



Selekcja cech osobniczych sygnału mowy z wykorzystaniem algorytmów genetycznych

KAMIL KAMIŃSKI, ANDRZEJ P. DOBROWOLSKI,
EWELINA MAJDA-ZDANCEWICZ

Wojskowa Akademia Techniczna, Wydział Elektroniki, Instytut Systemów Elektronicznych,
00-908 Warszawa, ul. gen. S. Kaliskiego 2, kkw.kaminski@gmail.com.

Streszczenie. W artykule przedstawiono system automatycznego rozpoznawania mowy zaimplementowany w środowisku *Matlab* oraz pokazano sposoby realizacji i optymalizacji poszczególnych elementów tego systemu. Główny nacisk położono na wyselekcjonowanie cech dystynktywnych głosu mówcy z wykorzystaniem algorytmu genetycznego, który pozwala na uwzględnienie synergii cech podczas selekcji. Pokazano również wyniki optymalizacji wybranych elementów klasyfikatora, m.in. liczby rozkładów Gaussa użytych do zamodelowania każdego z głosów. Ponadto, podczas tworzenia modeli poszczególnych głosów zastosowano uniwersalny model głosów.

Słowa kluczowe: biometria, automatyczne rozpoznawanie mowy, algorytmy genetyczne, selekcja cech

DOI: 10.5604/12345865.1197999

1. Wstęp

Ze względu na dużą redundancję przebiegu czasowego sygnału mowy, w systemach automatycznego rozpoznawania mowy oraz mowy (ASR — ang. *Automatic Speaker/Speech Recognition*) nie wykorzystuje się bezpośrednio sygnału w dziedzinie czasu. Projektanci systemów ASR poprzez złożone operacje na sygnałach mowy wydobywają z nich zestaw najistotniejszych — pod kątem określonego kryterium — cech dystynktywnych (deskryptorów). Otrzymane cechy tworzą niejednokrotnie liczny zbiór, a jednocześnie nie każda cecha skutecznie reprezentuje analizowany proces [1]. Stąd potrzeba wyselekcjonowania deskryptorów o największej zdolności dyskryminacyjnej. Istnieje wiele metod selekcji cech diagnostycznych, jednak

nie wszystkie z nich uwzględniają synergię¹ cech, która pozwala uzyskać większą zdolność dyskryminacji dla cech występujących w grupie niż dla poszczególnych cech z osobna. Niestety synergia cech w istotny sposób utrudnia ich selekcję, gdyż dobór najistotniejszych deskryptorów na podstawie indywidualnej oceny każdego z nich może prowadzić do uzyskania nieoptymalnego zbioru.

Prezentowana w artykule metoda selekcji cech z wykorzystaniem algorytmu genetycznego nie jest obciążona tą niedogodnością, gdyż w wyniku działania operacji genetycznych na pełnym zbiorze cech dystynktywnych możliwe jest otrzymanie optymalnego zbioru deskryptorów. W celu dokonania selekcji cech z wykorzystaniem algorytmu genetycznego autorzy posłużyli się dwiema niezależnymi komercyjnymi bazami głosów, dzięki czemu możliwe stało się sprawdzenie poprawności działania zastosowanego algorytmu selekcji cech oraz wyodrębnienie zbioru takich cech, które pozwalają na poprawną identyfikację mówcy niezależnie od zastosowanej bazy głosów. W kolejnych rozdziałach artykułu przedstawiono architekturę oraz sposób działania zaimplementowanego przez autorów systemu automatycznego rozpoznawania mówcy oraz wyniki badań przeprowadzonej selekcji cech i jej wpływu na wypadkową skuteczność identyfikacji mówców. W artykule znajdują się również wyniki optymalizacji zastosowanego w systemie klasyfikatora wykorzystującego tzw. modele mieszanin Gaussowskich GMMs (ang. *Gaussian Mixture Models*) oraz uniwersalny model głosów UBM (ang. *Universal Background Model*).

