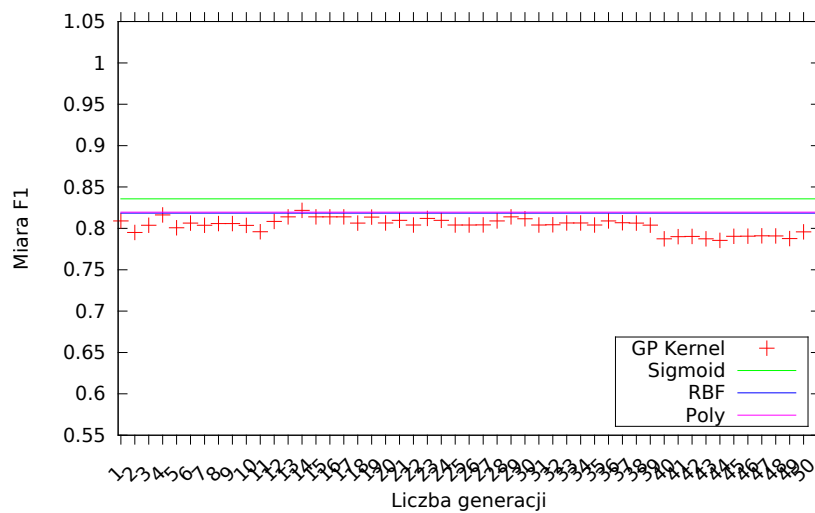
RYSUNEK 4.3: Trafność (ang. *accuracy*) klasyfikacji dla zbioru *breast* w funkcji czasu wykonania (ilości pokoleń).



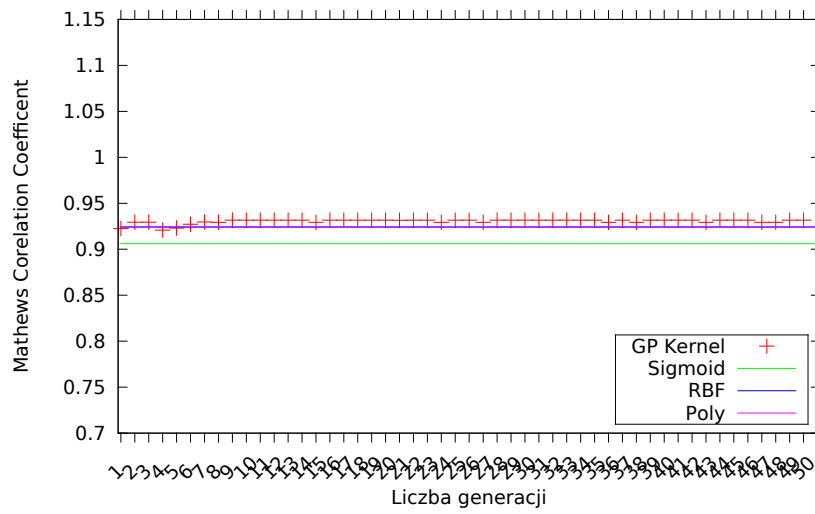
RYSUNEK 4.4: Trafność (ang. *accuracy*) klasyfikacji dla zbioru *heart* w funkcji czasu wykonania (ilości pokoleń).



