

**POLITECHNIKA WARSZAWSKA**

**WYDZIAŁ ELEKTRYCZNY**

Instytut Elektrotechniki Teoretycznej  
i Systemów Informacyjno – Pomiarowych

mgr inż. Tomasz Markiewicz

**SIECI NEURONOWE SVM W ZASTOSOWANIU DO  
KLASYFIKACJI OBRAZÓW KOMÓREK SZPIKU  
KOSTNEGO**

Rozprawa doktorska

Promotor: prof. dr hab. Stanisław Osowski

Warszawa 2005

## Spis treści

Wykaz najważniejszych symboli i oznaczeń	iv
1. Wprowadzenie	1
1.1 Stan aktualny wiedzy	1
1.2 Cel i teza pracy	2
1.3 Przegląd zawartości pracy	4
2. Klasyfikatory neuronowe typu SVM	6
2.1 Sieć SVM do problemu klasyfikacji	6
2.2 Rozpoznawanie wielu klas przy zastosowaniu sieci SVM	14
2.3 Algorytmy obliczeniowe SVM	15
2.3.1. <i>Algorytm LSVM Mangasariana</i>	15
2.3.2. <i>Algorytm ograniczeń aktywnych SVM<sup>Light</sup></i>	18
2.3.3. <i>Algorytm programowania sekwencyjnego Platta</i>	22
2.4 Program SVM_WIN	25
2.5 Analiza porównawcza algorytmów	26
3. Analiza morfologiczna obrazu	28
3.1 Tworzenie postaci dyskretnej obrazu	28
3.2 Progowanie i algorytm Otsu	29
3.3 Podstawowe operacje morfologiczne	30
3.3.1. <i>Erozja</i>	30
3.3.2. <i>Dylatacja</i>	32
3.3.3. <i>Otwarcie i zamknięcie</i>	33
3.4 Filtracja obrazu	34
3.5 Reprezentacja odległościowa obrazu binarnego	34
3.6 Segmentacja obrazu metodą działów wodnych	36
3.7 Transformacja geodezyjna i rekonstrukcja obrazu	37
4. Ekstrakcja obrazu komórek rakowych	39
4.1 Charakterystyka ogólna komórek krwiotwórczych	39
4.2 Algorytm ekstrakcji komórek	49
4.3 Przykłady wydzielonych komórek – wyniki segmentacji	56
5. Generacja i selekcja cech diagnostycznych do rozpoznania komórek	59
5.1 Cechy teksturalne	59

5.2	Cechy geometryczne	63
5.3	Cechy statystyczne	64
5.4	Cechy morfologiczne	65
5.5	Ocena jakości i selekcja cech	68
5.5.1.	<i>Analiza korelacyjna cech</i>	69
5.5.2.	<i>Selekcja cech na podstawie wartości średnich i wariancji danych</i>	71
5.5.3.	<i>Selekcja cech przy użyciu sieci neuronowej SVM o jądrze liniowym</i>	75
6.	Wyniki klasyfikacji	80
6.1	Wyniki rozpoznania 12 rodzajów komórek	81
6.2	Wyniki rozpoznania 17 rodzajów komórek	84
6.3	Wyniki rozpoznania 21 rodzajów komórek	86
6.4	Weryfikacja systemu na podstawie mielogramów wybranych pacjentów	89
7.	Podsumowanie i wnioski	94
	Literatura	97

## Wykaz najważniejszych symboli i oznaczeń

MLP	sieć neuronowa perceptronu wielowarstwowego
SVM	sieć neuronowa typu Support Vector Machine
$y$	sygnał wyjściowy sieci neuronowej
$\mathbf{x}$	wektor wejściowy cech
$\mathbf{w}$	wektor wag sieci neuronowej
$b$	stała polaryzacji sieci neuronowej
$d$	zadany sygnał
$p$	liczba danych uczących
$L(\mathbf{w}, \mathbf{b}, \boldsymbol{\alpha})$	funkcja Lagrange'a
$\boldsymbol{\alpha}$	wektor mnożników Lagrange'a
$\Phi(\mathbf{x})$	wektor funkcji nieliniowych odwzorowania
$N$	wymiar wektora wejściowego
$N_{SV}$	liczba wektorów podtrzymujących
$K(\mathbf{x}_i, \mathbf{x})$	skalarna funkcja jądra
$\xi$	nieujemne wartości dopełniające
$C$	parametr regularyzacyjny
LSVM	algorytm „Lagrangian SVM”
SMO	algorytm programowania sekwencyjnego Platt
BSVM	implementacja algorytmu Platt
$SVM^{Light}$	algorytm ograniczeń aktywnych
SE	element strukturujący
RGB	format zapisu obrazu

# 1. Wprowadzenie

Rozwój technologii informatycznych umożliwił zastosowanie maszyn w rozwiązaniu wielu problemów obliczeniowych. Powstały nowe narzędzia wydobywania i analizy informacji uzyskanej z pomiarów, dostarczające nowych możliwości, dotąd niedostępnych dla człowieka. Nastąpił znaczny postęp w technologii obliczeniowej, przetwarzaniu i gromadzeniu danych w czasie akceptowalnym przez użytkownika. Jedną z dziedzin, w której komputer znajduje coraz szersze zastosowanie, jest inżynieria biomedyczna [2,17,29,35,55,62]. Z jednej strony umożliwia ona automatyczną analizę wielu problemów biomedycznych, z drugiej staje się motorem dalszego rozwoju i postępów w rozumieniu zjawisk zachodzących w organizmach żywych. Główna tematyka pracy dotyczy problemu automatycznego rozpoznawania i klasyfikacji komórek krwiotwórczych szpiku kostnego u ludzi chorych na białaczkę [2,29,30,39,56] za pomocą współczesnych narzędzi matematycznych.

## 1.1 Stan aktualny wiedzy

Problem rozpoznawania komórek krwiotwórczych szpiku kostnego w ich cyklu rozwojowym jest kluczowy przy stawianiu diagnozy co do wystąpienia białaczki i jej rodzaju. W obecnej praktyce medycznej robi to ekspert ludzki (osoba przeszkolona o dużym doświadczeniu). Jest to zadanie bardzo odpowiedzialne a jednocześnie trudne i narażone na liczne błędy. Występuje wiele źródeł tych błędów: zmęczenie oczu, problemy z koncentracją związane z pogodą, trudność w rozpoznaniu dwu komórek sąsiadujących ze sobą, nietrafny (niereprezentatywny) wybór fragmentu rozmazu itp.

Dokładność oszacowania eksperckiego liczby występujących w rozmazie komórek różnego typu jest trudna do określenia. Można mówić jedynie o różnicach tego oszacowania w wykonaniu kilku ekspertów. Wg danych z Instytutu Hematologii w Warszawie, tolerowalne różnice mogą sięgać nawet 15%, co może być uznane za przybliżone oszacowanie błędu eksperta.

Nie istnieją dostępne na rynku urządzenia automatycznego rozpoznawania komórek. Bardzo ograniczone są również prace badawcze dostępne w czasopiśmie naukowych, konferencjach i Internecie. Po wnikliwej analizie literatury światowej poświęconej temu

tematowi znaleziono tylko kilka prac [2,19,56], z czego tylko jedna opublikowana w IEEE Transactions on Systems, Man & Cybernetic [56], czasopiśmie o najwyższym poziomie naukowym.

Metodyka stosowana przez większość autorów jest podobna. Po wydzieleniu pojedynczych komórek z obrazu poszukuje się cech diagnostycznych, możliwie najlepiej je opisujących. Zalicza się do nich cechy wynikające z opisu tekstury, geometrii komórek oraz rozkładu kolorów. Cechy te tworzą wektor  $x$  podlegający porównaniu z prototypem otrzymanym na etapie uczenia. Najczęściej stosowane klasyfikatory to klasyfikatory odległościowe mierzące dystans między aktualnym i prototypowym wektorem  $x$  lub klasyfikatory neuronowe, najczęściej sieć perceptronu wielowarstwowego (MLP).

Dokładności klasyfikacji komórek uzyskane takimi metodami pozostawiają wiele do życzenia. W pracy [2] uzyskano 61% dokładność (39% błędów) przy rozpoznawaniu 16 rodzajów komórek. Sohn [53] uzyskał dokładność 78% ( 22% błędów), ale przy rozpoznawaniu zaledwie 6 rodzajów komórek. W pracy [56] opublikowanej w IEEE Transactions on SMC stosując sieć MLP uzyskano dokładność ogólną rozpoznania rzędu 58% (42% błędów) dla 6 typów komórek. Wyniki prezentowane w publikacjach są więc dalekie od dokładności osiąganey przez eksperta ludzkiego. Jest wiele powodów tak małej dokładności:

- Komórki tego samego typu różnią się znacznie między sobą a jednocześnie są bardzo podobne do komórek innego typu.
- Obrazy rozmazu są bardzo zróżnicowane pod względem barwy i zależne od sposobu obróbki chemicznej i zastosowanych odczynników.
- Proces automatycznego przetwarzania obrazu rozmazu szpiku kostnego jest bardzo trudny i narażony na wiele błędów już w fazie wstępnego przetwarzania.
- Stosowane dotąd klasyfikatory odległościowe lub bazujące na sieciach neuronowych są bardzo niedoskonałe i bardzo wrażliwe na szумы powstałe w obrazach podczas obróbki chemicznej.

Praktycznie nie powstało również żadne szersze opracowanie omawiające problemy wydzielania takich komórek z obrazu, przetwarzania obrazu na cechy i poddawania ich rozpoznaniu i klasyfikacji.

## **1.2 Cel i teza pracy**

Głównym celem pracy jest opracowanie kompletnego systemu automatycznego rozpoznawania i klasyfikacji komórek krwiotwórczych, w którym udział człowieka będzie

ograniczony do minimum. Układ taki powinien charakteryzować się dokładnością zbliżoną do osiąganą w przypadku eksperta ludzkiego przy liczbach rodzajów komórek najczęściej występujących w typowym mielogramie.

Dla zrealizowania powyższego celu głównego konieczne jest rozwiązanie szeregu zadań pośrednich. Do najważniejszych należą:

- Opracowanie metod wstępnego przetwarzania obrazu cyfrowego rozmazu szpiku kostnego dla ekstrakcji pojedynczych komórek. Układ taki powinien charakteryzować się jak najwyższą sprawnością wydzielenia komórek, małą wrażliwością na zniekształcenia i szumy obrazu.
- Opracowanie skutecznych metod generacji cech diagnostycznych obrazu komórki, pozwalających na jak najlepsze zróżnicowanie wartości parametrów odpowiadających różnym typom komórek.
- Opracowanie metodyki selekcji cech diagnostycznych pozwalającej na różnicowanie jakości poszczególnych cech.
- Opracowanie układu klasyfikatora neuronowego charakteryzującego się najwyższą sprawnością i mało wrażliwego na zmienność wartości cech charakteryzujących komórki tego samego typu.
- Powiązanie wymienionych wyżej etapów w jeden projekt układu klasyfikacji komórek krwiotwórczych, spełniający wymagania zdefiniowane na wstępie.

Realizacja powyższych celów została przeprowadzona na bazie danych rzeczywistych otrzymanych przy współpracy z Instytutem Hematologii w Warszawie. Wszystkie próbki rozmazu szpiku kostnego pochodzą z bazy danych tego Instytutu. Poszczególne preparaty zostały dla potrzeb uczenia ocenione przez ekspertów Instytutu i użyte w pracy do trenowania układu klasyfikującego. Dane dla nowych preparatów były klasyfikowane przez opracowany system a następnie weryfikowane przez ekspertów Instytutu. Podejmując się rozwiązania postawionych sobie problemów autor sformułował następującą tezę pracy:

**Zastosowanie sieci SVM w połączeniu z metodami morfologicznymi przetwarzania obrazów i odpowiednią generacją cech diagnostycznych pozwala zbudować automatyczny klasyfikator komórek krwiotwórczych szpiku kostnego zapewniający dokładność zbliżoną do dokładności eksperta ludzkiego.**

### 1.3 Przegląd zawartości pracy

W pracy przedstawiono kompletny układ automatycznego rozpoznawania i klasyfikacji komórek występujących w rozmazie szpiku kostnego. Obrazy rozmazu pobrane z mikroskopu są poddawane wstępnemu przetwarzaniu a następnie ekstrakcji poszczególnych komórek za pomocą operacji morfologicznych [29,31,32,54,63]. Na podstawie barwnych obrazów pojedynczych komórek wyznaczane są cechy, które stanowią ich cyfrową reprezentację [13,29,30,39,50,58,60]. Z uwagi na dużą liczbę generowanych cech, są one analizowane i redukowane do optymalnego zestawu odpowiedniego dla rozróżnienia wybranych typów komórek [16,17,29]. Następnie tak stworzone wektory danych są podawane na wejście klasyfikatora neuronowego, którym w rozwiązaniu jest sieć SVM [4,8,9,18,40,58]. Praca bazuje na danych rzeczywistych uzyskanych dzięki współpracy z Instytutem Hematologii w Warszawie. Dane te zbierane były w okresie 3 ostatnich lat i dotyczą kilkudziesięciu pacjentów w różnym stadium rozwoju choroby.

Praca złożona jest z 7 rozdziałów, z których pierwszy stanowi wprowadzenie do tematyki rozprawy. Zawiera on krótki przegląd stosowanych dotąd rozwiązań, definiuje cele i tezę pracy, a także przedstawia jej zawartość. Rozdział drugi poświęcony jest przedstawieniu podstawowego narzędzia klasyfikacji zastosowanego w pracy, jakim jest sieć neuronowa typu Support Vector Machine (SVM). Zdefiniowano podstawowy problem uczenia charakterystyczny dla tych sieci, a następnie krótki przegląd współczesnych algorytmów rozwiązania tego zadania. Algorytmy uczenia i testowania sieci SVM zostały zaimplementowane w postaci programu SVM\_WIN na platformie Matlab.

Rozdział trzeci poświęcony jest analizie morfologicznej obrazu. Przedstawiono istotę podstawowych operacji morfologicznych zastosowanych przy ekstrakcji komórek z obrazu szpiku kostnego. Przedstawiono i zobrazowano na przykładach komórek działanie wszystkich etapów segmentacji, w tym operacji morfologicznych, filtracji filtrem Gaussa, tworzenia reprezentacji odległościowej obrazu binarnego a następnie końcowej segmentacji obrazu metodą działów wodnych.

Rozdział czwarty dotyczy ekstrakcji pojedynczych komórek z obrazu szpiku kostnego poprzez segmentację. Przedstawiono kompletny opis algorytmu oraz przykładowe wyniki dotyczące wszystkich komórek zawartych w preparacie szpiku kostnego. Omówiono występujące linie krwiotwórcze i ich schematy rozwojowe oraz cechy charakteryzujące poszczególne typy komórek. Całość wzbogacono przykładowymi obrazami ułatwiającymi zrozumienie problemów rozpoznania różnych typów komórek. Ważną cechą opracowanej



metody jest wysoka wydajność wydzielania komórek, szacowana na około 95% i działanie systemu nie wymagające ingerencji człowieka.

Rozdział piąty poświęcony jest generacji i selekcji cech diagnostycznych z obrazu pojedynczych komórek wydzielonych w etapie segmentacji. W tym celu wykorzystano cechy należące do grup statystycznych, geometrycznych, teksturalnych i morfologicznych. Przeprowadzono analizę wpływu poszczególnych cech na wynik klasyfikacji oraz metody selekcji cech najbardziej znaczących. Zamieszczone wyniki mają charakter uniwersalny i mogą znaleźć zastosowanie w innych problemach klasyfikacyjnych, nie związanych z rozpoznawaniem komórek.

Rozdział szósty zawiera wyniki automatycznej klasyfikacji dla różnej liczby typów komórek przy zróżnicowanej liczbie preparatów sporządzonych dla kilkudziesięciu pacjentów Instytutu Hematologii. Badania przeprowadzono dla 12, 17 i 21 typów rozpoznawanych komórek przy zmieniającej się liczbie pacjentów i testowanych preparatów. Badania potwierdziły dobrą skuteczność opracowanej metody o dokładności porównywalnej z dokładnością eksperta ludzkiego. Wytrenowany system został sprawdzony przy sporządzaniu kompletnego mielogramu nowych pacjentów Instytutu Hematologii i porównany z wynikami ekspertów. Porównania te pokazują wysoką zbieżność wyników i dają nadzieję (po wszechstronnym przetestowaniu na wielu pacjentach) na wdrożenie systemu w warunkach szpitalnych.

Rozdział siódmy stanowi podsumowanie pracy oraz wnioski wynikające z badań, w szczególności uwypuklono w nim najważniejsze oryginalne osiągnięcia autora rozprawy. Autor ma zamiar kontynuować badania, szczególnie w zakresie przetestowania urządzenia na dużo większej liczbie pacjentów, dalszego zwiększenia skuteczności systemu w diagnostyce pacjentów, optymalizacji uzysku i możliwości wdrożenia systemu wyposażonego w zautomatyzowany mikroskop, kamerę i komputer w praktyce szpitalnej.

## 2. Klasyfikatory neuronowe typu SVM

W ostatnich latach opracowano i udoskonalono nowy typ sieci neuronowej opartej na metodzie wektorów podtrzymujących, zwany Support Vector Machine (SVM). Znalazła ona duże zastosowanie w rozwiązaniu problemów klasyfikacji i regresji. Pierwsze numeryczne sformułowania metody zamieszczono w pracach [4,58,59] pochodzących z drugiej połowy lat dziewięćdziesiątych. Kolejne prace [8,9,11,24,28,46,48,52] rozwijały tę tematykę zarówno w zastosowaniach do bardziej złożonych problemów, jak i w optymalizacji algorytmów uczących. Stanowią one modyfikację i rozwinięcie pierwotnego sformułowania metody. Główna idea metody pozostała jednak nie zmieniona w stosunku do teorii opracowanej pierwotnie przez V. Vapnika [58].

### 2.1. Sieć SVM do problemu klasyfikacji

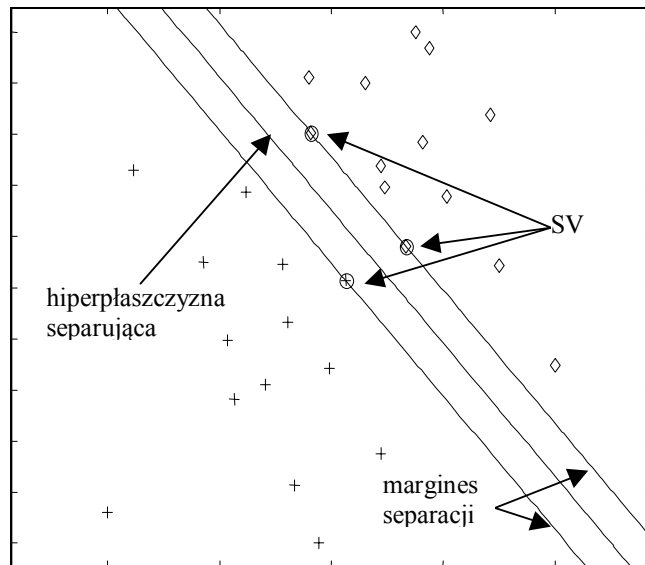
Istotą metody SVM w zastosowaniu do klasyfikacji jest maksymalizacja marginesu separacji pomiędzy dwoma klasami. Sieć SVM jest z definicji układem o jednym neuronie wyjściowym, stąd pojedyncza sieć może separować jedynie dwie klasy. Dla wzorców separowalnych liniowo należących do dwóch różnych klas definiuje się sygnał wyjściowy sieci równy  $y = y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ . Przy założonej tolerancji przypisania sygnału  $y(\mathbf{x})$  do określonej klasy warunki przynależności można zapisać w postaci:

- $y = \mathbf{w}^T \mathbf{x} + b \geq +1$  - przynależność do pierwszej klasy,
- $y = \mathbf{w}^T \mathbf{x} + b \leq -1$  - przynależność do drugiej klasy.

Przyjmuje się sygnał  $d$  zadany na wyjściu sieci równy jest 1 (klasa 1) lub  $-1$  (klasa 2). W zapisie łącznym oba warunki poprawnej klasyfikacji można przedstawić w postaci jednego ogólnego zapisu:

$$d(\mathbf{w}^T \mathbf{x} + b) \geq 1 \quad (2.1)$$

Zadaniem uczenia sieci SVM jest taki dobór wag, który maksymalizuje margines separacji między dwoma klasami (rys. 2.1).



Rys. 2.1 Klasyfikacja dwu klas separowalnych liniowo

Maksymalizacja marginesu separacji na danych uczących daje gwarancję dobrej generalizacji (klasyfikacja bezbłędna na danych testujących w trybie odtwarzania). Należy zauważyć, że odległość pomiędzy hiperpłaszczyzną separującą  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$  a dowolnym punktem przestrzeni  $\mathbf{x}$  jest definiowana w następującej postaci [58]:

$$\sigma = \frac{|g(\mathbf{x})|}{\|\mathbf{w}\|} \quad (2.2)$$

Przyjmując, że punkt  $\mathbf{x}=\mathbf{x}_0$  leży na prostej łączącej punkty najdalej wysunięte dla danej klasy (tzw. wektory podtrzymujące SV) otrzymuje się następujący wzór na odległość hiperpłaszczyzny od wektorów podtrzymujących:

$$\sigma(\mathbf{x}_0) = \frac{|\mathbf{w}^T \mathbf{x}_0 + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|} \quad (2.3)$$

Zakładając, że hiperpłaszczyzna leży dokładnie pośrodku przestrzeni separującej dwie klasy, otrzymuje się wzór na odległość między dwoma klasami w postaci:

$$r = \frac{2}{\|\mathbf{w}\|}$$

Niech  $p$  oznacza liczbę danych uczących. Problem maksymalizacji marginesu separacji między dwoma klasami sprowadza się do minimalizacji kwadratu normy wektora  $\mathbf{w}$ :

$$\min \|\mathbf{w}\|^2 = \min \mathbf{w}^T \mathbf{w} \quad (2.4)$$

przy ograniczeniach funkcyjnych w postaci:

$$d_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0 \quad (2.5)$$

definiowanych dla każdej pary danych uczących,  $i=1,2,\dots,p$ . Rozwiązanie tego problemu uzyskuje się za pośrednictwem tzw. funkcji Lagrange'a. Oznaczając przez  $\alpha_i$  mnożniki Lagrange'a odpowiadające ograniczeniom definiuje się funkcję Lagrange'a w następującej postaci [14]:

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^p \alpha_i [d_i(\mathbf{w}^T \mathbf{x}_i + b) - 1] \quad (2.6)$$

dla  $\alpha_i \geq 0$  ( $i=1,2,\dots,p$ ). Rozwiązanie problemu minimalizacji wartości funkcji Lagrange'a otrzymuje się poprzez minimalizację funkcji L względem wektora  $\mathbf{w}$  i wagi polaryzacji  $b$  oraz maksymalizację względem mnożników Lagrange'a  $\alpha_i$  dla  $i=1,2,\dots,p$ . Uwzględniając pierwszy warunek otrzymuje się:

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^p \alpha_i d_i \mathbf{x}_i = \mathbf{0}, \quad (2.7)$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^p \alpha_i d_i = 0, \quad (2.8)$$

przy spełnieniu ograniczeń  $\boldsymbol{\alpha} \geq \mathbf{0}$  dla  $i=1,2,\dots,p$ . Z równań (2.7) i (2.8) po zróżniczkowaniu funkcji Lagrange'a względem  $\mathbf{w}$  i  $b$  otrzymuje się:

$$\mathbf{w} = \sum_{i=1}^p \alpha_i d_i \mathbf{x}_i, \quad (2.9)$$

$$\sum_{i=1}^p \alpha_i d_i = 0. \quad (2.10)$$

Warunek (2.9) określa wektor wag  $\mathbf{w}$  jako funkcję mnożników Lagrange'a oraz danych uczących  $\mathbf{x}_i$  oraz  $d_i$ . Wstawiając wartości wektora wagowego  $\mathbf{w}$  do wzoru (2.6) przy uwzględnieniu warunku (2.10) otrzymuje się postać funkcji Lagrange'a w punkcie rozwiązania, którą oznaczmy jako  $Z(\boldsymbol{\alpha})$ :

$$Z(\boldsymbol{\alpha}) = \sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j \quad (2.11)$$

podlegającą maksymalizacji względem mnożników  $\alpha_i$  przy ograniczeniach (2.10). Jest to tzw. zadanie dualne. Rozwiązanie zadania maksymalizacji (2.11) przy warunku (2.10) pozwala wyznaczyć optymalne wartości mnożników Lagrange'a, a następnie poszukiwany wektor wagowy  $\mathbf{w}$  sieci przy wykorzystaniu wzoru (2.9). Należy zauważyć, że mnożniki Lagrange'a

odpowiadające znakowi większości ograniczenia  $d_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 > 0$  są z definicji równe zero. Jedynie mnożniki odpowiadające ograniczeniom aktywnym czyli przy znaku równości ograniczenia (2.5) są niezerowe. Stąd wzór określający wektor wagowy  $\mathbf{w}$  sieci można uprościć do postaci:

$$\mathbf{w} = \sum_{i=1}^{N_{SV}} \alpha_i d_i \mathbf{x}_i$$

gdzie  $N_{SV}$  oznacza liczbę wektorów podtrzymujących. Należy podkreślić, że zastosowanie sieci liniowej do problemu nieseparowalnego liniowo prowadzi zwykle do znanych błędów klasyfikacji.

Dla problemu nieseparowalnego liniowo wprowadza się nieliniowe odwzorowanie przestrzeni wielowymiarowej poprzez funkcję nieliniową  $\Phi(\mathbf{x})$ :

$$\Phi(\mathbf{x}) = \begin{bmatrix} \Phi_1(\mathbf{x}) \\ \Phi_2(\mathbf{x}) \\ \vdots \\ \Phi_k(\mathbf{x}) \end{bmatrix}$$

Funkcja  $\Phi(\mathbf{x})$  spełniająca rolę rzutowania danych oryginalnych nieseparowalnych liniowo w inną przestrzeń, o wymiarze  $k \geq N$ , w której są one separowalne na mocy twierdzenia Covera [18,58]. Wszystkie przedstawione wcześniej zależności pozostają słuszne przy zastąpieniu wektora  $\mathbf{x}$  przez wektor odwzorowań nieliniowych  $\Phi(\mathbf{x}) \in \mathbf{R}^k$ . Wektor  $\Phi(\mathbf{x})$  zastępuje zmienną  $\mathbf{x}$  w wyrażeniu (2.11), a zadanie dualne optymalizacji dotyczy wówczas maksymalizacji funkcji  $Z(\boldsymbol{\alpha})$  zdefiniowanej następująco [4]:

$$Z(\boldsymbol{\alpha}) = \sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j d_i d_j \Phi^T(\mathbf{x}_i) \Phi(\mathbf{x}_j) \quad (2.12)$$

Rozwiązanie problemu (2.12) odbywa się w identyczny sposób jak problemu (2.11), gdyż dane uczące przetransformowane nieliniowo nie wnoszą istotnej zmiany problemu optymalizacyjnego w stosunku do danych separowalnych liniowo. Rozwiązanie problemu (2.12) pozwala uzyskać optymalne wartości mnożników Lagrange'a, które z kolei wyznaczają wektor wagowy  $\mathbf{w}$  analogicznie jak w przypadku poprzednim, poprzez formalne zastąpienie  $\mathbf{x}$  wektorem  $\Phi(\mathbf{x})$ , tzn:

$$\mathbf{w} = \sum_{i=1}^{N_{SV}} \alpha_i d_i \Phi(\mathbf{x}_i) \quad (2.13)$$

gdzie podobnie jak poprzednio  $N_{SV}$  oznacza liczbę wektorów podtrzymujących. Wprowadźmy oznaczenie:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi^T(\mathbf{x}_i)\Phi(\mathbf{x}_j) \quad (2.14)$$

gdzie  $K(\mathbf{x}_i, \mathbf{x}_j)$  oznacza skalarną funkcję jądra. Transformacja przestrzeni  $\mathbf{x}$  poprzez funkcję jądra  $K(\mathbf{x}_i, \mathbf{x}_j)$  jest równoznaczna ze zrzutowaniem danych oryginalnych  $\mathbf{x}$  w inną przestrzeń wielowymiarową. Funkcja jądra  $K(\mathbf{x}_i, \mathbf{x}_j)$  jest z definicji symetryczna i wyrażana jako iloczyn skalarny dwu wektorów funkcyjnych  $\Phi(\mathbf{x}_i)$  i  $\Phi(\mathbf{x}_j)$ . Zostało pokazane [8,51,58], że funkcją jądra może być każda funkcja spełniająca warunki twierdzenia Mercera [58]. Zgodnie z tym twierdzeniem jądro  $K(\mathbf{x}_i, \mathbf{x}_j)$  jest symetryczne i daje się przedstawić jako iloczyn skalarny dwu wektorów funkcyjnych  $\Phi(\mathbf{x}_i)$  oraz  $\Phi(\mathbf{x}_j)$ , jeśli dla każdej funkcji ciągłej  $g(\mathbf{x}_i)$  spełniającej warunek:

$$\int g(\mathbf{x}_i)^2 dx_i < \infty \quad (2.15)$$

zachodzi:

$$\int K(\mathbf{x}_i, \mathbf{x}_j)g(\mathbf{x}_i)g(\mathbf{x}_j)dx_idx_j \geq 0 \quad (2.16)$$

Najczęściej stosowane w praktyce funkcje jądra to:

1. funkcja wielomianowa  $K(\mathbf{x}_i, \mathbf{x}) = (\mathbf{x}_i^T \mathbf{x} + \gamma)^p$  (2.17)

2. funkcja Gaussa (radialna)  $K(\mathbf{x}_i, \mathbf{x}) = e^{-\frac{|\mathbf{x}-\mathbf{x}_i|^2}{2\sigma^2}}$  (2.18)

3. funkcja liniowa  $K(\mathbf{x}_i, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}_i + \gamma)$  (2.19)

Zostało udowodnione, że wszystkie wymienione funkcje spełniają warunki Mercera bez żadnych ograniczeń [51,58].

Dla problemów nieseparowalnych liniowo nawet przy zastosowaniu rzutowania nieliniowego danych wprowadza się nieujemne wartości dopełniające  $\xi_i$  będące zmiennymi zmniejszającymi realny margines separacji. Oznaczmy przez:

$$\xi = \begin{bmatrix} \xi_1 \\ \xi_2 \\ \cdot \\ \xi_p \end{bmatrix}$$

wektor zmiennych dopełniających. Dodatkowo definiuje się wartość  $C$  będącą parametrem regularyzacyjnym, określającym wagę z jaką traktuje się możliwe błędy uczenia w stosunku do

wymagań wynikających z marginesu separacji. Funkcję celu podlegającą minimalizacji zapisać wówczas można jako tzw. problem pierwotny:

$$E(\mathbf{w}, \xi) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^p \xi_i \quad (2.20)$$

Ograniczenia w tym przypadku przybierają postać:

$$d_i (\mathbf{w}^T \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad (2.21)$$

$$\xi_i \geq 0. \quad (2.22)$$

dla  $i=1,2,\dots,p$ . Rozwiązanie zadania optymalizacji kwadratowej zdefiniowanego równaniami (2.20), (2.21) i (2.22) uzyskuje się podobnie jak poprzednio poprzez jego transformację do postaci funkcji Lagrange'a, a następnie minimalizację wartości tej funkcji:

$$\min L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^p \xi_i - \sum_{i=1}^p \alpha_i d_i [\mathbf{w}^T \Phi(\mathbf{x}_i) + b - 1 + \xi_i] \quad (2.23)$$

przy  $\xi_i \geq 0$  i mnożnikach Lagrange'a  $\alpha_i \geq 0$ .

Minimum wartości funkcji Lagrange'a  $L$  uzyskuje się minimalizując ją względem  $\mathbf{w}$  i  $\xi$  oraz maksymalizując względem mnożników Lagrange'a. Z warunku  $\frac{\partial L}{\partial \mathbf{w}} = 0$  wynika [58], że:

$$\mathbf{w} = \sum_{i=1}^p \alpha_i d_i \Phi(\mathbf{x}_i) \quad (2.24)$$

Podobnie z warunku  $\frac{\partial L}{\partial b} = 0$  wynika, że:

$$\sum_{i=1}^p \alpha_i d_i = 0 \quad (2.25)$$

Uwzględniając zależności (2.24) i (2.25) w sformułowaniu problemu pierwotnego otrzymuje się problem dualny, który można przedstawić następująco:

$$\max Z(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j d_i d_j K(\mathbf{x}_i, \mathbf{x}_j) \quad (2.26)$$

przy następujących ograniczeniach:

$$0 \leq \alpha_i \leq C, \quad (2.27)$$

$$\sum_{i=1}^p \alpha_i d_i = 0. \quad (2.28)$$

Rozwiązaniem zadania dualnego są optymalne wartości mnożników Lagrange'a  $\alpha_i$ . Ich znajomość pozwala wyznaczyć wektor wagowy  $\mathbf{w}$  sieci według zależności (2.24). Przy

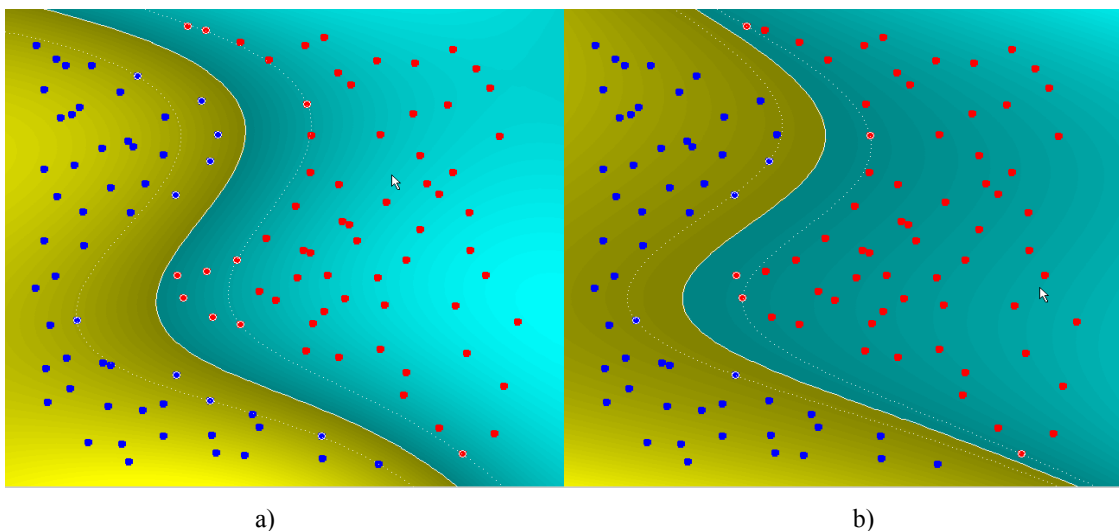
uwzględnieniu zerowych wartości  $\alpha_i$ , odpowiadających wektorom  $\mathbf{x}_i$  tworzącym ograniczenia nierównościowe, wzór (2.24) sprowadza się do postaci:

$$\mathbf{w} = \sum_{i=1}^{N_{sv}} \alpha_i d_i \Phi(\mathbf{x}_i)$$

Współczynnik  $b$  wyznacza się korzystając z dowolnego wektora podtrzymującego  $\mathbf{x}_{sv}$ , dla którego  $\mathbf{w}^T \Phi(\mathbf{x}_{sv}) + b = 1$ , skąd otrzymuje się:

$$b = 1 - \mathbf{w}^T \Phi(\mathbf{x}_{sv})$$

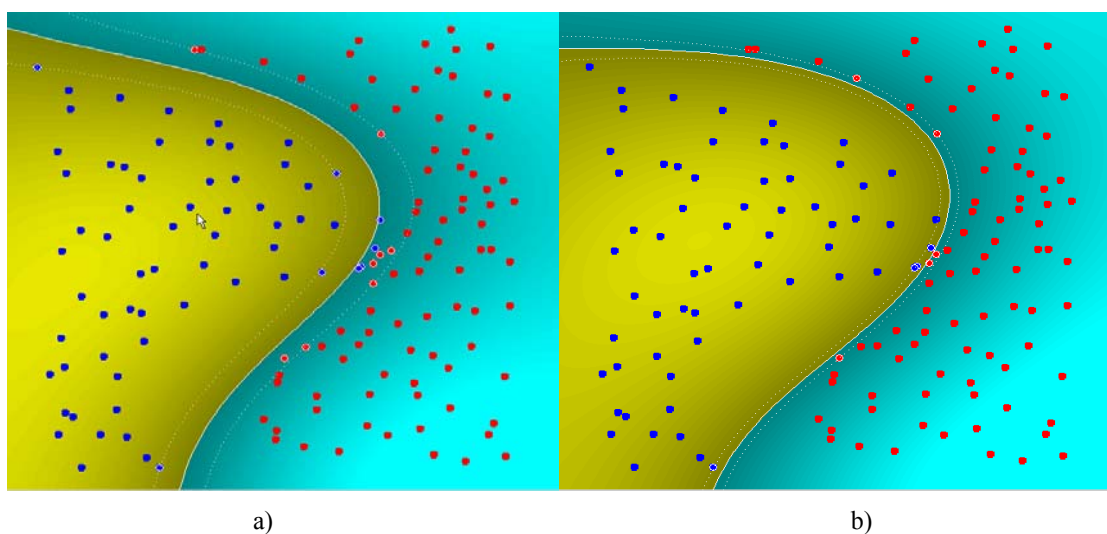
Ze sformułowania SVM problemu klasyfikacyjnego widać, że problem uczenia sieci neuronowej do klasyfikacji danych, niezależnie od separowalności klas, sprowadza się do zadania programowania kwadratowego względem mnożników Lagrange'a. Zadanie to należy do stosunkowo dobrze rozpoznanych w metodach optymalizacji i zawsze prowadzi do minimum globalnego.



Rys. 2.2 Klasyfikacja danych nieseparowalnych liniowo siecią SVM z jądrem gaussowskim (a) i wielomianowym (b)

Na rys. 2.2 a,b przedstawiono działanie klasyfikacyjne nieliniowej sieci SVM dla danych nieseparowalnych liniowo przy wartości  $C=30$ . Zastosowano jądro gaussowskie (rys. 2.2a) o  $\gamma=1$  oraz wielomianowe (rys. 2.2b) o rzędzie wielomianu równym 3. Oba rodzaje sieci dokonały bezbłędnej klasyfikacji, choć ułożenie hiperpłaszczyzny separującej i liczba wektorów podtrzymujących w obu przypadkach są różne. Liczba wektorów podtrzymujących była równa 19 (15.8%) dla jądra gaussowskiego i 9 (7.5%) dla jądra wielomianowego.