2. Baza głosów

Podczas prowadzonych badań autorzy wykorzystywali dwie powszechnie znane bazy głosów przeznaczone do badań związanych z rozpoznawaniem mowy i mówcy. Pierwszą z nich jest baza nagrań TIMIT stworzona przez MIT (ang. *Massachusetts Institute of Technology*), SRI (ang. *Stanford Research Institute*) oraz TI (ang. *Texas Instruments*) [2]. Baza zawiera nagrania głosowe 192 kobiet i 438 mężczyzn, zarejestrowane z szybkością próbkowania 16 kS/s, przy zapisie jednokanałowym z 16-bitową rozdzielczością amplitudową. Każdy z mówców reprezentowany jest przez 30-sekundowe nagranie utworzone przy użyciu dziesięciu niezależnych 3-sekundowych nagrań głosowych. W celu równomiernego wykorzystania głosów męskich i żeńskich dla potrzeb prezentowanych badań użyto 200 głosów kolejnych mówców. Nie były one w żaden sposób dobierane, wykorzystano 100 kolejnych męskich i 100 kolejnych żeńskich głosów z listy głosów dostępnych w bazie TIMIT. Nagrania użyte w prezentowanych badaniach zostały przepróbkowane do szybkości próbkowania 8 kS/s. Pozwala to na sprawdzenie skuteczności systemu w warunkach

¹ Synergia to współdziałanie różnych czynników, którego efekt jest większy niż suma poszczególnych oddzielnych działań.

zbliżonych do transmisji telefonicznej, co poszerza spektrum zastosowań prezentowanego systemu.

Kolejną bazą głosów użytą podczas badań była baza *2002 NIST Speaker Recognition Evaluation*. Jest to baza przeznaczona dla systemów rozpoznawania mówcy niezależnych od tekstu. Obejmuje ponad 156 godzin nagrań, rejestrowanych z szybkością próbkowania 8 kS/s oraz 16 kS/s [3]. Na nagrania składają się zarówno wypowiedzi pojedynczych mówców, jak i dialogi nagrywane w różnych warunkach. Dla potrzeb prowadzonych badań wykorzystano nagrania pojedynczych mówców rejestrowane z szybkością próbkowania 8 kS/s. W bazie znajduje się 330 nagrań spełniających powyższe kryterium (191 głosów żeńskich i 139 głosów męskich). Podobnie jak w przypadku bazy TIMIT, dla zrównoważenia liczby głosów męskich i żeńskich wykorzystano 200 głosów kolejnych mówców (100 głosów żeńskich i 100 męskich). Nagrania występujące w bazie *2002 NIST Speaker Recognition Evaluation* gromadzone są w formacie **.sph (sphere format)*, co powoduje konieczność właściwego wczytania plików lub ich konwersji. Autorzy do tego celu wykorzystali dostępny w środowisku *Matlab* pakiet narzędzi VOICEBOX przeznaczonych do przetwarzania sygnału mowy [4].

W celu wyodrębnienia części uczącej i części testowej nagrań zawartych w bazach TIMIT oraz *2002 NIST Speaker Recognition Evaluation* zdecydowano, aby przeznaczyć na część uczącą 25 s każdego z nagrań, a na część testową pozostałe niezależne 5 s nagrań. W przypadku bazy TIMIT 10 nagrań głosowych reprezentatywnych dla każdego z mówców zostało scalonych, a następnie podzielonych na dwie części (25 s i 5 s). Natomiast w bazie *2002 NIST Speaker Recognition Evaluation* zgromadzone nagrania przypadające na jednego mówcę były znacznie dłuższe, co również pozwoliło na wycięcie z każdego z nich 25 s segmentu uczącego i 5 s niezależnego segmentu testowego.

3. Wstępne przetwarzanie sygnałów mowy

Pierwszym i bardzo ważnym etapem działania systemu automatycznego rozpoznawania mówcy jest proces wstępnego przetwarzania sygnału, pozwalający na uniezależnienie zgromadzonych nagrań od ustawień sprzętu nagrywającego. Na tym etapie realizowana jest filtracja, normalizacja oraz selekcja ramek. Filtracja przeprowadzona została poprzez zastosowanie filtru pasmowo-przepustowego o skończonej odpowiedzi impulsowej, którego częstotliwość odcięcia oraz rząd poddano wcześniejszej optymalizacji [5].