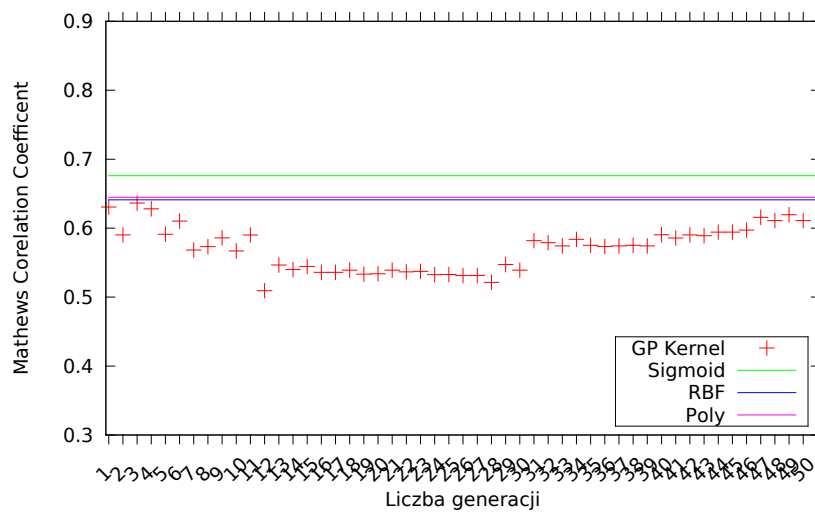
RYSUNEK 4.5: Wartość miary F1 dla wyników klasyfikacji zbioru *breast* w funkcji czasu wykonania (ilości pokoleń).



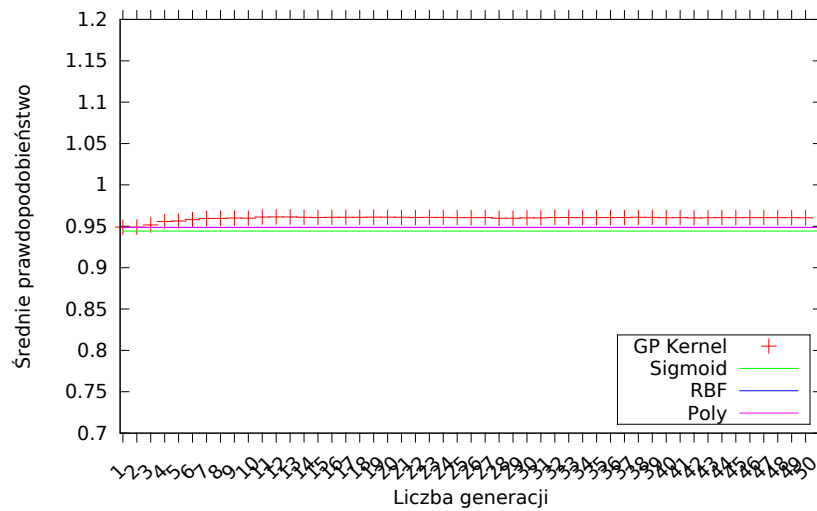
RYSUNEK 4.6: Wartość miary F1 dla wyników klasyfikacji zbioru *heart* w funkcji czasu wykonania (ilości pokoleń).



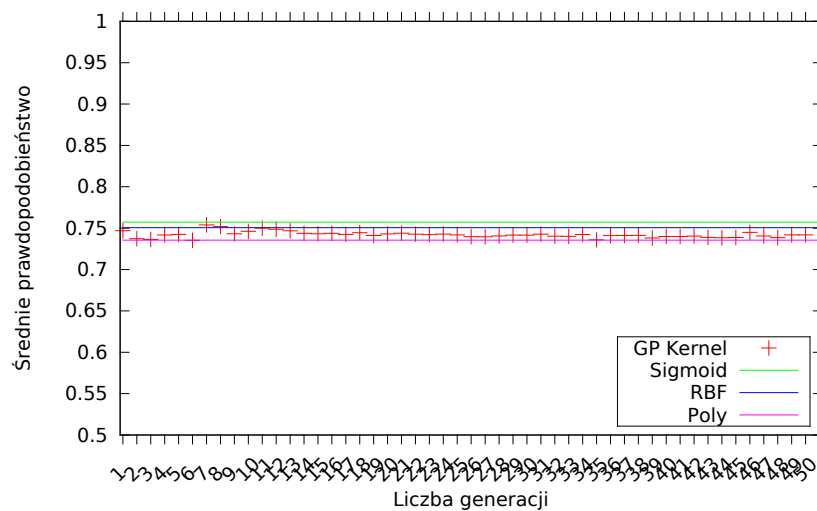
RYСУNEK 4.7: Wartość miary *Mathews Correlation Coefficient* (MCC) dla wyników klasyfikacji zbioru *breast* w funkcji czasu wykonania (ilości pokoleń).