Rys. 2.3 Klasyfikacja danych przy zastosowaniu sieci SVM o jądrze gaussowskim: a)  $C=50$ ,  
b)  $C=1000$

Na rys. 2.3 zilustrowano wpływ wartości stałej regularyzacyjnej  $C$  na ukształtowanie hiperpłaszczyzny separującej. Przy wartości  $\gamma=1$  oraz  $C=50$  (rys. 2.3a) margines separacji jest szeroki, liczba wektorów podtrzymujących równa 16 (10.5%). Algorytm uczący zakończył proces uczenia z czterema błędami. Rys. 2.3b pokazuje rozwiązanie bezbłędne problemu przy zastosowaniu 9 (5.9%) wektorów podtrzymujących i  $C=1000$ . Zwiększenie wartości parametru  $C$  umożliwiło bezbłędną klasyfikację danych uczących, ale jednocześnie zmniejszenie szerokości marginesu separacji, a więc pogorszenie własności generalizacyjnych sieci.

Interesująca jest interpretacja wartości optymalnych mnożników Lagrange'a otrzymanych w wyniku uczenia. Wartość  $\alpha_i=0$  oznacza, że warunek nierównościowy spełniany jest z nadmiarem. Zmienna  $\mathbf{x}_i$  nie tworzy wektora podtrzymującego i nie ma żadnego wpływu na wektor wagowy  $\mathbf{w}$  opisujący hiperpłaszczyznę. Usunięcie danej pary uczącej  $(\mathbf{x}_i, d_i)$  nie ma więc praktycznie żadnego wpływu na rozwiązanie problemu klasyfikacyjnego.

Wartość  $0 < \alpha_i < C$  oznacza, że zmienna  $\mathbf{x}_i$  odpowiada dokładnie ograniczeniu aktywnemu ( $\mathbf{w}^T \Phi(\mathbf{x}_i) + b = 1$ ). Wartość zmiennej dopełniającej jest równa  $\xi_i=0$ . Wektor  $\mathbf{x}_i$  tworzy zatem wektor podtrzymujący, mający wpływ na hiperpłaszczyznę. Jego usunięcie ze zbioru danych uczących miałyby więc wpływ na wynik rozwiązania optymalnego.

Wartość mnożnika Lagrange'a pozostająca na ograniczeniu górnym  $\alpha_i=C$  oznacza, że wektor  $\mathbf{x}_i$  znalazł się w zakresie marginesu separacji, a odpowiadająca mu zmienna dopełniająca  $\xi_i>0$ . Zmienna taka należy więc do zbioru wektorów podtrzymujących i ma również istotny wpływ na rozwiązanie problemu i ukształtowanie hiperpłaszczyzny separującej.

## 2.2. Rozpoznawanie wielu klas przy zastosowaniu sieci SVM

Sieci SVM z istoty swego działania dokonują rozdziału danych na dwie klasy. W odróżnieniu od sieci klasycznych, gdzie liczba klas poddanych klasyfikacji może być dowolna, rozpoznanie wielu klas przy pomocy tej techniki wymaga przeprowadzenia wielokrotnej klasyfikacji [21,61]. Do najbardziej znanych rozwiązań tego problemu należy metoda „jeden przeciw wszystkim”. W metodzie tej przy M klasach należy skonstruować M sieci, każda odpowiedzialna za rozpoznanie jednej klasy. Sieć i-ta jest trenowana na danych uczących, w których przykłady i-tej klasy są skojarzone z  $d_i=1$ , a pozostałe z  $d_i=-1$ . Po wytrenowaniu wszystkich sieci następuje etap odtwarzania, w którym ten sam wektor  $\mathbf{x}$  jest podawany na każdą sieć SVM. Określone są wartości wszystkich M sygnałów wyjściowych sieci (funkcji decyzyjnych):

$$\begin{aligned}y_1(\mathbf{x}) &= \mathbf{w}_1^T \boldsymbol{\phi}(\mathbf{x}) \\y_2(\mathbf{x}) &= \mathbf{w}_2^T \boldsymbol{\phi}(\mathbf{x}) \\&\dots \\y_M(\mathbf{x}) &= \mathbf{w}_M^T \boldsymbol{\phi}(\mathbf{x})\end{aligned}\tag{2.29}$$

Wektor  $\mathbf{x}$  jest następnie zaliczany do klasy o największej wartości funkcji decyzyjnej.

Innym rozwiązaniem jest zastosowanie metody „jeden przeciw jednemu” [21,61]. W metodzie tej konstruuje się  $M(M-1)/2$  klasyfikatorów typu SVM rozróżniających za każdym razem 2 klasy danych ze zbioru uczącego, kolejno parowanych ze sobą. Oznaczmy równanie decyzyjne sieci SVM rozróżniającej między klasą i-tą a j-tą w postaci:

$$y_{ij}(\mathbf{x}) = \mathbf{w}_{ij}^T \boldsymbol{\phi}(\mathbf{x})\tag{2.30}$$

Po wytrenowaniu wszystkich sieci można przystąpić do właściwej klasyfikacji przy założeniu konkretnej wartości wektora  $\mathbf{x}$ . Jeśli  $\text{sign}[\mathbf{w}_{ij}^T \boldsymbol{\phi}(\mathbf{x} + b)]$  wskazuje na i-tą klasę, należy zwiększyć o jeden sumę wskazań do tej klasy. W przeciwnym wypadku zwiększyć o jeden sumę wskazań do klasy j-tej. Klasa o największej sumie zwycięstw wśród  $M(M-1)/2$  sieci jest uważana za zwycięską (wektor  $\mathbf{x}$  zalicza się ostatecznie do tej klasy).

Są znane również inne podejścia do rozpoznawania wielu klas, pozwalające na rozwiązanie problemu w jednej strukturze SVM, przy uwzględnieniu wszystkich danych uczących na raz i z zastosowaniem dekompozycji zbioru uczącego na mniejsze podzbiory dla przyśpieszenia

procedury uczenia. Do takich metod należą między innymi metoda Cramera i Singera [11] oraz oryginalna metoda C. Hsu i C. Lina [21].

### 2.3 Algorytmy obliczeniowe SVM

Jak pokazano każde uczenie sieci SVM sprowadza się do rozwiązania zadania programowania kwadratowego, które może być rozwiązane przy użyciu wielu dostępnych aktualnie algorytmów optymalizacyjnych, takich jak MINOS [34] opracowany w Stanford Optimal Lab, pakiet OLS firmy IBM [22], LOQO uniwersytetu Princeton [57] czy Pakiet Optymalizacyjny Matlaba [33].

Pakiety ogólnie dostępne rozwiązania tego zadania nie są jednak dostosowane do zadań o dużej liczbie zmiennych optymalizowanych i danych uczących. W przypadku sieci neuronowych liczba danych uczących może sięgać nawet miliona. W pierwszej fazie optymalizacji, gdy każda dana ucząca tworzy wektor podtrzymujący liczbę zmiennych optymalizowanych jest równa liczbie danych uczących. To powoduje, że programy optymalizacyjne, ogólnie dostępne, są niewydolne i nie nadają się do zastosowań praktycznych w uczeniu sieci SVM.

W ostatnich latach powstały algorytmy specjalizowane do tego typu zadań, charakteryzujące się dużą szybkością działania i stosownością do rozwiązania zadań o ogromnej (mierzonej w milionach) liczbie danych uczących. Do najlepszych algorytmów specjalizowanych zaliczyć aktualnie można algorytm SMO Platt [48], SVM<sup>Light</sup> Joachimsa [24], oraz Lagrangian SVM (LSVM) Mangasariana [28]. Poniżej przedstawiono w zarysach główne aspekty tych algorytmów.

#### 2.3.1. Algorytm LSVM Mangasariana

Prezentując algorytm LSVM omówiono w pierwszej kolejności algorytm stosowany do sieci liniowej, a następnie jego uogólnienie w odniesieniu do sieci SVM stosującej jądro nieliniowe. Dualny problem uczenia w odniesieniu do sieci liniowej jest sformułowany następująco [58]:

$$\max Z(\boldsymbol{\alpha}) = \sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \cdot \mathbf{x}_j \quad (2.31)$$

przy ograniczeniach:

$$\sum_{i=1}^p \alpha_i d_i = 0 \quad (2.32)$$

$$0 \leq \alpha_i \leq C$$

Wprowadźmy oznaczenia:

$$\mathbf{D} = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_p \end{bmatrix}, \mathbf{A} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_N^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_N^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(p)} & x_2^{(p)} & \dots & x_N^{(p)} \end{bmatrix}, \mathbf{e} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_p \end{bmatrix}, \mathbf{1} = \text{diag}([1, 1, \dots, 1]).$$

Przy takich oznaczeniach równanie (2.31) można przekształcić do zadania minimalizacyjnego o następującej postaci:

$$\min \frac{1}{2} \boldsymbol{\alpha}^T \left( \frac{\mathbf{1}}{C} + \mathbf{D}(\mathbf{A}^T \mathbf{A} + \mathbf{e} \mathbf{e}^T) \mathbf{D} \right) \boldsymbol{\alpha} - \mathbf{e}^T \boldsymbol{\alpha} \quad (2.33)$$

Wprowadzając nowe oznaczenia:

$$\mathbf{H} = \mathbf{D} \cdot [\mathbf{A} \quad -\mathbf{e}] \quad (2.34)$$

$$\mathbf{Q} = \frac{1}{C} \mathbf{1} + \mathbf{H} \mathbf{H}^T \quad (2.35)$$

zadanie (2.33) można zapisać w sposób standardowy dla programowania kwadratowego [14]:

$$\min \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha} - \mathbf{e}^T \boldsymbol{\alpha} \quad (2.36)$$

Rozwiązanie zadania (2.36) musi spełniać warunki Kuhna-Tuckera, które tutaj można zapisać w postaci:

$$\begin{aligned} \boldsymbol{\alpha} &\geq \mathbf{0} \\ \mathbf{Q} \boldsymbol{\alpha} - \mathbf{e} &\geq \mathbf{0} \\ \boldsymbol{\alpha} &\perp \mathbf{Q} \boldsymbol{\alpha} - \mathbf{e} \end{aligned} \quad (2.37)$$

Dla dodatnich wektorów  $\mathbf{a}$  i  $\mathbf{b}$  ortogonalnych względem siebie  $\mathbf{0} < \mathbf{a} \perp \mathbf{b} > \mathbf{0}$  obowiązuje następująca relacja [28]:

$$\mathbf{a} = (\mathbf{a} - \beta \mathbf{b})_+ \quad (2.38)$$

przy stałej  $\beta > 0$ , gdzie oznaczenie  $\mathbf{x}_+$  oznacza wektor  $\mathbf{x}$  ze wszystkimi składowymi ujemnymi przyrównanymi do zera. Stosując tę identyczność do zależności (2.37) otrzymuje się:

$$\mathbf{Q} \boldsymbol{\alpha} - \mathbf{e} = ((\mathbf{Q} \boldsymbol{\alpha} - \mathbf{e}) - \beta \boldsymbol{\alpha})_+ \quad (2.39)$$

Warunek (2.39) prowadzi do prostego schematu iteracyjnego obliczeń:

$$\boldsymbol{\alpha}^{i+1} = \mathbf{Q}^{-1} (\mathbf{e} + ((\mathbf{Q} \boldsymbol{\alpha}^i - \mathbf{e}) - \beta \boldsymbol{\alpha}^i)_+) \quad (2.40)$$

gdzie  $i$  oznacza numer kolejnej iteracji. Udowodniono [28], że iteracje określone zależnością (2.40) są zbieżne liniowo do rozwiązania właściwego, o ile stała  $\beta$  spełnia warunek:

$$0 \leq \beta \leq \frac{2}{C} \quad (2.41)$$

Zauważmy, że w kolejnych krokach iteracyjnych wymagana jest inwersja macierzy  $\mathbf{Q}$ . W obliczeniu tej inwersji najwygodniej jest wykorzystać formułę Shermana-Morrisona-Woodbury'ego [14,15,28], zgodnie z którą:

$$\left(\frac{\mathbf{1}}{C} + \mathbf{H}\mathbf{H}^T\right)^{-1} = C\left(\mathbf{1} - \mathbf{H}\left(\frac{\mathbf{1}}{C} + \mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\right) \quad (2.42)$$

Biorąc pod uwagę, że  $\mathbf{H} \in \mathbf{R}^{p \times (N+1)}$ , iloczyn  $\mathbf{H}^T\mathbf{H}$  jest macierzą o wymiarze  $(N+1) \times (N+1)$ . Oznacza to, że odwracanie dotyczy macierzy o wymiarze  $(N+1) \times (N+1)$ , gdzie  $N$  jest wymiarem wektora  $\mathbf{x}$ . Zwykle  $N \ll p$ , co oznacza małą złożoność obliczeniową algorytmu.

W przypadku zastosowania jądra nieliniowego  $K(\mathbf{x}, \mathbf{x}_i)$ , można w prosty sposób zaadoptować powyższy algorytm uczenia. Zauważmy, że równania liniowe hiperpłaszczyzny  $\sum_{j=1}^N w_j x_j + b = 0$

w przypadku sieci nieliniowej zostaje zastąpione równaniem  $\sum_{j=1}^K w_j \Phi_j(\mathbf{x}) + b = 0$ . Rozwiązanie

problemu uzależnione jest od mnożników Lagrange'a, a równanie hiperpłaszczyzny można

przedstawić w postaci  $y(\mathbf{x}) = \sum_{i=1}^{Nsv} \alpha_i d_i K(\mathbf{x}, \mathbf{x}_i) = 0$ . Wyrażenie na  $\mathbf{Q}$  w standardowym

sformułowaniu (2.35) programowania kwadratowego można teraz zapisać w postaci:

$$\mathbf{Q} = \frac{\mathbf{1}}{C} + \mathbf{D}\mathbf{K}\left(\begin{bmatrix} \mathbf{A} - \mathbf{e} \\ \mathbf{A}^T \\ -\mathbf{e}^T \end{bmatrix}\right)\mathbf{D} \quad (2.43)$$

gdzie  $\mathbf{K}$  jest macierzą utworzoną przez funkcje jądra  $K(\mathbf{x}, \mathbf{x}_i)$ , dla  $i=1,2,\dots,p$ ,  $j=1,2,\dots,p$  analogicznej do wzorów (2.24) i (2.26) obowiązujących w sieciach liniowych. Zadanie minimalizacji dotyczy teraz problemu:

$$\min \left\{ \frac{1}{2} \boldsymbol{\alpha}^T \left( \frac{\mathbf{1}}{C} + \mathbf{D}\mathbf{K}\left(\begin{bmatrix} \mathbf{A} & -\mathbf{e} \\ \mathbf{A}^T \\ -\mathbf{e}^T \end{bmatrix}\right)\mathbf{D} \right) \boldsymbol{\alpha} - \mathbf{e}^T \boldsymbol{\alpha} \right\} \quad (2.44)$$

bezpośrednio odpowiadającego zależności (2.33) lub postaci ogólnej (2.36). Rozwiązanie problemu (2.44) musi spełniać warunki Kuhna-Tuckera, to znaczy:

$$\begin{aligned} \boldsymbol{\alpha} &\geq \mathbf{0} \\ \mathbf{Q}\boldsymbol{\alpha} - \mathbf{e} &\geq \mathbf{0} \\ \boldsymbol{\alpha} \perp \mathbf{Q}\boldsymbol{\alpha} - \mathbf{e} \end{aligned} \quad (2.45)$$

które są identyczne z (2.37), z jedyną różnicą dotyczącą sformułowania macierzy  $\mathbf{Q}$ . Rozwiązanie na wektor  $\boldsymbol{\alpha}$  otrzymuje się poprzez iteracje zdefiniowane w postaci (2.40). Jak udowodniono w [28] obowiązują tu te same warunki liniowej zbieżności do rozwiązania jak w

przypadku poprzednim. Jediną różnicą jest niemożność skorzystania z formuły Shermana-Morrisona-Woodbury'ego, ze względu na inną postać macierzy  $\mathbf{Q}$ . Oznacza to większą złożoność obliczeniową algorytmu w stosunku do przypadku liniowego.

### 2.3.2. Algorytm ograniczeń aktywnych SVM<sup>Light</sup>

Rozwiązanie zadania programowania kwadratowego w problemie klasyfikacji nieliniowej sprowadza się do rozwiązania zadania dualnego o postaci opisanej zależnościami (2.31) oraz (2.32). Definiując macierz kwadratową  $\mathbf{Q}$  o elementach  $Q_{ij} = d_i d_j K(\mathbf{x}_i, \mathbf{x}_j)$  oraz wprowadzając zapis wektorowy:

$$\mathbf{d} = \begin{bmatrix} d_1 \\ d_2 \\ \cdot \\ d_p \end{bmatrix}$$

zależność (2.31) przy ograniczeniach (2.32) można przedstawić w postaci zależności macierzowych, identycznych do (2.36), to znaczy:

$$\min \frac{1}{2} \mathbf{a}^T \mathbf{Q} \mathbf{a} - \mathbf{e}^T \mathbf{a} \quad (2.46)$$

przy ograniczeniach zapisanych w postaci wektorowej:

$$\begin{aligned} \mathbf{a}^T \mathbf{d} &= 0 \\ 0 &\leq \mathbf{a} \leq C \cdot \mathbf{e} \end{aligned} \quad (2.47)$$

Wymiar problemu optymalizacyjnego jest uzależniony od liczby danych uczących  $p$ . Przy dużych wymiarach ( $p > 10000$ ) operowanie macierzą  $\mathbf{Q}$  i przechowywanie jej w pamięci nawet współczesnych komputerów staje się bardzo trudne. Z drugiej strony ciągle odtwarzanie  $\mathbf{Q}$  powoduje ogromne spowolnienie algorytmu.

Radą na to jest podzielenie zbioru danych na mniejsze podzbiory i zmniejszenie w ten sposób wymiaru problemu optymalizacyjnego. Dekompozycja dzieli pełny zbiór danych na 2 podzbiory: aktywny (tzw. roboczy) dla którego poszukuje się rozwiązania optymalnego oraz nieaktywny, dla którego wszystkie warunki optymalności są z góry spełnione. Takie rozwiązanie problemu bazuje na tzw. strategii ograniczeń aktywnych [14,24]. Zgodnie z tą metodą w każdej iteracji zmienne optymalizowane  $\alpha_i$  są zaliczone do 2 kategorii:

- zbiór B zmiennych swobodnych, poddanych optymalizacji
- zbiór N zmiennych stałych o wartościach ustalonych przez ograniczenia.

Zmienne swobodne ze zbioru B podlegają aktualizacji, a ustalone pozostają chwilowo w ograniczeniach dolnych bądź górnych. W każdej iteracji sprawdzane są warunki optymalności rozwiązania. Zgodnie z metodyką funkcji Lagrange'a sprawdza się warunki Kuhna-Tuckera dla funkcji Lagrange'a odpowiadającej problemowi (2.46) i (2.47). Przy uwzględnieniu różnego typu ograniczeń nakładanych na zmienne można funkcję Lagrange'a zapisać w postaci [24]:

$$L(\mathbf{a}, \boldsymbol{\lambda}) = -\mathbf{a}^T \mathbf{e} + \frac{1}{2} \mathbf{a} \mathbf{Q}^T \mathbf{d} - \lambda_{eq} \mathbf{a}^T \mathbf{d} - \sum_{i=1}^p \lambda_i^{lo} \alpha_i - \sum_{i=1}^p \lambda_i^{up} (C - \alpha_i) \quad (2.48)$$

W równaniu tym  $\lambda_{eq}$ ,  $\lambda_i^{lo}$ ,  $\lambda_i^{up}$  ( $i=1,2,\dots,p$ ) oznaczają dodatkowe mnożniki Lagrange'a odpowiadające za ograniczenia, odpowiednio: równościowe ( $\mathbf{a}^T \mathbf{d}=0$ ), ograniczenia dolne kostkowe ( $\alpha_i \geq 0$ ) oraz ograniczenia górne kostkowe ( $C - \alpha_i \geq 0$ ). Optymalność rozwiązania wymaga, aby  $\frac{dL(\mathbf{a}, \boldsymbol{\lambda})}{d\mathbf{a}} = 0$  oraz iloczyn odpowiedniego mnożnika przez wartość funkcji ograniczenia był równy zeru, przy równoczesnym spełnieniu wszystkich ograniczeń. Oznacza to, że musi być spełniony następujący układ warunków:

$$\begin{aligned} -\mathbf{e} + \mathbf{Q}\mathbf{a} + \lambda_{eq} \mathbf{d} - \boldsymbol{\lambda}^{lo} - \boldsymbol{\lambda}^{up} &= \mathbf{0} \\ \lambda_i^{lo} \alpha_i &= 0 \\ \lambda_i^{up} (C - \alpha_i) &= 0 \\ \boldsymbol{\lambda}^{lo} &\geq \mathbf{0} \\ \boldsymbol{\lambda}^{up} &\geq \mathbf{0} \\ \mathbf{a}^T \mathbf{y} &= 0 \\ \mathbf{0} < \mathbf{a} < C \cdot \mathbf{e} \end{aligned} \quad (2.49)$$

dla  $i=1,2,\dots,p$  oraz  $\boldsymbol{\lambda}^{lo} = [\lambda_1^{lo} \quad \lambda_2^{lo} \quad \dots \quad \lambda_p^{lo}]^T$ ,  $\boldsymbol{\lambda}^{up} = [\lambda_1^{up} \quad \lambda_2^{up} \quad \dots \quad \lambda_p^{up}]^T$ . Jeśli wszystkie warunki są spełnione, rozwiązanie w danej iteracji jest optymalne. W przeciwnym przypadku algorytm dokonuje dekompozycji zmiennych  $\mathbf{a}$  i macierzy  $\mathbf{Q}$  odpowiadających podziałowi na zmienne aktywne (bazowe), oznaczone wskaźnikiem B i stałe oznaczone wskaźnikiem N, które mogą być przedstawione następująco:

$$\mathbf{a} = \begin{bmatrix} \mathbf{a}_B \\ \mathbf{a}_N \end{bmatrix} \quad \mathbf{d} = \begin{bmatrix} \mathbf{d}_B \\ \mathbf{d}_N \end{bmatrix} \quad \mathbf{Q} = \begin{bmatrix} \mathbf{Q}_{BB} & \mathbf{Q}_{BN} \\ \mathbf{Q}_{NB} & \mathbf{Q}_{NN} \end{bmatrix} \quad (2.50)$$

Ze względu na symetrię macierzy  $\mathbf{Q}$  mamy  $\mathbf{Q}_{BN} = \mathbf{Q}_{NB}^T$ . Przy takiej dekompozycji problem minimalizacyjny (2.46) można teraz zapisać w postaci:

$$\min \left\{ -\mathbf{a}_B^T (\mathbf{e} - \mathbf{Q}_{BN} \mathbf{a}_N) + \frac{1}{2} \mathbf{a}_B^T \mathbf{Q}_{BB} \mathbf{a}_B + \frac{1}{2} \mathbf{a}_N^T \mathbf{Q}_{NN} \mathbf{a}_N - \mathbf{a}_N^T \mathbf{e} \right\} \quad (2.51)$$

Biorąc pod uwagę, że zbiór N oznacza stałe wartości  $\mathbf{a}_N$ , nie mające wpływu na kierunek minimalizacji, można powyższe zadanie uprościć do:

$$\min \left\{ -\mathbf{a}_B^T (\mathbf{e} - \mathbf{Q}_{BN} \mathbf{a}_N) + \frac{1}{2} \mathbf{a}_B^T \mathbf{Q}_{BB} \mathbf{a}_B \right\} \quad (2.52)$$

przy ograniczeniach:

$$\begin{aligned} \mathbf{a}_B^T \mathbf{d}_B + \mathbf{a}_N^T \mathbf{d}_N &= 0 \\ \mathbf{0} < \mathbf{a} < C \cdot \mathbf{e} \end{aligned} \quad (2.53)$$

Postęp optymalizacji jest uzależniony od właściwego przypisania zmiennych do zbioru B i N, a następnie ograniczenia liczebności zbioru B.

W problemie optymalizacji część mnożników Lagrange'a przyjmuje wartość zerową. Odpowiadające im wektory uczące  $\mathbf{x}$  nie tworzą więc wektorów podtrzymujących SV. Nie wnoszą one zatem żadnej informacji w procesie uczenia i nie są brane pod uwagę w dalszym procesie uczenia. Część mnożników przyjmuje wartości górnego ograniczenia C. Odpowiadające im wektory podtrzymujące nazywać się będą wektorami SV z ograniczeniami (BSV). Gdyby z góry znane były wektory BSV, odpowiadające im wartości mnożników Lagrange'a byłyby równe C. Oznacza to dalszą redukcję liczby optymalizowanych zmiennych (mnożników). Pozostałe mnożniki stanowią zmienne podlegające optymalizacji. Oznaczmy je wskaźnikiem x. Wskaźnikiem y oznaczmy mnożniki odpowiadające BSV, a przez z zbiór mnożników zerowych odpowiadających danym uczącym nie tworzącym wektorów podtrzymujących. Mnożniki  $\mathbf{a}$ , wektor  $\mathbf{d}$  oraz macierz  $\mathbf{Q}$  można wówczas przedstawić w postaci:

$$\mathbf{a} = \begin{bmatrix} \mathbf{a}_x \\ \mathbf{a}_y \\ \mathbf{a}_z \end{bmatrix} \quad \mathbf{d} = \begin{bmatrix} \mathbf{d}_x \\ \mathbf{d}_y \\ \mathbf{d}_z \end{bmatrix} \quad \mathbf{Q} = \begin{bmatrix} \mathbf{Q}_{xx} & \mathbf{Q}_{xy} & \mathbf{Q}_{xz} \\ \mathbf{Q}_{yx} & \mathbf{Q}_{yy} & \mathbf{Q}_{yz} \\ \mathbf{Q}_{zx} & \mathbf{Q}_{zy} & \mathbf{Q}_{zz} \end{bmatrix} \quad (2.54)$$

Wszystkie składniki funkcji celu nie związane ze zbiorem x stanowią wartości stałe, nie podlegające minimalizacji. Problem optymalizacyjny może być wówczas zredukowany do:

$$\min \left\{ -\mathbf{a}_x^T (\mathbf{e} - \mathbf{Q}_{xx} \mathbf{e} C) + \frac{1}{2} \mathbf{a}_x^T \mathbf{Q}_{xx} \mathbf{a}_x \right\} \quad (2.55)$$

przy ograniczeniach:

$$\begin{aligned} \mathbf{a}_y^T \mathbf{d}_x + C \cdot \mathbf{e}^T \mathbf{d}_y &= 0 \\ 0 \leq \mathbf{a}_x \leq C \cdot \mathbf{e} \end{aligned} \quad (2.56)$$

Selekcja zbioru roboczego powinna zapewnić największy postęp w kierunku minimum funkcji celu. Oznaczmy przez  $\mathbf{p}$  kierunek największego spadku. Jeśli ten wektor ma q niezerowych



elementów, to te elementy tworzą aktualny zbiór roboczy  $B$ , oznaczony symbolem  $x$ , uczestniczący w definicji (2.55) zredukowanej funkcji celu. Istotnym problemem w procesie optymalizacji jest jak najwcześniejsze określenie, które mnożniki zdążają do swoich ograniczeń (dolnych bądź górnych). Pozwoli je wówczas wyeliminować i z góry ograniczyć liczbę zmiennych optymalizowanych.

T. Joachims skutecznie zaimplementował algorytm heurystyczny estymacji mnożników, pozwalający na szybką redukcję liczebności zmiennych optymalizowanych i znaczne przyspieszenie procesu uczenia przy dużej liczbie danych uczących. W pobliżu rozwiązania mnożnik Lagrange'a pozostający w ograniczeniu dolnym bądź górnym wskazuje jak dana zmienna przeciwdziała ograniczeniom, które reprezentuje mnożnik. Dodatnia (różna od zera) wartość mnożnika wskazuje, że zmienna przyjmuje wartość optymalną przy spełnieniu danego ograniczenia. W punkcie dalekim od optymalności używa się zwykle estymaty mnożnika. Dobrą estymatą  $\lambda^{eq}$  jest wyrażenie [24]:

$$\lambda^{eq} = \frac{1}{N_\alpha} \sum_{i=1}^{N_\alpha} \left( d_i - \sum_{j=1}^p \alpha_j d_j K(\mathbf{x}_i, \mathbf{x}_j) \right) \quad (2.57)$$

gdzie  $N_\alpha$  jest aktualnym rozmiarem wektora  $\alpha$ . Mnożniki  $\lambda_i^{lo}$  i  $\lambda_i^{up}$  mogą być estymowane wg następujących wzorów

$$\lambda_i^{lo} = d_i \left[ \left( \sum_{j=1}^p \alpha_j d_j K(\mathbf{x}_i, \mathbf{x}_j) + \lambda^{eq} \right) - 1 \right] \quad (2.58)$$

$$\lambda_i^{up} = -d_i \left[ \left( \sum_{j=1}^p \alpha_j d_j K(\mathbf{x}_i, \mathbf{x}_j) + \lambda^{eq} \right) + 1 \right] \quad (2.59)$$

Analizując zmiany wartości estymowanych mnożników  $\lambda_i^{lo}$  oraz  $\lambda_i^{up}$  w ostatnich  $k$  iteracjach można wnioskować o ich dalszych zmianach. Jeśli np. wartości te były ciągle dodatnie, można z dużą dozą prawdopodobieństwa przyjąć, że będą takie również w punkcie optymalnym. W takim przypadku można je z góry pominąć w optymalizacji i w ten sposób zredukować wymiar problemu optymalizacyjnego. Po wyselekcjonowaniu wektora zmiennych bazowych optymalizowanych rozwiązuje się w sposób standardowy problem optymalizacyjny (2.55) z ograniczeniami (2.56). Końcowe, pełne sprawdzenie wszystkich warunków optymalności dotyczy wyłącznie punktu ostatecznego rozwiązania. W przypadku ich niespełnienia, należy powtórzyć obliczenia, kontrolując w szczególny sposób te pary danych, dla których warunki optymalizacji zostały naruszone.

### 2.3.3. Algorytm programowania sekwencyjnego Platta

Algorytm programowania sekwencyjnego (Sequential Minimal Optimization) polega na dekompozycji problemu programowania kwadratowego na mniejsze podzadania, rozwiązywane sekwencyjnie aż do spełnienia wszystkich warunków optymalności Kuhna-Tuckera. Spośród wielu znanych podejść, takich jak "chunking" Vapnika [58,59], algorytmu Osuny i ich różnych odmian, za najlepszy uznaje się algorytm SMO Platta [48], w którym rozwiązywane podzadanie jest drugiego rzędu (dwa mnożniki Lagrange'a), przy sekwencyjnej wymianie par mnożników. Optymalizacja dwu mnożników na raz jest najmniejszym możliwym zadaniem do rozwiązania. Dzięki tej redukcji rozwiązanie problemu optymalizacyjnego dokonywane jest analitycznie, nie wymaga dużej pamięci a zbieżność algorytmu do rozwiązania jest zapewniona przez specjalny algorytm doboru par mnożników.

W zadaniu optymalizacji kwadratowej funkcji 2 zmiennych z ograniczeniami kostkowymi  $0 \leq \alpha_i \leq C$  oraz równościowymi typu  $\sum_{i=1}^2 \alpha_i d_i = const$ , najpierw sprawdzane są warunki ograniczeń kostkowych i równościowych, a następnie określone minimum funkcji przy spełnieniu tych warunków. Stąd rozwiązanie w każdym etapie znajduje się w obszarze dopuszczalnym, określonym przez ograniczenia. Warunki ograniczeń kostkowych wymuszają ich położenie wewnątrz kwadratu o boku równym  $C$ . Jednocześnie zmienne optymalnie dobrane muszą leżeć na prostej określonej równaniem  $\alpha_1 d_1 + \alpha_2 d_2 = const$ . Przy wartościach  $d_i = \pm 1$  ( $i=1,2$ ) równanie prostej odpowiada więc diagonalnej lub przeciwdiagonalnej tego kwadratu. Przy dochodzeniu do optymalnych wartości obu mnożników Lagrange'a dobiera się najpierw mnożnik  $\alpha_2$ , a następnie  $\alpha_1$ . Wartość  $\alpha_2$  dobiera się w dwu krokach [48]

- krok pierwszy:

$$\alpha_2 = \alpha_2 + \frac{d_2(E_1 - E_2)}{\eta} \quad (2.60)$$

gdzie:

$$\eta = K(\mathbf{x}_1, \mathbf{x}_1) + K(\mathbf{x}_2, \mathbf{x}_2) - 2K(\mathbf{x}_1, \mathbf{x}_2) \quad (2.61)$$

$$E_1 = y(\mathbf{x}_1) - d_1 \quad (2.62)$$

$$E_2 = y(\mathbf{x}_2) - d_2 \quad (2.63)$$

natomiast  $d_1$  i  $d_2$  oznaczają wartości zadane funkcji aproksymujących, przy czym:

$$\begin{aligned}
y(\mathbf{x}_1) &= \sum_{j=1}^2 d_j \alpha_j K(\mathbf{x}_j, \mathbf{x}_1) + b \\
y(\mathbf{x}_2) &= \sum_{j=1}^2 d_j \alpha_j K(\mathbf{x}_j, \mathbf{x}_2) + b
\end{aligned}
\tag{2.64}$$

odpowiadają mnożnikom  $\alpha_1$  i  $\alpha_2$  oraz parom uczącym  $(\mathbf{x}_1, d_1)$ ,  $(\mathbf{x}_2, d_2)$ , wybranym aktualnie w danym kroku optymalizacji.

- krok drugi

W kroku drugim następuje korekta  $\alpha_2$  w taki sposób, aby spełnić ograniczenia nakładane na mnożniki [47]:

$$\alpha_{2opt} = \begin{cases} H & \text{jesli } \alpha_2 \geq H \\ \alpha_2 & \text{jesli } L < \alpha_2 < H \\ L & \text{jesli } \alpha_2 < L \end{cases}
\tag{2.65}$$

Wartości L i H oznaczają ograniczenia odpowiednio dolne i górne nakładane na zmienne optymalizowane  $\alpha_1$  i  $\alpha_2$ . Wartości te są określone analitycznie i są równe [48]:

$$L = \begin{cases} \max \{0, \alpha_2 - \alpha_1\} & \text{jesli } d_1 \neq d_2 \\ \max \{0, \alpha_2 + \alpha_1 - C\} & \text{jesli } d_1 = d_2 \end{cases}
\tag{2.66}$$

$$H = \begin{cases} \min \{C, C + \alpha_2 - \alpha_1\} & \text{jesli } d_1 \neq d_2 \\ \min \{C, \alpha_2 + \alpha_1\} & \text{jesli } d_1 = d_2 \end{cases}
\tag{2.67}$$

Po wyznaczeniu wartości optymalnej  $\alpha_{2opt}$ , mnożnik  $\alpha_{1opt}$  określa się z prostej zależności liniowej:

$$\alpha_{1opt} = \alpha_1 + d_1 d_2 (\alpha_2 - \alpha_{2opt})
\tag{2.68}$$

Powyższe zależności (2.65) i (2.68) wyznaczają w sposób analityczny wartości optymalne mnożników Lagrange'a poddanych adaptacji na danym etapie. W algorytmie Platt'a po doborze wartości optymalnych dwu wybranych mnożników Lagrange'a następuje przejście do następnej pary. Zostało udowodnione przez Osunę [46], że taki sposób postępowania gwarantuje zbieżność do rozwiązania. Istotnym problemem w algorytmie jest sprawdzanie warunków optymalności rozwiązania. Wystąpić tu mogą trzy przypadki:

- $\alpha_i=0$

Przypadek taki odpowiada spełnieniu warunków ograniczeń z nadmiarem, tzn.  $d_1(\mathbf{w}^T \Phi(\mathbf{x}_i) + b) - 1 > 0$ . Oznacza to, że wektor  $\mathbf{x}_i$  odpowiadający  $\alpha_i=0$  nie tworzy wektora podtrzymującego i nie ma wpływu na przebieg hiperpłaszczyzny separacyjnej.

- $0 < \alpha_i < C$

Wobec niezerowej wartości  $\alpha_i$  zachodzi spełnienie warunku  $d_i(\mathbf{w}^T \Phi(\mathbf{x}_i) + b) - 1 = 0$ . Wektor  $\mathbf{x}_i$  odpowiadający  $\alpha_i$  tworzy więc wektor podtrzymujący.

- $\alpha_i = C$

Przypadek ten odpowiada spełnieniu warunku  $d_i(\mathbf{w}^T \Phi(\mathbf{x}_i) + b) + \xi_i - 1 = 0$  przy  $\xi_i > 0$ . Oznacza to, że przyjęcie przez  $\alpha_i$  wartości krańcowej ograniczenia zmniejsza margines separacji o wartość  $\xi_i$ . Przy  $\xi_i > 1$  jest to błąd klasyfikacji.

Oznaczmy przez  $R_i$  następującą wielkość:

$$R_i = d_i(\mathbf{w}^T \mathbf{x}_i + b) - d_i^2 = d_i(\mathbf{w}^T \mathbf{x}_i + b - d_i) = d_i E_i$$

gdzie  $E_i = \mathbf{w}^T \mathbf{x}_i + b - d_i$  jest predykcją błędu klasyfikacji. Przy spełnieniu warunków optymalności Kuhna-Tuckera spełnione są następujące przyporządkowania:

$$\alpha_i = 0 \Rightarrow R_i > 0$$

$$0 < \alpha_i < C \Rightarrow R_i = 0$$

$$\alpha_i = C \Rightarrow R_i < 0$$

Warunki Kuhna-Tuckera nie są spełnione w dwóch przypadkach:

$$1) \alpha_i < C \cap R_i < 0$$

$$2) \alpha_i > 0 \cap R_i > 0$$

Efektywność działania algorytmu SMO uzyskuje się przy optymalnym doborze kolejnej pary mnożników Lagrange'a. Platt zaproponował specjalną, heurystyczną procedurę wyboru dwu par danych i odpowiadających im mnożników Lagrange'a, gwarantującą redukcję wartości funkcji celu w każdym kroku optymalizacji.