Poziomy mocy sygnału mowy u różnych mówców są bardzo zróżnicowane, dlatego po filtracji przeprowadzono normalizację zgromadzonych sygnałów w stosunku do wartości maksymalnej dla każdego mówcy. Dzięki temu skalowaniu zachowane zostały relacje energetyczne pomiędzy poszczególnymi fragmentami zapisu.

W następnej kolejności realizowana była segmentacja sygnału mowy, która jest tożsama z operacją okienkowania. W prezentowanym systemie zastosowano okno Hamminga, którego widmo charakteryzuje się niskim poziomem listków bocznych, a tym samym znacznie redukuje zjawisko przecieku widma. Długość okna oraz jego przesunięcie również poddawane były procesowi optymalizacji w trakcie wcześniejszych badań.

Algorytm selekcji ramek realizowany jest trzyetapowo. W pierwszej kolejności eliminowane są ramki niespełniające wyznaczonego empirycznie kryterium mocy w ramce, stanowiącego iloczyn mocy najcichszej ramki oraz stałej wyznaczonej w procesie optymalizacji. Pozwala to na wyeliminowanie ciszy z analizowanych fragmentów mowy. W drugim etapie selekcjonowane są wyłącznie ramki zawierające głoski dźwięczne, ze względu na fakt, że w nich zawarta jest większość informacji o głosie mówcy istotna z punktu widzenia jego identyfikacji. Klasyfikacji na ramki zawierające głoski dźwięczne i bezdźwięczne dokonuje się na podstawie funkcji autokorelacji. Pierwsze maksimum funkcji autokorelacji, występujące dla zerowego przesunięcia, niesie informację o energii ramki, natomiast drugie pozwala określić jej dźwięczność. Wartość drugiego maksimum funkcji autokorelacji w ramce porównywana jest z wyznaczonym w procesie optymalizacji progiem dźwięczności.

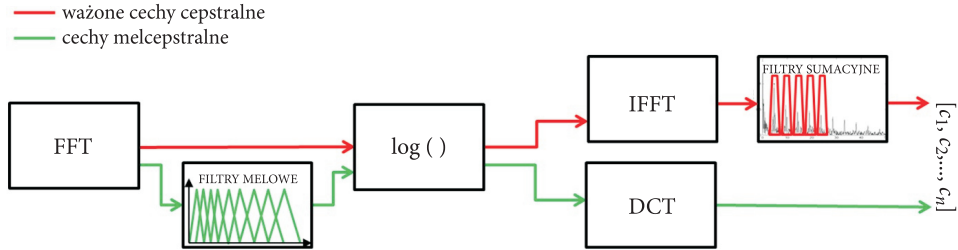
Ostatni etap eliminacji ramek polega na porównywaniu wartości częstotliwości podstawowej wyznaczonej przy użyciu dwóch metod. Wyznaczanie częstotliwości podstawowej metodą cepstralną jest mniej dokładne, ale bardziej niezawodne niż metodą autokorelacji, szczególnie przy silnym zaszumieniu. Dlatego porównanie tych wartości pozwala na określenie stopnia zaszumienia ramki sygnału mowy. Jeżeli różnica częstotliwości podstawowej wyznaczonej dwiema niezależnymi metodami będzie większa od wartości progowej wyznaczonej w procesie optymalizacji, to ramka zostanie odrzucona ze względu na zbytne zaszumienie sygnału, niekorzystnie wpływające na ostateczny wynik działania systemu ASR.

4. Generacja dystyngtywnych cech głosu

Generacja cech jest kluczowym elementem systemu automatycznego rozpoznawania mówcy. Błędy popełnione na tym etapie są już nie do nadrobienia w dalszych etapach działania systemu. Prezentowany w artykule system ASR wykorzystuje dwa rodzaje dystyngtywnych cech głosu. Są nimi ważone cechy cepstralne i cechy melcepstralne [6]. Sposób ich generacji przedstawiono na rysunku 1.