RYСУNEK 4.8: Wartość miary *Mathews Correlation Coefficient* dla wyników klasyfikacji zbioru *heart* w funkcji czasu wykonania (ilości pokoleń).



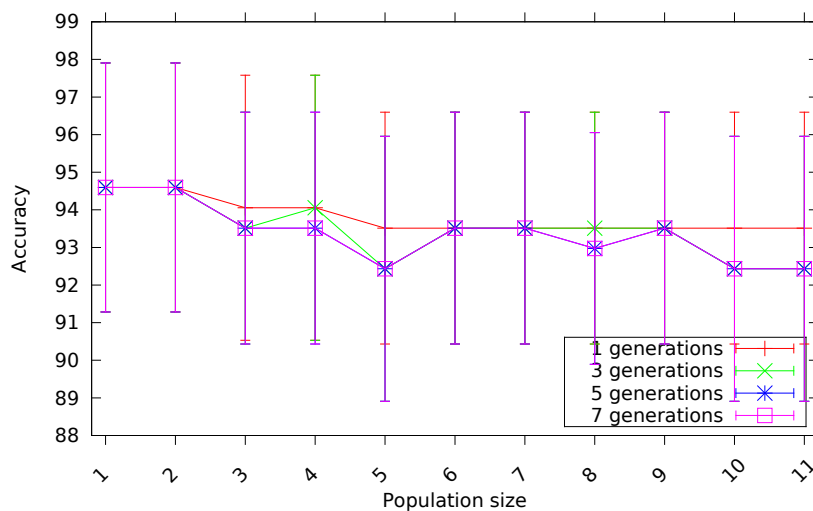
RYСУNEK 4.9: Średnia wartość prawdopodobieństwa przypisywanego przez SVM właściwej dla klasyfikowanego przykładu klasie w funkcji czasu wykonania (ilości pokoleń). Zbiór *breast*



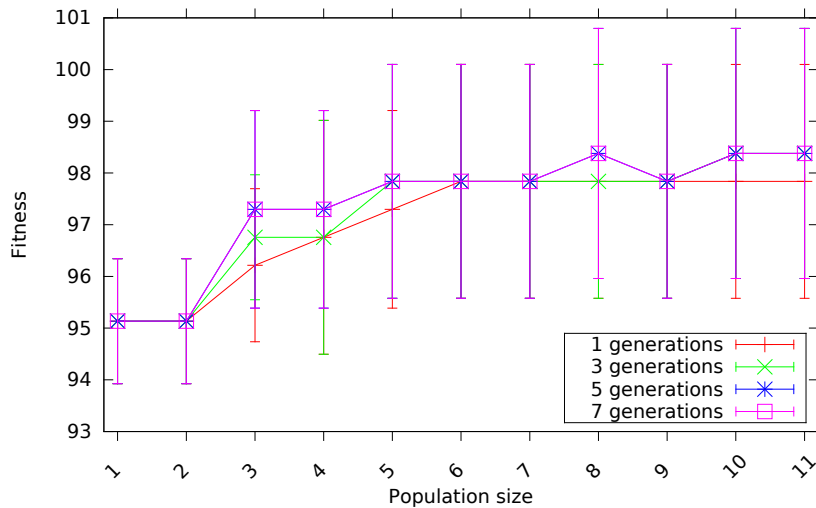
RYСУNEK 4.10: Średnia wartość prawdopodobieństwa przypisywanego przez SVM właściwej dla klasyfikowanego przykładu klasie w funkcji czasu wykonania (ilości pokoleń). Zbiór *heart*.

