

Žodžio pradžios ir galo nustatymas atpažįstant atskirai sakomus žodžius

G. Tamulevičius

Matematikos ir informatikos institutas,

A. Goštauto g. 12, LT-01108 Vilnius, Lietuva, el. p. g.tamulevicius@mch.mii.lt

A. Lipeika

Vilnius Gedimino technikos universitetas,

Naugarduko g. 41, LT-03227 Vilnius, Lietuva, el. p. lipeika@ktl.mii.lt

Įvadas

Pirmas automatinio atskirai sakomų žodžių atpažinimo etapas yra žodžio pradžios ir galo taškų nustatymas. Nuo žodžio pradžios ir galo taškų nustatymo tikslumo labai priklauso kalbos atpažinimo tikslumas. Mūsų atlikti eksperimentai parodė, kad apie du trečdalius atpažinimo klaidų padaroma dėl klaidingo žodžio pradžios ir galo taškų nustatymo.

Dažniausiai, atpažįstant kalbą, žodžių pradžios ir galo taškai yra nustatomi pagal kalbos signalo energiją [1]. Trumpalaikė signalo energija yra lyginama su iš anksto parinktu slenksčiu ir kai signalo energija peržengia slenksį, laikoma, kad prasidėjo kalbos signalas. Galo taškas nustatomas panašiu būdu, nagrinėjant kalbos signalą nuo galo. Slenksčio vertė priklauso nuo fono lygio. Tačiau šis algoritmas yra jautrus fono svyravimams ir įvairiems pašaliniams garsams. Klaidų dėl pašalinių garsų skaičiui sumažinti taikomi įvairūs euristiniai metodai, besiremiantys nesudėtinga logika. Vienas iš galimų sprendimų – išrinkti po keletą pradžios ir galo taškų kandidatų ir kalbos atpažinimo metu tikrinti visas galimas hipotezes [2,3]. Tačiau dėl to gerokai padidėja skaičiavimų apimtis, kadangi taip daroma lyginant su kiekvienu etalonu.

Pradžios ir galo taškams surasti taip pat naudojami ir kiti požymiai. Iš jų galima paminėti trumpalaikę autokoreliacinę [4] arba spektrinę [5] analizę. Žodžio galams nustatyti taip pat taikomi hipotezių tikrinimo metodai [6]. Tačiau dėl įvairių priežasčių šie metodai nepaplito ir žodžio ribų suradimo uždavinys tebėra aktualus.

Žodžio ribų nustatymas iš kalbos signalo

Į kalbos signalą galima žiūrėti kaip į atsitiktinį besikeičiančių savybių signalą. Remdamiesi šia prielaida, žodžio galams nustatyti pabandėme taikyti atsitiktinių sekų savybių pasikeitimo momentų nustatymo metodą [7], laikydami, kad nagrinėjamame signale yra du pasikeitimo momentai – žodžio galai. Nagrinėjamas signalas yra

vaizduojamas autoregresijos (AR) modeliu. Tikėtumo funkcijai maksimizuoti taikomas dinaminio programavimo metodas. Kadangi nagrinėjamo signalo modelio parametrai nežinomi, kaip pradinius parametrų įverčius naudojome parametrų įverčius iš fiksuoto ilgio atkarpų signalo pradžioje ir gale, o kaip žodžio modelį – visą likusį signalą. Su šiais parametrais gaunami pradiniai žodžio galo taškų įverčiai. Toliau skaičiuojama iteratyviai. Panaudojus pradinius galo taškų įverčius, iš naujo įvertinami AR modelio parametrai, nustatomi galo taškai ir tikrinama, ar galo taškų įverčiai pasikeitė. Iteratyvus procesas tęsiasi tol, kol galo taškų įverčiai nustoja keistis.

Pirmiausia nagrinėkime žodžio galų nustatymo uždavinį, kai modelio parametrai yra žinomi. Paskui jį pritaikysime nežinomiems modelio parametrams.