Obie zaprezentowane techniki generacji cech w pierwszej fazie działania realizują transformację sygnału do dziedziny częstotliwości, ze względu na fakt, że analiza sygnału czasowego obarczona jest nadmierną redundancją (proces ten jest inspirowany działaniem ludzkich organów komunikacji głosowej). Otrzymane widmo amplitudowe znacznie bardziej uwydatnia różnice treści nagrywanych wypowiedzi

niż osobnicze atrybuty związane m.in. z tonem krtaniowym. Dlatego konieczne są dalsze operacje na sygnale pozwalające wydobyć osobnicze cechy głosu mówcy.



Rys. 1. Schemat procesu generacji cech

W przypadku generacji ważonych cech cepstralnych kolejnym procesem jest logarytmowanie widma amplitudowego, dzięki któremu multiplikatywny związek między składową wolnozmienną i amplitudami poszczególnych impulsów pochodzących od pobudzenia zamienia się w związek addytywny. Poddanie takiego sygnału odwrotnej transformacji Fouriera powoduje, że wolnozmiennne przebiegi związane z transmitancją traktu głosowego zostają usytuowane blisko zera na osi czasu cepstralnego zwanego pseudoczasem, natomiast impulsy związane z dźwiękiem krtaniowym zaczynają się mniej więcej w obrębie okresu sygnału krtaniowego i powtarzają się co ten okres. Ostatnim etapem generacji ważonych cech cepstralnych jest wymnożenie otrzymanego sygnału w dziedzinie pseudoczasu przez bank filtrów sumacyjnych, które uwzględniają nie tylko maksymalne amplitudy prążków w cepstrum, lecz także wartości je otaczające, które również niosą informację osobniczą o głosie mówcy.

Podczas generacji cech melcepstralnych otrzymane po transformacji amplitudowe widmo Fouriera wymnażane jest przez bank filtrów melowych, które naśladując ludzki organ słuchu i jego nieliniową wrażliwość na pobudzenia z różnych zakresów częstotliwości, powodują poprawę percepcji. Następnie otrzymane cechy są logarytmowane, podobnie jak w przypadku cech cepstralnych. Ostatnim etapem generacji cech melcepstralnych jest poddanie sygnału transformacji cosinusowej w celu dekorrelacji cech.

5. Selekcja cech dystynktywnych z wykorzystaniem algorytmu genetycznego

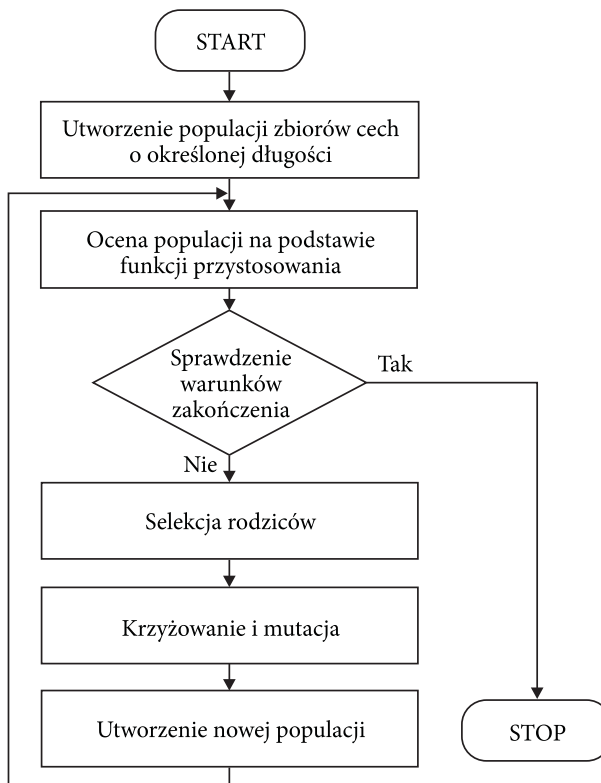
Na etapie generacji cech wygenerowano maksymalnie liczny zbiór cech dystynktywnych, które mogą być wykorzystane w systemie automatycznego rozpoznawania mówcy. Prowadzone na świecie badania dowodzą, że nie zawsze użycie maksymalnego

zestawu cech pozwala uzyskać najlepsze wyniki [1]. Selekcja cech daje niejednokrotnie możliwość otrzymania większej lub takiej samej dokładności klasyfikacji dla zredukowanego wektora cech, co z kolei przekłada się na skrócenie czasu obliczeń.