Algorytm Platta zakłada, że optymalizację  $\alpha_1$  i  $\alpha_2$  przeprowadza się na podstawie aktualnych wartości mnożników Lagrange'a, spełniających warunki ograniczeń (przy inicjalizacji można założyć  $\alpha = \mathbf{0}$ ). Przy poszukiwaniu właściwych na danym etapie uczenia mnożników  $\alpha_1$  i  $\alpha_2$  spośród wszystkich mnożników Lagrange'a Platt zaproponował dwie pętle przeszukujące dane: pętlę zewnętrzną i wewnętrzną.

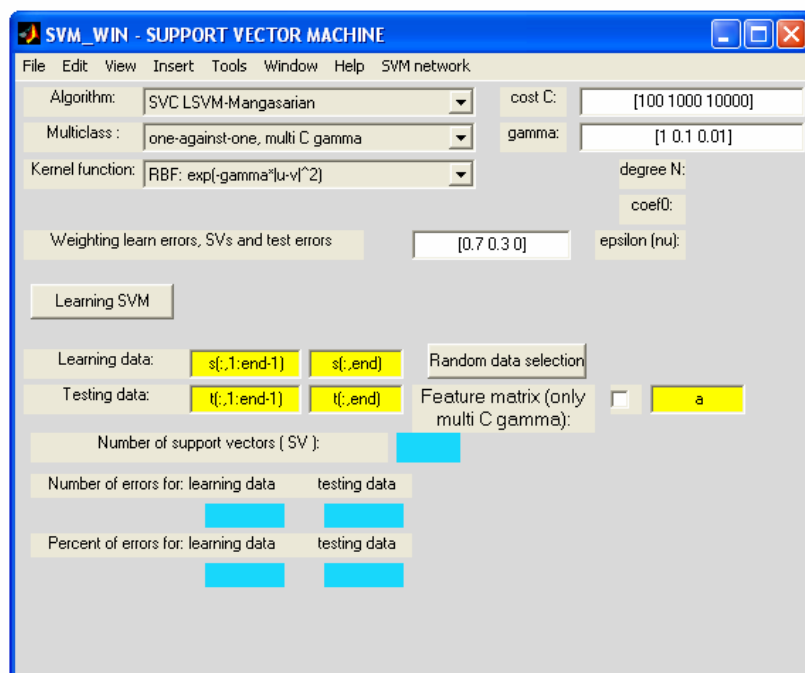
Pętla zewnętrzna selekcjonuje pierwszy mnożnik  $\alpha_1$  a pętla wewnętrzna drugi mnożnik  $\alpha_2$  w taki sposób, że następuje proces maksymalizacji funkcji celu w problemie dualnym. Pętla zewnętrzna przebiega przez wszystkie pary danych uczących selekcjonując te, dla których  $0 < \alpha_i < C$ . Następnie sprawdza spełnienie warunków Kuhna-Tuckera dla wyselekcjonowanych

przypadków i przyjmuje pierwszy dla którego te warunki są naruszone. Mnożnik Lagrange'a odpowiadający tak wyselekcjonowanemu przypadkowi staje się równy  $\alpha_1$ .

Po wyselekcjonowaniu mnożników  $\alpha_1$  pętla wewnętrzna przeszukuje wszystkie przypadki danych dla których  $0 < \alpha_i < C$  poszukując takiego, dla którego  $|E_2 - E_1|$  przyjmuje wartość maksymalną. Na wyselekcjonowanych 2 przypadkach następuje optymalizacja obu mnożników  $\alpha_1$  i  $\alpha_2$  zgodnie ze wzorami (2.65) i (2.68). Jeśli taka optymalizacja nie prowadzi do satysfakcjonującego uzysku funkcji celu, następuje ponowny, tym razem losowy wybór pierwszego mnożnika  $\alpha_1$  spośród przypadków odpowiadających  $0 < \alpha_i < C$  i powtórzenie procedury wyboru  $\alpha_2$ , a następnie optymalizacja. Jeśli to również zawiedzie, następuje powtórzenie losowego wyboru  $\alpha_1$ , ale spośród wszystkich danych uczących, połączone z dalszą procedurą wyboru  $\alpha_2$  i optymalizacją wartości funkcji celu.

## 2.4 Program SVM\_WIN

Dla celów pracy stworzono implementację wybranych algorytmów w środowisku Matlaba.



Rys. 2.4 Widok graficznego interfejsu użytkownika programu SVM\_WIN

Komunikacja z programem dokonuje się poprzez graficzny interfejs użytkownika pokazany na rys. 2.4. Program został wyposażony w następujące opcje:

- wybór algorytmu rozwiązywania problemu programowania kwadratowego (*Algorithm*)
- wybór rodzaju funkcji jądra (*Kernel function*)
- wybór strategii wyłaniania zwycięskiej klasy ("jeden-przeciw-wszystkim", "jeden-przeciw-jednemu", drzewo decyzyjne - DAG) (*Multiclass*)
- zadawanie wag przy zastosowaniu kryterium doboru optymalnych hiper-parametrów ( $C, \gamma$ ) sieci (*Weighting learn errors, SVs and test errors*)
- możliwość doboru parametrów sieci niezależnie dla każdej pary klas (opcja *multi C gamma* w polu *Multiclass*)
- możliwość wybierania wyselekcjonowanych dla każdej pary klas zestawu cech poprzez macierz  $\mathbf{a}$  (*Feature matrix*). Kolejne kolumny macierzy odpowiadają kolejnym parom klas i zawierają indeksy używanych cech.

## 2.5 Analiza porównawcza algorytmów

Do analizy porównawczej algorytmów wybrano typowe zadania benchmarkowe (10 zadań) znajdujące się na stronie internetowej [64]. Są to problemy klasyfikacji binarnej. Porównano algorytmy: LSVM, SVM<sup>Light</sup> i SMO Platt. Porównanie dotyczyło czasu działania, liczby wektorów podtrzymujących oraz liczby błędów uczących i testujących. Wyniki porównania zgromadzono w tabeli 2.1.

Tabela 2.1 Porównanie działania algorytmów na problemach testowych

Nazwa danych	Liczba danych uczących/ testujących	Rozmiar wektora $x$	LSVM				SVM <sup>Light</sup>				SMO Platt			
			$N_{sv}$	czas [s]	błąd ucz. [%]	błąd test. [%]	$N_{sv}$	czas [s]	błąd ucz. [%]	błąd test. [%]	$N_{sv}$	czas [s]	błąd ucz. [%]	błąd test. [%]
banana	400 / 4900	2	178	0.375	6.75	10.51	90	0.375	7.00	10.41	91	0.047	7.00	10.43
breast cancer	200 / 77	9	186	0.078	15.50	22.08	113	0.156	16.50	24.68	113	0.016	16.50	24.68
diabetis	468 / 300	8	426	0.578	19.02	22.67	242	0.313	20.09	22.33	243	0.062	20.30	22.33
flare solar	666 / 400	9	579	1.438	30.48	32.75	464	1.203	31.83	32.00	460	0.359	31.83	32.00
german	700 / 300	20	611	1.969	0	24.33	517	0.359	0	24.33	518	0.266	0	24.33
heart	170 / 100	13	153	0.062	0	18.00	70	0.016	11.76	17.00	70	0.015	11.76	17.00
image	1300 / 1010	18	1294	50.51	0	3.76	202	2.172	0.15	3.07	201	0.36	0.15	3.07
ringnorms	400 / 7000	20	217	0.453	0	6.52	207	0.031	0	1.57	203	0.063	0	1.59
thyroid	140 / 75	5	78	0.031	1.43	4.00	29	0.063	2.14	4.00	31	0.001	2.14	4.00
titanic	150 / 2051	3	139	0.032	17.33	23.31	91	0.203	19.33	22.62	91	0.015	19.33	22.62

Z otrzymanych rezultatów wynika, że ze względu na dokładność klasyfikacji algorytm LSVM Mangasariana jest zwykle nieco gorszy (tylko w problemie breast cancer okazał się lepszy),

pozostałe dawały praktycznie taki sam błąd. Ponadto czas działania LSVM znacznie rośnie wraz ze wzrostem liczby danych uczących. Również liczba wektorów podtrzymujących generowana przez LSVM jest znacznie wyższa. W przypadku algorytmów SVM<sup>Light</sup> i SMO Platta liczby wektorów podtrzymujących były prawie identyczne, ale ten drugi był średnio około pięć razy szybszy. Również w innych próbach przeprowadzonych przez autora (np. sztuczny nos [5,6]) algorytm Platta okazał się zdecydowanie najlepszy.

### 3. Analiza morfologiczna obrazu

Morfologia matematyczna jest dziedziną zajmującą się analizą struktur przestrzennych. W sensie matematycznym bazuje na geometrii i algebrze, tworząc technikę pełnej analizy obrazu poprzez złożenie wielu transformacji elementarnych.

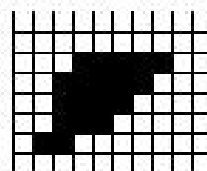
#### 3.1. Tworzenie postaci dyskretnej obrazu

Wiele operacji morfologicznych obrazu jest zdefiniowanych dla przestrzeni ciągłej. Jednakże na komputerze możliwa jest tylko dyskretna postać danych odpowiadających podziałowi obrazu na piksele. Dlatego pierwszym etapem jest transformacja danych z dwuwymiarowej przestrzeni Euklidesowej  $\mathbf{R}^2$  na przestrzeń dyskretną  $\mathbf{Z}^2$ . Proces ten, będący próbkowaniem obrazu ciągłego, definiuje przyjęta dyskretyzacja. Piksele przyjmują wartości zależne od lokalnej zmienności w obrazie i ich rozkład zależy od przyjętej sieci ich położenia. W praktyce najczęściej przyjmuje się sieć prostokątną ze względu na łatwość zapisu macierzowego. W takim przypadku wartość piksela jest definiowana przez zmienność obrazu w położonym względem niego centralnie polu prostokątnym o wymiarach zgodnych z odległościami w przyjętej sieci.

Reprezentacja dyskretna obrazu wprowadza pewien błąd, którego wielkość jest związana z rozdzielczością dyskretyzacji obszaru. Dla zbyt dużego rozmiaru jednostkowego obszaru, któremu odpowiada jeden piksel, mogą zniknąć niektóre cechy obrazu, jak to przedstawiono na rys. 3.1 W najprostszym przypadku transformacji obrazu binarnego przyjmuje się wartość piksela 1, jeżeli 50% odpowiadającego jemu obszaru lub więcej jest ciemna, stąd wąski obszar jasny z rys. 3.1a po dyskretyzacji zniknął (rys. 3.1b).



a)



b)

Rys. 3.1 Odwzorowanie dyskretne obrazu binarnego ilustrujące problem doboru rozdzielczości dyskretyzacji, a) – obraz wejściowy, b) – reprezentacja dyskretna



Drugą płaszczyzną uproszczeń związanych z transformacją do postaci dyskretnej jest reprezentacja kolorów. Pierwszym uproszczeniem jest przyjęcie np. średniej wartości nasycenia barwy na obszarze przyporządkowanym jednemu pikselowi. Błąd wynikający z tego uproszczenia zależy od rozdzielczości dyskretyzacji obszaru. Drugim uproszczeniem jest skończona liczba wartości reprezentującej natężenia barwy (np. w skali 8-bitowej szarości jest to liczba całkowita od 0 do 255). Błąd wprowadzony przez to uproszczenie nie zależy od rozdzielczości dyskretyzacji, a od liczby bitów użytych do kodowania stopnia szarości.

### 3.2. Progowanie i algorytm Otsu

Jedną z podstawowych transformacji obrazów reprezentowanych w skali szarości jest progowanie, oznaczane przez  $T_{[t_1, t_2]}$ , gdzie  $t_1$  jest progiem dolnym, a  $t_2$  – górnym. Operatory progów dla wszystkich pikseli obrazu wejściowego jednoznacznie określają transformację do postaci binarnej. Operacja ta może być przeprowadzona albo dla zakresu zdefiniowanego dwoma progami albo przy założeniu tylko jednego progów. W pierwszym przypadku dla punktu  $x$  obrazu  $f$  określa ją zależność:

$$T_{[t_1, t_2]}[f(x)] = \begin{cases} 1 & \text{if } t_1 \leq f(x) \leq t_2 \\ 0 & \text{else} \end{cases} \quad (3.1)$$

Przy założeniu jednej wartości progów przyjmuje się automatycznie  $t_2 = t_{\max}$ . Należy w tym miejscu zauważyć, że przy jednej wartości progów wynik operacji dla progów wyższego zawiera się w wyniku dla progów niższego, tj.:

$$T_{[t_{\max}, t_{\max}]}[f(x)] \subseteq T_{[t_{\max-1}, t_{\max}]}[f(x)] \subseteq \dots \subseteq T_{[t_0, t_{\max}]}[f(x)] \quad (3.2)$$

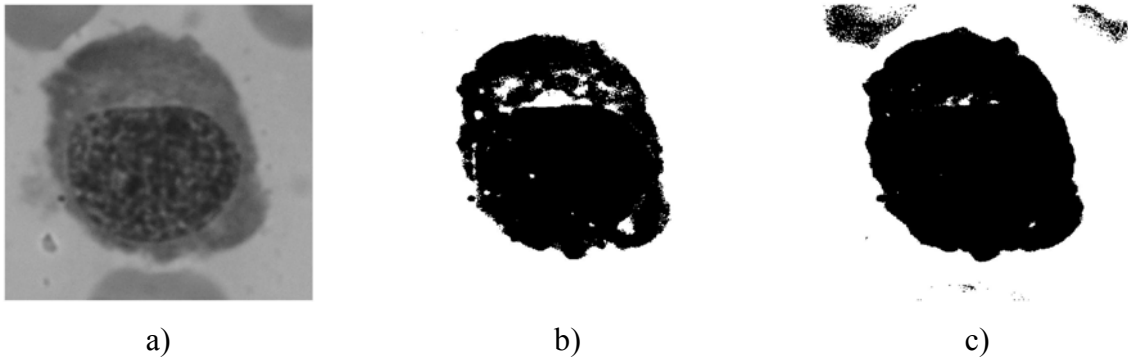
W roku 1979 Nobuyuki Otsu [47] zaproponował algorytm automatycznego wyznaczania progów na podstawie histogramu obrazu. Celem algorytmu jest wyznaczenie progów optymalnie separujących obiekt od tła obrazu. Dla każdej możliwej wartości progowania wyznaczone są wariancje klas  $\sigma_1^2, \sigma_2^2$  i wariancja międzyklasowa  $\sigma_B^2$ , gdzie za klasy 1 i 2 przyjmuje się odpowiednio obiekt i tło. Ich obliczenie przeprowadza się na podstawie znormalizowanego histogramu poprzez określenie prawdopodobieństwa przynależności piksela do klasy 1 i 2, oznaczonego odpowiednio jako  $\omega_1$  i  $\omega_2$ . Następnie oblicza się średnie poziomy szarości obydwu klas i całego obrazu, oznaczone jako  $\mu_1, \mu_2, \mu_T$ . Jeżeli  $p_i$  oznacza  $i$ -tą wartość histogramu, to wariancje dla wartości progów  $k$  są określone wzorami:

$$\sigma_1^2 = \sum_{i=0}^{k-1} \frac{(i - \mu_1)^2 p_i}{\omega_1} \quad (3.3)$$

$$\sigma_2^2 = \sum_{i=k}^{t_{\max}} \frac{(i - \mu_2)^2 p_i}{\omega_2}$$

$$\sigma_B^2 = \omega_1(\mu_1 - \mu_T)^2 + \omega_2(\mu_2 - \mu_T)^2 = \omega_1\omega_2(\mu_2 - \mu_1)^2 \quad (3.4)$$

Próg  $k$ , dla którego tak wyznaczona wariancja przyjmuje wartość maksymalną, jest progiem optymalnym. Na przykład dla obrazu komórki z rys. 3.2a wyznaczona w ten sposób wartość progu  $k$  wynosi 115. Obraz binarny komórki odpowiadający optymalnemu progowaniu wg Otsu przedstawia rys. 3.2c. Dla porównania na rys. 3.2b przedstawiono obraz binarny otrzymany przy zastosowaniu progu ( $t=100$ ) ustalonego arbitralnie przez użytkownika.



Rys. 3.2 Wynik operacji progowania  $T_{[t_1, t_{\max}]}$ : a) obraz wejściowy, b) wynik progowania dla wartości progu  $t=100$ , c) wynik optymalnego progowania metodą Otsu

### 3.3. Podstawowe operacje morfologiczne

#### 3.3.1. Erozja

Podstawową operacją morfologiczną jest erozja. Do jej zdefiniowania konieczny jest wybór elementu strukturującego (SE). Element taki musi mieć określony obszar i punkt bazowy  $x$ . Operacja erozji na zbiorze binarnym  $X$ , elementem strukturującym SE, jest zdefiniowana [54,63] jako zbiór punktów bazowych  $x$  takich elementów SE, które zawarte są całkowicie w  $X$ . Operację erozji elementem SE, oznaczaną dalej przez  $\varepsilon_{SE}$ , można przedstawić wzorem:

$$\varepsilon_{SE}(X) = \{x \mid SE_x \subseteq X\} \quad (3.5)$$

Rozszerzenie tej definicji na skalę szarości obrazu  $f$  można przedstawić jako minimum translacji obrazu  $f$  o wszystkie wektory  $-s$  zbioru wektorów zawartych w SE o początkach w

punkcie bazowym. Wartość erozji odpowiadająca danemu pikselowi  $\mathbf{x}$  obrazu  $f$  jest równa minimum funkcji:

$$\varepsilon_{SE}(\mathbf{x})|_f = \min_{\mathbf{s} \in SE} f(\mathbf{x} + \mathbf{s}) \quad (3.6)$$

Istnieje możliwość użycia bardziej złożonego elementu strukturującego (przestrzennego), któremu przypisuje się pewną wysokość. Przykładem może być dysk w skali szarości, którego trzeci wymiar określa wysokość części sfery rozpiętej nad nim z określoną w centrum maksymalną wysokością wybieraną przez użytkownika. Erozję takim elementem  $SE_v$  określa wzór:

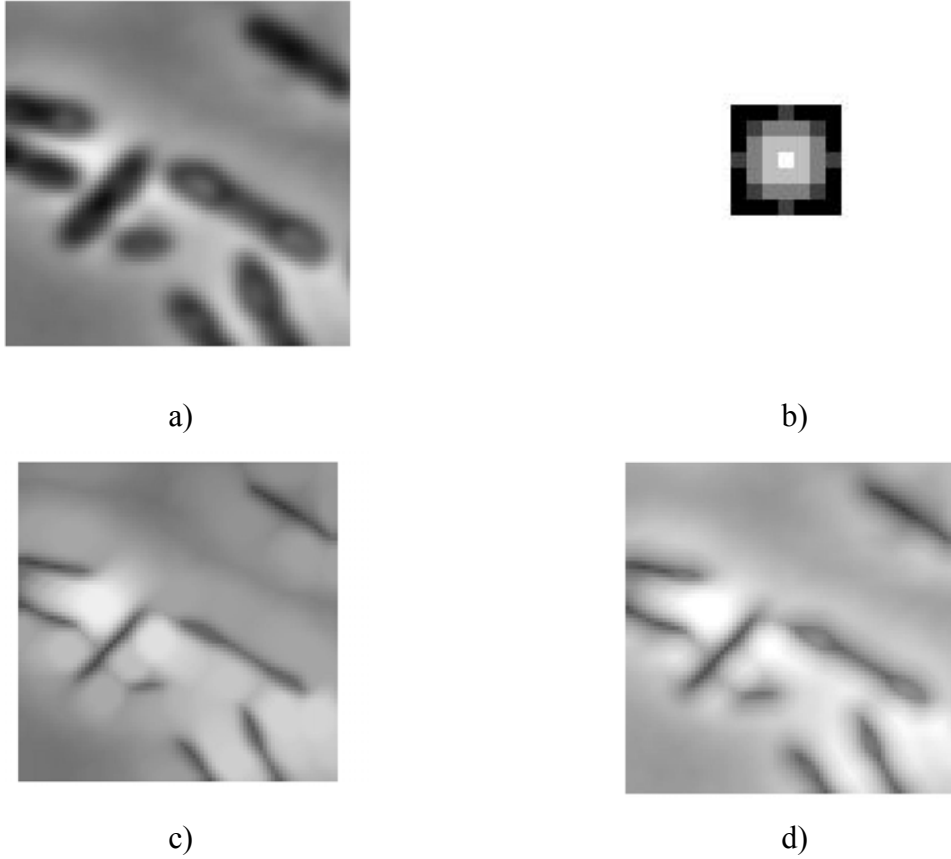
$$\varepsilon_{SE_v}(\mathbf{x})|_f = \min_{\mathbf{s} \in SE_v} \{f(\mathbf{x} + \mathbf{s}) - SE_v(\mathbf{s})\} \quad (3.7)$$



Rys. 3.3 Elementy strukturujące: a) płaski w kształcie dysku o promieniu 6 pikseli, b) wypukły w kształcie dysku o promieniu 6 i wysokości 250

Dla tak zdefiniowanej erozji piksel  $\mathbf{x}$  przyjmuje wartość najmniejszą z wartości obrazu poddanego translacji o każdy wektor  $\mathbf{s}$  należący do  $SE_v$  pomniejszoną o wartość  $SE_v$  dla tego wektora. Elementy strukturujące w kształcie dysku płaskiego bądź wypukłego przedstawione są na rys. 3.3.

Na rys. 3.4 przedstawiono wynik erozji obrazu (a) wykonanej przy pomocy elementu strukturującego w kształcie dysku prostokątnego o rozmiarze 7 pikseli przedstawionego na rys. 3.4b. Wynik przedstawiony na rys.3.4c odpowiada dyskowi płaskiemu, natomiast na rys. 3.4d – dyskowi wypukłemu o wysokości 40 pikseli. Erozja w obu wypadkach spowodowała zwężenie szczegółów tworzących krawędzie obrazu, ale widoczne są istotne różnice obrazów zerodowanych w zależności od rodzaju użytego elementu strukturującego.



Rys. 3.4 Obraz a) wejściowy, b) element strukturujący w kształcie dysku, c) obraz po erozji elementem płaskim w kształcie przedstawionym na rys. b), d) wynik erozji elementem wypukłym w kształcie dysku o wysokości 40 pikseli.

### 3.3.2. Dylatacja

Operacja dylatacji zbioru  $X$  elementem strukturującym  $SE$  jest zdefiniowana [54,63] jako zbiór takich położeń punktów  $\mathbf{x}$  odpowiadających punktom bazowym  $SE$ , że zbiór punktów wspólnych  $SE$  i  $X$  jest niepusty:

$$\delta_{SE}(X) = \{\mathbf{x} \mid SE_{\mathbf{x}} \cap X \neq \emptyset\} \quad (3.8)$$

Przy rozszerzeniu na skalę szarości można tę operację zapisać w postaci:

$$[\delta_{SE}(f)](x) = \max_{s \in SE} f(x + s) \quad (3.9)$$

Jeżeli przez  $C$  oznaczymy zbiór komplementarny do  $f$ , zdefiniowany jako transformacja:

$$C = f^c(x) = t_{\max} - f(x) \quad (3.10)$$

to zachodzą następujące relacje pomiędzy erozją i dylatacją [54,63]:

$$\varepsilon_{SE} = C \delta_{SE} C \quad (3.11)$$

$$\delta_{SE} = C \varepsilon_{SE} C$$

$$\delta_{SE} \varepsilon_{SE} \leq f \leq \varepsilon_{SE} \delta_{SE} \quad (3.12)$$

### 3.3.3. Otwarcie i zamknięcie

Zdefiniowane dotąd operacje mogą tworzyć nowe przekształcenia obrazu. Jednym z możliwych postępowań jest wykonanie na obrazie kolejno operacji erozji i dylatacji. Szczególny przypadek, który zachodzi dla erozji elementem SE i dylatacji jego odbiciem lustrzanym  $\underline{SE}$  względem środka  $x$ , nazywany jest otwarciem. Tak zdefiniowana operacja otwarcia  $\gamma$  jest równoważna sumie takich położeń elementów SE, dla których SE jest zawarte w zbiorze  $X$ :

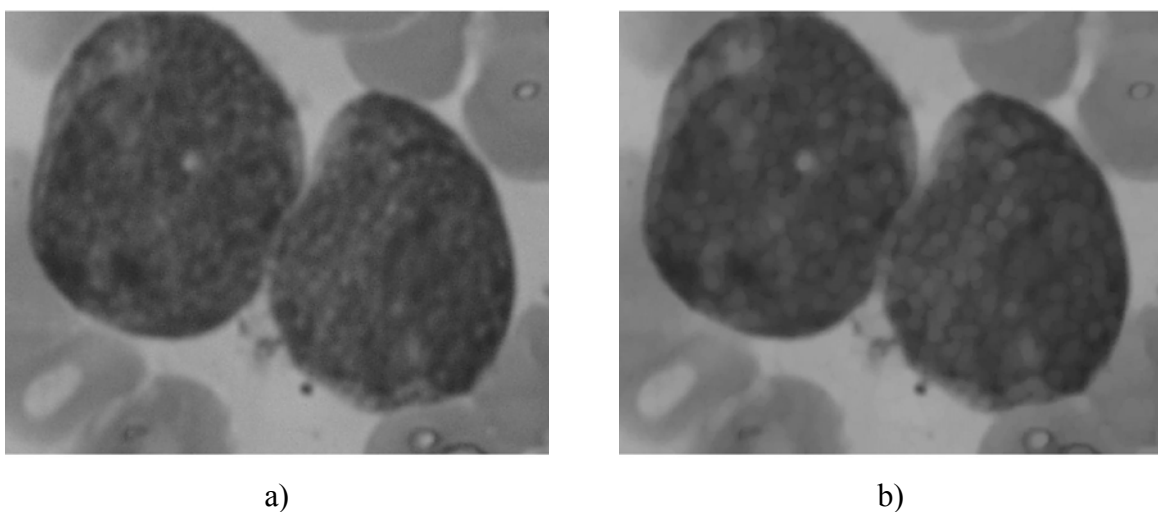
$$\gamma_{SE}(X) = \bigcup_x \{SE_x \mid SE_x \subseteq X\} \quad (3.13)$$

Efektom wykonania operacji otwarcia jest rozdzielenie elementów połączonych cienkimi liniami, wygładzanie konturów i usuwanie niewielkich elementów.

Operację przeprowadzoną w odwrotnej kolejności do otwarcia, czyli najpierw dylatację a następnie erozję, nazywamy zamknięciem. Jak łatwo zauważyć, zamknięcie zbioru  $X$  jest równoważne otwarciu zbioru komplementarnego do niego, oznaczonego jako  $X^C$ . Można więc zdefiniować operację zamknięcia  $\phi$  jako:

$$\phi_{SE}(X) = \left[ \bigcap_x \{SE_x \mid SE_x \subseteq X^C\} \right]^C \quad (3.14)$$

W wyniku operacji zamknięcia obszary leżące blisko siebie są łączone, mniejsze braki w obszarach i wklęsłości konturów zapełniane. Operację otwarcia obrazu  $f$  komórki (będącą jednocześnie operacją zamknięcia jego dopełnienia) przedstawia rys. 3.5.



Rys. 3.5 Ilustracja operacji otwarcia elementem w kształcie dysku o promieniu 2: a) obraz wejściowy, b) wynik operacji.

Ważną cechą operacji otwarcia i zamknięcia jest idempotentność - wykonanie ich kolejny raz przy użyciu takiego samego elementu strukturującego nie zmienia poprzedniego zbioru  $X$ .

### 3.4. Filtracja dolnoprzepustowa obrazu

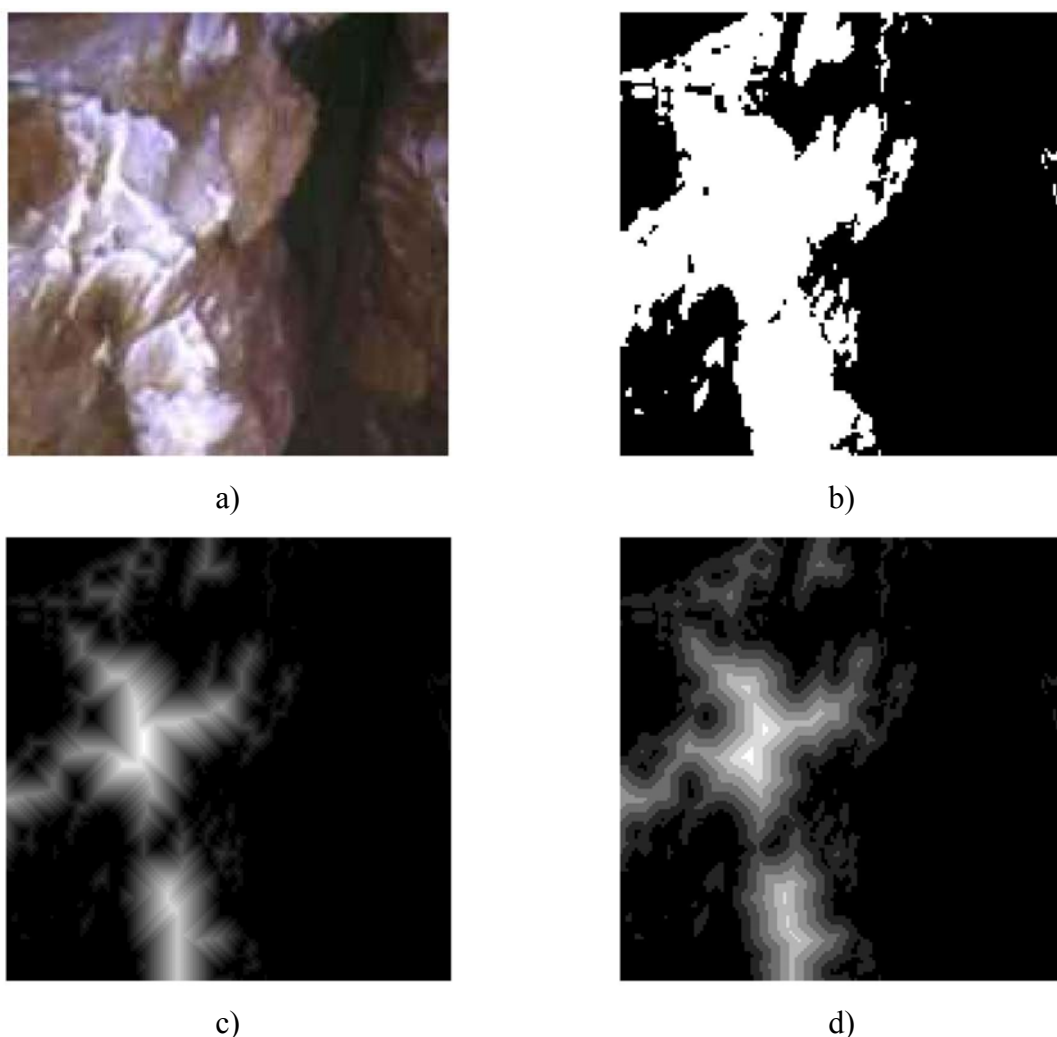
Jednym z elementów występujących przy przetwarzaniu obrazu wejściowego jest filtracja dolnoprzepustowa. W wyniku otrzymuje się obraz o zredukowanych wartościach drobnych zakłóceń, odpowiadających wysokim częstotliwościom. Filtr działa w sposób uśredniający poszczególne poziomy szarości pikseli. Efektem takiej filtracji poza eliminacją zakłóceń dla pojedynczych pikseli jest także rozmycie krawędzi. Dla przyjętej maski jasność wyjściową transformowanego piksela obrazu oblicza się jako sumę iloczynów współczynników maski i wartości jasności pikseli w obszarze równym wielkości maski. Dla redukcji zakłóceń w pojedynczych pikselach bez nadmiernego wpływu na rozmycie krawędzi dobrze nadaje się maska o kształcie funkcji Gaussa i małym rozmiarze np.: 3x3, 4x4 lub 5x5 pikseli. Drugą wielkością określającą filtr jest odchylenie standardowe charakteryzujące szybkość redukcji wartości współczynników filtru względem elementu centralnego maski. Przykład maski gaussowskiej  $\mathbf{W}_G$  o wymiarze 5x5 przedstawiono poniżej [32]

$$\mathbf{W}_G = \begin{bmatrix} 0.0369 & 0.0392 & 0.0400 & 0.0392 & 0.0369 \\ 0.0392 & 0.0416 & 0.0424 & 0.0416 & 0.0392 \\ 0.0400 & 0.0424 & 0.0433 & 0.0424 & 0.0400 \\ 0.0392 & 0.0416 & 0.0424 & 0.0416 & 0.0392 \\ 0.0369 & 0.0392 & 0.0400 & 0.0392 & 0.0369 \end{bmatrix}$$

### 3.5. Reprezentacja odległościowa obrazu binarnego

Kluczową operacją dla procesu segmentacji obrazu jest utworzenie reprezentacji odległościowej dla obrazu binarnego otrzymanego poprzez progowanie dzielące obszar interesujących obiektów od tła. Polega to na wygenerowaniu macierzy odległości poszczególnych pikseli segmentowanego obszaru od tła. Dla odległości, określanej jako najmniejsza liczba pikseli oddzielających dany piksel od tła, duży wpływ mają nawet niewielkie nierówności brzegów obszarów. Dla takiej reprezentacji odległościowej segmentacja zwracałaby zbyt dużą liczbę wydzielonych regionów, których źródłem mogą być niewielkie zakłócenia lub dyskretyzacja. Aby tego uniknąć należy wygenerować macierz odległości przy zastosowaniu elementu strukturującego SE. Rozmiar takiego elementu określa kolejne strefy odległości, które można otrzymać np. poprzez wykonanie erozji takim elementem. Piksele usuwane w kolejnej iteracji wskutek erozji otrzymują wartość równą tej

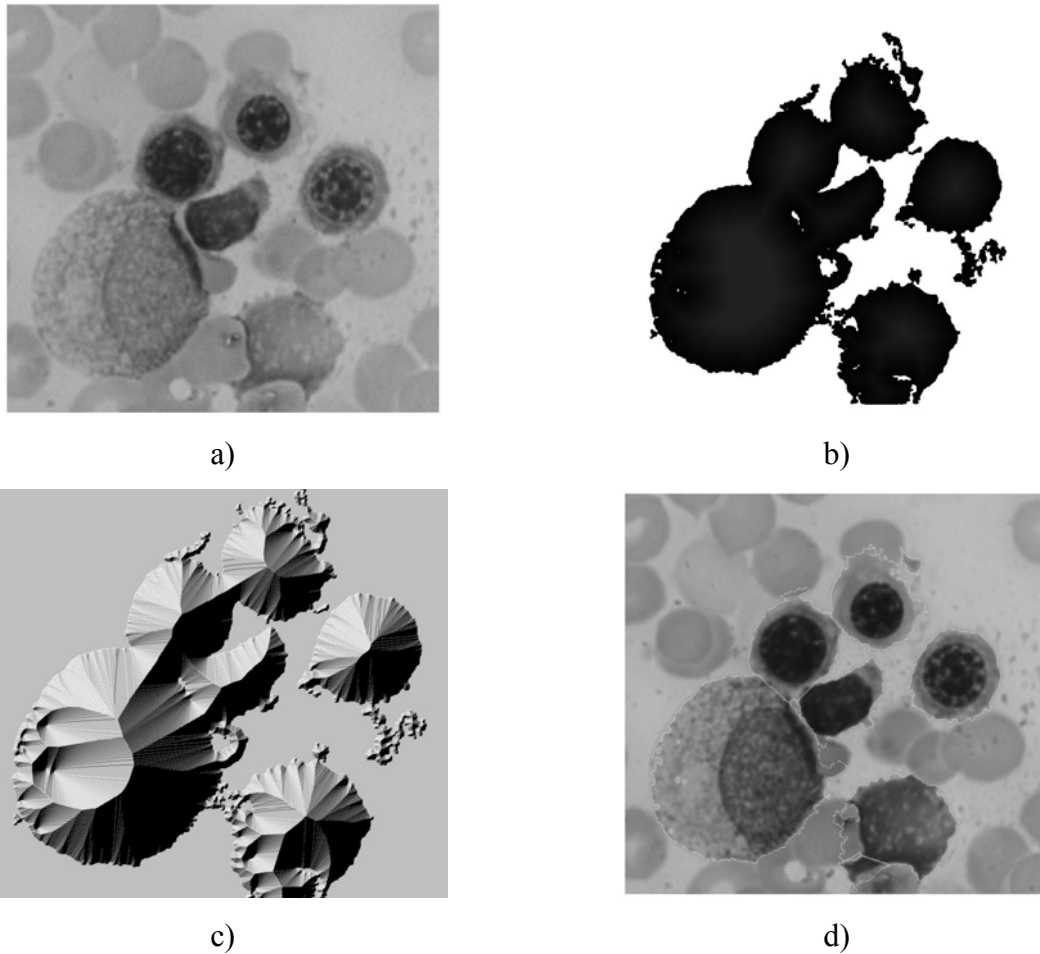
iteracji. Taka metoda ma cechy operacji otwarcia, czyli wygładza brzegi poszczególnych stref odległości, ale zachowuje nawet niewielkie elementy tła wewnątrz segmentowanego obszaru, istotne przy segmentacji, szczególnie elementów licznie stykających się ze sobą. Na rys. 3.6 przedstawione są przykłady zastosowania różnych reprezentacji odległościowych obrazu. W reprezentacji przy zastosowaniu funkcji ciągłej (3.6c) widoczne są wąskie ciemne linie na jasnym tle. Odpowiadają one niższym wartościom odległości i powstają z wklęsłości konturu obszaru złożonych nawet z jednego piksela, a otrzymana reprezentacja charakteryzuje się dużą zmiennością wyników dla poszczególnych położzeń pikseli. Zastosowanie erozji elementem SE w kształcie np. rombu o promieniu 3-ch pikseli generuje wynik bardziej stabilny o mniejszej liczbie różnych poziomów szarości.



Rys. 3.6 Przykłady reprezentacji odległościowej obrazu: a) obraz oryginalny, b) przetransformowany do postaci binarnej, c) mapa odległościowa obrazu przy zastosowaniu ciągłej funkcji odległości, d) mapa odległościowa obrazu przy wykorzystaniu erozji elementem w kształcie rombu i promieniu 3-ch pikseli.