Uždavinio formulavimas, kai modelio parametrai žinomi

Nagrinėkime atsitiktinę seką $x = \{x(1), x(2), \dots, x(N)\}$, kuri yra tiesinės, diskretinės, laikui bėgant kintančios sistemos išėjimas. Sistema yra vaizduojama autoregresine skirtumine kintamų parametrų lygtimi:

$$x(n) = -a_1(n)x(n-1) - \dots - a_p(n)x(n-p) + b(n)v(n). \quad (1)$$

Pradžioje tarsime, kad parametrai yra žinomi ir tenkina sąlygą

$$A(n) = \begin{cases} A_1, & n = \dots, 1, 2, \dots, u_1, \\ A_2, & n = u_1 + 1, \dots, u_2, \\ A_3, & n = u_2 + 1, \dots, N; \end{cases} \quad (2)$$

čia $A(n) = [a_1(n), a_2(n), \dots, a_p(n), b(n)]$ yra AR modelio parametrai; $u = [u_1, u_2]$ – šuoliško parametrų pasikeitimo momentai, kurie tenkina sąlygą $p < u_1 < u_2 < N$. Mūsų uždavinys – rasti pasikeitimo momentų įverčius $\hat{u} = [\hat{u}_1, \hat{u}_2]$. Ieškosime labiausiai tikėtino parametrų pasikeitimo momentų įverčio.

Tikėtinumo funkcijos maksimizavimas

Pasikeitimo momentus $\hat{u} = [\hat{u}_1, \hat{u}_2]$ rasime maksimizuodami tikėtinumo funkcijos logaritmą

$$\hat{u} = \arg \max_u \log p(x|u), \quad (3)$$

kur [7]

$$\begin{aligned} \log p(x|u) = & \log p(x(1), \dots, x(p)) - \\ & - (N-p)/2 \log(2\pi) - (u_1-p) \log b(1) - \\ & - (u_2-u_1) \log b(2) - (N-u_2) \log b(3) - \\ & - \frac{1}{2b^2(1)} \sum_{n=p+1}^{u_1} \left[\sum_{j=0}^p a_j(1)x(n-j) \right]^2 - \\ & - \frac{1}{2b^2(2)} \sum_{n=u_1+1}^{u_2} \left[\sum_{j=0}^p a_j(2)x(n-j) \right]^2 - \\ & - \frac{1}{2b^2(3)} \sum_{n=u_2+1}^N \left[\sum_{j=0}^p a_j(3)x(n-j) \right]^2. \end{aligned} \quad (4)$$

Tiesiogiai maksimizuoti (4) išraiškos beveik neįmanoma dėl didelės skaičiavimų apimtys. Galima įrodyti, kad jos maksimumo vieta sutampa su kitos tikslo funkcijos $\Theta(u|x) = L_1(u_1|x) + L_2(u_2|x)$ maksimumo vieta [7], kur

$$\begin{aligned} L_i(k|x) = & -(k-p) \log b(i) - (N-k) \log b(i+1) - \\ & - \frac{1}{2b^2(i)} \sum_{n=p+1}^k \left[\sum_{j=0}^p a_j(i)x(n-j) \right]^2 - \\ & - \frac{1}{2b^2(i+1)} \sum_{n=k+1}^N \left[\sum_{j=0}^p a_j(i+1)x(n-j) \right]^2. \end{aligned} \quad (5)$$

Čia $i = 1, 2; k = p+1, 2, \dots, N$.

(5) išraišką galima skaičiuoti rekurentiškai:

$$\begin{aligned} L_i(k|x) = & L_i(k-1|x) - \log b(i) + \log b(i+1) - \\ & - \frac{1}{2b^2(i)} \left[\sum_{j=0}^p a_j(i)x(n-j) \right]^2 + \\ & + \frac{1}{2b^2(i+1)} \left[\sum_{j=0}^p a_j(i+1)x(n-j) \right]^2. \end{aligned} \quad (6)$$

Čia $i = 1, 2; k = 2, \dots, N$.

Kadangi pradinės sąlygos nepriklauso nuo pasikeitimo momentų, (6) išraiškai spręsti galima naudoti nulines pradines sąlygas.

Funkcija $\Theta(u|x)$ yra suma dviejų funkcijų, kurių kiekviena priklauso tik nuo vieno kintamojo. Todėl funkcijai $\Theta(u|x)$ maksimizuoti galima taikyti dinaminio programavimo metodą. Apskaičiuokime Belmano funkcijas [7]:

$$g_1(u_2|x) = \max_{\substack{u_1 \\ p < u_1 < u_2}} L_1(u_1|x) \quad u_2 = p+2, \dots, N, \quad (7)$$

$$\begin{aligned} g_2(u_3|x) = & \max_{\substack{u_2 \\ p+1 < u_2 < u_3}} [L_2(u_2|x) + g_1(u_2|x)], \\ & u_3 = p+3, \dots, N. \end{aligned} \quad (8)$$

Ieškomus pasikeitimo momentus surandame iš Belmano funkcijų

$$\hat{u}_k = \min [\arg \max_n g_k(n|x)], \quad k = 2, 1; \quad (9)$$

čia $\hat{u}_3 = N$.