4.4.1 Monotoniczność funkcji trafności

Miejscami funkcja trafności nie jest monotoniczna, a ściślej niemalejąca, względem liczby pokoleń oraz wielkości populacji (co widać np. na wykresach 4.11 i 4.13). Wydawałoby się, że tak być nie powinno (algorytm genetyczny zwraca najlepszego osobnika z całego swojego przebiegu, więc wszystkie osobniki, które pojawiły się w pierwszych 5 pokoleniach pojawią się w pierwszych 7 pokoleniach, więc trafność dla po 7 pokoleniach powinna być co najmniej tak dobra jak po 5). Jednak może się tak zdarzyć ze względu na to, że trafność pokazana na wykresach to trafność klasyfikacji zbioru walidującego, natomiast trafność użyta przez algorytm genetyczny jako miara dostosowania (ang. *fitness*) to trafność klasyfikacji zbioru testującego. Widać to na wykresie 4.12, który przedstawia wartość przystosowania dla tych samych danych, dla których na wykresie 4.11 jest pokazana trafność klasyfikacji na zbiorze walidującym - tutaj funkcja wykazuje mniej braku monotoniczności.

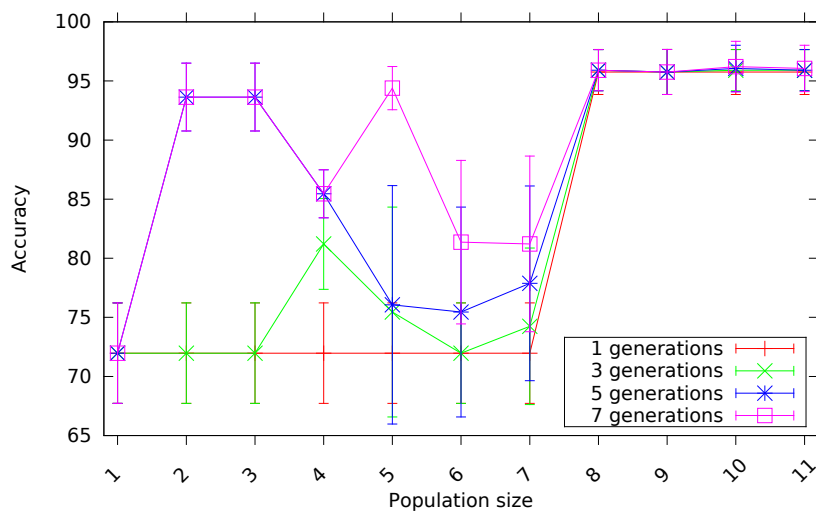


RYСУNEK 4.11: Trafność klasyfikacji dla zbioru *iris* w funkcji rozmiaru populacji dla różnych ilości pokoleń, dla małych populacji.

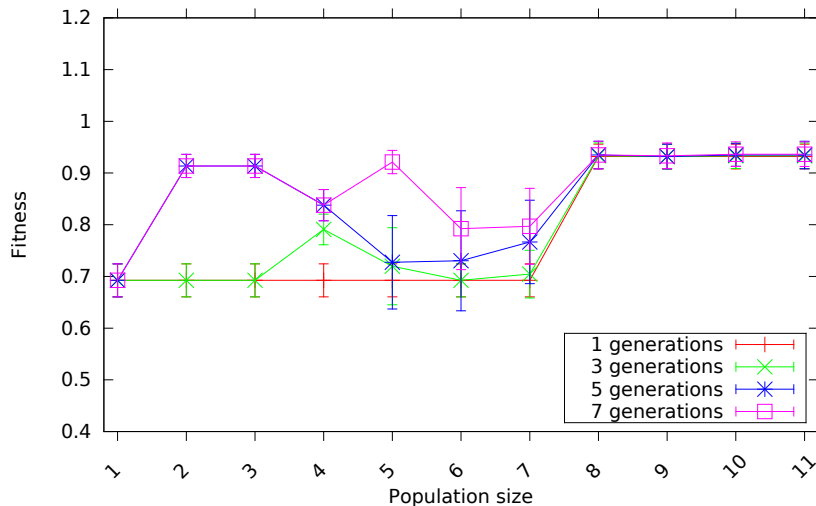


RYSUNEK 4.12: Najlepsza wartość funkcji przystosowania (ang. *fitness*) *iris* w funkcji rozmiaru populacji dla różnych ilości pokoleń, dla małych populacji.