### 3.6. Segmentacja obrazu metodą działów wodnych

Danymi wejściowymi dla segmentacji obrazu jest macierz reprezentacji odległościowej. Polecaną metodą segmentacji jest metoda działów wodnych [54]. Intuicyjnie można ją zinterpretować w postaci opadu deszczu na pasmo górskie, odwzorowane poprzez macierz odległości. Poszczególnymi regionami podlegającymi segmentacji są strefy zbierania się wody, a liniami granic obiektów - granie górskie.



Rys 3.7 Przykład segmentacji obrazu w skali szarości (a) poprzez postać binarną podobszarów (b), reprezentację odległościową do metody działów wodnych (c), obraz z liniami podziału (d).

Dla zformalizowania opisu metody należy wprowadzić pewne oznaczenia. Niech  $h_{\min}$  i  $h_{\max}$  będą odpowiednio najmniejszą i największą wartością macierzy odległościowej danych. Niech  $M_i$  oznacza region i-tego minimum lokalnego a  $CB(M_i)$  - region zbierania się wody odpowiadający danemu minimum. Dla każdego  $h$  z przedziału  $[h_{\min}, h_{\max}]$   $CB_h(M_i)$  oznacza podzbiór punktów  $CB(M_i)$  o wartościach mniejszych lub równych  $h$ , natomiast  $X_h$  podzbiór



wszystkich regionów zbierania się wody, które zawierają punkty o wartościach odległości mniejszych lub równych  $h$ . Należy teraz stopniowo przeprowadzić proces "zatapiania", zaczynając od wartości  $h=h_{\min}$ .

Podczas przejścia z poziomu  $h$  na  $h+1$  można wyróżnić dla regionów  $Y=\{CB_{h+1}(M_i)\}$  trzy przypadki:

1. Region  $Y$  jest rozłączny z regionami o poziomie  $h$ . Jest on zatem minimum lokalnym i tworzy nowy region zbierania wody  $CB_{h+1}(M_i)$ .
2. Region  $Y$  ma część wspólną z jednym regionem  $CB_h(M_i)$ . W takim przypadku  $CB_{h+1}(M_i):=Y$ .
3. Region  $Y$  ma część wspólną z kilkoma regionami  $CB_h(M_i)$ . Jest on dzielony liniami równoodległymi od konturów  $CB_h(M_i)$  i odpowiednie strefy są przypisywane do odpowiadających im  $CB_h(M_i)$ .

W wyniku takiego postępowania otrzymuje się podział zupełny obrazu na strefy przynależne do określonego minimum lokalnego. Zaletą metody jest to, że podział obrazu jest niezależny od położenia środków ciężkości poszczególnych obiektów, a jedynie od zmian odległości zawartych w wejściowej macierzy odległościowej charakteryzującej obszar. Przykład segmentacji ilustruje rys. 3.7. Wynikiem segmentacji obrazu metodą działów wodnych jest wydzielenie tylu obszarów ile było minimów lokalnych. Może zatem zajść potrzeba usunięcia najmniejszych z nich. Jedną z metod umożliwiających taką operację jest transformacja geodezyjna i rekonstrukcja.

### 3.7. Transformacja geodezyjna i rekonstrukcja

Transformacjami geodezyjnymi nazywa się osobny rodzaj operacji morfologicznych wykonywanych na obrazie wejściowym  $f$  z użyciem maski  $g$  pełniącej rolę ograniczeń. Istotą operacji geodezyjnych są transformacje morfologiczne obrazu, spełniające jednocześnie ograniczenia nakładane przez maskę. Przykładem może być dylatacja geodezyjna  $\delta_g$  obrazu  $f$  zdefiniowana w taki sposób, że spełniony jest warunek:

$$\delta_g(f) = \delta(f) \wedge g \quad (3.15)$$

Wynikiem takiej operacji jest obraz monotoniczny względem znaczników zawartych w obrazie. Oznacza to, że wszystkie pary przyległych pikseli zrekonstruowanego obrazu muszą spełniać taką samą relację, jaką spełniają ich odpowiedniki zawarte w masce. Analogicznie otrzymuje się erozję geodezyjną  $\varepsilon_g$  zdefiniowaną jako:

$$\varepsilon_g(f) = \varepsilon(f) \vee g \quad (3.16)$$

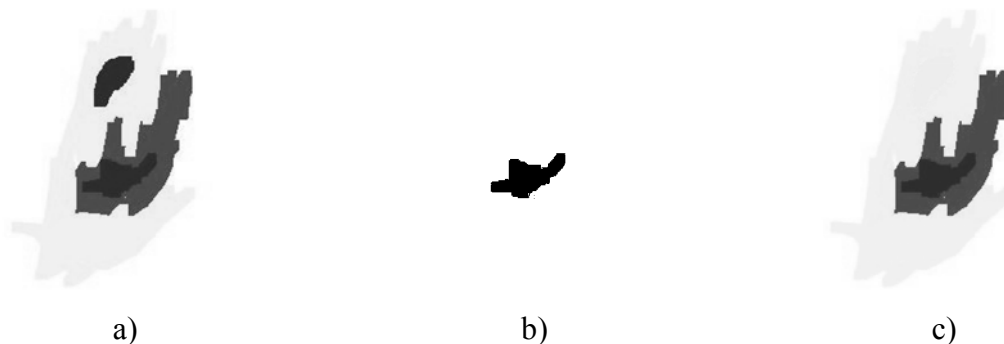
Zgodnie z tym wzorem wynik erozji geodezyjnej  $f$  nie musi zawierać się w  $f$ . Na podstawie tych dwóch operacji (3.15) i (3.16) można zdefiniować transformację geodezyjną  $v_g$  obrazu  $f$  w punkcie  $\mathbf{x}$  określoną wzorem:

$$[v_g(f)](\mathbf{x}) = \begin{cases} [\delta_g(f)](\mathbf{x}) & \text{if } f(\mathbf{x}) \leq g(\mathbf{x}) \\ [\varepsilon_g(f)](\mathbf{x}) & \text{else} \end{cases} \quad (3.17)$$

Omówione operacje geodezyjne mogą być wykorzystane do rekonstrukcji obrazu. Jeżeli w obrazie  $f$  są zawarte centra poszczególnych elementów obrazu, a w masce  $g$  obraz wejściowy, to za pomocą np. dylatacji geodezyjnej obrazu  $f$  z maską  $g$  można odtworzyć zaznaczone obiekty. Cechą tej metody jest monotoniczność odtworzonych obiektów względem centrów zawartych w  $f$ .

Przykład rekonstrukcji obrazu ilustruje rys. 3.8. Wynikiem jej może być obraz nie zawierający mniejszego ciemniejszego obszaru jak to przedstawiono na rys. 3.8c ponieważ odtwarzając kształt obiektu operacja rekonstrukcji usunęła wszystkie lokalne maksima pozostawiając tylko jedno w punkcie początkowym zawartym w  $f$ .

Zauważmy, że rekonstrukcja geodezyjna może być przeprowadzona w taki sposób, że zachowane zostaną jedynie obszary zwarte o kształcie zbliżonym do kołowego, podczas gdy zwykle kryterium obszaru zachowałoby wszystkie obszary o odpowiedniej wielkości niezależnie od ich kształtu, np. obszary długie i wąskie.



Rys. 3.8 Rekonstrukcja obrazu: a) maska  $g$ , b) obraz  $f$  spełniająca rolę znacznika, c) wynik rekonstrukcji.

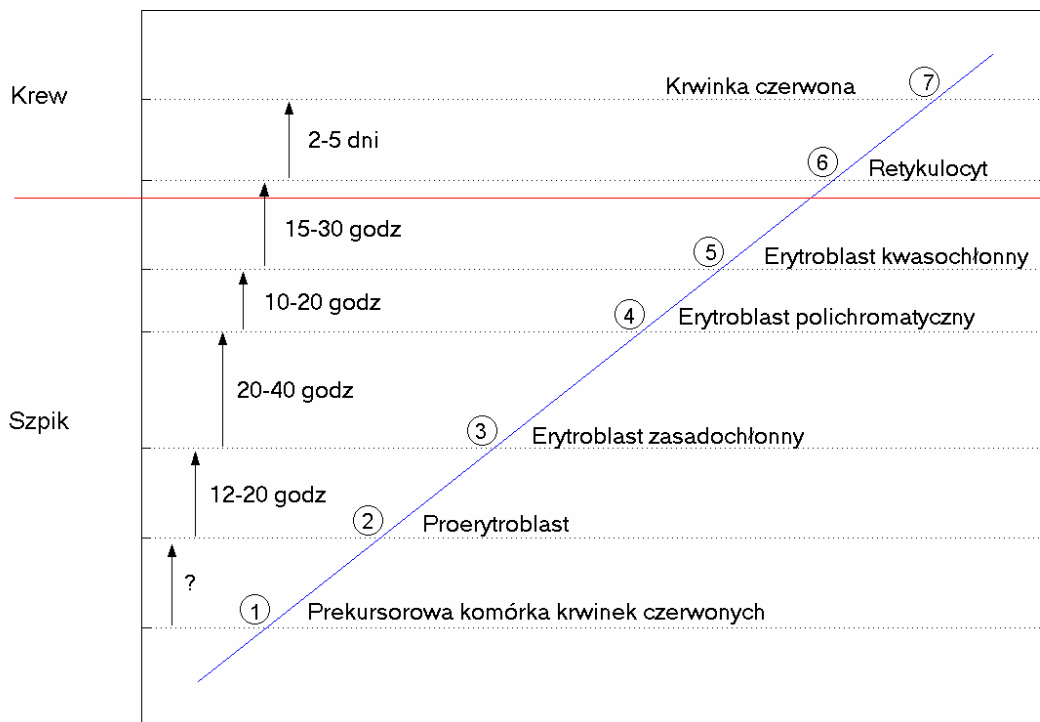
## **4. Ekstrakcja obrazu komórek rakowych**

Podstawowymi elementami obrazu rozmazu szpiku kostnego podlegającymi analizie są komórki jądrzaste należące do różnych linii rozwojowych i będące na różnym stadium rozwoju. Ich skład procentowy pełni istotną rolę w diagnostyce medycznej chorych m.in. z podejrzeniem białaczki szpikowej [3,20,23,26]. Dla zautomatyzowania procesu określania składu procentowego tych komórek w rozmazie konieczne jest wydzielenie każdej z nich z obrazu mikroskopowego (tzw. ekstrakcja), a następnie przypisanie do konkretnej grupy (typu komórki).

### **4.1. Charakterystyka ogólna komórek krwiotwórczych**

W układzie krwiotwórczym człowieka istnieje wiele typów komórek odpowiedzialnych za prawidłowe funkcjonowanie organizmu. Część z nich jest rozróżnialna w obrazach rozmazu przy pomocy specjalnych barwień, inne nie są rozróżnialne. Aktualnie w ramach mielogramu, będącego procentowym składem szpiku, oznacza się komórki występujące w poszczególnych fazach rozwoju w ramach trzech linii rozwojowych. Pierwszą tworzą komórki układu czerwonokrwinkowego, rozwijające się w ramach erytropoezy normoblastycznej i erytropoezy megaloblastycznej. Drugą grupę tworzy układ białokrwinkowy, a trzecią utkanie chłonne. Jeżeli w rozmazie pojawiają się komórki, które nie występują normalnie u zdrowego człowieka, są one dodatkowo wymieniane. Wszystkie wymienione komórki podlegające liczeniu na etapie diagnozy są komórkami jądrzastymi. Ponadto w szpiku występują krwinki czerwone, nie posiadające barwiącego się jądra, płytki krwi, pojedyncze komórki łatwo rozróżnialne oraz cienie komórkowe, będące pozostałościami po ich rozpadzie. W ramach poszczególnych linii rozwojowych każdemu etapowi odpowiadają określone typy komórek, przypisane im. Ponieważ cykl rozwojowy komórki jest procesem ciągłym, trudno jest określić jednoznacznie granicę pomiędzy kolejnymi stadiami rozwoju w ramach danej linii i jednoznacznie sklasyfikować daną komórkę. Różnice w wynikach klasyfikacji mogą wystąpić zarówno w wyniku oceny obrazu przez kilku laborantów, jak tego samego laboranta w wyniku zmęczenia, czy sugestii poprzednimi obrazami lub komórkami położonymi obok.

Schemat cytogenezy układu erytrocytowego (czerwonokrwinkowego) obrazujący poszczególne typy komórek i czas ich dojrzewania przedstawiono na rys. 4.1 [23].

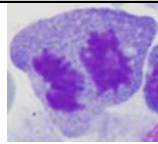
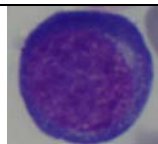
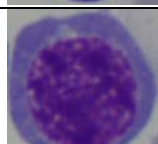
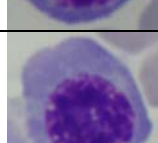
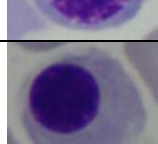
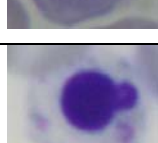


Rys. 4.1 Schemat cytogenezy układu erytrocytowego (dane wg [23])

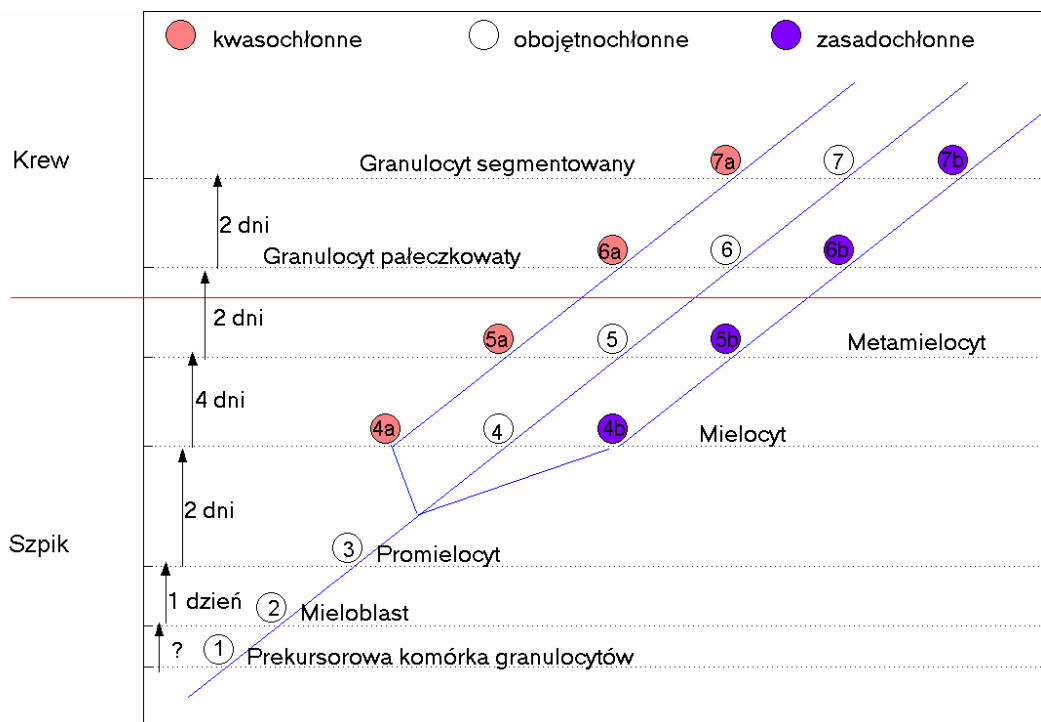
W diagnostyce hematologicznej prekursorowa komórka krwinek czerwonych nie jest rozróżnialna i nie podlega liczeniu. Najwcześniejszym rozróżnialnym stadium rozwoju jest proerytroblast przechodzący kolejne fazy od erytroblasta zasadochłonnego, polichromatycznego, aż do kwasochłonnego (ortochromatycznego). W postaci retykulocytu i dojrzałej krwinki czerwonej może występować zarówno we krwi jak i szpiku, ale nie zawiera już jądra, które uległo ekspulsji, i w związku z tym nie podlega klasyfikacji. Tabela 4.1 przedstawia przykładowe obrazy komórek tego układu i ich cechy charakterystyczne ułatwiające identyfikację. W zestawieniu komórek przedstawionych w tabeli 4.1 ujęto również postać podziałową erytropoezy, będącą komórką w trakcie podziału oraz paraerytroblast występujący w stanach chorobowych.

W układzie czerwonokrwinkowym rozróżnia się również erytropoezę megaloblastyczną, która może prawidłowo występować w okresie ciąży i u noworodków, natomiast w innych przypadkach jest symptomem chorobowym. W jej ramach rozróżnia się kolejne fazy rozwoju: promegaloblasty, megaloblasty zasadochłonne, polichromatyczne i ortochromatyczne. Podstawową cechą różnicującą te komórki w stosunku do komórek erytropoezy normoblastycznej jest ich wielkość - są one około dwukrotnie większe od swoich odpowiedników erytropoezy normoblastycznej.

Tabela 4.1 Zestawienie komórek układu erytrocytowego i erytropoezy normoblastycznej oraz ich cech (dane wg. [23,26])

Komórka	Obraz	Rozmiar	Kształt komórki	Kształt jądra	Chromatyna	Jąderka	Strefa przejścia	Cytoplazma	Ziarnistość
Postacie podziałowe				dwa jądra nieregularne	fioletowawa	brak	brak przejaśnienia	niebieska	brak
Proerytroblast		20-25 $\mu\text{m}$	okrągły, rzadko owalny	okrągłe	ciemna, grubogrudkowa	czasem ciemniejsze, nie zawsze widoczne	przejaśnienie przyjądrowe	ciemnoniebieska, niewiele	brak
Erytroblast zasadochłonny		13-18 $\mu\text{m}$	okrągły, zniekształcony przez sąsiadów	okrągłe	zbita, promieniste pasma	brak	brak przejaśnienia	ciemnoniebieska, więcej niż w proerytroblast	brak
Erytroblast polichromatyczny		10-15 $\mu\text{m}$	okrągły, zniekształcony przez sąsiadów	okrągłe	ciemna, wyraźna kondensacja	brak	brak przejaśnienia	sina, więcej niż erytroblast zasadochłonny	brak
Erytroblast ortochromatyczny		8-12 $\mu\text{m}$	okrągły, zniekształcony przez sąsiadów	okrągłe	bardzo ciemna, prawie jednolita	brak	brak przejaśnienia	różowa, jak w erytrocytach, ponad połowa	brak
Paraerytroblast		8-12 $\mu\text{m}$	okrągły	okrągłe, czasem z wypustką	jednolita, niebieska	brak	brak przejaśnienia	jasnoniebieska	drobna, przy brzegach komórki

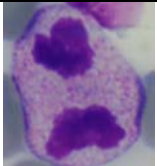
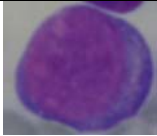
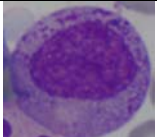
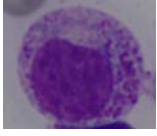
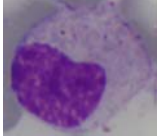

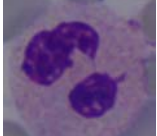
W układzie białokrwinkowym występuje większe zróżnicowanie linii dojrzewania komórek, niż w innych układach. Po pierwszych fazach rozwoju w postaci nierozróżnialnej komórki prekursorowej, rozróżnialnych fazach mieloblastu i promielocytu, w kolejnych stadiach rozwoju komórki zalicza się do trzech odrębnych linii komórkowych: neutrofilii, eozynofili i bazofili. Ponieważ znacznie mniej wiadomo na temat kinetyki wytwarzania, różnicowania, krążenia i migracji komórek linii eozynofili i bazofili, trudno jednoznacznie powiedzieć, czy powstają one przy różnicowaniu promielocytu, czy też mają osobne komórki prekursorowe. Równoległe do nich z osobnej komórki prekursorowej rozwijają się monoblasty i promonocyty, w standardowym barwieniu nierozróżnialne od mieloblastów, przechodzące w ostateczną formę monocytu, który już jest rozróżnialny od innych komórek. Rys 4.2 przedstawia poszczególne linie rozwoju komórek układu granulocytowego (tzn. białokrwinkowego).

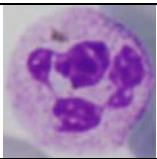
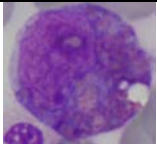
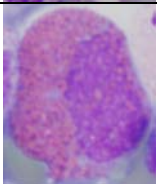
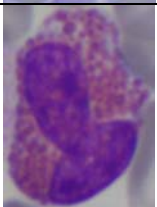
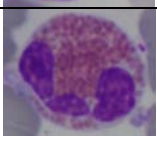
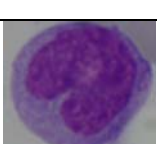


Rys. 4.2 Schemat cytoogenezy układu granulocytowego (granulopoeza) (dane wg [23])

Na rys. 4.2 symbolami 4,5,6,7 oznaczone są komórki obojętnochłonne, czyli neutrofile, symbolami 4a,5a,6a,7a ich odpowiedniki linii kwasochłonnej, czyli eozynofile, oraz symbolami 4b,5b,6b,7b zasadochłonne, czyli bazofile. W tabeli 4.2 zestawiono przykładowe obrazy komórek układu granulocytowego i cechy charakterystyczne ułatwiające ich identyfikację.

Tabela 4.2 Zestawienie wybranych komórek układu granulocytowego i ich cech (dane wg. [23,26])

Komórka	Obraz	Rozmiar	Kształt komórki	Kształt jądra	Chromatyna	Jąderka	Strefa przejścia	Cytoplazma	Ziarnistość
Postacie podziałowe									
Mieloblast		15-25 $\mu\text{m}$	owalny, niekiedy okrągły	owalny, nieco nieregularny rzadko okrągły	delikatny, drobnogrudkowy	1-4 średnie, jaśniejsze	wydłużone przejaśnienie, nie zawsze	niebieska, ciemniejsza na brzegach, niewiele	brak lub azurochłonna dość gruba
Promielocyt		10-30 $\mu\text{m}$	owalny lub okrągły	owalny, położone niecentralnie	początki kondensacji	1-2 średnie lub duże, jaśniejsze, często niewidoczne	wyraźne przejaśnienie	jasna niebieska, strefa przejaśnienia, więcej do dużo	obfita lub bardzo obfita, azurochłonna gruba
Mielocyt		15-25 $\mu\text{m}$	owalny lub okrągły	owalne lub nieco nerkowate	częściowo skondensowana	niewidoczne	brak	bladoniebieska lub różowawa, połowa powierzchni komórki, zanika	obfita, gruba azurochłonna lub obojętnochłonna
Metamielocyt		14-20 $\mu\text{m}$	owalny lub okrągły	wydłużony, podkowiały leży blisko krawędzi komórki	skondensowana	brak	brak	różowa	nieliczna azurochłonna, obojętnochłonna nieregularna
Granulocyt pałeczkowaty			owalny lub okrągły	w kształcie pałeczki, wewnątrz komórki	skondensowana	brak	brak	jasnoróżowa	
Granulocyt segmentowaty			owalny lub okrągły	2-4 segmentowe zwężone z min. 1/3 długości	skondensowana	brak	brak	jasnoróżowa	

Granulocyt segmentowany wielopłatowy			owalny lub okrągły	co najmniej 5 segmentowe	skondensowana	brak	brak	jasnoróżowa	
Mielocyt kwasochłonny			owalny lub okrągły	owalne lub nieco nerkowate	częściowo skondensowana	brak	brak	niewidoczna pod ziarnistością	liczna średnia azurochłonna
Metamielocyt kwasochłonny			owalny lub okrągły	wydłużony, podkowiasty leży blisko krawędzi komórki	skondensowana	brak	brak	niewidoczna pod ziarnistością	liczna średnia azurochłonna
Granulocyt pałeczkowaty kwasochłonny			owalny lub okrągły	w kształcie pałeczki, wewnątrz komórki	skondensowana	brak	brak	niewidoczna pod ziarnistością	liczna gruba azurochłonna
Granulocyt segmentowany kwasochłonny			owalny lub okrągły	2-4 segmentowe zwężone z min. 1/3 długości	skondensowana	brak	brak	niewidoczna pod ziarnistością	liczna gruba azurochłonna
Monocyt		15-25 μm	okrągły, owalny lub nieregularny	nieregularny	grubogrudkowa lub marmurkowa	niewidoczne	brak	niebieskosina, jaśniejsza, więcej	niewiele azurochłonna bardzo drobna



W szpiku kostnym zdrowego człowieka i w większości chorób hematologicznych niektóre komórki nie występują lub występują bardzo rzadko. Dotyczy to m.in. linii bazofili, czyli części zasadochłonnej układu granulocytowego. Ponadto w stosowanym barwieniu nierozróżnialne od linii mieloidalnej są monoblasty i promonocyty. Typowe obrazy takich komórek można znaleźć w pracy [21,27].

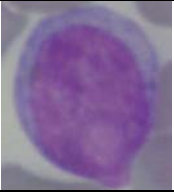
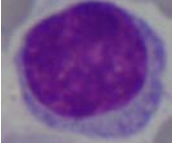
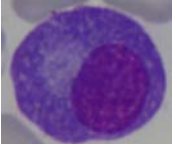
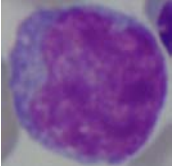
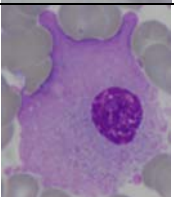
Trzecią grupą komórek jest tzw. układ siateczkowo-śródbłonkowy, określane również jako utkanie chłonne. W ramach tego układu występują nierozróżnialne limfoblasty, rozróżnialne prolimfocyty, limfocyty, plazmoblasty, proplazmoblasty, plazmocyty i paraplazmocyty. Z tej grupy najczęściej oznaczeniu podlegają prolimfocyty, limfocyty i plazmocyty. Ponadto podobne do limfocytów są komórki chłoniaka, których liczba (w przypadku występowania) też jest oznaczana. Dodatkowo określa się jakość szpiku (czy jest bogatokomórkowy czy ubogokomórkowy, ilość tkanki tłuszczowej), obecność megakariocytów, płytkowatość oraz grudkowatość. Najczęściej występujące komórki utkanie chłonne i pozostałe, otrzymane w preparatach, przedstawia tabela 4.3.

Dla sporządzenia mielogramu opisującego skład procentowy szpiku kostnego należy zliczyć od 200 do 400 komórek aby wynik procentowy udziału poszczególnych typów był obiektywny. Czasami w przypadku szpiku o wzmożonej komórkowości w obrazach jego rozmazu komórki są ułożone tak gęsto, że policzenie ich jest bardzo trudne. Występuje wówczas potrzeba zatrudnienia na raz dwóch laborantów, z których jeden nazywa widziane pod mikroskopem komórki, a drugi je sumuje. Czasem dla weryfikacji wyniku konieczne jest oznaczenie składu przez zespoły pracujące niezależnie. Często wykonuje się np. 2 rozmazy i skład procentowy szacuje się częściowo z jednego, częściowo z drugiego preparatu. Końcowy wynik jest zawsze procentowy i dodatkowo zawiera udział wymienionych wcześniej głównych układów w szpiku kostnym. Wyniki zliczania komórek przez 2 niezależnych ekspertów mogą różnić się nawet o 15% i taką dokładność uznaje się za zadawalającą.

Należy zauważyć, iż pomimo pozornie identycznego zabarwienia widzianego pod mikroskopem, na wykonanych zdjęciach widać czasami wyraźną różnicę w barwie komórek, czerwonych krwinek i tła oraz jasności obrazu. Powodem tego może być kilka czynników:

- inne ustawienia mikroskopu
- inne ustawienia aparatu (np. temperatura barwowa dobrana automatycznie lub ustalona ręcznie)
- wpływ czynników z otoczenia (np. światło dzienne, z oświetlenia)
- inna partia odczynników użyta w procesie przygotowania

Tabela 4.3 Komórki utkane chłonne oraz pozostałe i ich cechy (dane wg. [23,26])

Komórka	Obraz	Rozmiar	Kształt komórki	Kształt jądra	Chromatyna	Jąderka	Strefa przejścia	Cytoplazma	Ziarnistość
Prolimfocyt		12-15 $\mu\text{m}$	okrągły, niekiedy owalny	okrągły	homogenna	małe lub średnie, jaśniejsze, 1-2 sztuki	brak	niebieska, zwykle ciemna, niewiele	brak
Limfocyt		10-15 $\mu\text{m}$	okrągły, niekiedy owalny	okrągły, lekko owalny	homogenna, skondensowana	czasem słabo widoczne małe jądro	brak	niebieska, niewiele	prawidłowo brak
Plazmocyty		15-20 $\mu\text{m}$	owalna	okrągły	skondensowana	niewidoczne	przejaśnienie przyjądrowe	granatowa, niekiedy pojedyncze wodniczki	brak
Chłoniak									
Megakariocyt		do 100 $\mu\text{m}$	owalny, niekiedy okrągły	wielopłatowe, nieregularne	skondensowana	niewidoczne	brak	różowa, dużo	drobna różowa

- inny odczyn wody destylowanej użytej do przyrządzenia odczynników
- różny czas barwienia.

Niektóre czynniki wpływające na parametry obrazu można wyeliminować (np. ustawiania mikroskopu czy wpływ czynników otoczenia), inne są nie do wyeliminowania. Dlatego też w następnym rozdziale omówiono również wpływ ustawień aparatu (temperatury barwowej) jako jedyne go czynnika możliwego do regulacji w systemie pomiarowym.

Wszystkie zdjęcia wykorzystane w niniejszej pracy wykonano na tym samym urządzeniu, w większości przy takich samych ustawieniach, z powiększeniem 1000×, przy rozdzielczości 1712×1368 pikseli, zapisane w formacie RGB. Wszystkie obrazy poddane segmentacji były wykonane przy takich samych ustawieniach mikroskopu i aparatu (przy maksymalnym naświetleniu, czułości ASO 100, automatycznie dobieranej temperaturze barwowej itp.) w celu eliminacji wpływu regulacji manualnej. Rozmazy były barwione przy użyciu standardowego barwienia panoptycznego metodą Maya-Grunwalda-Giemsy (MGG).

Zdjęcia te nie podlegały dodatkowemu przetwarzaniu przy użyciu jakiegokolwiek programu graficznego. Wszystkie zarejestrowane obrazy komórek były klasyfikowane przez autora i weryfikowane przez jednego, czasem dwóch, pracowników laboratorium Instytutu Hematologii w Warszawie. Dzięki temu zminimalizowano różnice subiektywne wynikające z indywidualnych ocen różnych laborantów.

Jak już wspomniano, wyniki zliczania poszczególnych komórek ujmowane są w formie mielogramów. Prawidłowy mielogram zdrowego człowieka wg danych z Instytutu Hematologii objęty jest normą przedstawioną w tabeli 4.4. Należy podkreślić, że normy przyjmowane na świecie niekoniecznie pokrywają się z normami polskimi. Przykładowo w tabeli 4.5 przedstawiono zakresy wartości procentowego udziału poszczególnych komórek, uznane za prawidłowe wg prof. A. V. Hoffbrandta i J. Pettita [20].

Porównując obie tabele można zauważyć, że niektóre komórki są oceniane na świecie łącznie. Jest to związane zarówno z trudnością ich rozróżnienia, jak i mniejszym znaczeniem medycznym podziału pomiędzy te typy. Praktycznie poza podziałem erytropoezy na 4 rodzaje komórek, większość pozostałych rodzajów komórek wybranych do badań pokrywa się z normą światową podaną przez Hoffbrandta.

Do dalszych badań w pracy wybrano tylko najczęściej występujące typy komórek. Jest to związane z koniecznością zapewnienia wystarczająco dużej liczby przykładów każdej rozważanej klasy reprezentującej określony typ komórki.

Tabela. 4.4 Norma dla szpiku dorosłych wg Instytutu Hematologii

Udział procentowy komórek różnych typów w rozmazie szpiku kostnego	
Komórki	Obserwowany zakres [%]
<b>A) Układ czerwonokrwinkowy</b>	10 - 30
Postacie podzielowe	0 – 0.3
Proerytroblasty	0 – 1.5
Erytroblasty zasadochłonne	0.5 – 5
Erytroblasty polichromatyczne	5 – 15
Erytroblasty ortochromatyczne	5 – 15
<b>B) Układ białokrwinkowy</b>	62 – 77
Postacie podzielowe	0 – 0.2
Mieloblasty	0.5 – 3
Promielocyty	0.5 – 5
Mielocyty obojętnochłonne	5 – 18
Mielocyty kwasochłonne	0.5 – 2.5
Mielocyty zasadochłonne	0 - .03
Metamielocyty obojętnochłonne	8 – 25
Metamielocyty kwasochłonne	0 – 2
Metamielocyty zasadochłonne	0 – 0.1
Pałeczki obojętnochłonne	10 – 30
Pałeczki kwasochłonne	0.4 – 1
Pałeczki zasadochłonne	0 – 0.1
Segmenty obojętnochłonne	11 – 30
Segmenty kwasochłonne	0.4 – 3
Segmenty zasadochłonne	0 – 0.5
<b>C) Utkanie chłonne</b>	
Limfocyty	3 – 12
Plazmocyty	0.2 - 2

Tabela 4.5 Norma dla szpiku dorosłych wg Hoffbrandt, Pettit [20].

Udział procentowy komórek różnych typów w rozmazie szpiku kostnego			
Komórki		Obserwowany zakres	95% zakresu (średnia)
Blasty		0.0 – 3.2	0.0 – 3.0 (1.4)
Promielocyty		3.6 – 13.2	3.2 – 12.4 (7.8)
Mielocyt obojętnochłonny		4.0 – 21.4	3.7 – 10.0 (7.6)
Mielocyt kwasochłonny		0.0 – 5.0	0.0 – 2.8 (1.3)
Metamielocyt		1.0 – 7.0	2.3 – 5.9 (4.1)
Pałeczka i segment	mężczyzna	21.0 – 45.6 *	21.9 – 42.9 (32.1)
	kobieta	29.6 – 46.6 *	28.8 – 45.9 (37.4)
Eozynofile (kwasochłonne)		0.4 – 4.2	0.3 – 4.2 (2.2)
Eozynofile i mielocyt kwasochłonny		0.9 – 7.4	0.7 – 6.3 (3.5)
Bazofile (zasadochłonne)		0.0 – 0.8	0.0 – 0.4 (0.1)
Erytroblasty	mężczyzna	18.0 – 39.4 *	16.2 – 40.1 (28.1)
	kobieta	14.0 – 31.8 *	13.0 – 32.0 (22.5)
Limfocyty		4.6 – 22.6	6.0 – 20.0 (13.1)
Plazmocyty		0.0 – 1.4	0.0 – 1.2 (0.6)
Monocyty		0.0 – 3.2	0.0 – 2.6 (1.3)
Makrofagi		0.0 – 1.8	0.0 – 1.3 (0.4)
Stosunek granulopoezy do erytropoezy	mężczyzna	1.1 – 4.0 +	1.1 – 4.1 (2.1)
	kobieta	1.6 – 5.4 +	1.6 – 5.2 (2.8)
Oznaczona różnica pomiędzy mężczyzną i kobietą: *P<0.001; +P<0.01.			

#### 4.2. Algorytm ekstrakcji komórek

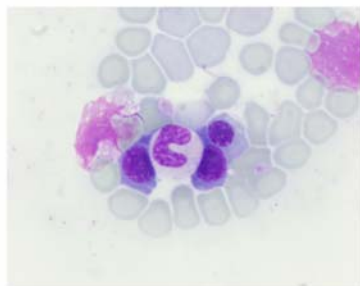
Podstawowym etapem automatycznej klasyfikacji komórek jest ich prawidłowa ekstrakcja z obrazu całego rozmazu. Istotne jest, aby użyty algorytm mógł być stosowany bez udziału parametrów ustawianych manualnie dla różnych zdjęć czy preparatów. Pierwszym krokiem jest skalowanie obrazu odpowiednio do przyjętych wartości referencyjnych. Istnieje kilka możliwości takiego skalowania:

- rozciąganie histogramu
- przesunięcie histogramu

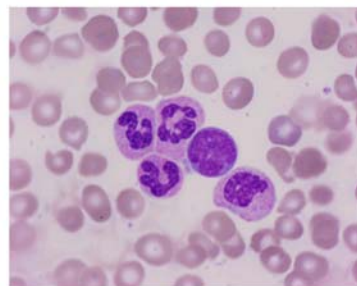
- przeskalowanie histogramu za pomocą jednej wartości referencyjnej
- przeskalowanie histogramu za pomocą dwóch wartości referencyjnych.

Ze względu na to, że w danym obrazie mogą występować komórki o różnym nasyceniu barw, zrezygnowano z pierwszej możliwości. Również przesunięcie histogramu wydaje się nieodpowiednie, gdyż może spowodować znaczne rozjaśnienie wydzielanych komórek i mieć istotny wpływ na ich cechy. Dlatego też rozważono tylko dwie ostatnie możliwości. W przypadku skalowania obrazu za pomocą dwóch wartości referencyjnych wybiera się jako jedną z nich tło (osobna kwestia czy to będzie najjaśniejszy punkt tła, czy np. wartość średnia – wtedy powstaje problem określenia co jest tłem, a co już nie), a jako drugą - barwę czerwonych krwinek, która u różnych osobników jest zbliżona do siebie. W czasie badań okazało się, że takie skalowanie może tworzyć obrazy o bardzo różnym zabarwieniu (należy pamiętać, że operacje są wykonywane na każdej składowej RGB równocześnie i niezależnie), jak pokazano na rys. 4.3 c i d. Dlatego zdecydowano się na stosowanie skalowania przy pomocy jednej wartości referencyjnej. Wybrano najjaśniejszy punkt tła, który transformowano do barwy białej. Najciemniejszy punkt obrazu pozostawał niezmienny, natomiast pozostałe były transformowane liniowo. Zaletą tego podejścia jest jednoznaczność wartości minimalnej i maksymalnej jasności obrazu. Jak już wspomniano w przypadku przyjęcia wartości średniej tła pozostaje problem uprzedniego wyznaczenia jego granic. Jest to szczególnie trudne jeżeli chodzi o rozdzielenie stref tłuszczowych i krwinek czerwonych. Podejście to wiąże się z dwoma problemami: niektóre obrazy zawierają jasne plamki będące efektem nierównomiernego barwienia (pęcherzyki powietrza, zatłuszczenia) oraz ciemne plamki będące zanieczyszczeniami. Tego typu czynniki mogą mieć niekorzystny wpływ na proces skalowania, jednakże prawie zawsze algorytm segmentacji jest na nie odporny. Rysunek 4.3 przedstawia przykłady skalowania obrazów oryginalnych (rys. 4.3a,b) przy użyciu dwóch wartości referencyjnych (rys. 4.3c,d) oraz przy użyciu jednej wartości referencyjnej (rys. 4.3e,f). Skalowanie wg 2 wartości referencyjnych może całkowicie zniekształcić obraz (rys. 4.3d), natomiast przyjęcie jednej wartości referencyjnej poprawia jakość obrazu, a jednocześnie daje lepsze zrównoważenie składowych barw, szczególnie w zakresie tła i krwinek czerwonych. Po wykonaniu skalowania obraz poddawany jest dalszemu przetworzeniu i segmentacji dla oddzielenia poszczególnych komórek, dla których będą generowane cechy diagnostyczne.