Iki šiol laikėme, kad AR modelio parametrai yra žinomi. Mums reikia pritaikyti šį algoritmą esant nežinomiems modelio parametrams.

Tikėtinumo funkcijos maksimizavimas, kai modelio parametrai nežinomi

Tikėtinumo funkcijai maksimizuoti, kai AR modelio parametrai nežinomi, naudojame apibendrintą matematinės vilties maksimizavimo algoritmą [8].

Naudojant apibendrintą matematinės vilties maksimizavimo algoritmą, reikalingi pradiniai nežinomų parametru įverčiai. Todėl fono pradiniais parametru įverčiais ėmėme nežinomų parametru įverčius, apskaičiuotus iš fiksuoto ilgio atkarpų signalo pradžioje ir gale. Žodžio pradiniais parametru įverčiais laikėme nežinomų parametru įverčius, apskaičiuotus iš viso likusio signalo. Panaudojus šiuos parametrus, gaunami pradiniai žodžio galo taškų įverčiai. Toliau skaičiuojama iteratyviai. Panaudojus pradinius galo taškų įverčius, iš naujo įvertinami AR modelio parametrai, anksčiau aprašytu būdu nustatomi galo taškai (pasikeitimo momentai) ir tikrinama, ar galo taškų įverčiai pasikeitė. Iteratyvus procesas tęsiasi tol, kol galo taškų įverčiai nustoja keistis. Taip iteratyviai yra maksimizuojamas tikėtinumas, kol tikėtinumo funkcija palaipsniui konverguoja į kritinį tašką.

Žodžio galų nustatymas iš segmentuoto kalbos signalo energijos

Iki šiol nagrinėjome situaciją, kai žodžio galams nustatyti tiesiogiai naudojamas pats kalbos signalas. Alternatyvus ir plačiausiai taikomas žodžio galų nustatymo būdas [1] remiasi kalbos signalo skaidymu į segmentus (kadrus) ir signalo energijos kadruose matavimu. Pasiūlytą žodžio galų nustatymo būdą galima taikyti ir skaidant kalbos signalą į kadrus.

Tarkim, $x = \{ x(1), x(2), \dots, x(N) \}$ yra signalo energija nuosekliai sunumeruotuose signalo kadruose ir energijos vertės kadruose yra nepriklausomi normalieji atsitiktiniai dydžiai. Tada (2) išraišką galime perrašyti taip:

$$A(n) = \begin{cases} A_1 = N(\mu_1, \sigma_1^2), & n = \dots, 1, 2, \dots, u_1, \\ A_2 = N(\mu_2, \sigma_2^2), & n = u_1 + 1, \dots, u_2, \\ A_3 = N(\mu_3, \sigma_3^2), & n = u_2 + 1, \dots, N; \end{cases} \quad (10)$$

čia μ_i, σ_i^2 , kur $i=1,2,3$ yra signalo energijos kadruose vidurkiai ir dispersijos; A_1 ir A_3 reiškia foninio triukšmo energijos parametrus kadruose prieš žodį ir po žodžio; A_2 reiškia žodžio energijos parametrus kadruose.

Taip formuluojant uždavinį, (4) išraišką galima perrašyti:

$$\begin{aligned} \log p(x|u) = \log \prod_{n=1}^N p(x(n)) = & -\frac{N}{2} \log(2\pi) - \\ & -u_1 \log \sigma_1 - (u_2 - u_1) \log \sigma_2 - (N - u_2) \log \sigma_3 - \\ & -\frac{1}{2\sigma_1^2} \sum_{n=1}^{u_1} [x(n) - \mu_1]^2 - \frac{1}{2\sigma_2^2} \sum_{n=u_1+1}^{u_2} [x(n) - \mu_2]^2 - \\ & -\frac{1}{2\sigma_3^2} \sum_{n=u_2+1}^N [x(n) - \mu_3]^2. \end{aligned} \quad (11)$$

Tada dalinių tikėtinumo funkcijų (5) išraiška įgauna tokią formą:

$$\begin{aligned} L_i(k|x) = & -k \log \sigma_i - (N - k) \log \sigma_{i+1} - \\ & -\frac{1}{2\sigma_i^2} \sum_{n=1}^k [x(n) - \mu_i]^2 - \\ & -\frac{1}{2\sigma_{i+1}^2} \sum_{n=k+1}^N [x(n) - \mu_{i+1}]^2; \end{aligned} \quad (12)$$

čia $i=1,2$; $k=1,2,\dots,N$.