Zatem przynajmniej część braku monotoniczności funkcji trafności na zbiorze walidującym wynika z przeuczenia algorytmu - znaleziona przez algorytm genetyczny funkcja jądrowa lepiej sprawdza się przy klasyfikacji zbioru testującego niż walidującego. Nie jest to jednak jedyna przyczyna braku monotoniczności - widać to na wykresie 4.14 przedstawiającym wartość funkcji przystosowania dla zbioru *vowel* - jej przebieg jest bardzo podobny do przebiegu ukazanej na rys. 4.13 funkcji trafności klasyfikacji zbioru walidującego na tym samym zbiorze. Co więc jest przyczyną braku monotoniczności? Warto zauważyć, że funkcja jest niemalejąca ze względu na ilość pokoleń oraz że dla jednego pokolenia funkcja jest monotoniczna. Sugeruje to, że "winnym" może być selekcja - w praktyce nie zachodzi ono w przypadku gdy algorytm genetyczny działa przez jedno pokolenie.



RYSUNEK 4.13: Trafność klasyfikacji dla zbioru *vowel* w funkcji rozmiaru populacji dla różnych ilości pokoleń, dla małych populacji.



RYSUNEK 4.14: Trafność klasyfikacji dla zbioru *vowel* w funkcji rozmiaru populacji dla różnych ilości pokoleń, dla małych populacji.

Gdy przyjrzeć się dokładnie przebiegowi ewolucji widać, że rzeczywiście tak jest. Dodatkowy osobnik (rys. 4.15), który odróżnia w pokoleniu pierwszym populację o wielkości 3 i 4 jest przodkiem innego osobnika (rys. 4.16), który w trzecim pokoleniu osiąga fitness większy niż osobnik (rys. 4.18), który w przypadku populacji wielkości 3 był przodkiem osobnika (rys. 4.19), który to okazał się najlepszym podczas przebiegu całego algorytmu. W rezultacie tego "geny" potencjalnego zwycięzcy nie przetrwały w przebiegu algorytmu z populacją liczącą 4 osobników. Jak widać osobnik najlepszy we wszystkich pokoleniach nie musi być wcale potomkiem osobników najlepszych w poszczególnych pokoleniach - czasem połączenie dwóch osobników przeciętnych może dać osobnika bardzo dobrego.

■

RYSUNEK 4.15: Funkcja z pierwszego pokolenia, która w przebiegu z wielkością populacji 4 osiągnęła fitness 0.4242424. Przodek funkcji z rys.4.16

■

RYSUNEK 4.16: Funkcja z trzeciego pokolenia, która w przebiegu z wielkością populacji 4 osiągnęła fitness 0.78030306. Potomek funkcji z rys.4.15, przodek zwycięskiej funkcji (rys.4.17) z przebiegu z populacją o wielkości 4.



RYSUNEK 4.17: Funkcja z ostatniego pokolenia w przebiegu z wielkością populacji 4 osiągnęła fitness 0.8333333. Przodek zwycięskiej funkcji z rys.4.19



RYSUNEK 4.18: Funkcja z trzeciego pokolenia, która w przebiegu z wielkością populacji 3 i 4 osiągnęła fitness 0.6515151. Przodek zwycięskiej funkcji z rys.4.19



RYSUNEK 4.19: Zwycięska funkcja w przebiegu z populacją wielkości 3, potomek funkcji z rys.4.18. Osiągnęła fitness 0.9015151.