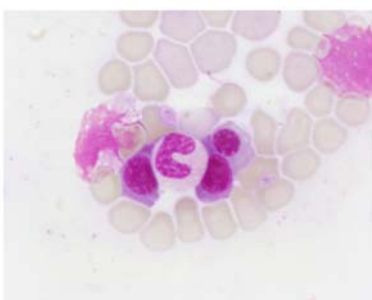
Aby segmentacja wносиła jak najmniejszy błąd, obraz cyfrowy musi być wcześniej odpowiednio przygotowany za pomocą operacji morfologicznych. Do takich operacji zalicza się m.in. erozję, dylatację, otwarcie i zamknięcie. Stosuje się również inne złożone operacje



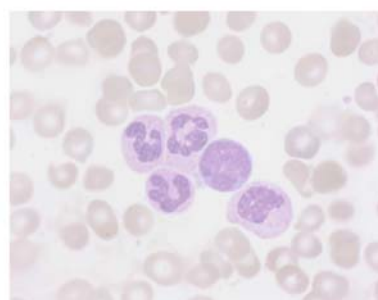
a)



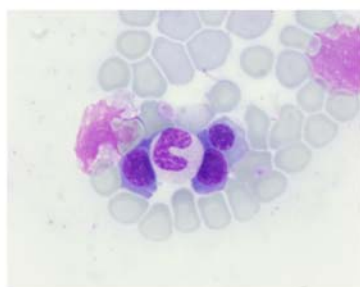
b)



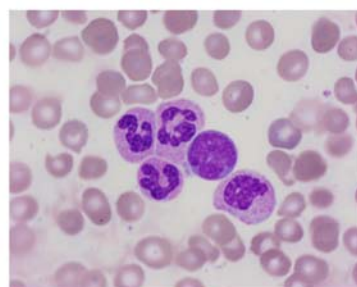
c)



d)



e)



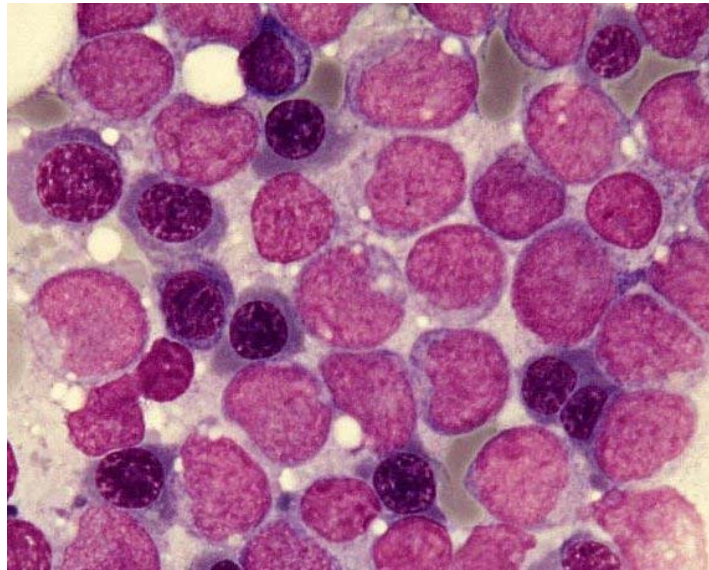
f)

Rys. 4.3 Przykład skalowania obrazów (a,b) przy użyciu dwóch wartości referencyjnych (c,d) i przy użyciu jednej wartości referencyjnej (e,f)

np.: "zamiatanie", czy filtrację przy użyciu filtru o charakterystyce typu "top hat" (transformacja kapelusza). Należy zaznaczyć, że operacje morfologiczne działają tylko na obrazach binarnych lub w skali szarości. Stąd dokonuje się wstępnie przetransformowanie obrazu na jeden z tych formatów.

Przykładowym obrazem poddanym analizie w pracy jest rozmaz szpiku kostnego chorego na ostrą białaczkę mieloblastyczną, zarejestrowany z powiększeniem 1000× w barwieniu Maya-Grunwalda-Giemsy (MGG) (rys 4.4). Poszczególne komórki wypełniają

prawie cały widoczny obszar, co stanowi dodatkowe utrudnienie w procesie segmentacji. Barwa cytoplazmy mieloblastów jest zbliżona do zabarwienia otoczenia, natomiast chromatyna ma złożoną strukturę.

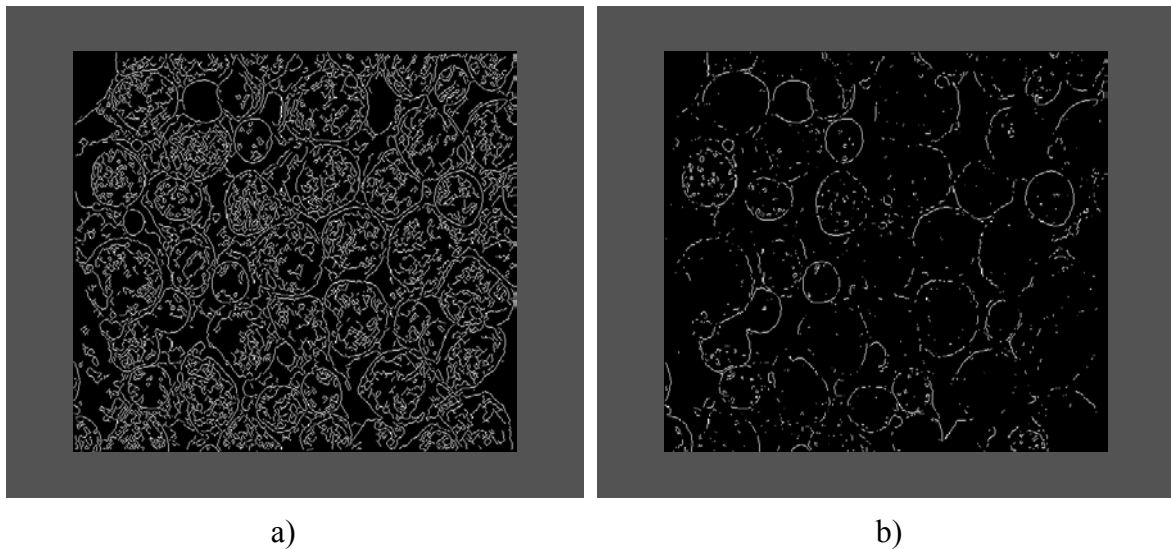


Rys 4.4 Obraz szpiku kostnego z przewagą mieloblastów (komórki o jaśniejszych jądrach), widoczne również erytroblasty wielobarwliwe (komórki o ciemniejszych jądrach).

Użycie standardowych procedur wykrywania kształtów dla obrazu kolorowego przetransformowanego do skali szarości z jednakowymi wagami dla poszczególnych składników barw nie przynosi zadowalających rezultatów. Widać to wyraźnie na rys 4.5, gdzie zastosowano algorytm wykrywania krawędzi Canny (rys. 4.5a) i Sobela (rys. 4.5b). Przy zastosowaniu standardowej procedury otrzymuje się niepełne kontury komórek (brakuje odcinków na styku komórki z innymi) oraz kontury będące pochodną złożoności tekstur cytoplazmu i chromatyny. W efekcie przy zastosowaniu metody filtru Sobela otrzymuje się szereg oddzielnych fragmentów krawędzi i zbiór punktów heterogenności chromatyny, a dla algorytmu Canny ogromny nadmiar krawędzi.

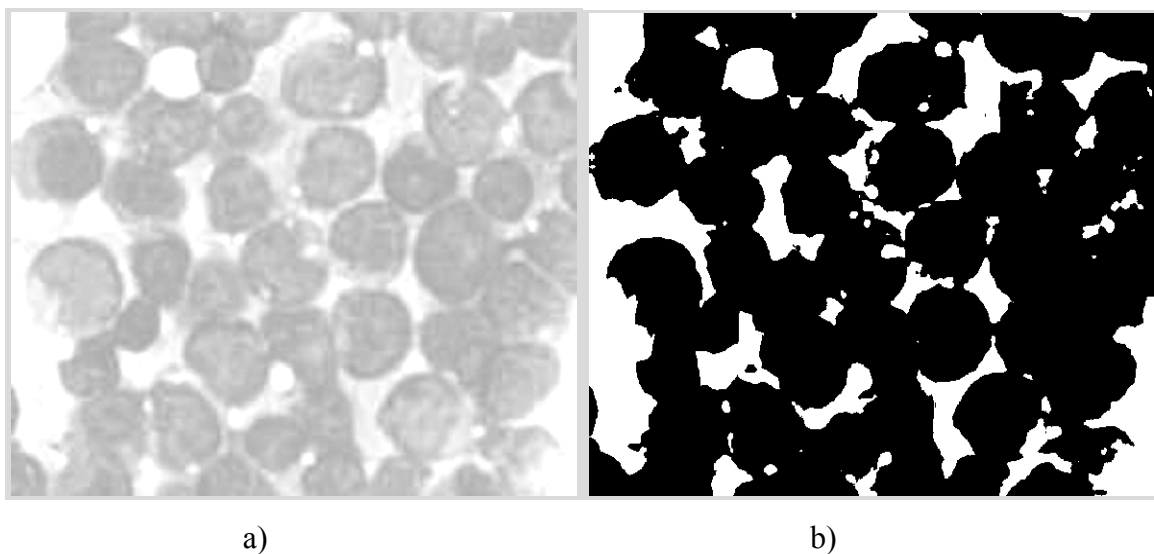
Jak już wspomniano, podstawowym krokiem w procesie przygotowania obrazu do segmentacji jest przetransformowanie go do skali szarości. Z przeprowadzonych badań wynika, że najbardziej uwidaczniającą brzegi komórek jest różnica pomiędzy barwą niebieską i zieloną. W wyniku takiej operacji otrzymuje się wysokie prawdopodobieństwo wyeliminowania szarych fragmentów obrazu nie należących do komórek jądrzastych, czyli tłuszczu i czerwonych krwinek. Rysunek 4.6a przedstawia dopełnienie wyniku takiej operacji.





Rys 4.5 Efekt wykrywania krawędzi za pomocą a) algorytmu Canny; b) filtru Sobela

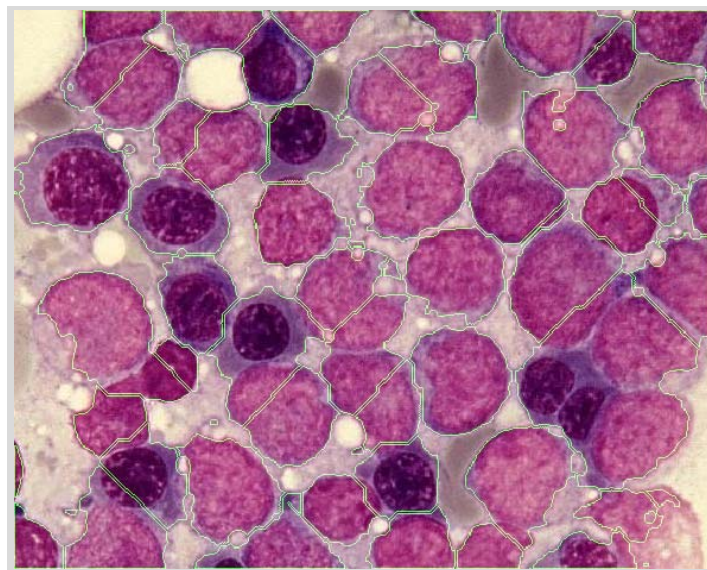
Tak otrzymany obraz w skali szarości należy przekształcić na obraz binarny, który w następnych etapach zostanie poddany dalszym operacjom morfologicznym. Skuteczna i powtarzalna dla wielu obrazów jest metoda segmentacji przez progowanie. Najlepsze efekty otrzymuje się dla wartości progu z przedziału (0.8 – 0.98) w skali znormalizowanej (0 - 1). Na rys 4.6b przedstawiono efekt takiej segmentacji dla progu 0.94, uznanego za najbardziej uniwersalny dla wielu obrazów. Tak powstały obraz binarny można poddać operacjom morfologicznym np. zamykania. Operacja taka jest przeprowadzona na kolorze białym i usuwa wszelkie drobne pozostałości nie należące do komórek, wygładzając ich kontury. Element strukturujący może mieć postać kwadratu, linii, dysku, koła lub sześciokąta.



Rys 4.6 Obraz w skali szarości otrzymany przez odjęcie składowej zielonej od niebieskiej (a) i po segmentacji przez progowanie (b)

Ponieważ naturalne kształty obiektów występujących na obrazie są zwykle nieliniowe zalecane jest stosowanie elementu o brzegach nieliniowych, np. dysku o wymiarze 3. Użycie małego elementu (wymiar elementu 1 lub 2) daje zwykle niewielki efekt, natomiast zbyt dużego (wymiar powyżej 5) powoduje wyraźną utratę powierzchni komórek. Zdaniem autora optymalny jest wymiar 3. Alternatywą jest metoda w postaci kilkakrotnej erozji i dylatacji, ale wadą tego postępowania jest możliwość usunięcia z tła niewielkich obszarów pomiędzy blisko ułożonymi komórkami.

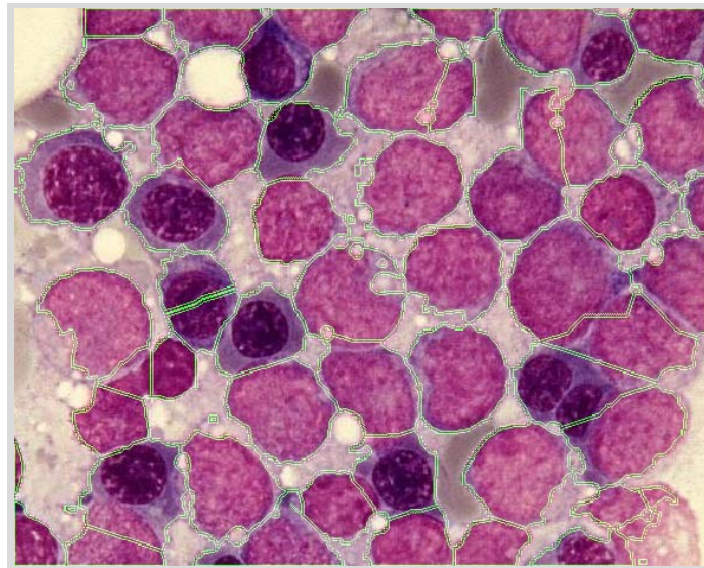
Kolejnym etapem jest wygenerowanie macierzy odległości poszczególnych pikseli komórek od tła i zastosowanie segmentacji metodą działów wodnych. Standardowa funkcja w Matlabie zwraca macierz o odległościach określonych liczbą pikseli dzielących punkt



Rys 4.7 Obraz po segmentacji metodą działów wodnych dla standardowej macierzy odległości

obliczeniowy od tła. Taka transformacja obrazu tworzy silne spadki w postaci linii, będące efektem niewielkich nierówności brzegu komórek. W efekcie po użyciu segmentacji metodą działów wodnych otrzymuje się niektóre komórki podzielone, co jest błędem trudnym do naprawienia przy automatycznym (pozbawionym ingerencji człowieka) przetwarzaniu obrazu. Przykład takiego sposobu segmentacji przedstawiono na rys 4.7, na którym wiele komórek zostało podzielonych na 2 części. Zdecydowanie lepsze efekty daje generacja macierzy odległości przy użyciu określonego kształtu elementu strukturującego. Macierz taka jest generowana przy wielokrotnym powtórzeniu operacji erozji wybranym elementem SE. Wszystkie piksele usunięte z powierzchni komórek w kolejnym kroku otrzymują odległość

równą temu krokowi. Efektem takiej modyfikacji będzie zmniejszenie złych podziałów komórek na części.



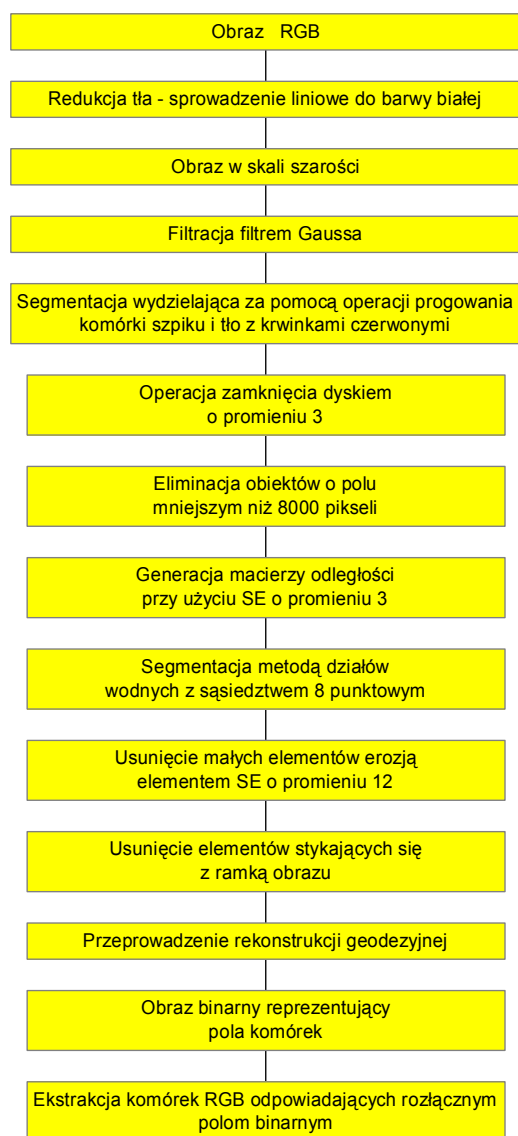
Rys 4.8 Obraz po segmentacji metodą działów wodnych dla zmodyfikowanej macierzy odległości

Ponadto aby poprawić segmentację komórek stykających się z komórkami leżącymi na brzegu obrazu (niepełnymi nie podlegającymi klasyfikacji) obraz jest przedłużany we wszystkie cztery strony o skrajny wiersz lub kolumnę, dzięki czemu niepełne komórki są uzupełniane i nie wpływają na macierz odległości w analizowanym obszarze. Rysunek 4.8 przedstawia wynik segmentacji dla zmodyfikowanej macierzy odległości. Widać wyraźną poprawę jakości segmentacji (mniej komórek podzielonych na części).

Ostatnim etapem jest eliminacja niewielkich obszarów nie będących komórkami jądrazystymi lub będących częścią ich cytoplazmy. Do tego celu można zastosować metodę rekonstrukcji za pomocą erozji geodezyjnej. Dla analizowanych obrazów dobrano promień oktagonu równy 12 pikseli, przy pomocy którego wykonano erozję tworząc obraz znacznika rekonstrukcji.

Dodatkowo w obrazie znacznika wskazane jest usunięcie elementów dotykających brzegu obszaru obrazu. Zaletą zastosowania tej operacji na obrazie znacznika jest zachowanie komórek w niewielkim stopniu stykających się z krawędzią obszaru, a leżących na analizowanym obszarze. Jedynie komórki podzielone przez granice obszaru zostają usunięte.

Rys 4.9 przedstawia podsumowanie zastosowanego algorytmu segmentacji. Uwzględniono w nim podstawowe operacje wykonywane na obrazie, prowadzące do wydzielenia poszczególnych komórek.

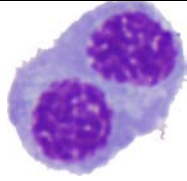
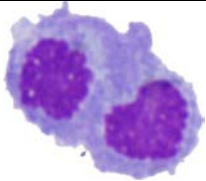
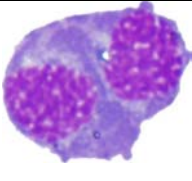
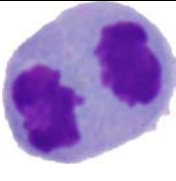
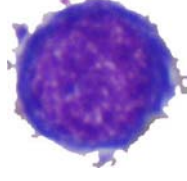
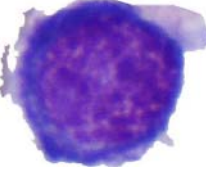
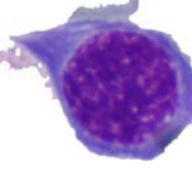
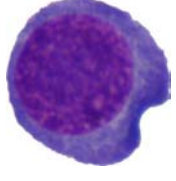
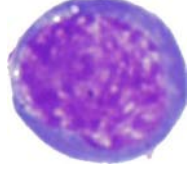
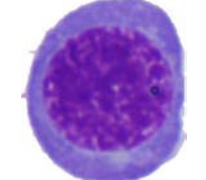
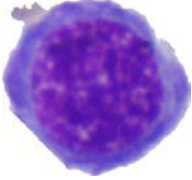
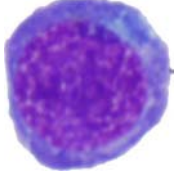
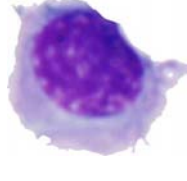

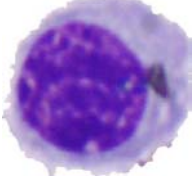
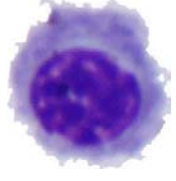


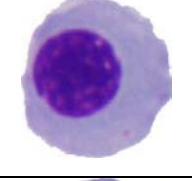

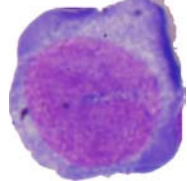

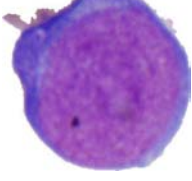
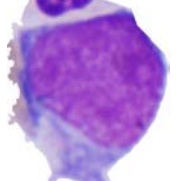
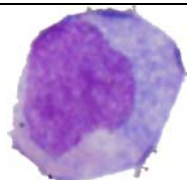
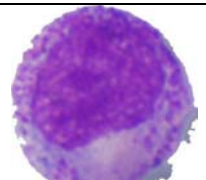
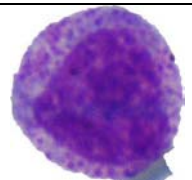
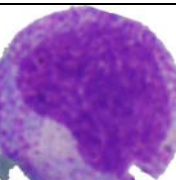
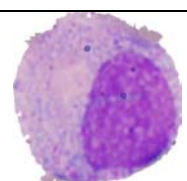
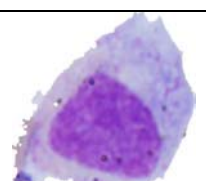
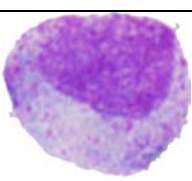
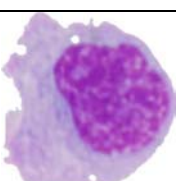


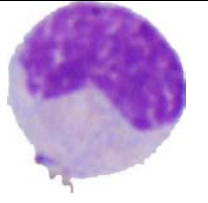
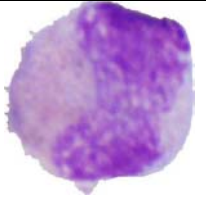
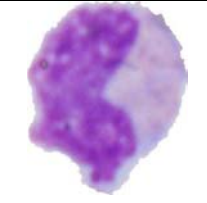
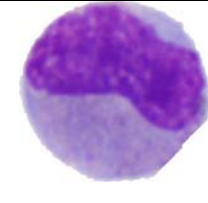
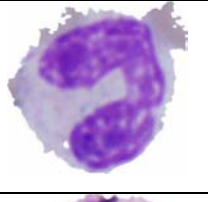

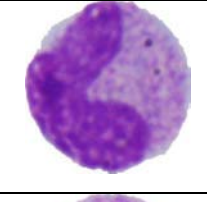
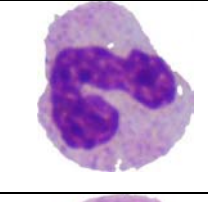


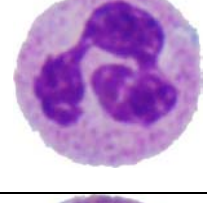
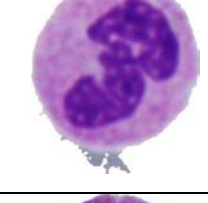
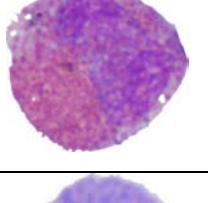
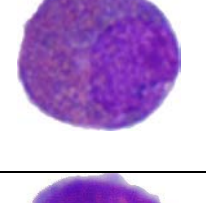
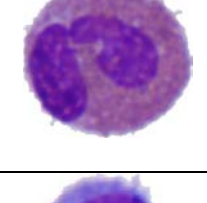
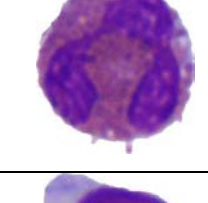
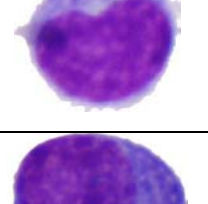
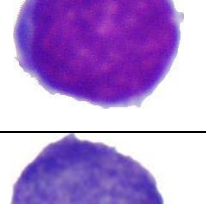
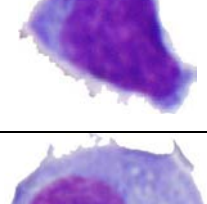
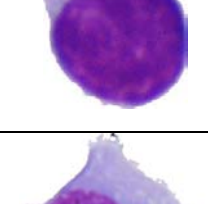




Rys 4.9 Schemat algorytmu segmentacji

### 4.3. Przykłady wydzielonych komórek – wyniki segmentacji

W tabeli 4.6 przedstawiono wybrane przykłady komórek krwiotwórczych, po 4 dla każdego rodzaju. Komórki pochodzą z rozmazów różnych pacjentów Instytutu Hematologii w Warszawie. Zostały one wykonane w różnym czasie w okresie 3 ostatnich lat.

Tabela 4.6 Przykłady wydzielonych komórek należących do różnych typów.

Nazwa	1	2	3	4
Postacie podziałowe erytropoezy				
Proerytroblasty				
Erytroblasty zasadochłonny				
Erytroblasty polichromatyczny				
Erytroblasty ortochromatyczny				
Blasty				
Promielocyty				
Mielocyty obojętnochłonne				

Metamielocyty obojętnochłonne				
Pałeczkowane obojętnochłonne				
Segmenty obojętnochłonne				
Eozynofile				
Limfocyty				
Plazmocyty				

## 5. Generacja i selekcja cech diagnostycznych do rozpoznania komórek

Rozpoznawanie obrazów jest zawsze związane z postrzeganiem i analizą cech charakteryzujących te obrazy. W rozwiązaniu proponowanym w pracy przy automatycznym rozpoznawaniu komórek użyto wielu różnych cech, które w ogólności można zaliczyć do następujących grup:

- cechy teksturalne – opisujące tekstury, a więc rozkłady statystyczne odcieni szarości lub kolorów, bez dogłębnej analizy szczegółów obrazu,
- cechy geometryczne – charakteryzujące właściwości geometryczne obrazu, takie jak pole powierzchni, obwód, wypukłości, symetria i inne parametry opisujące kształt obiektu,
- cechy statystyczne – określające statystyczną zmienność obiektu przy pomocy momentów statystycznych, np. wartość średnia, wariancja, skośność, kurtoza,
- cechy morfologiczne – parametry charakteryzujące obraz, powstałe bezpośrednio w wyniku zastosowania operacji morfologicznych na badanym obrazie, np. porównanie powierzchni obiektu przed i po wykonaniu określonej operacji morfologicznej.

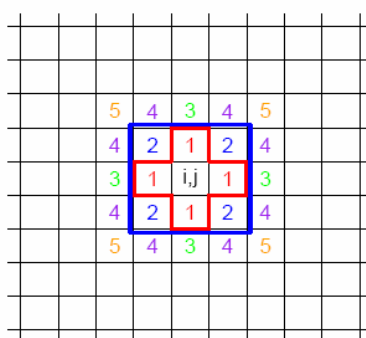
Tak wygenerowane cechy obrazu komórek lepiej lub gorzej charakteryzują je pod względem rozróżniania różnych typów. Służą one stworzeniu wektora wejściowego  $x$  dla klasyfikatora, przy czym można wykorzystać je wszystkie naraz bądź jedynie te, które najlepiej różnicują różne typy komórek. Należy przy tym pamiętać o normalizacji, czyli sprowadzeniu ich wartości do zbliżonych poziomów. Normalizacja taka została przeprowadzona w pracy przez podzielenie wartości rzeczywistych cech przez wartość maksymalną odpowiadającą cechom tworzącym bazę danych komórek na etapie uczenia. W pracy, w sposób ciągły (na przestrzeni 2 lat), dokonywano akwizycji obrazów kolejnych komórek stosując wartości normalizujące przyjęte na wstępie. Wskutek tego niektóre cechy mogą mieć wartości nieznacznie przekraczające zakres (0-1), co nie ma jednak większego wpływu na jakość działania klasyfikatora.

### 5.1 Cechy teksturalne

Istnieje wiele metod generacji cech tekstury. Różnią się one liczbą wydzielanych cech i podstawami matematycznymi ich generacji. Do klasyfikacji komórek po wielu próbach wstępnych wybrano cechy odpowiadające algorytmom Markova i Unsera. Wybór ten był

uzasadniony wcześniejszymi badaniami nad jakością otrzymywanych cech tekstur obrazów komórek przy użyciu (poza wymienionymi) również metody Haralicka, Gabora i opisu fraktalnego. Do celów pracy zmodyfikowano i zaadoptowano w środowisku Matlab'a program generacji cech teksturalnych, opracowany w [1].

Algorytm Markowa [1,13,50,60] opisuje tekstury wyrażając wartości odcieni szarości obrazu za pomocą funkcji odcieni sąsiedztwa zdefiniowanej maską. Cechą charakterystyczną każdego piksela jest jego odległość od piksela centralnego maski. W pierwszym kroku określa się rozmiar maski sąsiedztwa definiujący liczbę generowanych cech. Przykładowa maska M o rozmiarze 5×5 jest pokazana na rys. 5.1.



Rys. 5.1 Maska sąsiedztwa M o rozmiarze 5×5 dla piksela o indeksie (i,j)

Dla przyjętej maski sąsiedztwa M formułowane są równania opisujące stopnie szarości  $f(\mathbf{x})$  pikseli zawartych w masce:

$$f(\mathbf{x}) = \sum_{\mathbf{x}+\mathbf{d} \in M} \Theta_{\mathbf{d}} [f(\mathbf{x}+\mathbf{d}) + f(\mathbf{x}-\mathbf{d})] + \varepsilon(\mathbf{x}) \quad (5.1)$$

gdzie  $\mathbf{d}$  oznacza odległość mierzona w pikselach na masce sąsiedztwa,  $\Theta_{\mathbf{d}}$  – wyznaczony współczynnik zależności a  $\varepsilon(\mathbf{x})$  – różnicę pomiędzy wartością rzeczywistą stopnia szarości piksela i obliczoną z maski sąsiedztwa. Z powyższych równań estymuje się wektor rozwiązań  $\Theta$  odpowiadający równaniu:

$$\mathbf{g}_v = \Theta \cdot \mathbf{g}_m$$

gdzie

$$\begin{aligned} g_v(i) &= \sum_{\mathbf{x} \in \Omega} [f(\mathbf{x}+i) + f(\mathbf{x}-i)] \cdot f(\mathbf{x}) \\ g_m(i, j) &= \sum_{\mathbf{x} \in \Omega} [f(\mathbf{x}+i) + f(\mathbf{x}-i)] \cdot [f(\mathbf{x}+j) + f(\mathbf{x}-j)] \end{aligned} \quad (5.2)$$

oraz

$$\varepsilon = \frac{1}{N_{\mathbf{x} \in \Omega}} \sum_{\mathbf{x} \in \Omega} \sum_d [f(\mathbf{x}) \cdot \Theta_{\mathbf{d}}] \cdot [f(\mathbf{x}) \cdot \Theta_{\mathbf{d}}]$$



$\Omega$  oznacza zbiór pikseli  $\mathbf{x}$  posiadających sąsiadów w odległości  $d$ ,  $N$  - liczbę takich pikseli  $a$  i  $i$  oraz  $j$  - odległości od piksela centralnego. Dla zastosowanej maski  $M$  o rozmiarze  $10 \times 10$  otrzymuje się w sumie 11 cech. W klasyfikacji komórek wykorzystano jedynie ostatnią cechę, generowaną dla jądra i cytoplazmy w skali szarości, gdyż rozkłady pozostałych cech dla badanych obrazów nie wnosily istotnej wartości diagnostycznej przy rozpoznaniu obrazów różnych typów komórek.

Algorytm Unsera [1,13,50,60] oparty jest na histogramach sum i różnic poziomów szarości pikseli położonych w odległości  $d$  w wybranych kierunkach. Jest uproszczoną wersją metody GLCM Haralicka. W wykorzystanym wariancie przyjęto  $d=1$  oraz jeden kierunek  $0^\circ$  z uwagi na brak ukierunkowania analizowanych obiektów obrazu (komórek, w szczególności ich jąder lub cytoplazmy). Przyjmując  $f(\mathbf{x})$  jako oznaczenie stopnia szarości piksela o położeniu  $\mathbf{x}$ , wyznaczane są sumy (funkcja  $sum(\mathbf{x})$ ) i różnice (funkcja  $diff(\mathbf{x})$ ):

$$\begin{aligned} sum(\mathbf{x}) &= f(\mathbf{x}) + f(\mathbf{x} + d) \\ diff(\mathbf{x}) &= f(\mathbf{x}) - f(\mathbf{x} + d) \end{aligned} \quad (5.3)$$

Następnie generowane są histogramy sum i różnic, obydwa w zakresach dwa razy większych od zakresu skali szarości obrazu:

$$\begin{aligned} H_{sum}(i) &= card\{p \mid sum(p) = (i_1, i_2, \dots, i_m)\} \\ H_{diff}(j) &= card\{p \mid diff(p) = (j_1, j_2, \dots, j_m)\} \end{aligned} \quad (5.4)$$

przy czym  $i_1=0$ ,  $i_m=511$ ,  $j_1=-255$ ,  $j_m=255$  w przypadku skali szarości kodowanej 8 - bitowo.

Wykorzystano następujące cechy generowane na podstawie histogramów:

- wartość średnia histogramu sum:

$$M_{sum} = \frac{\sum_{i=0}^{511} i \cdot H_{sum}(i)}{\sum_{i=0}^{511} H_{sum}(i)} \quad (5.5)$$

- kątowny moment drugiego rzędu histogramu sum:

$$M_{2sum} = \frac{\sum_{i=0}^{511} (H_{sum}(i))^2}{\left(\sum_{i=0}^{511} H_{sum}(i)\right)^2} \quad (5.6)$$

- kątowny moment drugiego rzędu histogramu różnic:

$$M_{2_{diff}} = \frac{\sum_{i=-255}^{255} (H_{diff}(i))^2}{\left( \sum_{i=-255}^{255} H_{diff}(i) \right)^2} \quad (5.7)$$

- kontrast histogramu sum:

$$K_{sum} = \frac{\sum_{i=0}^{511} (i - M_{sum})^2 \cdot H_{sum}(i)}{\sum_{i=0}^{511} H_{sum}(i)} \quad (5.8)$$

- kontrast histogramu różnic:

$$K_{diff} = \frac{\sum_{i=-255}^{255} (i - M_{diff})^2 \cdot H_{diff}(i)}{\sum_{i=-255}^{255} H_{diff}(i)} \quad (5.9)$$

- entropia histogramu sum:

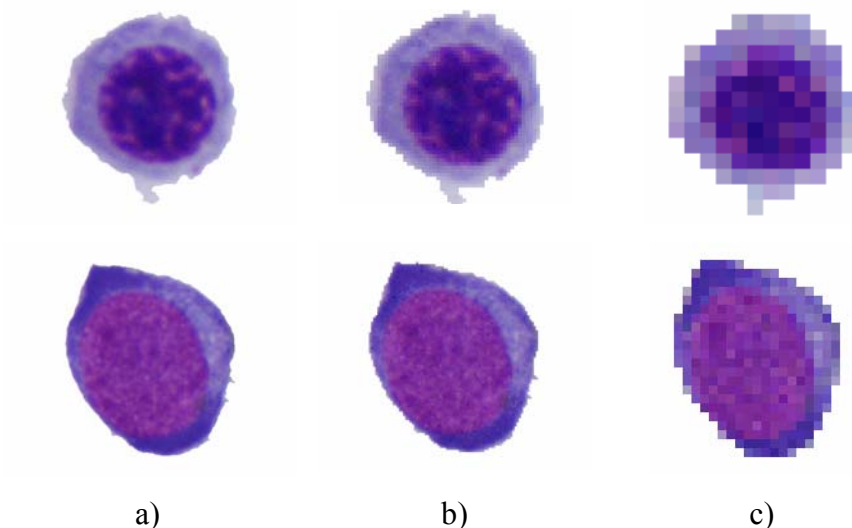
$$E_{sum} = - \sum_{i=0, H_{sum}(i)>0}^{511} H_{sum}(i) \cdot \log \frac{H_{sum}(i)}{\sum_{j=0}^{511} H_{sum}(j)} \quad (5.10)$$

- entropia histogramu różnic:

$$E_{diff} = - \sum_{i=-255, H_{diff}(i)>0}^{255} H_{diff}(i) \cdot \log \frac{H_{diff}(i)}{\sum_{j=-255}^{255} H_{diff}(j)} \quad (5.11)$$

Powyższe siedem cech (5.5 – 5.11) wygenerowano dla każdej składowej barw (RGB) dla jądra oraz cytoplazmy oddzielnie. W sumie wygenerowano w ten sposób 42 cechy dla każdego oryginalnego obrazu komórek.