(12) kaip ir (6) išraišką galima skaičiuoti rekurentiškai, imant nulines pradines sąlygas:

$$\begin{aligned} L_i(k|x) = & L_i(k-1|x) - \log \sigma_i + \log \sigma_{i+1} - \\ & -\frac{1}{2\sigma_i^2} [x(n) - \mu_i]^2 + \frac{1}{2\sigma_{i+1}^2} [x(n) - \mu_{i+1}]^2. \end{aligned} \quad (13)$$

Čia $i=1,2$; $k=2,\dots,N$.

Skaičiavimai atliekami tokia pat tvarka kaip ir nustatant žodžio galus iš signalo, tik čia $x(n)$ yra signalo energija kadruose ir vietoj (2), (4), (6) išraiškų naudojamos (10) – (13).

Ekspirimentai

Preliminarūs eksperimentai buvo atlikti naudojant 50 skirtingų žodžių ištarių. Pirmas ištariškas buvo naudojamas žodžių etalonams sudaryti, kiti trys – kaip testiniai. Ieškant žodžio galo taškų tradiciniu metodu, paremtu signalo trumpalaikės energijos lyginimu su slenksčiu, buvo gauta 12 % atpažinimo klaidų.

Surandant žodžio galų taškus mūsų pasiūlytu metodu (žodžio galų nustatymas iš kalbos signalo), kai buvo naudojamas nulinės eilės AR modelis, kalbos signalo vertės nėra prognozuojamos, klaidų skaičius buvo 4 %. Naudojant aukštesnės eilės AR modelį, klaidų skaičius labai padidėja. Pavyzdžiui, naudojant 4-os eilės modelį, klaidų skaičius padidėja iki 26 %.

Ieškant žodžio galų taškų kitu mūsų pasiūlytu metodu (žodžio galų nustatymas iš segmentuoto kalbos signalo

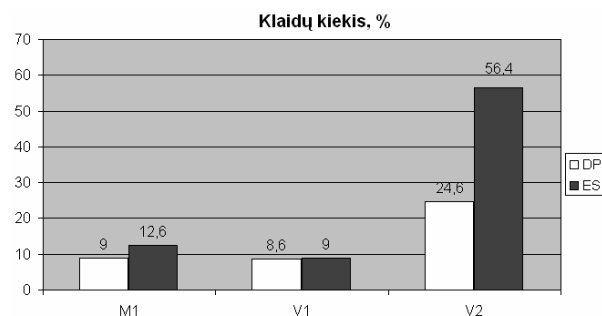
energijos) klaidų, taip pat buvo 4 %, tačiau naudojant segmentuotą kalbos signalą, labai sumažėja skaičiavimų apimtis, todėl šį metodą patogiaus taikyti. Kadangi nulinės eilės AR modelis nepagerino žodžio galų nustatymo rezultatų ir praktiškai žodžio galus nustatyti yra lengviau taikant dinaminį programavimą segmentuotam kalbos signalui, tolesniuose eksperimentuose lyginome šio metodo darbingumą su tradiciniu metodu, paremtu signalo trumpalaikės energijos lyginimu su slenksčiu.

Šiame eksperimente buvo naudojama 50 skirtingų lietuvių kalbos žodžių balsų bazė. Viena moteris ir du vyrai ištarė šiuos žodžius po 11 kartų. Atskirai sakomų žodžių atpažinimo sistema [9], kurioje taikomi abu šie žodžio galų nustatymo metodai, buvo naudojama eksperimentams. Kiekvieno kalbėtojo pirmos sesijos balso įrašai buvo naudojami žodžių etalonams sudaryti, o 10 sesijų (500 žodžių) – testui. Eksperimento rezultatai pateikti 1 lentelėje.

1 lentelė. Žodžio galų nustatymo eksperimento rezultatai. DP – dinaminio programavimo metodas, ES – slenksčio metodas

Kalbėtojas	Klaidų skaičius, %					
	M1 (mot.)		V1 (vyr.)		V2 (vyr.)	
Metodas	DP	ES	DP	ES	DP	ES
I sesija	10	18	2	6	26	56
II sesija	10	12	2	8	24	56
III sesija	6	4	4	6	18	52
IV sesija	8	24	2	4	34	56
V sesija	8	16	10	6	18	56
VI sesija	6	12	16	14	28	50
VII sesija	6	14	14	16	22	60
VIII sesija	10	8	12	10	22	60
IX sesija	14	14	12	12	26	60
X sesija	12	4	12	8	28	58
Vidutinė klaida	9	12.6	8.6	9	24.6	56.4

Žodžio galų nustatymo eksperimento rezultatai, suvidurkinti per 10 sesijų, pateikti 1 pav.