4.5 Czas wykonania

4.6 Użycie pamięci

4.7 Posdumowanie wyników

Rozdział 5

Case study - klasyfikacja danych ADHD 200

ADHD 200 był międzynarodowym konkursem, który zakończył się we wrześniu 2011 roku. Dzięki współpracy ośmiu szpitali i ośrodków naukowych z całego świata w ramach konkursu udostępniono zbiór zawierający dane medyczne 776 dzieci, z czego 285 z ADHD. Zadaniem uczestników konkursu było skonstruowanie klasyfikatora diagnozującego ADHD na podstawie tych danych. Zbiór testowy zawierał dane xxx dzieci, których danych nie było w zbiorze uczącym, bez podanej diagnozy. Celem skonstruowanego klasyfikatora było przypisanie diagnozy do przykładów ze zbioru testującego.

Wyniki konkursu pokazały, że zbiór danych był trudny w klasyfikacji. Największa osiągnięta trafność klasyfikacji wyniosła 60.51% (szczegółowe wyniki dostępne na stronie konkursu: [B]).

5.1 Opis zbioru danych

5.1.1 Surowe dane

Dane dostarczone przez organizatorów konkursu składają się z:

- Danych klinicznych:
 - Płeć
 - Wiek
 - Współczynnik IQ
 - Prawo/lewo ręczność
- Danych obrazowych:
 - Strukturalnych - dane pochodzące z obrazowania rezonansu magnetycznego (*MRI*, ang. *Magnetic Resonance Imaging*). Są to trójwymiarowe obrazy o rozdzielczości ok. 256x254x160 punktów, obrazujące strukturę mózgu osoby badanej
 - Funkcjonalnych - dane pochodzące z obrazowania funkcjonalnego rezonansu magnetycznego (*FMRI*, ang. *functional Magnetic Resonance Imaging*) będące sekwencją ok 120 trójwymiarowych obrazów o rozdzielczości ok 250x250x250 punktów, obrazującą aktywność mózgu osoby badanej rejestrowaną przez ok 6 minut.