Ze względu na wysoką rozdzielczość obrazu i przyjętą odległość  $d=1$ , cechy Unsera wygenerowane dla obrazu oryginalnego opisują głównie charakter mikroskopowy tekstury, pomijając jej charakter makroskopowy. Dlatego też zdecydowano się wygenerować takie same cechy dla obrazu przekształconego, o niższej rozdzielczości. Cechy te wygenerowano tylko dla jądra, dla każdej składowej barw oddzielnie, przy czterokrotnym i 16-krotnym zmniejszaniu rozdzielczości. Uzyskano w ten sposób dodatkowo 42 cechy, wykorzystane w badaniach dla większej liczby rozpoznawanych klas. Przykłady obrazów komórek o zmniejszonej rozdzielczości przedstawiono na rysunku 5.2.



Rys. 5.2 Przykłady przekształcania obrazu oryginalnego (a), na obraz o zmniejszonej rozdzielczości 4-krotnie (b) i 16-krotnie (c)

## 5.2 Cechy geometryczne

Istotnymi cechami różniącymi poszczególne typy komórek są cechy geometryczne, odnoszące się zarówno do jądra jak i całej komórki. Ponieważ kształt geometryczny analizowanych obiektów nie zależy od składowej barwy, cechy te były generowane przy użyciu masek binarnych komórki i jej jądra. Należy zaznaczyć, że największy wpływ na jakość cech geometrycznych ma poprawność przeprowadzonej segmentacji obrazu (dla cech całej komórki) oraz segmentacji części składowych pojedynczej komórki (jądra i cytoplazmy). Zdefiniowano następujące cechy geometryczne:

1. pole powierzchni jądra komórki mierzone w pikselach,
2. stosunek pola powierzchni jądra do pola powierzchni całej komórki,
3. najdłuższa oś jądra,
4. najkrótsza oś jądra,
5. ekscentryczność jądra,
6. powierzchnia wielokąta wypukłego opisanego na jądrze,
7. powierzchnia jądra po usunięciu wewnętrznych dziur,
8. promień koła o takiej samej powierzchni jak jądro,
9. stosunek powierzchni jądra do powierzchni wielokąta wypukłego opisanego na jądrze,
10. stosunek powierzchni jądra do powierzchni najmniejszego prostokąta opisanego na jądrze,
11. obwód jądra,

12. średnia odległość pikselowa od centralnego piksela jądra,
13. zwartość mierzona jako stosunek kwadratu obwodu jądra do powierzchni jądra,
14. symetria osiowa jądra,
15. symetria pól względem najdłuższej osi jądra,
16. obwód wielokąta wypukłego opisanego na jądrze,
17. różnica obwodu wielokąta wypukłego opisanego na jądrze i obwodu samego jądra,
18. liczba wklęsłości występujących w jądrze,
19. stosunek powierzchni wielokąta wypukłego opisanego na jądrze do powierzchni jądra po usunięciu wewnętrznych dziur.

### 5.3 Cechy statystyczne

Do pomiaru stopnia zmienności nasycenia barw w całej komórce, w cytoplazmie i jądrze można zastosować cechy statystyczne, wykorzystując również momenty statystyczne wyższych rzędów określone na podstawie histogramu. Rozwinięciem tego opisu jest generacja takich cech dodatkowo dla macierzy gradientu utworzonej z obrazu wejściowego.

Wygenerowane zostały następujące cechy statystyczne:

1. wartość średnia dla obrazu oryginalnego,
2. wariancja dla obrazu oryginalnego,
3. skośność dla obrazu oryginalnego,
4. kurtoza dla obrazu oryginalnego,
5. wartość średnia elementów macierzy modułów gradientu obrazu,
6. wariancja wartości elementów macierzy modułów gradientu obrazu,
7. skośność wartości elementów macierzy modułów gradientu obrazu,
8. kurtoza wartości elementów macierzy modułów gradientu obrazu.

Skośność i kurtoza obrazu oraz skośność gradientu i kurtoza gradientu zostały wygenerowane dla całych komórek, pozostałe cechy osobno dla jądra i cytoplazmy. Wszystkie te cechy wyznaczono dla dwóch składowych – czerwonej i zielonej. Zrezygnowano z wykorzystania składowej niebieskiej ze względu na jej dużą zależność (prawie liniową) od pozostałych barw. W ten sposób otrzymano dla całych komórek 4 cechy, dla ich jądra 12 i dla ich cytoplazmy 12 cech, co dało w sumie 24 cechy statystyczne.

## 5.4 Cechy morfologiczne

Mianem cech morfologicznych określono różne cechy obiektów obrazów, zmieniające się w zależności od przeprowadzonych na nich operacjach morfologicznych. W pierwszym kroku z jądra wydzielany jest ciemniejszy fragment przy użyciu operacji progowania. Wartość progu była określana dwoma metodami: metodą Otsu i jako 2/5 długości histogramu obrazu. Następnie na tym fragmencie przeprowadzano operacje erozji dla zbadania szybkości zmian pola powierzchni i liczby rozłącznych jego części, a więc zmienności jego cech geometrycznych. Wygenerowano następujące cechy dla obydwu wartości progowania:

1. pole powierzchni fragmentu jądra przed erozją,
2. pole powierzchni fragmentu jądra po jednokrotnej erozji dyskiem o promieniu 4 piksele,
3. liczba rozłącznych fragmentów jądra po jednokrotnej erozji dyskiem o promieniu 4 piksele,
4. najmniejsza liczba erozji jądra dyskiem o promieniu 4 piksele usuwająca całkowicie jego powierzchnię.

W ten sposób otrzymano 8 cech morfologicznych. Parametry dysku (4 piksele) dobrano w wyniku wielu wstępnych eksperymentów, mających ustalić najlepsze jego parametry dla uzyskania dobrego różnicowania klas komórek.

W tabeli 5.1 przedstawiono pełny zestaw cech uszeregowanych i ponumerowanych według sposobu ich generacji: cechy statystyczne, teksturalne, geometryczne i morfologiczne. Podano również numery przyporządkowane poszczególnym cechom, ułatwiające posługiwanie się nimi w dalszych badaniach.

Tabela 5.1 Maksymalny zbiór cech użytych przy definiowaniu wektora wejściowego  $x$

Nr cechy	Nazwa	Nr cechy	Nazwa
<b>CECHY STATYSTYCZNE</b>			
1	skośność histogramu komórki barwy czerwonej	2	skośność histogramu komórki barwy zielonej
3	kurtoza histogramu komórki barwy czerwonej	4	kurtoza histogramu komórki barwy zielonej
5	wartość średnia jądra barwy czerwonej	6	wartość średnia jądra barwy zielonej
7	wartość średnia cytoplazmy barwy czerwonej	8	wartość średnia cytoplazmy barwy zielonej
9	wariancja jądra barwy czerwonej	10	wariancja jądra barwy zielonej
11	wariancja cytoplazmy barwy czerwonej	12	wariancja cytoplazmy barwy zielonej
13	wartość średnia gradientu jądra barwy czerwonej	14	wartość średnia gradientu jądra barwy zielonej
15	wartość średnia gradientu cytoplazmy barwy czerwonej	16	wartość średnia gradientu cytoplazmy barwy zielonej
17	wariancja gradientu jądra barwy czerwonej	18	wariancja gradientu jądra barwy zielonej
19	wariancja gradientu cytoplazmy barwy czerwonej	20	wariancja gradientu cytoplazmy barwy zielonej
21	skośność histogramu gradientu komórki barwy czerwonej	22	skośność histogramu gradientu komórki barwy zielonej
23	kurtoza histogramu gradientu komórki barwy czerwonej	24	kurtoza histogramu gradientu komórki barwy zielonej

CECHY TEKSTURALNE			
25	11-ta cecha Markova dla jądra w skali szarości	26	11-ta cecha Markova dla cytoplazmy w skali szarości
27	wartość średnia histogramu sum dla cytoplazmy barwy czerwonej	28	moment kątowy histogramu sum dla cytoplazmy barwy czerwonej
29	moment kątowy histogramu różnic dla cytoplazmy barwy czerwonej	30	kontrast histogramu sum dla cytoplazmy barwy czerwonej
31	kontrast histogramu różnic dla cytoplazmy barwy czerwonej	32	entropia histogramu sum dla cytoplazmy barwy czerwonej
33	entropia histogramu różnic dla cytoplazmy barwy czerwonej	34	wartość średnia histogramu sum dla cytoplazmy barwy zielonej
35	moment kątowy histogramu sum dla cytoplazmy barwy zielonej	36	moment kątowy histogramu różnic dla cytoplazmy barwy zielonej
37	kontrast histogramu sum dla cytoplazmy barwy zielonej	38	kontrast histogramu różnic dla cytoplazmy barwy zielonej
39	entropia histogramu sum dla cytoplazmy barwy zielonej	40	entropia histogramu różnic dla cytoplazmy barwy zielonej
41	wartość średnia histogramu sum dla cytoplazmy barwy niebieskiej	42	moment kątowy histogramu sum dla cytoplazmy barwy niebieskiej
43	moment kątowy histogramu różnic dla cytoplazmy barwy niebieskiej	44	kontrast histogramu sum dla cytoplazmy barwy niebieskiej
45	kontrast histogramu różnic dla cytoplazmy barwy niebieskiej	46	entropia histogramu sum dla cytoplazmy barwy niebieskiej
47	entropia histogramu różnic dla cytoplazmy barwy niebieskiej	48	wartość średnia histogramu sum dla jądra barwy czerwonej
49	moment kątowy histogramu sum dla jądra barwy czerwonej	50	moment kątowy histogramu różnic dla jądra barwy czerwonej
51	kontrast histogramu sum dla jądra barwy czerwonej	52	kontrast histogramu różnic dla jądra barwy czerwonej
53	entropia histogramu sum dla jądra barwy czerwonej	54	entropia histogramu różnic dla jądra barwy czerwonej
55	wartość średnia histogramu sum dla jądra barwy zielonej	56	moment kątowy histogramu sum dla jądra barwy zielonej
57	moment kątowy histogramu różnic dla jądra barwy zielonej	58	kontrast histogramu sum dla jądra barwy zielonej
59	kontrast histogramu różnic dla jądra barwy zielonej	60	entropia histogramu sum dla jądra barwy zielonej
61	entropia histogramu różnic dla jądra barwy zielonej	62	wartość średnia histogramu sum dla jądra barwy niebieskiej
63	moment kątowy histogramu sum dla jądra barwy niebieskiej	64	moment kątowy histogramu różnic dla jądra barwy niebieskiej
65	kontrast histogramu sum dla jądra barwy niebieskiej	66	kontrast histogramu różnic dla jądra barwy niebieskiej
67	entropia histogramu sum dla jądra barwy niebieskiej	68	entropia histogramu różnic dla jądra barwy niebieskiej
CECHY GEOMETRYCZNE			
69	pole powierzchni jądra komórki mierzone w pikselach	70	stosunek pola powierzchni jądra do pola powierzchni całej komórki
71	najdłuższa oś jądra	72	najkrótsza oś jądra
73	ekscentryczność jądra	74	powierzchnia wielokąta wypukłego opisanego na jądrze
75	powierzchnia jądra po usunięciu dziur	76	promień koła o takiej samej powierzchni jak jądro
77	stosunek powierzchni jądra do powierzchni wypukłej	78	stosunek powierzchni jądra do powierzchni najmniejszego prostokąta
79	obwód	80	średnia odległość od centralnego piksela
81	zwartość mierzona jako stosunek kwadratu obwodu jądra do powierzchni jądra	82	symetria osiowa
83	symetria pól względem najdłuższej osi jądra	84	obwód wielokąta wypukłego opisanego na jądrze
85	różnica obwodu wielokąta wypukłego i jądra	86	liczba wklęsłości

87	stosunek powierzchni wielokąta wypukłego do powierzchni jądra po usunięciu dziur		
<b>CECHY MORFOLOGICZNE</b>			
88	pole powierzchni fragmentu jądra przed erozją dla progu 2/5 długości histogramu	89	pole powierzchni fragmentu jądra po jednokrotnej erozji dyskiem o promieniu 4 pikseli dla progu 2/5 długości histogramu
90	liczba rozłącznych fragmentów jądra po jednokrotnej erozji dyskiem o promieniu 4 pikseli dla progu 2/5 długości histogramu	91	najmniejsza liczba erozji jądra dyskiem o promieniu 4 pikseli usuwająca jego powierzchnię dla progu 2/5 długości histogramu
92	pole powierzchni fragmentu jądra przed erozją dla progu ustalonego metodą Otsu	93	pole powierzchni fragmentu jądra po jednokrotnej erozji dyskiem o promieniu 4 pikseli dla progu ustalonego metodą Otsu
94	liczba rozłącznych fragmentów jądra po jednokrotnej erozji dyskiem o promieniu 4 pikseli dla progu ustalonego metodą Otsu	95	najmniejsza liczba erozji jądra dyskiem o promieniu 4 pikseli usuwająca jego powierzchnię dla progu ustalonego metodą Otsu
<b>CECHY TEKSTURALNE DLA OBRAZU O CZTERO-KROTNEJ REDUKCJI ROZDZIELCZOŚCI</b>			
96	wartość średnia histogramu sum dla jądra barwy czerwonej	97	moment kątowy histogramu sum dla jądra barwy czerwonej
98	moment kątowy histogramu różnic dla jądra barwy czerwonej	99	kontrast histogramu sum dla jądra barwy czerwonej
100	kontrast histogramu różnic dla jądra barwy czerwonej	101	entropia histogramu sum dla jądra barwy czerwonej
102	entropia histogramu różnic dla jądra barwy czerwonej	103	wartość średnia histogramu sum dla jądra barwy zielonej
104	moment kątowy histogramu sum dla jądra barwy zielonej	105	moment kątowy histogramu różnic dla jądra barwy zielonej
106	kontrast histogramu sum dla jądra barwy zielonej	107	kontrast histogramu różnic dla jądra barwy zielonej
108	entropia histogramu sum dla jądra barwy zielonej	109	entropia histogramu różnic dla jądra barwy zielonej
110	wartość średnia histogramu sum dla jądra barwy niebieskiej	111	moment kątowy histogramu sum dla jądra barwy niebieskiej
112	moment kątowy histogramu różnic dla jądra barwy niebieskiej	113	kontrast histogramu sum dla jądra barwy niebieskiej
114	kontrast histogramu różnic dla jądra barwy niebieskiej	115	entropia histogramu sum dla jądra barwy niebieskiej
116	entropia histogramu różnic dla jądra barwy niebieskiej		
<b>CECHY TEKSTURALNE DLA OBRAZU O 16-KROTNEJ REDUKCJI ROZDZIELCZOŚCI</b>			
117	wartość średnia histogramu sum dla jądra barwy czerwonej	118	moment kątowy histogramu sum dla jądra barwy czerwonej
119	moment kątowy histogramu różnic dla jądra barwy czerwonej	120	kontrast histogramu sum dla jądra barwy czerwonej
121	kontrast histogramu różnic dla jądra barwy czerwonej	122	entropia histogramu sum dla jądra barwy czerwonej
123	entropia histogramu różnic dla jądra barwy czerwonej	124	wartość średnia histogramu sum dla jądra barwy zielonej
125	moment kątowy histogramu sum dla jądra barwy zielonej	126	moment kątowy histogramu różnic dla jądra barwy zielonej
127	kontrast histogramu sum dla jądra barwy zielonej	128	kontrast histogramu różnic dla jądra barwy zielonej
129	entropia histogramu sum dla jądra barwy zielonej	130	entropia histogramu różnic dla jądra barwy zielonej
131	wartość średnia histogramu sum dla jądra barwy niebieskiej	132	moment kątowy histogramu sum dla jądra barwy niebieskiej
133	moment kątowy histogramu różnic dla jądra barwy niebieskiej	134	kontrast histogramu sum dla jądra barwy niebieskiej
135	kontrast histogramu różnic dla jądra barwy niebieskiej	136	entropia histogramu sum dla jądra barwy niebieskiej
137	entropia histogramu różnic dla jądra barwy niebieskiej		

## 5.5 Ocena jakości i selekcja cech

Przedstawione w tabeli 5.1 zestawy cech są zmaksymalizowanym zbiorem cech ustalonym w ciągu prawie trzech lat współpracy z Instytutem Hematologii w Warszawie nad automatycznym rozpoznawaniem komórek. Były one stopniowo definiowane i badane, w miarę rozszerzania bazy danych różnych typów komórek. W miarę wzrostu liczby typów, trudność zadania rośnie i okazuje się, że satysfakcjonujący do tej pory zbiór cech wymaga rozszerzenia. Badania rozpoczęto z zestawem 87 cech teksturalnych, geometrycznych i statystycznych, które dobrze sprawdzały się przy rozpoznawaniu do 16 typów komórek. Przy zwiększeniu tej liczby do 21 zaszła potrzeba ich dalszego rozszerzenia – w ten sposób powstały następne cechy (morfologiczne i dodatkowe teksturalne) od numeru 88 aż do 137.

Należy podkreślić, że podane wyżej liczby cech stanowią maksymalny rozmiar wektora  $x$ . Badania pokazały, że użycie maksymalnego zestawu cech nie prowadzi do najlepszych wyników, gdyż nie są one jednakowo ważne w procesie rozpoznania komórek. Pewne cechy w procesie rozpoznania pełnią funkcję szumu pomiarowego, pogarszając możliwość rozpoznania różnych typów komórek. Ważnym elementem procesu staje się zatem ocena jakości cech i opracowanie metod ich selekcji przy tworzeniu wektora wejściowego  $x$ . Zdolność generalizacji sieci neuronowej nie jest bezpośrednio powiązana z rozmiarem wektorów wejściowych, a raczej z ich składem i takim doбором, który najlepiej różnicuje różne typy komórek. W wyniku badań stwierdzono, że cechy wysoce skorelowane mają zwykle niekorzystny wpływ na jakość klasyfikacji, dominując nad innymi i tłumiąc w ten sposób ich korzystne działanie.

W badaniu jakości cech można zastosować dwie strategie. W pierwszej bada się każdą cechę indywidualnie, oceniając jej jakość pod kątem różnicowania klas. Cechy o najgorszym wskaźniku różnicowania są odrzucane, a pozostałe tworzą wektor wejściowy  $x$ . Taki sposób postępowania dobrze różnicuje jakość cech działających indywidualnie. W praktyce okazuje się, że włączenie równoległego działania wielu cech na raz może zmienić „jakość” danej cechy. Pewne cechy (nawet te gorsze) współpracując ze sobą wzbogacają się nawzajem, podnosząc wzajemnie ich wartość diagnostyczną. Stąd w praktyce ważniejsza jest metoda druga, polegająca na wartościowaniu poszczególnych cech działających jednocześnie. W tym podrozdziale zostanie przedstawiona analiza jakości cech w obu ujęciach.

W przypadku oceny każdej cechy indywidualnie, jej jakość będzie mierzona poprzez badanie stopnia korelacji zachodzącej pomiędzy wybraną cechą a pozostałymi oraz analizę rozkładu jej wartości średnich i wariancję w ramach poszczególnych klas. W drugim ujęciu

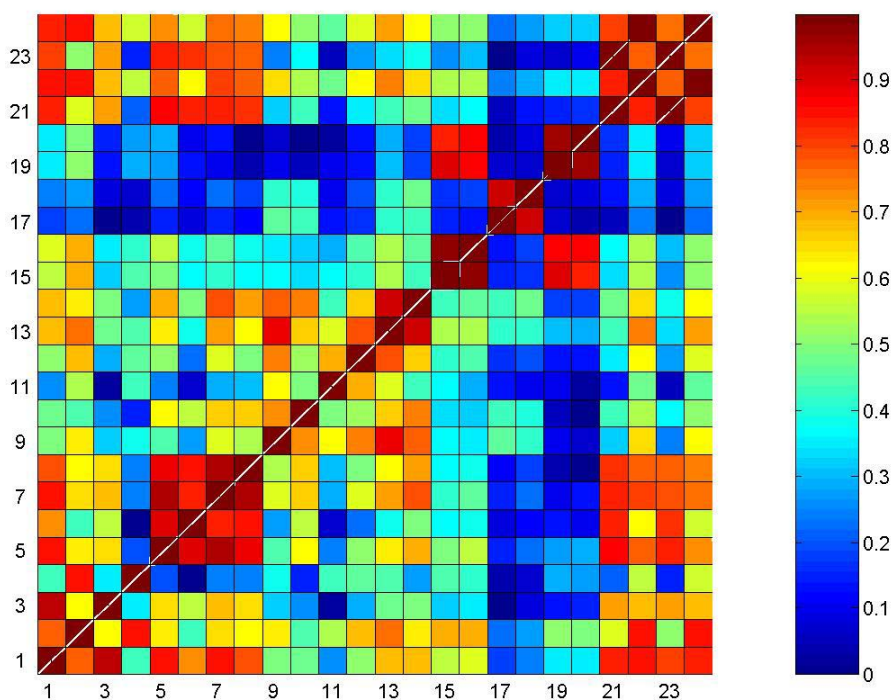


wartościowania, cechy działające łącznie badane będą przy zastosowaniu liniowej sieci SVM. Wyniki przedstawione w dalszych punktach rozdziału dotyczyć będą 3 podstawowych grup cech: statystycznych, geometrycznych i teksturalnych ( w sumie 87 cech).

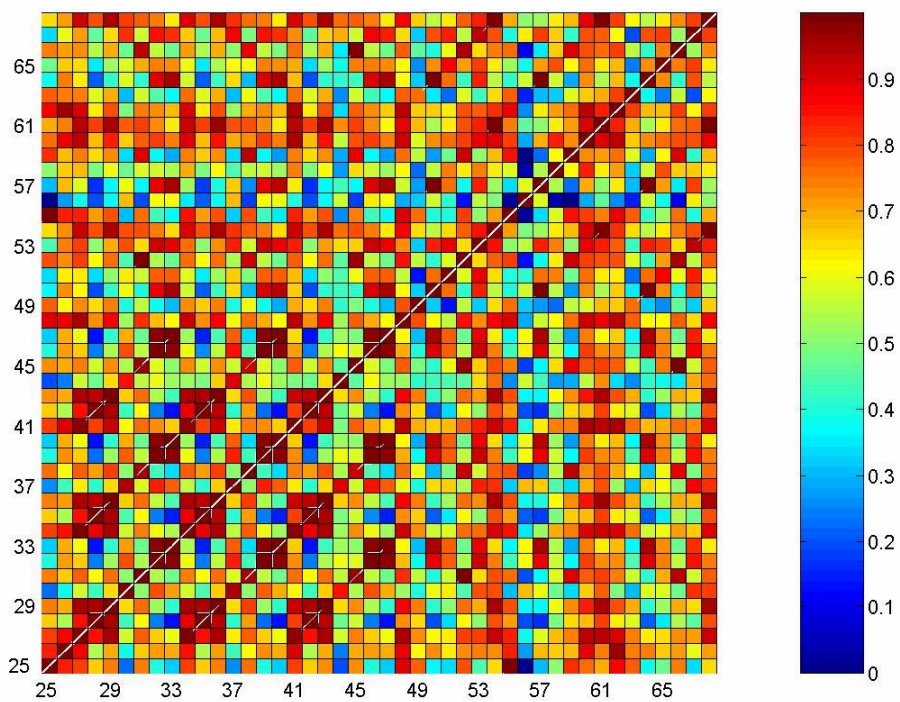
### 5.5.1. Analiza korelacyjna cech

Pierwsza redukcja wielowymiarowej przestrzeni cech jest możliwa przez eliminowanie tych cech, które są z innymi silnie skorelowane. Analizowano korelacje między składnikami pełnego zbioru cech. Dla uproszczenia analizy prezentację wyników pokazano jedynie dla trzech wybranych grup cech: statystycznych, geometrycznych i teksturalnych, przy czym ograniczono się do korelacji wewnątrz każdej grupy oddzielnie. Jest to uzasadnione tym, że poszczególne grupy są zwykle od siebie niezależne, a korelacja między nimi jest stosunkowo niewielka. Na rys. 5.3, 5.4 i 5.5 przedstawiono korelacje występujące między cechami statystycznymi, teksturalnymi i geometrycznymi.

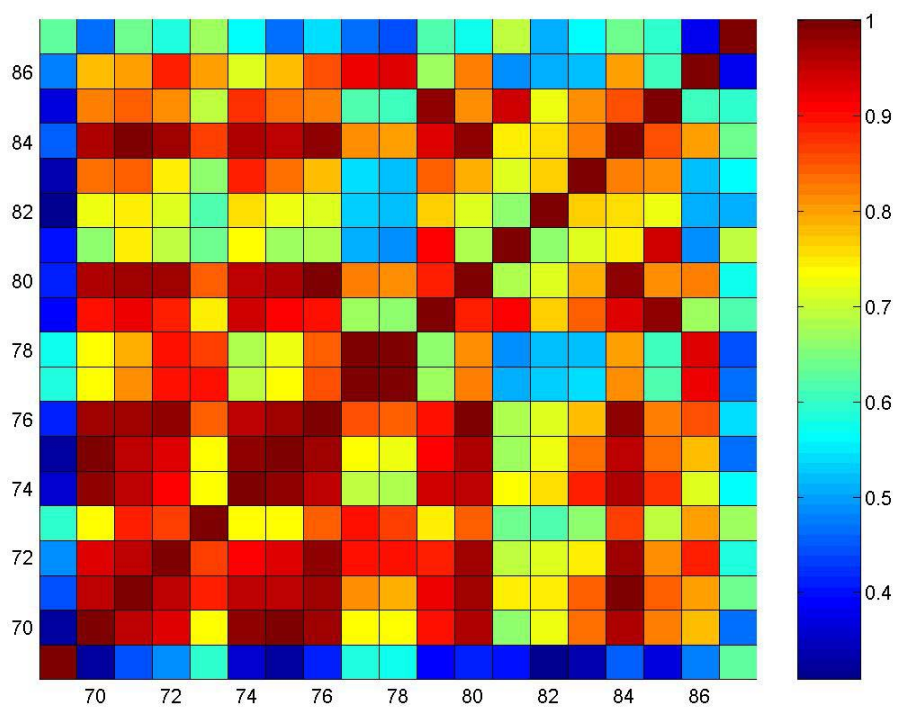
Największa zmienność wartości korelacji wystąpiła między cechami teksturalnymi, ale nawet dla nich można łatwo wybrać pewne cechy wysoce skorelowane. Największa procentowo liczba cech skorelowanych występuje wśród cech geometrycznych. Cechy statystyczne wydają się być stosunkowo mało skorelowane, co jest zgodne z przewidywaniami. Na podstawie analizy korelacyjnej można zauważyć, że wiele cech



Rys. 5.3 Obraz korelacji cech statystycznych



Rys. 5.4 Obraz korelacji cech teksturalnych



Rys. 5.5 Obraz korelacji cech geometrycznych

skorelowanych ze sobą (aż do 1/3) może zostać wyeliminowanych z zestawu tworzącego oryginalny wektor  $x$ , stanowiący wejście dla klasyfikatora. Wg powszechnie przyjętej opinii [50] spośród cech wysoce skorelowanych można wyeliminować większość z nich bez uszczerbku dla jakości działania klasyfikatora.

### 5.5.2. Selekcja cech na podstawie wartości średnich i wariancji danych

Innym sposobem selekcji indywidualnej cech jest analiza relacji zachodzących pomiędzy położeniami centrów i wariancją cech dla każdej z 2 klas. Miara użyteczności  $W_k(i,j)$  cechy  $k$  w rozróżnianiu klasy  $i$ -tej oraz  $j$ -tej jest zdefiniowana w postaci:

$$W_k(i, j) = \frac{|\mu_k(j) - \mu_k(i)|}{\sigma_k(j) + \sigma_k(i)} \quad (5.12)$$

gdzie  $\mu_k(i), \mu_k(j)$  oznaczają wartość średnią cechy  $k$  dla komórek odpowiednio klasy  $i$ -tej oraz  $j$ -tej, a  $\sigma_k(i), \sigma_k(j)$  odpowiednie odchylenia standardowe. Preferowane są cechy, dla których odległości centrów dwu różnych klas są jak największe, a wariancja wewnątrz danej klasy jak najmniejsza. W rezultacie im większa jest wartość  $W_k(i,j)$ , tym cecha  $k$  ma większy wpływ na rozróżnialność tych klas. Cechy powinny być oceniane dla każdej pary klas oddzielnie, gdyż dobre różnicowanie jednej pary klas nie oznacza tego samego dla innej pary. Należy zaznaczyć, że to kryterium nie wykrywa automatycznie cech skorelowanych, jest więc dobrym uzupełnieniem analizy korelacyjnej. Kryterium to jest pracochłonne i generuje ogromne ilości danych. Przy  $N$  cechach należy zbadać wszystkie możliwe kombinacje cech branych po dwie, czyli  $\binom{N}{2} = \frac{N!}{2!(N-2)!} = \frac{N(N-1)}{2}$ . Przy  $N=158$  daje to liczbę kombinacji równą 12403, trudną do pełnej analizy. W praktyce eliminuje się najpierw cechy skorelowane i dopiero po takiej redukcji cech przeprowadza się analizę wartości średnich i wariancji. Tabela 5.2 przedstawia przykładowe wartości wariancji dla 12 wybranych klas i 20 wybranych cech.

W tabeli 5.3 przedstawiono z kolei przykładowe odległości pomiędzy centrami poszczególnych klas. Biorąc pod uwagę małe wariancje cech (wartość od  $1E-4$  do  $1E-2$ ) odległości między centrami wydają się być raczej dobre. Jednakże dla wybranej cechy (na przykład nr 10) niektóre odległości pomiędzy parami klas nie są wystarczająco duże. Taka cecha powinna zostać rozważona jako kandydat do redukcji ze zbioru danych tworzących wektor wejściowy  $x$  klasyfikatora.

Tabela 5.2 Wariancje 20 wybranych cech dla 12 typów komórek

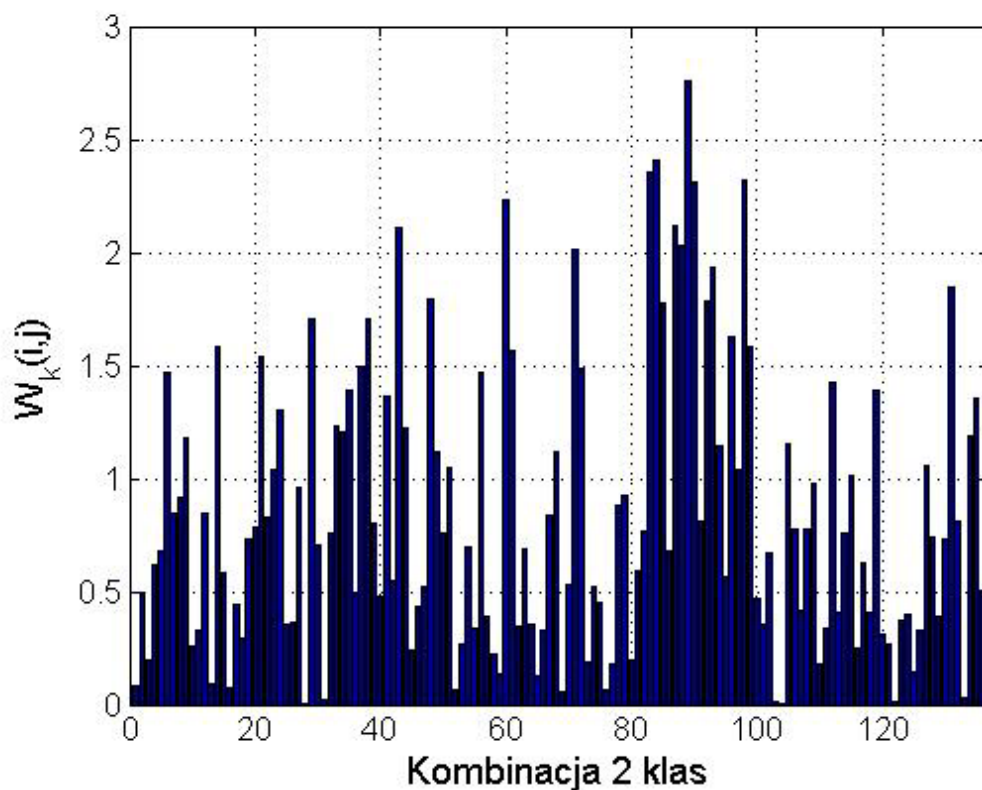
1	2	3	4	5	6	7	8	9	10	11	12
1.37E-2	3.53E-2	6.34E-3	2.16E-2	1.31E-2	5.44E-3	3.98E-2	2.13E-2	9.56E-3	8.98E-3	1.02E-2	8.81E-3
9.98E-3	6.64E-3	1.95E-2	6.25E-3	9.97E-3	5.54E-3	2.15E-2	1.40E-2	5.17E-3	7.86E-3	1.09E-2	4.80E-3
4.42E-3	3.85E-2	3.69E-3	1.99E-2	1.52E-2	3.36E-3	3.79E-2	1.27E-2	6.57E-3	9.14E-3	3.30E-3	7.17E-3
2.86E-3	2.44E-3	1.64E-2	1.45E-3	3.80E-3	1.41E-3	2.75E-2	8.75E-3	1.36E-3	2.88E-3	4.16E-3	1.58E-3
1.41E-2	1.53E-2	1.47E-2	1.45E-2	3.08E-2	1.31E-2	1.87E-2	1.68E-2	1.81E-2	1.00E-2	1.28E-2	5.22E-3
1.24E-2	9.53E-3	1.16E-2	1.19E-2	4.54E-2	1.81E-2	4.87E-3	6.83E-3	2.13E-2	1.70E-2	9.86E-3	5.10E-3
1.44E-2	1.41E-2	2.92E-2	1.44E-2	3.36E-2	1.94E-2	4.38E-2	2.17E-2	3.43E-2	1.68E-2	1.27E-2	1.15E-2
2.13E-2	2.44E-2	4.74E-2	1.99E-2	3.14E-2	2.18E-2	4.08E-2	2.12E-2	3.79E-2	1.93E-2	1.54E-2	1.18E-2
5.08E-3	8.47E-3	5.59E-3	7.36E-3	9.83E-3	1.10E-2	1.21E-2	7.75E-3	8.10E-3	8.16E-3	2.86E-2	1.18E-3
1.08E-2	1.13E-2	6.04E-3	1.27E-2	5.04E-3	1.11E-2	1.73E-2	1.31E-2	9.13E-3	2.74E-2	1.99E-2	7.02E-4
2.25E-2	2.05E-2	2.72E-2	4.47E-2	9.68E-3	1.41E-2	2.76E-2	2.36E-2	1.85E-2	1.56E-2	4.17E-3	2.18E-3
2.13E-2	1.33E-2	2.27E-2	8.58E-3	7.74E-3	6.92E-3	1.51E-2	1.40E-2	1.55E-2	5.51E-3	5.10E-3	1.55E-3
9.93E-3	2.05E-2	1.15E-2	2.50E-2	2.94E-2	2.18E-2	6.49E-3	7.59E-3	2.98E-2	2.12E-2	9.30E-3	5.90E-3
9.54E-3	1.55E-2	7.72E-3	6.62E-3	1.02E-2	8.43E-3	1.35E-2	1.28E-2	1.45E-2	1.23E-2	4.67E-3	3.94E-3
9.75E-3	2.01E-3	1.78E-2	4.22E-3	2.71E-2	2.74E-3	1.10E-2	5.09E-3	1.51E-2	1.84E-3	3.04E-2	6.88E-3
9.15E-3	3.64E-3	2.39E-2	7.42E-3	2.37E-2	5.32E-3	1.24E-2	5.92E-3	1.58E-2	2.50E-3	2.85E-2	8.76E-3
1.05E-3	8.41E-3	1.41E-2	7.31E-3	1.37E-2	2.28E-2	1.68E-2	1.44E-2	1.05E-2	3.59E-2	2.74E-3	6.14E-3
1.95E-3	1.36E-2	2.05E-2	5.35E-3	1.13E-2	1.42E-2	1.79E-2	1.44E-2	5.88E-3	1.72E-2	3.21E-3	1.19E-2
4.23E-3	1.00E-3	2.65E-2	2.09E-3	2.02E-2	2.03E-3	9.52E-3	2.29E-3	1.33E-2	1.63E-3	1.64E-2	9.92E-3
5.54E-3	2.59E-3	5.76E-2	4.30E-3	2.15E-2	7.30E-3	1.55E-2	2.63E-3	1.59E-2	3.10E-3	1.86E-2	1.77E-2

Obie miary: wariancja (a ściślej odchylenie standardowe) danych tworzących klasę oraz odległości między centrami dwu klas, powiązane są ze sobą miarą użyteczności cechy  $Wk(i,j)$  określoną wzorem (5.12). Wartość tej miary, niezależnie od wymienionych wyżej w tabelach 5.2 i 5.3, stanowić może następne kryterium selekcji cech.