Podczas oceny jakości cech niektóre z nich mogą mieć postać szumu pomiarowego pogarszającego możliwość rozpoznawania danego wzorca, natomiast inne mogą być ze sobą silnie skorelowane, co powoduje dominację tych cech nad pozostałymi i zwykle niekorzystnie wpływa na jakość klasyfikacji.

Istotnym elementem jest wybór metody selekcji cech. W literaturze dostępnych jest wiele różnych metod selekcji, poczynając od szybkich metod rankingowych, a kończąc na czasochłonnych, zawierających złożone klasyfikatory. Do najbardziej znanych miar jakości i metod selekcji cech należą: współczynnik Fishera, t-statystyki, korelacja wzajemna, sekwencyjna selekcja w przód, algorytmy genetyczne, liniowa analiza dyskryminacyjna.

Autorzy w prezentowanym systemie zastosowali algorytm genetyczny do selekcji najistotniejszych z punktu widzenia systemu ASR cech osobniczych. Metoda ta



Rys. 2. Schemat blokowy selekcji cech przy użyciu algorytmu genetycznego [1]

uwzględnia synergię cech i umożliwia otrzymanie ich optymalnego zbioru, jednakże jest dość czasochłonna. Zasada działania zaimplementowanego klasyfikatora cech wykorzystującego algorytm genetyczny została przedstawiona na rysunku 2.

W pierwszej kolejności w maksymalnym potencjalnym zbiorze cech dys-
tynktywnych obliczana jest informacja wzajemna (3) pomiędzy poszczególnymi
cechami oraz pomiędzy cechami a wektorem przynależności klasowej. Pozwala ona
na określenie, jak bardzo znajomość zmiennej Y ułatwia przewidywanie wartości
zmiennej X . Innymi słowy, jak bardzo poznanie jednej ze zmiennych zmniejsza
niepewność drugiej zmiennej. W celu wyznaczenia informacji wzajemnej konieczne
jest określenie stopnia nieuporządkowania poszczególnych zmiennych zwanego
entropią. Niska entropia zmiennej zapewnia jej wysoką przewidywalność, co jest
efektem pożądanym. Entropię każdej zmiennej oraz entropię łączną wyznacza się
wg zależności (1) i (2) [7]:

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x) \quad (1)$$

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y) \quad (2)$$

$$I(Y; X) = H(X) + H(Y) - H(X, Y). \quad (3)$$

W rezultacie utworzona zostaje macierz informacji wzajemnych pomiędzy
wszystkimi potencjalnymi cechami $I(yh_i; y)$ oraz wektor informacji wzajemnej
pomiędzy każdą z cech a wektorem przynależności klasowych $I(yh_i; yh_j)$.

Niestety obliczanie informacji wzajemnej dla dużego zbioru cech jest bardzo
czasochłonne. Dlatego powtarzanie wspomnianych obliczeń dla każdej możliwej
kombinacji byłoby niewykonalne w rozsądnym czasie, stosując deterministyczne
formy rozwiązań. W przypadku algorytmu genetycznego bazowe informacje wza-
jemnie maksymalnego zbioru cech wykorzystywane są każdorazowo przy ocenie
przystosowania populacji osobników. Populacja początkowa tworzona jest w sposób
pseudolosowy poprzez generację wektorów (chromosomów) zawierających numery
cech. Następnie każdy chromosom w populacji jest oceniany na podstawie funkcji
przystosowania (6), z wykorzystaniem uśrednionych wartości informacji wzajemnej
pomiędzy cechami a wektorem przynależności klasowej (4) oraz pomiędzy cechami
wzajemnie między sobą (5) dla występujących w populacji zbiorów cech [7]:

$$V = \frac{1}{N_n} \sum_{i=1}^{N_n} I(yh_i; y) \quad (4)$$

$$P = \frac{1}{N_n^2} \sum_{i=1}^{N_n} \sum_{j=1}^{N_n} I(yh_i; yh_j) \quad (5)$$

$$S = V - P. \quad (6)$$

Jeżeli różnica pomiędzy maksymalną wartością funkcji przystosowania w danej populacji a wartością średnią funkcji przystosowania w populacji nie przekroczy wyznaczonego empirycznie poziomu, to dalsze obliczenia są przerywane. Przerwanie obliczeń następuje również w przypadku osiągnięcia granicznej liczby pokoleń, która stanowi wartość stałą.