5.1.2 Preprocessing

5.2 Konstrukcja i selekcja cech

5.3 Wyniki klasyfikacji

5.3.1 Kernel GP

5.3.2 Porównanie z innymi algorytmami

SVM

Inne klasyfikatory

Rozdział 6

Podsumowanie

Literatura

- [BGV92] Bernhard E. Boser, Isabelle M. Guyon, Vladimir N. Vapnik. A training algorithm for optimal margin classifiers. *Proceedings of the fifth annual workshop on Computational learning theory, COLT '92*, strony 144–152, New York, NY, USA, 1992. ACM.
- [BHB02] Claus Bahlmann, Bernard Haasdonk, Hans Burkhardt. On-line handwriting recognition with support vector machines a kernel approach. *Proceedings of the Eighth International Workshop on Frontiers in Handwriting Recognition (IWFHR'02)*, IWFHR '02, strony 49–, Washington, DC, USA, 2002. IEEE Computer Society.
- [BTB05] S. Boughorbel, J. P Tarel, Nozha Boujemaa. Conditionally positive definite kernels for svm based image recognition. *Multimedia and Expo, 2005. ICME 2005. IEEE International Conference on*, strony 113–116, 2005.
- [CL11] Chih-Chung Chang, Chih-Jen Lin. LIBSVM: a library for support vector machines. *ACM Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27, Maj 2011.
- [GBNT04] Ran Gilad-Bachrach, Amir Navot, Naftali Tishby. Margin based feature selection - theory and algorithms. *Proceedings of the twenty-first international conference on Machine learning, ICML '04*, strony 43–, New York, NY, USA, 2004. ACM.
- [GSST06] Christian Gagné, Marc Schoenauer, Michèle Sebag, Marco Tomassini. Genetic programming for kernel-based learning with co-evolving subsets selection. *Proceedings of the 9th international conference on Parallel Problem Solving from Nature, PPSN'06*, strona 1008–1017, Berlin, Heidelberg, 2006. Springer-Verlag.
- [HCL03] Chih-Wei Hsu, Chih-Chung Chang, Chih-Jen Lin. *A practical guide to support vector classification*. 2003.
- [HM05] Tom Howley, Michael G. Madden. The genetic kernel support vector machine: Description and evaluation. *Artificial Intelligence Review*, 24(3-4):379–395, 2005.
- [HM06] Tom Howley, Michael G. Madden. An evolutionary approach to automatic kernel construction. Stefanos Kollias, Andreas Stafylopatis, Włodzisław Duch, Erkki Oja, redaktorzy, *Artificial Neural Networks – ICANN 2006*, wolumen 4132 serii *Lecture Notes in Computer Science*, strony 417–426. Springer Berlin Heidelberg, 2006.
- [KS03] K. Krawiec, J. Stefanowski. *Uczenie maszynowe i sieci neuronowe*. Wydaw. Politechniki Poznańskiej, 2003.
- [Luk09] Sean Luke. *Essentials of Metaheuristics*. Sean Luke, 2009. available at <http://cs.gmu.edu/~sean/book/metaheuristics/>.
- [Mit97] Thomas M. Mitchell. *Machine Learning*. McGraw-Hill, Inc., New York, NY, USA, wydanie 1, 1997.
- [PLM08] Riccardo Poli, William B. Langdon, Nicholas Freitag McPhee. *A Field Guide to Genetic Programming*. Lulu Enterprises, UK Ltd, 2008.
- [RS04] Thomas Philip Runarsson, Sven Sigurdsson. Asynchronous parallel evolutionary model selection for support vector machines. *Neural Information Processing–Letters and Reviews*, 3(3):59–68, 2004.

- [Sea10] Luke Sean. *The ECJ Owner's Manual*. Yale Univ Pr, Pa/xdziernik 2010.
- [SL07] Keith M. Sullivan, Sean Luke. Evolving kernels for support vector machine classification. *Proceedings of the 9th annual conference on Genetic and evolutionary computation*, GECCO '07, strona 1702–1707, New York, NY, USA, 2007. ACM.
- [SS02] B. Schölkopf, A.J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. Adaptive computation and machine learning. MIT Press, 2002.
- [Sta03] Carl Staelin. Parameter selection for support vector machines. *Hewlett-Packard Company, Tech. Rep. HPL-2002-354R1*, 2003.
- [STC04] John Shawe-Taylor, Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, New York, NY, USA, 2004.
- [YJZ02] Kai Yu, Liang Ji, Xuegong Zhang. Kernel nearest-neighbor algorithm. *Neural Process. Lett.*, 15(2):147–156, Kwiecie/n 2002.

Zasoby internetowe

[A] ECJ

<http://cs.gmu.edu/~eclab/projects/ecj/>

[B] ADHD 200

http://fcon_1000.projects.nitrc.org/indi/adhd200/

[C] Libsvm datasets

<http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

[D] UCI Machine Learning Repository

<http://archive.ics.uci.edu/ml/datasets/>

[E] Hyperopt - Distributed Asynchronous Hyperparameter Optimization in Python

<http://jaberg.github.io/hyperopt/>

[F] <http://crsouza.blogspot.com> - lista funkcji jądrowych <http://crsouza.blogspot.com/2010/03/kernel-functions-for-ma.html>



© 2013 Tomasz Ziętkiewicz

Instytut Informatyki, Wydział Informatyki i Zarządzania
Politechnika Poznańska

Skład przy użyciu systemu \LaTeX .

Bib \TeX :

```
@mastersthesis{ mnowak-masterthesis,  
  author = "Tomasz Ziętkiewicz",  
  title = "{Optymalizacja klasyfikatora SVM za pomocą programowania genetycznego}",  
  school = "Poznan University of Technology",  
  address = "Pozna{\n}, Poland",  
  year = "2013",  
}
```