Tabela 5.3 Odległości pomiędzy centrami klas dla wybranej cechy (nr 10) dla 12 typów komórek

	1	2	3	4	5	6	7	8	9	10	11	12
1	0.000	1.113	1.010	0.813	1.789	0.687	1.154	0.728	0.870	0.766	0.768	1.300
2	1.113	0.000	1.803	0.427	2.420	1.166	1.025	0.831	1.594	0.978	1.415	2.010
3	1.010	1.803	0.000	1.479	1.384	1.062	1.602	1.444	0.873	1.205	1.101	0.525
4	0.813	0.427	1.479	0.000	2.095	0.781	1.020	0.736	1.234	0.601	1.067	1.672
5	1.789	2.420	1.384	2.095	0.000	1.420	2.461	2.307	1.113	1.605	1.276	1.464
6	0.687	1.166	1.062	0.781	1.420	0.000	1.440	1.135	0.525	0.299	0.489	1.267
7	1.154	1.025	1.602	1.020	2.461	1.440	0.000	0.529	1.756	1.320	1.566	1.814
8	0.728	0.831	1.444	0.736	2.307	1.135	0.529	0.000	1.464	1.043	1.267	1.694
9	0.870	1.594	0.873	1.234	1.113	0.525	1.756	1.464	0.000	0.755	0.552	1.120
10	0.766	0.978	1.205	0.601	1.605	0.299	1.320	1.043	0.755	0.000	0.629	1.411
11	0.768	1.415	1.101	1.067	1.276	0.489	1.566	1.267	0.552	0.629	0.000	1.369
12	1.300	2.010	0.525	1.672	1.464	1.267	1.814	1.694	1.120	1.411	1.369	0.000

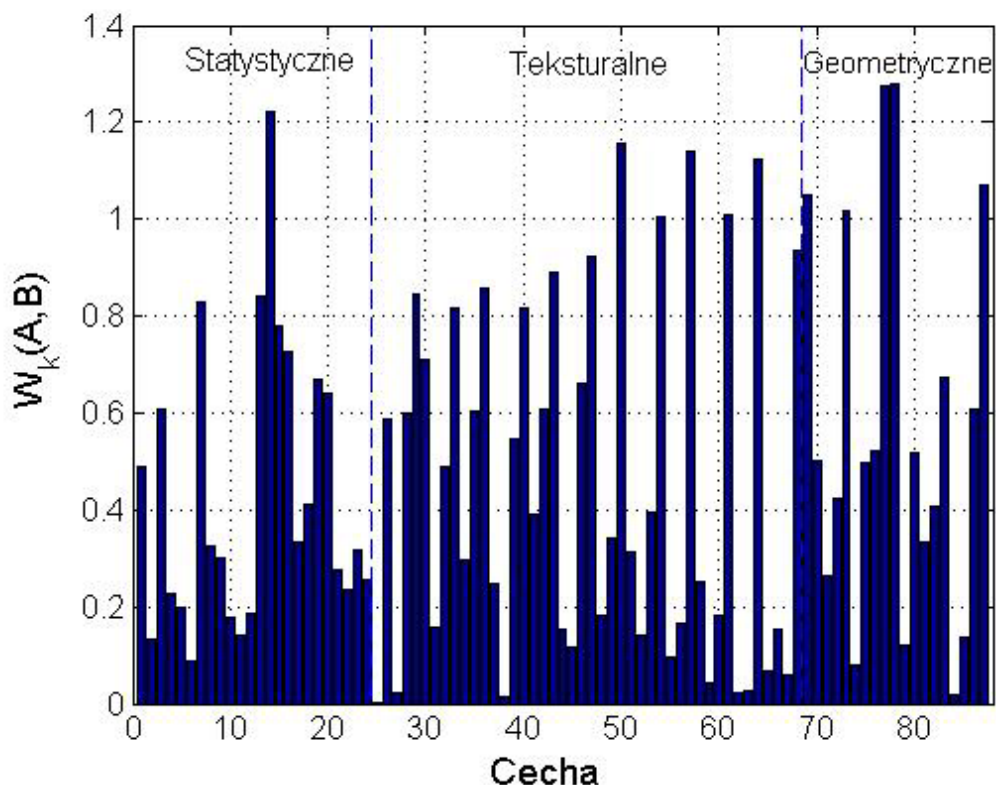
Ze względu na zastosowaną w pracy metodę klasyfikacji „jeden przeciw jednemu” przeprowadzono selekcję cech niezależnie dla każdej pary klas. Przykładowy rozkład miary  $W_k(i,j)$  określającej użyteczność jednej cechy (cecha nr 10) dla wybranych kombinacji par klas przedstawia rysunek 5.6.



Rys. 5.6 Wykres ilustrujący zmianę wartości miary  $W_k(i,j)$  cechy Nr 10 dla wybranych kombinacji par klas

Wartość miary  $W_k(i,j)$  cechy dziesiątej zmienia się od 0.04 do 2.8 zależnie od pary klas. Oznacza to, że ta sama cecha dla jednej pary klas może mieć istotny wpływ na jakość klasyfikacji (duża wartość  $W_k$ ), dla innych niewielki (mała wartość  $W_k$ ). Podobna analiza przeprowadzona dla wszystkich cech i par klas pozwala określić przydatność cechy dla rozpoznania klas. Dla każdej pary klas można przedstawić wykres miary  $W_k(i,j)$  dla kolejnych cech i na tej podstawie wyselekcjonować cechy najlepsze, o najwyższej wartości  $W_k(i,j)$ . Rysunek 5.7 przedstawia w sposób graficzny ocenę przydatności kolejnych cech dla rozpoznania 2 wybranych klas A i B. Jak widać występuje znaczna zmienność użyteczności cech nawet w ramach poszczególnych grup. Cechy o największych wartościach  $W_k(i,j)$  stanowią z punktu widzenia danej miary najlepsze kandydaty do wektora  $x$ . Jest oczywiste, że każda grupa cech powinna mieć swoich reprezentantów w wektorze cech  $x$ , stanowiących

wejście dla klasyfikatora neuronowego. Przykładowo, przy selekcji 35 najlepszych cech dla przedstawionej na rys. 5.7 pary klas A i B, reprezentację wektora  $x$  stanowiłyby 8 cech statystycznych, 19 cech teksturalnych i 8 geometrycznych, czyli w przybliżeniu od jednej trzeciej do połowy cech każdej grupy. Oznacza to, że każda grupa cech jest ważna w rozpoznawaniu komórek.

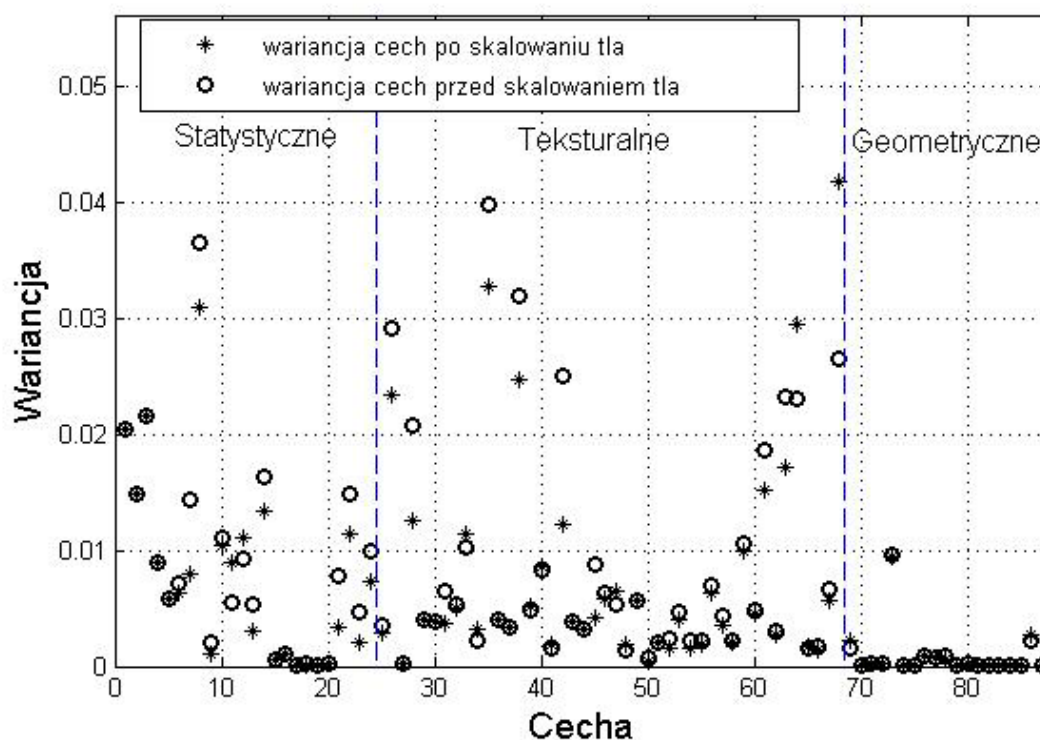


Rys. 5.7 Wartości miary  $W_k$  kolejnych cech dla klasyfikacji wybranej pary klas A i B

Interesująca jest ocena jakości cech skorelowanych ze sobą przy użyciu zastosowanej miary  $W_k(i,j)$ . Dla rozważanej pary klas A i B otrzymano następujące wyniki:

- dla cech statystycznych o najwyższej korelacji: (15 i 16), (19 i 20), (21 i 23) i (22 i 24) otrzymano bardzo zbliżone wartości miary  $W_k(A,B)$
- dla cech teksturalnych o najwyższej korelacji: (28, 35, 42), (29, 36, 43), (31 i 52), (32 i 39), (33 i 40), (44 i 66), (54, 61, 68) i (50, 57, 60) otrzymano również zbliżone wartości miary. Dla zbiorów cech (25 i 55) oraz (27, 34, 41) wartości miary  $W_k(A,B)$  znacznie różnią się między sobą.
- dla cech geometrycznych o najwyższej korelacji: (70 i 75), (76 i 80) oraz (77 i 78) otrzymano również zbliżone wartości  $W_k(A,B)$ . Wyjątek stanowiły cechy 71 i 84 (skorelowane ze sobą), dla których wartości miary  $W_k(A,B)$  znacznie się różniły.

Analiza otrzymanych wyników pokazuje, że w większości przypadków cechy skorelowane charakteryzują się podobnymi wartościami miary jakości  $W_k(i,j)$ . Nie można uznać tego za regułę, gdyż niektóre, skorelowane ze sobą cechy, mają indywidualnie różny wpływ na rozróżnienie dwu klas między sobą.



Rys. 5.8 Wpływ skalowania na wariancje cech komórek erytroblastów polichromatycznych

Kolejnym istotnym problemem w selekcji cech jest zbadanie wpływu skalowania tła obrazu (opisanego w punkcie 5.2) na rozkład wartości cech. Wraz ze zmianą tła zmianie ulega wariancja cech w ramach klas, najczęściej malejąc. Jest to korzystny efekt skalowania i dlatego do generacji cech użyto obrazów przeskalowanych. Rysunek 5.8 przedstawia wpływ skalowania na zmianę wartości wariancji 87 cech przykładowej klasy komórek (erytroblastów polichromatycznych).

W przypadku 27 cech wariancje uległy znacznemu zmniejszeniu, z czego 24 uzyskały wartości bliskie zeru, co dobrze świadczy o stabilności tych cech. Tylko w przypadku 4 cech wartości wariancji uległy zwiększeniu (zwykle nieznacznemu).

### 5.5.3. Selekcja cech przy użyciu sieci neuronowej SVM o jądrze liniowym

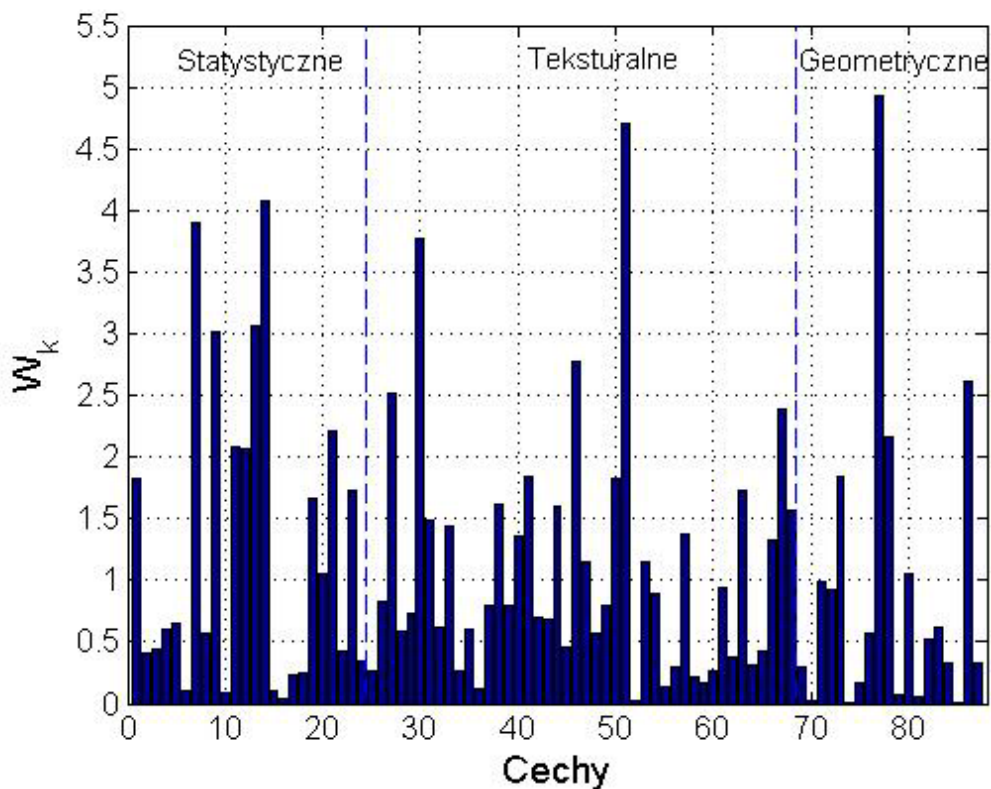
Selekcja cech metodami przedstawionymi w punktach poprzednich ma to ograniczenie, że nie uwzględnia wzajemnego wpływu cech działających równocześnie. Cecha nawet nie najlepsza indywidualnie, może zmienić miarę swojego oddziaływania na wynik

klasyfikacji przy współdziałaniu ze zbiorem innych cech. Stąd pożądanym jest badanie wpływu poszczególnych cech działających razem, gdyż tylko wówczas jest pewność, że ocena znaczenia danej cechy jest obiektywna.

W pracy wykorzystano do tego celu liniową sieć SVM. Idea zastosowania sieci SVM do badania wpływu poszczególnych składowych wektora  $\mathbf{x}$  na wynik klasyfikacji została zaczerpnięta z pracy [16,17], gdzie stosowano ją w analizie genetycznej. Dla sieci liniowej SVM wynik klasyfikacji można przedstawić w postaci:

$$y(\mathbf{x}) = \text{sign}[\mathbf{w}^T \mathbf{x} + b] = \text{sign}\left[\sum_{k=1}^N w_k x_k + b\right] \quad (5.13)$$

gdzie  $\mathbf{x}$  jest N-wymiarowym wektorem wejściowym,  $\mathbf{w}$  – macierzą wag a  $b$  – wagą polaryzacji. W wyniku uczenia liniowej sieci SVM następuje taka adaptacja wag, aby uzyskać możliwie najlepszy wynik klasyfikacji. Im wyższa jest wartość wagi  $w_k$  tym silniejszy wpływ składnika  $x_k$  wektora  $\mathbf{x}$  na działanie klasyfikatora. Jest to zatem naturalna selekcja cech, współdziałających ze sobą w procesie klasyfikacji.



Rys. 5.9 Wartości wag  $w_k$  kolejnych cech dla wybranej pary klas A i B

Tylko cechy wejściowe odpowiadające największym wartościom wag  $w_k$  są istotne w klasyfikacji, gdyż mają większy wpływ na wartość sygnału wyjściowego sieci, i tylko one powinny być wybierane do składu wektora  $\mathbf{x}$ . Oczywiście taką ocenę cech należy



przeprowadzić dla każdej pary klas oddzielnie. Rysunek 5.9 przedstawia wynik takiej oceny cech przy klasyfikacji dla tej samej pary klas, co na rysunku 5.7.

Jest widoczna istotna różnica znaczenia poszczególnych cech. O ile w poprzednim przypadku dominowały cechy teksturalne, tym razem oddziaływanie poszczególnych rodzajów cech jest bardziej zrównoważone. Ponadto dokonując analogicznej jak poprzednio analizy cech skorelowanych otrzymano:

- dla cech statystycznych o najwyższej korelacji: (19, 20) i (21, 23) otrzymano różniące się wartości natomiast dla cech (15, 16) i (22, 24) zbliżone.
- dla cech teksturalnych o najwyższej korelacji: (27, 34, 41), (29, 36, 43), (31, 52), (32, 39), (25, 52), (44, 66), (54, 61, 68) i (50, 57, 60) otrzymano różniące się wartości, tylko dla cech (33, 40) i (28, 35, 42) wartości były zbliżone.
- dla cech geometrycznych o najwyższej korelacji: (70, 75), (76, 80), (77, 78) i (71, 84) wartości te znacznie się różniły.

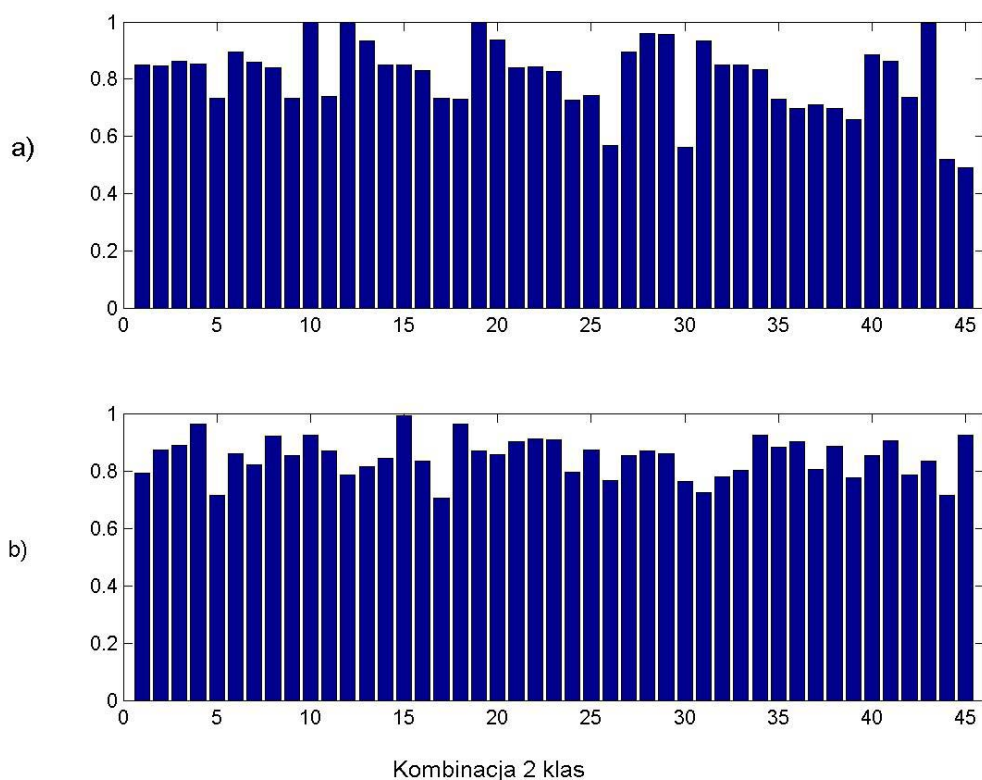
Jak pokazano prawie wszystkie cechy skorelowane charakteryzowały się inną wartością  $w_k$  co oznacza ich różny wpływ na wynik klasyfikacji. Jest to bardzo interesująca właściwość tego sposobu wyłaniania istotnych cech procesu.

W przypadku ograniczenia się do 35 najważniejszych cech do zbioru tworzącego wektora  $x$  weszłyby 11 cech statystycznych, 18 teksturalnych i 6 geometrycznych, ale w zupełnie innym zestawie niż przy selekcji indywidualnej. Z cech statystycznych pozostało tylko 5, poprzednio wybranych: 7, 13, 14, 19 i 20. Z cech teksturalnych swoją wartość potwierdziły cechy: 30, 46, 47, 50, 57, 61 i 68, w sumie 7, a z cech geometrycznych tylko cztery: 73, 77, 78 i 86. Oznacza to, że tylko 16 z 35 poprzednio wybranych indywidualnych cech zostało wyselekcjonowanych tą metodą do zbioru wspólnego.

Dla oceny wpływu korelacji cech przeprowadzono analizę korelacyjną wybranych cech w obydwu metodach. Współczynniki korelacji 10 najlepszych cech o najwyższych wskaźnikach jakości dla każdej pary klas obrazują wykresy na rys. 5.10.

Jak można zauważyć przy selekcji cech metodą pierwszą cztery pary cech mają wskaźnik korelacji bliski 1, natomiast przy selekcji z użyciem sieci SVM o jądrze liniowym, jest tylko jedna taka para. W drugim przypadku wyższa jest wartość średniej korelacji między cechami i wynosi 0.8475 wobec 0.8096 w pierwszym przypadku. Analiza wyników w obu przypadkach dowodzi, że generalnie korelacja między cechami jest stosunkowo wysoka i nie ona decyduje o ich jakości. Co więcej zdarzają się nawet cechy bardzo wysoko skorelowane, które współdziałając z innymi, zapewniają najlepsze działanie klasyfikatora. Można to

tłumaczyć bardzo wysoką nieseparowalnością poszczególnych klas mającą bezpośredni wpływ na wskaźnik korelacji pomiędzy cechami.



Rys. 5.10 Współczynniki korelacji 10 najlepszych cech dla pary klas przy selekcji na podstawie wartości średnich i wariancji danych a) i przy użyciu sieci neuronowej SVM o jądrze liniowym b).

Praktycznym sprawdzianem jakości wyselekcjonowanych cech jest porównanie wyników klasyfikacji komórek z uwzględnieniem wszystkich cech z odpowiadającymi im wynikami przy zastosowaniu ograniczonej liczby cech. Do porównania wskazana jest analiza wyników klasyfikacji przy zastosowaniu wszystkich cech, ograniczonej liczby najlepszych oraz ograniczonej liczby najgorszych cech wg każdej z metod selekcji. Przykładowe wyniki porównawcze błędów klasyfikacji dla 16 typów komórek zamieszczono w tabeli 5.4.

W przypadku selekcji z zastosowaniem liniowej sieci SVM wynik dla 30 najlepszych cech był znacznie lepszy od wyniku uzyskanego z wykorzystaniem wszystkich 87 cech, natomiast znacznie gorszy przy zastosowaniu 30 cech najgorszych. Jest to zgodne z oczekiwaniami i potwierdza prawidłowość metody selekcji.

W przypadku klasyfikatora wykorzystującego ocenę indywidualną cech opartą na wartości średniej i odchyleniu standardowym (miara  $W_k(i,j)$ ) wszystkie trzy wyniki były do

Tabela 5.4 Wyniki klasyfikacji komórek należących do 16 klas przy użyciu 87 cech i ograniczonych zestawów cech wyselekcjonowanych za pomocą obydwu metod selekcji

Błąd % klasyfikacji	Wszystkie 87 cech	30 najlepszych cech	30 najgorszych cech
Selekcja cech na podstawie wartości średnich i wariancji danych	21.13 %	24.54 %	27.36 %
Selekcja cech przy użyciu sieci neuronowej SVM o jądrze liniowym	21.13 %	18.71 %	33.60 %

siebie zbliżone. Pomimo wyselekcjonowania cech najlepszych w sensie indywidualnym nie udało się poprawić wyników klasyfikacji, a wynik ostateczny był nawet gorszy. Oznacza to, że jakość cech ulega istotnej zmianie przy współdziałaniu z innymi. Ich ocena ma więc sens jedynie przy uwzględnianiu współdziałania wszystkich cech na raz. Jest to właściwość bardzo charakterystyczna dla komórek rakowych, która niekoniecznie musi się potwierdzić w innych zadaniach klasyfikatorów, np. klasyfikacji genów [17], gdzie jedną z uznanych metod oceny jakości cech jest analiza ich wariancji i wartości średnich.

## 6. Wyniki klasyfikacji

Układy rozpoznawania i klasyfikacji komórek uzyskanych z rozmazu szpiku kostnego poddane zostały wielu testom, mającym na celu ustalenie ich skuteczności przy różnej liczbie typów komórek, zmieniającej się liczbie pacjentów i zróżnicowanych warunkach wytwarzania rozmazu. Należy podkreślić, że pozyskiwane obrazy komórek pochodziły z przestrzeni kilku lat. Były wytwarzane przez różnych laborantów Instytutu Hematologii, przy zastosowaniu odczynników pochodzących z różnych serii produkcyjnych. To powodowało duże zróżnicowanie barwienia preparatów mające znaczny wpływ na końcowy wynik przetwarzania i klasyfikacji. Choroby zdiagnozowane u pacjentów były bardzo różnorodne, poczynając od białaczek ostrych poprzez przewlekłe, aż do zaburzeń rzadziej występujących, jak np. plazmocytoma czy choroba Addisona-Bermyera. Wskutek tego występowała duża różnorodność mielogramów poszczególnych preparatów jak i liczebności pacjentów, od których pochodziły poszczególne typy komórek. Komórki występujące prawidłowo w hematopoezie (np. układu czerwokrwińkowego i białokrwińkowego obojętnochłonne) występowały w większej liczbie preparatów, natomiast komórki normalnie nie występujące lub występujące w bardzo małym procencie, pochodziły od ograniczonej liczby pacjentów. W badaniach klasyfikacyjnych wykorzystano tylko te typy komórek, których znacząca liczebność pozwalała prawidłowo nauczyć sieć neuronową i oczekiwać właściwej generalizacji klasyfikatora. Typy komórek o bardzo małej liczebności (poniżej 20) nie były brane pod uwagę. Dla zrównoważenia różnych typów komórek i zapobieżenia dominacji komórek najbardziej licznych pewne typy komórek występujące najczęściej były częściowo eliminowane ze zbioru uczącego. Również w testowaniu zastosowano górny pułap liczbowy komórek, aby uzyskać wiarygodne wyniki procentowego udziału błędów (komórki o dużej liczebności pozwalały zwykle na uzyskanie najlepszych wyników klasyfikacji).

Zadanie automatycznej klasyfikacji przeprowadzono dla trzech różnych zestawów komórek, zawierających odpowiednio dwanaście, siedemnaście i dwadzieścia jeden typów komórek włączając w to wspólną klasę cieni i krwinek czerwonych. W każdym przypadku zastosowano technikę tzw. "cross-validation". Dane były losowo dzielone na 5 równych (w przybliżeniu) części i przeprowadzono klasyfikacje dla wszystkich możliwych kombinacji, przy czterech częściach tworzących zbiór danych uczących i jednej tworzącej zbiór testujący.

W procesie uczenia parametry klasyfikatora były dobierane niezależnie od danych testujących. Zbiór danych uczących był dzielony na 4 części, klasyfikator był uczony czterokrotnie na trzech z nich a czwarta służyła do weryfikacji parametrów uczących. Przeprowadzając powyższe czynności dla każdej możliwej kombinacji hiperparametrów  $C$  i  $\gamma$ , z listy możliwych wartości i stosując kryterium wagowe błędu  $E=0.9*N_{\text{wer}} + 0.1*N_{\text{SV}}$ , (gdzie  $N_{\text{wer}}$  oznacza liczbę błędów dla danych weryfikujących a  $N_{\text{SV}}$  – liczbę wektorów podtrzymujących), wybierano optymalne wartości  $C$  i  $\gamma$  klasyfikatora dla każdej pary klas niezależnie. Dopiero dla tak znalezionych parametrów przeprowadzono uczenie na całym zbiorze danych uczących i właściwe testowanie na danych testujących.

### 6.1 Wyniki rozpoznania 12 rodzajów komórek

Dane poddane klasyfikacji pochodziły z preparatów od 14 pacjentów, dla których wykonano średnio po 32 zdjęcia fragmentów obrazów rozmazu szpiku kostnego. Liczebności komórek należących do poszczególnych klas podane są w kolumnie trzeciej tabeli 6.1. Proces uczenia sieci SVM o jądrze radialnym przeplatał się z doбором optymalnego zestawu cech diagnostycznych, tworzących wektor  $x$ . Uczenie rozpoczęto z pełnym zestawem cech (87), po którym dokonano określenia błędów zarówno uczących jak i testujących. Następnie zestaw cech został poddany redukcji poprzez wyeliminowanie pewnej liczby cech najmniej znaczących. Uszeregowanie cech odbywało się zgodnie z metodyką liniowej sieci SVM, omówioną w rozdziale 5. Aby ustalić optymalną liczbę cech przeprowadzono wiele prób uczenia sieci z różną liczbą cech, porównując za każdym razem błędy uczenia i weryfikacji.

Tabela 6.1 Wynik rozpoznania 12 rodzajów komórek

Klasa	Typ komórki	Liczba komórek	Średnia liczba błędów dla danych uczących	Błąd procentowy dla danych uczących	Liczba błędów dla danych testujących	Błąd procentowy dla danych testujących
1	erytroblast zasadochłonny	96	5.5	5.73%	9	9.38%
2	erytroblast polichromatyczny	138	8.5	6.16%	15	10.87%
3	erytroblast kwasochłonny	74	12	16.22%	19	25.68%
4	blast	130	0.75	0.58%	3	2.31%
5	promielocyt	23	10	43.48%	11	47.83%
6	mielocyt	55	4.25	8.02%	11	20.75%
7	metamielocyt	30	7.5	25%	10	33.33%
8	granulocyt pałeczkowaty	22	9.25	42.05%	12	54.55%
9	granulocyt segmentowany	38	0.75	1.97%	4	10.53%
10	prolimfocyt	21	7.75	36.90%	10	47.62%
11	limfocyt	148	1.5	1.01%	5	3.38%
12	plazmocyt	23	6.25	27.17%	4	17.39%
Razem		796	75.5	9.49%	113	14.20%

Za optymalną liczbę cech uznano tę, która gwarantowała najmniejszy błąd na danych weryfikujących (30). Dopiero dla takiej sieci przeprowadzono pełny proces uczenia i testowania na danych testujących (nie uczestniczących w uczeniu i weryfikacji). Wyniki w formie błędów rozpoznania danych uczących i testujących przedstawione są w tabeli 6.1. Są to wartości uśrednione, wynikające z przyjętej w badaniach strategii "cross-validation". Stąd wynikają wartości nie całkowito-liczbowe błędów dla danych uczących.

Błąd klasyfikacji dla danych uczących wyniósł średnio 9.49% natomiast dla danych testujących 14.2%, co sugeruje stosunkowo dobrą generalizację klasyfikatora. Zwraca uwagę fakt, że klasy nielicznie reprezentowane, np. 5, 7, 8 i 10 uzyskały znacznie gorszy wynik od pozostałych. Można przypuszczać, że ta liczba danych nie była wystarczająca do ich prawidłowej reprezentacji w procesie uczenia. W tabeli 6.2 i 6.3 zamieszczono pełne wyniki rozkładu błędów dla danych uczących i testujących. W kolejnych wierszach tabel podano aktualne wyniki klasyfikacji danych w ramach jednej klasy. Każdy niezerowy element macierzy występujący poza diagonalną oznacza błąd, czyli przypisanie komórki do niewłaściwej klasy. Element  $a_{i,j}$  macierzy oznacza przypisanie i-tego typu komórki do j-tej klasy.

Tabela 6.2 Macierz błędów klasyfikacji dla danych uczących (wartość średnia)

Klasa	1	2	3	4	5	6	7	8	9	10	11	12
1	90.5	5.25										0.25
2	6.25	129.5	2								0.25	
3	1	10.25	62								0.75	
4				129.25			0.25				0.5	
5	1			6.25	13	2.75						
6	1.5			0.25		48.75	0.25	0.25	1		1	
7	0.25			1	0.5	5	22.5	0.75				
8	1.5			0.25			1.75	12.75	5.75			
9						0.5	0.25		37.5			
10				0.25						13.25	7.5	
11	0.25	1	0.25							1.5	145	
12	3.5	0.25	0.25		0.75	0.25					1.25	16.75

W tabelach kolorami oznaczono klasy w ramach poszczególnych układów krwiotwórczych występujące bezpośrednio po sobie w procesie dojrzewania. Błędy powstałe na styku takich komórek mogą zdarzyć się nawet najbardziej doświadczonemu ekspertowi, gdyż bardzo trudno jest przyporządkować jednoznacznie komórki na etapie przejścia z jednego stadium do drugiego.

Tabela 6.3 Macierz błędów klasyfikacji dla danych testujących

Klasa	1	2	3	4	5	6	7	8	9	10	11	12
1	87	9										
2	8	123	6								1	
3	1	16	55								2	
4				127	1		1				1	
5	1			6	12	3					1	
6	1			5	1	42	2	1	1			
7		1		2	1	4	20	1	1			
8	3	1				1	2	10	5			
9	1						2	1	34			
10				1						11	9	
11		1		1					1	2	143	
12				1				1			2	19

Błędów takich nie można traktować na równi z pozostałymi. W związku z tym w tabeli 6.4 przedstawiono wynik klasyfikacji przy pominięciu błędów rozpoznawania pomiędzy sąsiednimi komórkami w tej samej linii ich dojrzwania.

Tabela 6.4 Wyniki klasyfikacji dla 12 klas przy pominięciu błędów pomiędzy kolejnymi fazami rozwoju komórek

Klasa	Średnia liczba błędów dla danych uczących	Błąd procentowy dla danych uczących	Liczba błędów dla danych testujących	Błąd procentowy dla danych testujących
1	0.25	0.26%	0	0
2	0.25	0.18%	1	0.72%
3	0.75	1.01%	1	1.35%
4	0.75	0.58%	2	1.54%
5	1	4.35%	2	8.70%
6	4	7.55%	8	15.09%
7	1.75	5.83%	5	16.67%
8	1.75	7.95%	5	22.73%
9	0.75	1.97%	3	7.89%
10	0.25	1.19%	1	4.76%
11	1.5	1.01%	3	2.03%
12	6.25	27.17%	4	17.39%
Razem	19.25	2.42%	35	4.40%

Widoczne jest znaczące, bo nawet trzykrotne zmniejszenie błędu rozpoznawania poszczególnych typów komórek. Błąd klasyfikacji dla danych testujących (nie uczestniczących w uczeniu) przy pominięciu błędów pomiędzy kolejnymi typami komórek

wyniósł zaledwie 4.4 %, co jest bardzo dobrym wynikiem, całkowicie akceptowalnym w praktyce medycznej.

## 6.2 Wyniki rozpoznania 17 rodzajów komórek

Zwiększenie liczby klas komórek poddanych rozpoznaniu utrudnia proces klasyfikacji. Wynika to bezpośrednio z wprowadzenia nowych klas komórek, dla których dane nie tworzą zbiorów idealnie odseparowanych od już istniejących. Innym powodem jest zwiększenie liczby pacjentów, dla których komórki mogą mieć bardziej zróżnicowany wygląd, a odpowiadające im zbiory parametrów większą wariancję.

Przy 17 rodzajach komórek dane poddane klasyfikacji pochodziły z preparatów od 19 pacjentów. Dla nowych 5 pacjentów wykonano średnio po 230 zdjęć fragmentów obrazów rozmazu szpiku kostnego. Liczebności poszczególnych typów komórek podane zostały w kolumnie trzeciej tabeli 6.5. Podobnie jak poprzednio, zastosowano strategię "cross-validation" dzieląc zbiór danych na 5 części i używając 4 zbiorów do uczenia, a jednego do testowania. Wyniki w postaci średniej liczby błędów i błędu procentowego na danych uczących i testujących przedstawiono w tabeli 6.5. Przy 17 rodzajach komórek średni błąd uczenia był równy 13.45%, a błąd testowania 18.71%.

Tabela 6.5 Wyniki klasyfikacji dla 17 klas

Klasa	Typ komórki	Liczba komórek	Średnia liczba błędów dla danych uczących	Błąd procentowy dla danych uczących	Liczba błędów dla danych testujących	Błąd procentowy dla danych testujących
1	postać podziałowa erytropoezy	15	.25	1.7%	1	6.7%
2	proerytroblast	25	9.5	38.0%	9	36.0%
3	erytroblast zasadochłonny	138	25	18.1%	32	23.2%
4	erytroblast polichromatyczny	403	24.25	6.0%	54	13.4%
5	erytroblast kwasochłonny	238	14.5	6.1%	34	14.3%
6	promegaloblast	23	2.25	9.8%	5	21.7%
7	blast	153	14	9.2%	18	11.8%
8	promielocyt	99	25.25	25.5%	34	34.3%
9	mielocyt	195	26.25	13.5%	41	21.0%
10	metamielocyt	172	41.25	24.0%	58	33.7%
11	granulocyt pałeczkowaty	165	71.5	43.3%	77	46.7%
12	granulocyt segmentowany	250	23.75	9.5%	28	11.2%
13	granulocyt segmentowany kwasochłonny	39	3.5	9.0%	6	15.4%
14	prolimfocyt	19	8.5	44.7%	12	63.4%
15	limfocyt	295	13.75	4.7%	19	6.4%
16	plazmocyt	40	8.25	20.7%	6	15.00%
17	cienie i krwinki czerwone	334	30.5	9.1%	52	15.6%
Razem		2603	342.25	13.15%	487	18.71%



Niektóre klasy nielicznie reprezentowane, np. 2 i 14 zostały rozpoznane ze znacznie większym błędem niż pozostałe. Może to świadczyć o ciągle niewystarczającej reprezentacji tych komórek w procesie uczenia. Należy jednak zauważyć, że klasy 8, 10 i 11 pomimo znaczącej liczebności nie uzyskały zadawalającego wyniku klasyfikacji. Tym razem głównym powodem jest ich ogromne podobieństwo do sąsiadów z tej samej linii rozwojowej. W tabeli 6.6 zamieszczono rozkłady błędów klasyfikacji dla danych testujących. Podobnie jak poprzednio, wiersze macierzy oznaczają wyniki klasyfikacji danych w ramach jednej klasy, czyli aktualne przypisanie komórek danej klasy przez układ klasyfikatora.

Tabela 6.6 Macierz błędów klasyfikacji dla danych testujących

Klasa	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
1	14	1															
2		16	3				2	2	1								
3		1	106	28				1	1								1
4			15	349	36										3		
5				33	204							1					
6			1	2	2	18											
7	1				1		135	3	2	4					5		2
8							9	65	20	1	1						3
9		1	2	1			2	8	154	17	3	3			1		3
10			1				2		33	114	10	8			1		3
11							1		4	21	88	43			1		7
12			1						1	4	19	222					3
13									1		1	1	33		1		2
14							1	1						6	11		
15				5	1		1		3	1				1	276		7
16			1			1			3						1	34	
17					1		4	3	14	2	3	4	4		17		282

Jest widoczne, że większość spośród znaczących błędów klas 8, 10, 11 i 14 dotyczyła komórek sąsiadujących ze sobą na etapie rozwoju. Trudności w rozpoznaniu sąsiednich typów komórek potęguje zmienność cech wynikająca z różnic w barwieniu preparatów lub / i ustawień mikroskopu podczas badania. Z uwagi na znaczne zwiększenie liczby preparatów z których pochodziły obrazy komórek zwiększyła się również w sposób naturalny zmienność ich cech.

W tabeli 6.7 przedstawiono wyniki klasyfikacji dla danych testujących przy pominięciu błędów rozpoznawania pomiędzy sąsiednimi klasami w rozwoju komórek. Średni błąd rozpoznania jest równy 6.76%, co można uznać za wynik zadawalający.