Jeżeli maksymalna wartość funkcji przystosowania nie spełnia warunków zakończenia i wykazuje wzrost przystosowania względem średniej wartości tej funkcji, to kolejnym krokiem algorytmu jest selekcja osobników najlepiej przystosowanych. Selekcja odbywa się w sposób pseudolosowy, jednak faworyzuje osobniki lepiej przystosowane wg zależności (7) [7]

$$\text{round} \left(N \frac{e^{ar} - 1}{e^a - 1} \right), \quad (7)$$

gdzie N stanowi wielkość populacji, i jest wartością pseudolosową z przedziału $[0; 1]$, natomiast a jest wartością stałą.

Następnie dwa wylosowane według powyższego wzoru chromosomy podlegają krzyżowaniu. Odbywa się ono wielopunktowo w ten sposób, że każda cecha dla nowego chromosomu dobierana jest pseudolosowo spośród dwóch cech występujących na danym indeksie w wektorach cech tzw. chromosomów rodziców. Procedura ta jest powtarzana aż do otrzymania nowego wektora o tej samej długości co wektory rodzicielskie. W przypadku wystąpienia duplikacji cech w nowym wektorze następuje mutacja wektora poprzez podmienienie cechy przez cechę niewystępującą w wektorze. W wyniku operacji genetycznych powstaje nowa populacja, która poddawana jest sprawdzeniu funkcją przystosowania. Powyższe operacje są powtarzane aż do spełnienia warunków stopu. Ostatecznym wynikiem jest wektor cech najlepiej przystosowanego chromosomu z ostatniego pokolenia.

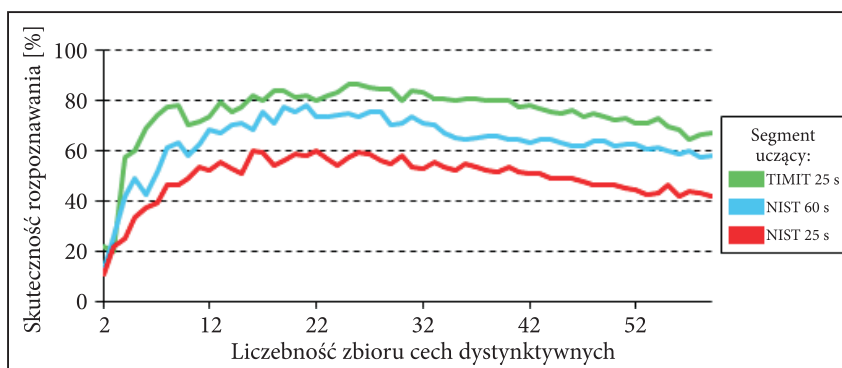
6. Klasyfikator GMM-UBM

W procesie klasyfikacji zastosowano liniową kombinację rozkładów normalnych GMMs, które umożliwiają generację modeli głosów zawierających dużą liczbę cennych informacji, a zarazem oszczędnych z punktu widzenia wymaganej pamięci [8]. Model GMM tworzony dla każdego mówcy, modeluje wielowymiarowy rozkład gęstości prawdopodobieństwa na podstawie wyekstrahowanych cech

dystynktywnych wygenerowanych w oparciu o dane uczące. Parametry startowe modelu mówcy, tj. wartości oczekiwane, macierze kowariancji oraz wagi rozkładów, mogą być dobierane w sposób pseudolosowy lub zdeterminowany wg algorytmu GMM-UBM [8]. Metoda ta polega na utworzeniu głosu uniwersalnego na podstawie głosów wielu osób. Pozwala to na dopasowanie się modelu określonego mówcy w mniejszej liczbie iteracji, co wpływa korzystnie na skuteczność identyfikacji oraz na przyspieszenie procesu tworzenia modeli głosów.

7. Wyniki eksperymentów

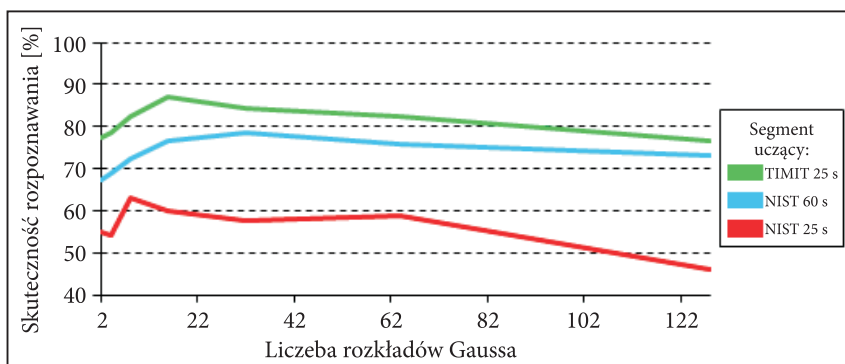
Głównym celem prezentowanych badań było zastosowanie skutecznego algorytmu selekcji dystynktywnych cech głosów mówców. Wykorzystano do tego celu algorytm genetyczny, dla którego danymi wejściowymi był maksymalny potencjalny zbiór cech dystynktywnych, wektor przynależności klasowych oraz określona długość pożądanego wektora cech po jego redukcji. Długość ta podlegała optymalizacji, której wyniki zilustrowano na rysunku 3. Do testów wykorzystano dane pochodzące z baz TIMIT i NIST. Długość segmentu uczącego dla każdego mówcy wynosiła 25 s, a długość odcinka testowego 5 s. W związku z tym, że system optymalizowany był we wcześniejszych badaniach [5] dla znacznie dłuższego segmentu uczącego, postanowiono wykonać test również dla 60 s segmentu uczącego z wykorzystaniem nagrań z bazy NIST.



Rys. 3. Skuteczność rozpoznawania w zależności od liczebności zbioru cech dystynktywnych

Wyselekcjonowane wektory cech zostały porównane, następnie wykonano kilka powtórzeń selekcji dla wektorów cech charakteryzujących się największą skutecznością. W rezultacie wybrano zbiór 22 cech pozwalający uzyskać najwyższą skuteczność identyfikacji.

W kolejnym etapie badań dokonano sprawdzenia, dla jakiej liczby rozkładów Gaussa następuje najlepsza generalizacja danych, a tym samym najwyższa skuteczność identyfikacji. Testy przeprowadzono również dla trzech różnych długości segmentów uczących.



Rys. 4. Skuteczność rozpoznawania w zależności od liczby rozkładów Gaussa użytych do zamodelowania każdego z głosów mówców

Na podstawie otrzymanych wyników widać, że dla różnych długości segmentów uczących liczba GMMs dająca najwyższą skuteczność identyfikacji jest niejednakowa. Postanowiono wybrać liczbę 16 jako kompromisową wartość liczby rozkładów Gaussa, pozwalającą na najlepszą generalizację.

TABELA 1

Wyniki skuteczności prawidłowej identyfikacji mówcy dla zbioru 22 cech dystynktywnych

| Segment uczący | 50 osób (25 kobiet + 25 mężczyzn) | | 100 osób (50 kobiet + 50 mężczyzn) | | 200 osób (100 kobiet + 100 mężczyzn) | |
|----------------|--------------------------------------|------------------------|---------------------------------------|------------------------|---|------------------------|
| | prawidłowo rozpoznane | procentowa skuteczność | prawidłowo rozpoznane | procentowa skuteczność | prawidłowo rozpoznane | procentowa skuteczność |
| TIMIT | 49 | 98,0% | 91 | 91,0% | 176 | 88,0% |
| NIST 60 | 48 | 96,0% | 82 | 82,0% | 154 | 77,0% |
| NIST 25 | 41 | 82,0% | 65 | 65,0% | 123 | 61,5% |

Ostatnim etapem prowadzonych badań było zastosowanie uniwersalnego modelu głosów do utworzenia poszczególnych modeli głosów powstałych z wyselekcjonowanych cech dystynktywnych w procesie uczenia. Do jego utworzenia

wykorzystano głosy pochodzące z bazy TIMIT niebiorące udziału w testowaniu (92 głosy żeńskie i 338 głosów męskich). Dokonano również próby wykorzystania głosów pochodzących z bazy NIST do utworzenia uniwersalnego modelu głosów, jednak otrzymano znacznie słabszą skuteczność identyfikacji, spowodowaną najprawdopodobniej znacznie niższą jakością nagrań w porównaniu do bazy TIMIT. Ostateczna skuteczność identyfikacji dla prezentowanego systemu ASR przedstawiona została w tabeli 1.