Tabela 6.7 Wyniki klasyfikacji dla 17 klas przy pominięciu błędów pomiędzy kolejnymi fazami rozwoju komórek

Klasa	Liczba błędów dla danych testujących	Błąd procentowy dla danych testujących
1	1	6.7%
2	6	24%
3	3	2.2%
4	3	0.7%
5	1	0.4%
6	5	21.7%
7	15	9.8%
8	5	5.1%
9	16	8.2%
10	15	8.7%
11	13	7.9%
12	9	3.6%
13	6	15.4%
14	2	10.5%
15	18	6.1%
16	6	15.0%
17	52	15.6%
Razem	176	6.76%

Jak można zauważyć tylko dla klasy 2 i 6 błąd rozpoznania przewyższa 20%. Obie te klasy miały stosunkowo nieliczną reprezentację w zbiorze danych (odpowiednio 25 i 23 komórki). Błędy klasyfikacji pozostałych rodzajów komórek są na akceptowalnym poziomie, z reguły nie przekraczając 10%. Porównując wyniki przedstawione w tabeli 6.5 i 6.7 można zauważyć, że praktycznie tylko 1/3 błędów (6.76% w stosunku do 18.71%) jest znacząca, gdyż dotyczy komórek istotnie różniących się od siebie.

### 6.3 Wyniki rozpoznania 21 rodzajów komórek

Dane poddane klasyfikacji pochodziły tym razem z preparatów od 26 pacjentów. Dla nowych 7 pacjentów wykonano po około 200 zdjęć fragmentów obrazów rozmazu szpiku kostnego, które dodano do istniejącej bazy danych 19 pacjentów. W uczeniu zastosowano 2666 wektorów uczących  $x_i$  i stowarzyszonych z nimi wartości zadanych  $d_i$ . Zastosowano rozszerzoną reprezentację cech, dołączając do zbioru standardowego 87-wymiarowego cechy morfologiczne i teksturalne dotyczące zredukowanej rozdzielczości (patrz tabela 5.1). Dla uzyskania najlepszych rezultatów rozpoznania konieczna była redukcja wymiaru wektora  $x$ , przeprowadzona przy użyciu selekcji cech opartej na liniowej sieci SVM. Po wykonaniu

wielu prób wstępnych za optymalną uznano liczbę 35 cech wyselekcjonowanych przez sieć liniową SVM. Tak utworzone wektory wejściowe  $x_i$  posłużyły do uczenia sieci SVM o radialnej funkcji jądra. Sieć wytrenowaną poddano testowaniu na danych nie uczestniczących w uczeniu z zastosowaniem strategii "cross-validation". W tabeli 6.8 przedstawiono podstawowe dane zbiorcze dotyczące tego zadania w trybie testowania. Są to: aktualne typy komórek poddanych rozpoznaniu, liczebność poszczególnych klas, liczby błędnych klasyfikacji oraz procentowy błąd rozpoznania.

Tabela 6.8 Wyniki klasyfikacji dla 21 klas

Klasa	Typ komórki	Liczba komórek	Liczba błędów dla danych testujących	Błąd procentowy dla danych testujących	Liczba błędów przy pominięciu błędów między sąsiadami	Błąd procentowy przy pominięciu błędów między sąsiadami
1	postać podziałowa erytropoezy	23	6	26.1%	6	26.1%
2	proerytroblast	39	7	18.0%	6	15.4%
3	erytroblast zasadochłonny	139	38	27.3%	6	4.3%
4	erytroblast polichromatyczny	350	54	15.4%	5	1.4%
5	erytroblast kwasochłonny	227	31	13.7%	0	0%
6	promegaloblast	30	4	13.3%	2	6.7%
7	megaloblast zasadochłonny	64	10	15.6%	2	3.1%
8	megaloblast polichromatyczny	46	18	39.1%	4	8.7%
9	blast	115	16	13.9%	12	10.4%
10	promielocyt	94	31	33.0%	8	8.5%
11	mielocyt	155	41	26.5%	21	13.5%
12	metamielocyt	128	54	42.2%	12	9.4%
13	granulocyt pałeczkowaty	140	61	43.6%	7	5.0%
14	granulocyt segmentowany	225	23	10.2%	2	0.9%
15	eozynofile	127	20	15.8%	20	15.8%
16	monocyt	55	24	43.6%	24	43.6%
17	prolimfocyt	45	12	26.7%	1	2.2%
18	limfocyt	281	17	6.1%	14	5.0%
19	proplazmocyt	84	16	19.1%	1	1.2%
20	plazmocyt	93	24	25.8%	11	11.8%
21	cienie i krwinki czerwone	209	19	9.2%	19	9.2%
Razem		2666	526	19.73%	183	6.86%

Wyniki dotyczą wyłącznie danych testujących nie uczestniczących w uczeniu. Ostatnie dwie kolumny tabeli pokazują wyniki przy pominięciu błędów wynikających z sąsiedztwa dwu komórek tej samej linii rozwojowej. Średni błąd rozpoznania danych uczących był równy 13.1%, a więc porównywalny z błędem rozpoznania 17 typów komórek. Błąd rozpoznania komórek dla danych testujących (19.73%) tylko nieznacznie odbiega od wyniku uzyskanego

dla 17 typów (18.71%). Szczegółowa analiza wyników wykazała, że prawie dwie trzecie błędów dotyczyło komórek należących do sąsiednich faz rozwoju. Szczegółowe wyniki dla danych testujących pokazano w tabeli 6.9, w formie macierzy przypisań komórek do odpowiednich klas. Przy pominięciu błędów odpowiadających sąsiednim komórkom niezadawalające rezultaty klasyfikacji otrzymano dla jednej klasy (monocyty). Monocyty ze względu na wysokie podobieństwo z mielocytami zostały z nimi pomyłone aż 12-krotnie. Jest to znany problem również dla eksperta ludzkiego. Zwykle dla rozwiązania tego problemu stosuje się inne barwienie drugiego preparatu wykonanego dla tego samego pacjenta.

Tabela 6.9 Macierz błędów klasyfikacji dla danych testujących

Klasa	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
1	17	1	0	2	0	0	0	0	0	2	0	1	0	0	0	0	0	0	0	0	0
2	0	32	1	0	0	2	0	0	2	1	0	0	0	0	0	0	0	0	0	1	0
3	0	1	101	31	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	3	0
4	0	0	20	296	29	0	0	0	0	0	0	0	2	0	0	0	0	2	0	1	0
5	0	0	0	31	196	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	1	0	0	0	26	2	0	0	0	0	0	0	0	0	0	0	0	0	1	0
7	0	1	0	1	0	2	54	6	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	14	28	0	0	2	0	0	0	0	0	0	0	1	1	0
9	0	0	0	0	1	1	0	0	99	4	3	1	0	0	1	0	3	0	0	0	2
10	0	1	0	0	1	0	0	0	6	63	17	1	0	0	1	0	0	0	0	2	1
11	0	0	1	1	0	0	0	0	2	6	114	14	6	1	0	4	0	2	0	1	2
12	0	0	0	0	0	0	0	0	0	1	25	74	17	7	0	2	0	0	0	0	2
13	0	0	0	0	0	0	0	0	0	1	4	17	79	37	0	0	0	1	0	0	1
14	0	0	0	0	0	0	0	0	0	0	1	1	21	202	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	1	0	0	1	17	0	0	0	0	0	2
16	0	0	1	0	0	0	0	0	0	1	12	4	0	0	0	31	1	1	0	4	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	33	11	0	0	1
18	0	0	2	0	0	0	0	0	0	0	2	0	0	0	0	0	3	264	1	1	8
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	68	15	0
20	1	0	2	0	0	0	0	0	0	0	5	0	0	0	0	0	0	3	13	69	0
21	0	0	0	0	1	0	0	0	0	0	6	1	0	1	0	1	0	6	0	0	187

Wyniki uzyskane dla 21 klas komórek dają przegląd skali trudności zadania automatycznej klasyfikacji komórek szpiku kostnego. Przy tej liczbie komórek występują typowe problemy rozpoznania i klasyfikacji, wynikające z częściowego pokrywania się przestrzeni parametrów (braku separacji klas), dużej różnorodności danych będącej wynikiem dużej liczebności pacjentów, jak również problemów obliczeniowych dla algorytmów uczenia będących

rezultatem ograniczonej liczby danych uczących (przy 21 klasach liczba danych przekroczyła 2500, przy dużym (od 30 do 137) rozmiarze wektorów wejściowych  $x$ ).

#### 6.4 Weryfikacja systemu na podstawie mielogramów wybranych pacjentów

Badania systemu rozpoznania i klasyfikacji komórek przedstawione dotąd miały na celu wszechstronne przetestowanie dokładności działania układu. Stąd duże zróżnicowanie rodzajów komórek oraz zastosowana metodologia "cross-validation" uważana na świecie za najbardziej obiektywny sposób oceny dokładności metody. W praktyce szpitalnej pożądany jest nieco inny sposób oceny, polegający na porównaniu mielogramów sporządzonych przez eksperta ludzkiego i układ automatycznej klasyfikacji dla każdego pacjenta oddzielnie. Mielogram stanowi ocenę składu szpiku kostnego, niezbędną do postawienia właściwej diagnozy rozpoznania choroby pacjenta przyjmowanego do szpitala. W ramach tej oceny najistotniejszymi informacjami otrzymywanymi z rozmazu szpiku kostnego są:

- udział poszczególnych układów krwiotwórczych w szpiku
- zachowanie bądź brak dojrzewania komórek w ramach poszczególnych linii
- komórkowatość szpiku

Komórkowatość szpiku jest oceniana w niewielkim powiększeniu (np. 40 $\times$ ), na podstawie stosunku liczby komórek do powierzchni rozmazu i jest tym mniejsza im większe są plamy tłuszczowe. Na podstawie sporządzonego mielogramu ocenia się udział linii krwiotwórczych i prawidłowość ich dojrzewania. W tabelach 6.10 – 6.12 zamieszczono przykładowe wyniki klasyfikacji komórek szpiku kostnego dla trzech różnych pacjentów w formie sporządzanej w laboratorium medycznym. Kolumny trzecia i czwarta przedstawiają wyniki podane przez eksperta ludzkiego, natomiast piąta i szósta – wyniki uzyskane przez opracowany układ automatycznej klasyfikacji. Ostatnia kolumna pokazuje różnicę (w punktach procentowych) między wskazaniami eksperta ludzkiego i układu automatycznego.

Tabela 6.10 Porównanie rzeczywistego składu mielogramu pierwszego pacjenta i składu otrzymanego w procesie automatycznej klasyfikacji

Oznaczenie układu / klasa	Typ komórki	Rzeczywista liczba komórek (wg eksperta)	Udział procentowy (wg eksperta)	Liczba komórek (wg systemu autora)	Udział procentowy (wg systemu autora)	Błąd (w punktach procentowych)
A1	<b>ERYTROPOEZA NORMOBLASTYCZNA</b>	173	49.71%	172	48.04%	1.67%
1	postać podziałowa erytropoezy	3	0.86%	3	0.83%	0.03%
2	proerytroblast	7	2.01%	7	1.95%	0.06%
3	erytroblast zasadochłonny	25	7.18%	24	6.70%	0.48%
4	erytroblast polichromatyczny	86	24.71%	86	24.02%	0.69%
5	erytroblast kwasochłonny	52	14.94%	52	14.52%	0.42%
A2	<b>ERYTROPOEZA MEGALOBLASTYCZNA</b>	0	0%	0	0%	0%
6	promegaloblast	0	0%	0	0%	0%
7	megaloblast zasadochłonny	0	0%	0	0%	0%
8	megaloblast polichromatyczny	0	0%	0	0%	0%
B	<b>UKŁAD BIAŁOKRWINKOWY</b>	148	42.53%	155	43.30%	0.77%
9	blast	3	0.86%	2	0.55%	0.31%
10	promielocyt	19	5.45%	19	5.30%	0.15%
11	mielocyt	37	10.63%	48	13.40%	2.77%
12	metamielocyt	20	5.74%	13	3.63%	2.11%
13	granulocyt pałeczkowaty	31	8.90%	28	7.82%	1.08%
14	granulocyt segmentowany	26	7.47%	33	9.21%	2.26%
15	eozynofile	12	3.44%	12	3.35%	0.09%
16	monocyt	10	2.87%	9	2.51%	0.36%
C	<b>UTKANIE CHŁONNE</b>	17	4.89%	18	6.15%	1.26%
17	prolimfocyt	0	0%	0	0%	0%
18	limfocyt	16	4.59%	20	5.58%	0.99%
19	proplazmocyt	0	0%	0	0%	0%
20	plazmocyt	1	0.28%	2	0.55%	0.27%
	<b>POZA MIELOGRAMEM</b>					
21	cienie i krwinki czerwone	72		62		
Razem	suma rozpoznanych komórek mielogramu	348		358		

Tabela 6.11 Porównanie rzeczywistego składu mielogramu drugiego pacjenta i otrzymanego w procesie automatycznej klasyfikacji

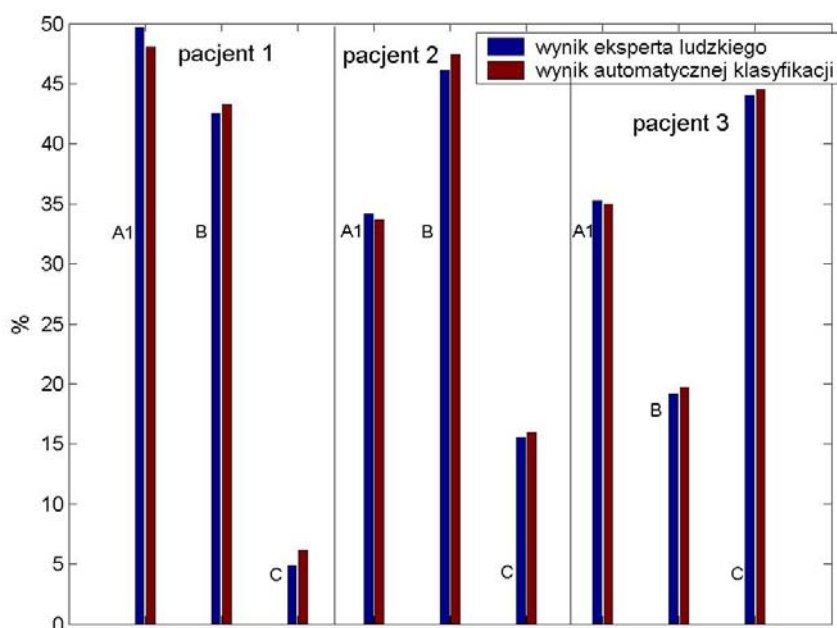
Oznaczenie układu / klasa	Typ komórki	Rzeczywista liczba komórek (wg eksperta)	Udział procentowy (wg eksperta)	Liczba komórek (wg systemu autora)	Udział procentowy (wg systemu autora)	Błąd (w punktach procentowych)
A1	<b>ERYTROPOEZA NORMOBLASTYCZNA</b>	57	34.13%	59	33.71%	0.43%
1	postać podziałowa erytropoezy	1	0.59%	1	0.57%	0.02%
2	proerytroblast	0	0%	0	0%	0%
3	erytroblast zasadochłonny	5	2.99%	6	3.42%	0.43%
4	erytroblast polichromatyczny	25	14.97%	25	14.28%	0.69%
5	erytroblast kwasochłonny	26	15.56%	27	15.42%	0.14%
A2	<b>ERYTROPOEZA MEGALOBLASTYCZNA</b>	0	0%	0	0%	0%
6	promegaloblast	0	0%	0	0%	0%
7	megaloblast zasadochłonny	0	0%	0	0%	0%
8	megaloblast polichromatyczny	0	0%	0	0%	0%
B	<b>UKŁAD BIAŁOKRWINKOWY</b>	77	46.11%	83	47.43%	1.32%
9	blast	2	1.19%	2	1.14%	0.05%
10	promielocyt	8	4.79%	9	5.14%	0.35%
11	mielocyt	19	11.37%	15	8.57%	2.80%
12	metamielocyt	13	7.78%	19	10.85%	3.07%
13	granulocyt pałeczkowaty	12	7.18%	8	4.57%	2.61%
14	granulocyt segmentowany	16	9.58%	22	12.57%	2.99%
15	eozynofile	7	4.19%	8	4.57%	0.38%
16	monocyt	7	4.19%	5	2.85%	1.34%
C	<b>UTKANIE CHŁONNE</b>	26	15.57%	28	16.00%	0.43%
17	prolimfocyt	0	0%	0	0%	0%
18	limfocyt	22	13.17%	25	14.28%	1.11%
19	proplazmocyt	0	0%	0	0%	0%
20	plazmocyt	4	2.39%	3	1.71%	0.68%
	<b>POZA MIELOGRAMEM</b>					
21	cienie i krwinki czerwone	139		131		
Razem	suma rozpoznanych komórek mielogramu	167		175		

Tabela 6.12 Porównanie rzeczywistego składu mielogramu trzeciego pacjenta i otrzymanego w procesie automatycznej klasyfikacji

Oznaczenie układu / klasa	Typ komórki	Rzeczywista liczba komórek (wg eksperta)	Udział procentowy (wg eksperta)	Liczba komórek (wg systemu autora)	Udział procentowy (wg systemu autora)	Błąd (w punktach procentowych)
A1	<b>ERYTROPOEZA NORMOBLASTYCZNA</b>	68	35.23%	70	35%	0.23%
1	postać podziałowa erytropoezy	2	1.03%	2	1%	0.03%
2	proerytroblast	4	2.07%	6	3%	0.93%
3	erytroblast zasadochłonny	17	8.80%	11	5.5%	3.3%
4	erytroblast polichromatyczny	23	11.91%	25	12.5%	0.59%
5	erytroblast kwasochłonny	22	11.39%	26	13%	1.61%
A2	<b>ERYTROPOEZA MEGALOBLASTYCZNA</b>	0	0%	0	0%	0%
6	promegaloblast	0	0%	0	0%	0%
7	megaloblast zasadochłonny	0	0%	0	0%	0%
8	megaloblast polichromatyczny	0	0%	0	0%	0%
B	<b>UKŁAD BIAŁOKRWINKOWY</b>	34	19.17%	39	19.7%	0.53%
9	blast	2	1.03%	5	2.5%	1.47%
10	promielocyt	8	4.14%	6	3%	1.14%
11	mielocyt	4	2.07%	5	2.5%	0.43%
12	metamielocyt	3	3.01%	4	2%	1.01%
13	granulocyt pałeczkowaty	1	0.51%	1	0.5%	0.01%
14	granulocyt segmentowany	9	4.66%	10	5%	0.34%
15	eozynofile	7	3.61%	8	4%	0.39%
16	monocyt	3	1.55%	2	1%	0.55%
C	<b>UTKANIE CHŁONNE</b>	85	44.04%	89	44.5%	0.46%
17	prolimfocyt	0	0%	3	1.5%	1.5%
18	limfocyt	83	43%	81	40.5%	2.5%
19	proplazmocyt	0	0%	0	0%	0%
20	plazmocyt	2	1.03%	5	2.5%	1.47%
	<b>POZA MIELOGRAMEM</b>					
21	cienie i krwinki czerwone	79		72		
Razem	suma rozpoznanych komórek mielogramu	207		200		



Na rys. 6.1 pokazano wykresy przedstawiające udziały procentowe głównych linii rozwojowych komórek dla poszczególnych pacjentów. Potwierdzają one bardzo dobrą dokładność wyników uzyskanych przy zastosowaniu systemu opracowanego w pracy. Przedstawione wyniki dla trzech przykładowych pacjentów świadczą o wysokiej dokładności automatycznego ustalania składu szpiku kostnego. Błąd bezwzględny oceny udziału procentowego danej linii rozwojowej w szpiku wyniósł maksymalnie 1.67 %, co przy średnim



Rys. 6.1 Porównanie oceny składu procentowego linii rozwojowych w mielogramie określone przez eksperta ludzkiego i automatyczną klasyfikację

udziale układu równym 33% (praktycznie w większości przypadków występują tylko trzy układy) daje maksymalny błąd względny poniżej 5%. W praktyce laboratoryjnej przyjmuje się za akceptowalny poziom błędów równy 5%.

Drugim wskaźnikiem oceny dokładności automatycznego wyznaczania mielogramu może być obliczenie błędów średniego estymacji udziału poszczególnych typów komórek (ostatnia kolumna tabel 6.10 – 6.12 w punktach procentowych). Dla pierwszego pacjenta otrzymano 0.84% dla komórek występujących w preparacie, dla drugiego pacjenta 1.19%, a dla trzeciego 1.08%. Taki poziom błędów jest znacznie niższy od zmienności wyników oceny laboratoryjnej związanej z losowym wyborem analizowanego miejsca rozmazu, dla rozmazów pochodzących od jednego pacjenta. W efekcie wyniki otrzymane w automatycznej klasyfikacji są bardzo dokładne i dają właściwy obraz proporcji zachodzących pomiędzy poszczególnymi układami krwiotwórczymi w szpiku kostnym.

## 7. Podsumowanie i wnioski końcowe

Problem automatycznego rozpoznania i klasyfikacji komórek krwiotwórczych zawartych w szpiku kostnym jest zadaniem bardzo złożonym. Do jego realizacji niezbędne jest wykonanie wielu zadań pośrednich, w tym automatyczne pobieranie obrazów z mikroskopu, przetwarzanie wstępne obrazu i ekstrakcja pojedynczych komórek, generacja i selekcja cech tworzących dane wejściowe dla klasyfikatora i wreszcie końcowa klasyfikacja. Aby otrzymać dobry system klasyfikacyjny komórek należało każdy z tych etapów opracować w taki sposób, aby wprowadzał jak najmniej błędów, niezależnie od biologicznej i chemicznej zmienności badanych preparatów. Przedstawione w pracy wyniki dowodzą, że automatyczna klasyfikacja wybranych typów komórek szpiku kostnego jest możliwa z błędem zbliżonym do akceptowalnego poziomu w praktyce diagnostyki laboratoryjnej stosowanej w szpitalach.

Pierwszym wstępnym zadaniem było ustalenie optymalnych parametrów ustawień mikroskopu i kamery, służących do akwizycji obrazów cyfrowych. Powinny być one praktycznie niezależne od badanego preparatu. W rozwiązaniu zaproponowanym w pracy udział człowieka sprowadza się do wybrania odpowiedniego miejsca w preparacie przy zastosowaniu mikroskopu ze sterowanym stolikiem i automatyczną regulacją ostrości. Kolejnym etapem jest preprocesing obrazu, a następnie jego segmentacja przeprowadzona przy zastosowaniu metod morfologicznych. Zadanie to było bardzo trudne ze względu na złożony charakter struktur obrazu. Komórki należące do tej samej klasy są mocno zróżnicowane. Dotyczy to zarówno jądra o bardzo różnych barwach, wielkościach, kształtach i teksturze, jak i cytoplazmy o zróżnicowanej barwie i różnej ziarnistości. Dodatkowym utrudnieniem jest gęste upakowanie komórek w preparacie, powodujące zlewanie się kilku komórek w jedną. Utrudnieniem jest występowanie w obrazie krwinek czerwonych, które nie podlegają analizie, oraz cieni komórkowych powstałych z rozpadu komórek, trudnych do rozróżnienia dla systemu działającego w trybie automatycznym. W wyniku przeprowadzenia wielu prób stworzono skuteczny algorytm przetwarzania wstępnego stosujący skalowanie, szereg operacji morfologicznych, filtrację i wygładzanie brzegów. Zastosowanie segmentacji obrazu metodą działów wodnych pozwoliło na uzyskanie precyzyjnej ekstrakcji komórek przy braku ingerencji człowieka. Badania na rozległej bazie danych (ponad 4000 obrazów)

potwierdziły prawie 95% skuteczność wydzielenia zarówno komórek jądrzastych jak i cieni komórkowych oraz nielicznych krwinek czerwonych.

Kolejnym etapem jest przetworzenie obrazu poszczególnych komórek na właściwy zestaw cech tworzących dane wejściowe dla klasyfikatora. Cechy zastosowane w pracy należą do grup: teksturalnych, statystycznych, geometrycznych i morfologicznych. Wydają się stanowić wyczerpujący zestaw parametrów umożliwiających rozróżnienie różnych klas komórek od siebie. Aby uzyskać najwyższą sprawność klasyfikatora przeprowadzono wnikliwą analizę cech mającą na celu wyselekcjonowanie tych najlepszych, gwarantujących najwyższą sprawność rozpoznania. W wyniku przeprowadzonych badań stwierdzono, że selekcja na bazie analizy jakości pojedynczych cech nie jest właściwym podejściem w tym trudnym zadaniu. Dobre wyniki uzyskano dopiero przy ocenie poszczególnych cech działających jednocześnie. Możliwość takie stwarzała sieć liniowa SVM, w sposób znakomity oceniająca jakość poszczególnych cech w procesie rozpoznania i klasyfikacji komórek.

Jako klasyfikator wybrano sieć neuronową typu SVM o jądrze nieliniowym typu gaussowskiego. Jest to sieć stosunkowo nowa, powstała w latach dziewięćdziesiątych i stosowana obecnie na coraz szerszą skalę. Przeprowadzone testy [5,6,35,36,37,38,40,43] wykazały jej bezwzględną przewagę w większości typowych zadań klasyfikacyjnych. Szczególnie dobrą skuteczność sieci SVM obserwuje się w problemach klasyfikacji wieloklasowej. W zadaniu rozpoznania komórek jej największą zaletą jest możliwość sterowania szerokością marginesu separacji, dzięki której można stworzyć sieć klasyfikującą odporną na dużą zmienność kształtu i barwy komórek należących do tej samej klasy. Stworzony od podstaw układ automatycznego rozpoznania i klasyfikacji komórek krwiotwórczych pozwolił na uzyskanie dokładności rozpoznania porównywalnej z ekspertem ludzkim. Przedstawione wyniki badań dają realną nadzieję na stworzenie automatycznego układu rozpoznawania i klasyfikacji komórek szpiku kostnego zastępującego pracę tego eksperta.

Za najważniejsze osiągnięcia oryginalne pracy autor uważa:

1. opracowanie układu wstępnego przetwarzania i segmentacji obrazu rozmazu szpiku kostnego umożliwiającego precyzyjną ekstrakcję obrazów poszczególnych komórek bez aktywnego udziału człowieka;
2. stworzenie bazy danych zawierającej ponad 3000 opisanych obrazów rozmazu szpiku kostnego i ponad 8000 wydzielonych z nich komórek;

3. zaproponowanie skutecznych algorytmów generacji cech diagnostycznych komórek wykorzystujących opis tekstury, geometrii oraz parametrów statystycznych i morfologicznych;
4. przeprowadzenie dogłębnej analizy jakościowej pojedynczych cech diagnostycznych i zaproponowanie skutecznego rozwiązania problemu ich selekcji na bazie liniowej sieci SVM;
5. opracowanie programu uczącego sieci klasyfikującej SVM, implementującego najskuteczniejsze obecnie algorytmy uczenia: Platta, SVM<sup>Light</sup>, LSVM i potwierdzenie ich skuteczności w rozwiązywaniu zadania klasyfikacji komórek krwiotwórczych;
6. powiązanie wymienionych wcześniej operacji w jeden kompletny system klasyfikacyjny i jego wszechstronne przetestowanie dla różnej liczby typów komórek pochodzących od wielu pacjentów;
7. Wykonanie bardzo dużej liczby eksperymentów numerycznych dla potwierdzenia skuteczności działania opisanego systemu.

Autor zamierza w przyszłości kontynuować badania w tej dziedzinie. Kierunki dalszych badań dotyczyć będą:

- optymalizacji metod obliczeniowych prowadzących do zmniejszenia czasu przetwarzania wstępnego danych;
- stworzenia interfejsu użytkownika łatwego w obsłudze dla personelu szpitalnego;
- dalszej poprawy skuteczności klasyfikacji komórek sąsiadujących ze sobą;
- przetestowania systemu na reprezentowalnej liczbie nowych pacjentów;
- wdrożenia systemu w praktyce szpitalnej.

## Literatura

1. M. Basa, Klasyfikacja tekstur przy pomocy sieci neuronowych, Praca dyplomowa magisterska, Politechnika Warszawska, Warszawa, 2001
2. M. Beksac, M. S. Beksac, V. B. Tippi, H. A. Duru, M. U. Karakas, A. Nurcarak, An artificial intelligent diagnostic system on differential recognition of hematopoietic cells from microscopic images, *Cytometry*, 1997, vol. 30, pp. 145-150
3. J. M. Bennett, D. Catovsky, M. T. Daniel, G. Flandrin, D. A. Galton, H. R. Gralnick, C. Sultan, Proposals for the classification of the acute leukaemias. French-American-British (FAB) co-operative group, *Br. J. Haematol.*, 1976, vol. 33, pp. 451-458
4. B. Boser, I. Guyon, V. Vapnik, An training algorithm for optimal margin classifier, *Proc. V Annual Workshop on Comp. Learn. Theory*, Pittsburg, 1992, pp. 144-152
5. K. Brudzewski, S. Osowski, T. Markiewicz, J. Ulaczyk, Classification of gasoline with supplement of bio-products by means of an electronic nose and SVM neural network, *Sensors and Actuators*, 2005
6. K. Brudzewski, S. Osowski, T. Markiewicz, Classification of milk by means of an electronic nose and SVM neural network, *Sensors and Actuators B-Chem.* 98 (2-3), 2004, pp. 291-298
7. K. Brudzewski, S. Osowski, T. Markiewicz, J. Ulaczyk, Support Vector Machine for Recognition of Bio-products in Gasoline, *Lecture Notes on Computer Science*, 2005, vol. 3697, pp. 899-904
8. C. Burges, A tutorial on support vector machines for pattern recognition, (in *Knowledge discovery and data mining*, ed. Usama Fayyad, Kluwer, 2000), pp. 1-43
9. C. Chang, C. Lin, LIBSVM: a library for Support Vector Machines, <http://www.csie.ntu.edu.tw/~cjlin/libsvm>, 2003
10. A. Cichocki, R. Unbehauen, *Neural networks for optimization and signal processing*, Wiley, N. Y., 1993
11. K. Crammer, Y. Singer, On the Learnability and Design of Output Codes for Multiclass Problems, *Proc. of COLT*, 2000.
12. W. Duch, J. Korbicz, L. Rutkowski, R. Tadeusiewicz, (red.) *Sieci neuronowe*, Akademicka Oficyna Wydawnicza, Warszawa, 2000
13. R. O. Duda, P. E. Hart, P. Stork, *Pattern classification and scene analysis*, Wiley, N.Y., 2003
14. P. Gill, W. Murray, M. Wright, *Practical optimization*, Academic Press, New York, 1981
15. G. Golub, C. Van Loan, *Matrix computation*, North Oxford Academic, 1990
16. I. Guyon, A. Elisseeff, An introduction to variable and feature selection, *Journal of Machine Learning Research*, 2003, vol. 3, pp. 1158 – 1182
17. I. Guyon, J. Weston, S. Barnhill, V. Vapnik, Gene selection for cancer classification using Support Vector Machines, *Machine Learning*, 2002, vol. 46, pp. 389-422
18. S. Haykin, *Neural networks, comprehensive foundation*, Prentice Hall, New Jersey, 1999
19. H. Hengen, S. Spoor, M. Pandit, *Analysis of blood & bone marrow smears*, SPIE Med. Imag., San Diego, 2002
20. A.V. Hoffbrand, J.E. Pettit, *Atlas hematologii klinicznej*, Czelej, Lublin 2003
21. C. W. Hsu, C. J. Lin, A comparison methods for multi class support vector machines, *IEEE Trans. Neural Networks*, 2002, vol. 13, pp. 415-425
22. IBM Corporation, *IBM optimization subroutine guide and reference*, IBM Systems Journal, vol. 31, 1992, SC23-0519
23. K. Janicki, *Hematologia kliniczna*, PZWL, Warszawa, 1992
24. T. Joachims. Making large-scale SVM learning practical. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods --- Support Vector Learning*, pages 169-184, Cambridge, MIT Press, MA, 1999
25. W. Kwiatkowski, *Metody automatycznego rozpoznawania wzorców*, WAT, Warszawa, 2001
26. K. Lewandowski, A. Hellmann, *Haematology atlas*, Multimedia Medical Publisher, Gdansk, 2001
27. O. Lezoray, H. Cardot, Cooperation of color pixel classification schemes and color watershed, *IEEE Trans. Image Processing*, vol. 11, 2002, pp. 783-789
28. O. L. Mangasarian, D. R. Musicant, Lagrangian support vector machines, *Journal of Machine Learning Research*, 2001, pp. 161-177

29. T. Markiewicz, S. Osowski, B. Marianska, L. Moszczyński, Automatic Recognition of the Blood Cells of Myelogenous Leukemia Using SVM , IJCNN Montreal, 2005
30. T. Markiewicz, L. Moszczyński, Analysis of features for blood cell recognition, VI CPEE, Zakopane, 2004
31. T. Markiewicz, S. Osowski, L. Moszczyński, Myelogenous leukemia cell image preprocessing for feature generation, V Int. Workshop Computational Methods in Electrical Engineering, Jazlowiec, Ukraina, 2003, pp. 70-73
32. Matlab Image Processing Toolbox User's, MathWorks, Natick, 2005
33. Matlab Optimization Toolbox User's Guide, MathWorks, Natick, 2005
34. B. Murtagh, M. Saunders, MINOS 5.1 user guide, Technical Report SOL 83-20R, Stanford University, 1983, revised 1987
35. S. Osowski, L. Tran Hoai, T. Markiewicz, Support Vector Machine based expert system for reliable heart beat recognition, IEEE Trans. on Biomedical Engineering, 2004, vol. 51 , pp. 582-589
36. S. Osowski, T. Markiewicz, Analiza porównawcza algorytmów uczących sieci neuronowych RBF, Software 2.0, 2004, luty, pp. 56-62
37. S. Osowski, T. Markiewicz, OLS Versus SVM Approach to Learning of RBF Networks, IJCNN Montreal, 2005
38. S. Osowski, L. Tran Hoai, T. Markiewicz, Recognition of Heartbeats Using Support Vector Machine Networks – a Comparative Study, Lecture Notes on Computer Science, 2005, vol. 3697, pp. 637-642
39. S. Osowski, T. Markiewicz, B. Mariańska, L. Moszczyński, Feature generation for the cell image recognition of myelogenous leukemia, IEEE Int. Conf. EUSIPCO, Vienna, 2004, pp. 753-756
40. S. Osowski, K. Siwek, T. Markiewicz, MLP and SVM Networks – a Comparative Study, IEEE Int. Conf. NORSIG, 2004, Helsinki, pp. 37-41
41. S. Osowski, K. Siwek, T. Markiewicz, Comparative analysis of learning algorithms of MLP network, Int. Conference on Fundamentals of Electrotechnics and Circuit Theory (SPETO), 2004, pp. 465-468
42. S. Osowski, T. Markiewicz, B. Świdorski, Analiza porównawcza algorytmów uczących sieci RBF, XXVI IC-SPETO, 2003, pp. 475-480
43. S. Osowski, L. Tran Hoai, T. Markiewicz, K. Siwek, Support vector machine and neuro-fuzzy network for heart beat recognition, IV CPEE Workshop, Zakopane, 2002, pp. 187-190
44. S. Osowski, T. Markiewicz, M. Basa, L. Tran Hoai, SVM network for texture recognition, Int. Conf. in Signals and Electronic Systems, Wrocław, 2002, pp. 267-300
45. S. Osowski, T. Markiewicz, Algorytm ortogonalizacji w uczeniu wielowyjściowej sieci RBF, XXV SPETO, Ustroń, 2002, pp. 529 – 532
46. E. Osuna, R. Freund, F. Girosi, An improved training algorithm for SVM, (in "Neural networks for signal processing VII", J. Principe, L. Gile, N. Morgan, E. Wilson eds), IEEE Press, New York, 1997, pp. 276-285
47. N. Otsu, A Threshold Selection Method from Gray-Level Histograms, IEEE Transactions on Systems, Man, and Cybernetics, vol. 9, no. 1, 1979, pp. 62-66
48. J. Platt, Fast training of SVM using sequential optimization, (in "Advances in kernel methods – support vector learning", B. Scholkopf, C. Burges, A. Smola eds, MIT Press, Cambridge), 1998, pp. 185-208
49. R. Sałat, T. Markiewicz, A novel approach to fault location in transmission line of power system using Support Vector Machine, Proc. IV Int. Workshop Computational Problems of Electrical Engineering, Zakopane, 2002
50. J. Schurmann, Pattern classification, a unified view of statistical and neural approaches, Wiley, N. Y., 1996
51. B. Schölkopf, A. Smola, Learning with Kernels, MIT Press, Cambridge, MA, 2002
52. A. Smola, B. Schölkopf, A tutorial on support vector regression, NeuroColt Technical Report NV2-TR-1998-030 (1998) <http://www.neurocolt.com>
53. S. Sohn, Bone marrow white blood cell classification, M.S. thesis, Univ. Missouri, Columbia, 2000
54. P. Soille, Morphological image analysis, principles and applications, Springer, Berlin, 2003
55. A. Statnikov, C.F. Aliferis, I. Tsamardinos, D. Hardin, S. Levy, A Comprehensive Evaluation of Multicategory Classification Methods for Microarray Gene Expression Cancer Diagnosis, Bioinformatics, 2005 Mar 1;21(5):631-43
56. N. Theera-Umpon, P. Gader, System-level training of neural networks for counting white blood cells, IEEE Trans. SMS-C, 2002, vol. 32, pp. 48-53
57. R. J. Vanderbei, LOQO: an interior point code for quadratic programming, TRSOR-94-15, Statistics and Operations Research, Princeton University, N.J., 1994
58. V. Vapnik, Statistical learning theory, Wiley, N.Y., 1998

59. V. Vapnik, An overview of statistical learning theory, IEEE Trans. Neural Networks, 1999, vol. 10, pp. 988-999
60. T. Wagner, Texture analysis ( in "Handbook of Computer Vision and Application", Jahne, B., Haussecker, H., and Geisser P. eds.), Academic Press, 1999, pp. 275-309
61. J. Weston, C. Watkins, Multiclass support vector machines, Technical report CSD-TR-98-04, 1998, University of London
62. W. Wolberg, W. N. Street, O. L. Mangasarian, Machine learning to diagnose breast cancer from image-processed features, Rep. of Uni. Wisconsin, 1994
63. K. W. Zieliński, M. Strzelecki, Komputerowa analiza obrazu biomedycznego, Wstęp do morfometrii i patologii ilościowej, PWN, Warszawa 2001
64. Benchmarks, <http://users.rsise.anu.edu.au/~raetsch/data/>