8. Wnioski

Przeprowadzone badania pozwoliły wyselekcjonować zbiór cech dystynktywnych pozwalających na skuteczną identyfikację mówców. Zredukowano wymiar wektora z 60 do 22 cech, co stanowi znaczący zysk pamięciowy oraz zwiększa szybkość obliczeń. Dzięki zastosowaniu algorytmu genetycznego podczas selekcji uwzględniono synergię cech, co znacząco przełożyło się na wypadkową skuteczność systemu względem wcześniejszych prób selekcji cech z wykorzystaniem miary Fishera i techniki PCA (ang. *Principal Component Analysis*) [5]. Dalsze prace wiążą się z koniecznością wykonania stosownej normalizacji wyników klasyfikatora GMM, a następnie utworzenia układu decyzyjnego umożliwiającego odrzucenie mówcy nieznajdującego się w zamkniętej bazie głosów.

Artykuł opracowany na podstawie referatu wygłoszonego na XXVIII Międzynarodowej Konferencji Naukowo-Technicznej „Ekomilitaris 2014”, Zakopane, 9-12.09.2014 r.

Wyniki przedstawione w artykule są efektem badań finansowanych w ramach działalności statutowej Wydziału Elektroniki Wojskowej Akademii Technicznej.

Artykuł wpłynął do redakcji 5.11.2014 r. Zweryfikowaną wersję po recenzjach otrzymano 20.10.2015 r.

LITERATURA

- [1] OSOWSKI S., *Metody i narzędzia eksploracji danych*, BTC, Legionowo, 2013.
- [2] GAROFOLO J.S. et al., *TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1*, Linguistic Data Consortium, Philadelphia, 1993.
- [3] MARTIN A., PRZYBOCKI M., *2002 NIST Speaker Recognition Evaluation LDC2004S04*, Linguistic Data Consortium, Philadelphia, 2004.
- [4] BROOKES M., *VOICEBOX: Speech Processing Toolbox for MATLAB*, <http://www.ee.ic.ac.uk/np/staff/dmb/voicebox/voicebox.html>, 2002.
- [5] KAMIŃSKI K., MAJDA E., DOBROWOLSKI A.P., *Automatic speaker recognition using Gaussian Mixture Models*, 17th IEEE SPA Conference, 2013, 220-225.
- [6] DOBROWOLSKI A.P., MAJDA E., *Cepstral analysis in the speakers recognition systems*, 15th IEEE SPA Conference, 2011, 85-90.

- [7] LUDWIG O., NUNES U., *Novel Maximum-Margin Training Algorithms for Supervised Neural Networks*, IEEE Transactions on Neural Networks, t. 21, nr 6, 2010, 972-984.
- [8] REYNOLDS D.A., QUATIERI T.F., DUNN R.B., *Speaker Verification Using Adapted Gaussian Mixture Models*, Digital Signal Processing, nr 10, 2000, 19-41.
- [9] GOLDBERG D.E., *Algorytmy genetyczne i ich zastosowanie*, WNT, Warszawa, 2003.

K. KAMIŃSKI, A.P. DOBROWOLSKI, E. MAJDA-ZDANCEWICZ

Selection of individual features of a speech signal using genetic algorithms

Abstract. The paper presents an automatic speaker's recognition system, implemented in the Matlab environment, and demonstrates how to achieve and optimize various elements of the system. The main emphasis was put on features selection of a speech signal using a genetic algorithm which takes into account synergy of features. The results of optimization of selected elements of a classifier have been also shown, including the number of Gaussian distributions used to model each of the voices. In addition, for creating voice models, a universal voice model has been used.

Keywords: biometrics, automatic speaker recognition, genetic algorithms, feature selection

DOI: 10.5604/12345865.1197999