1 pav. Vidutinis atpažinimo klaidų skaičius procentais lyginant pasiūlytą žodžio galų nustatymo metodą su slenksciniu žodžio galų nustatymo metodu. Kalbėtojai M1, V1 ir V2. DP – dinaminio programavimo metodas, ES – slenksčio metodas

Kalbėtojų M1 ir V1 balso įrašai yra geros kokybės ir turi žemą triukšmo lygį. Kalbėtojo V2 balso įrašų triukšmo lygis aukštas ir juose yra daug pašalinių garsų. Tai pagrindinė priežastis, kodėl V2 balso įrašuose buvo daug klaidų. Tačiau taikant dinaminio programavimo metodą,

klaidų visuose balso įrašuose, ypač daug triukšmo turinčiuose V2 balso įrašuose, gauta mažiau.

Išvados

Buvo tiriamas dinaminio programavimo metodo taikymas žodžio galams nustatyti atpažįstant atskirai sakomus žodžius. Pirmame tyrimų etape kalbos ir foninio triukšmo spektrams modeliuoti buvo naudojamas tiesinės prognozės modelis, bet eksperimentinis tyrimas parodė, kad tiesinės prognozės modelis netinka žodžio galams nustatyti.

Atliekant tolesnius tyrimus, kalbos signalas buvo skaidomas į segmentus (kadrus) ir žodžio galams nustatyti buvo naudojama signalo energija segmentuose. Žodžio galų nustatymas maksimizuojant tikėtimumo funkciją taip pat rėmėsi dinaminio programavimu. Sukurto algoritmo darbingumas buvo vertinamas pagal žodžių atpažinimo klaidų skaičių. Algoritmo darbingumo tyrimai parodė, kad:

- esant žemo triukšmo lygio signalams, žodžio galai nustatomi dinaminio programavimo metodu truputį tiksliau negu slenksčio metodu;
- pasiūlyto metodo pranašumas labiau išryškėja esant aukšto triukšmo lygio signalams.

Be to, taikant pasiūlytą metodą nereikia iš anksto parinkti jokio slenksčio, nes automatiškai prisitaikoma prie įrašymo aplinkos triukšmo lygio. Šis žodžio galų nustatymo metodas pritaikytas kuriant atskirai sakomų žodžių atpažinimo sistemą.

G. Tamulevičius, A. Lipeika. Žodžio pradžios ir galo nustatymas atpažįstant atskirai sakomus žodžius // Elektronika ir elektrotechnika. Kaunas: Technologija, 2005. – Nr. 2(58). – P. 61–64.

Nagrinėjamas dinaminio programavimo metodo taikymas žodžio galams nustatyti atpažįstant atskirai sakomus žodžius. Žodžio galų nustatymas remiasi tikėtimumo funkcijos maksimizavimu. Nežinomų parametrų problemai spręsti taikomas matematinės vilties maksimizavimo principas. Sprendimų priėmimas remiasi kalbos signalo ir foninio triukšmo energija. Pasiūlyto metodo darbingumas tiriamas naudojant atskirai sakomų lietuvių kalbos žodžių duomenų bazę. Algoritmo darbingumo tyrimai parodė, kad esant žemo triukšmo lygio signalams žodžio galai nustatomi dinaminio programavimo metodu truputį tiksliau negu slenksčio metodu. Tačiau pasiūlyto metodo pranašumas labiau išryškėja esant aukšto triukšmo lygio signalams. Be to, taikant pasiūlytą metodą, nereikia iš anksto parinkti jokio slenksčio, nes automatiškai prisitaikoma prie įrašymo aplinkos triukšmo lygio. Šis žodžio galų nustatymo metodas pritaikytas kuriant atskirai sakomų žodžių atpažinimo sistemą. Il. 1, bibl. 9 (lietuvių kalba; santraukos lietuvių, anglų ir rusų k.).

G. Tamulevičius, A. Lipeika. On Endpoint Detection in Isolated Word Recognition // Electronics and Electrical Engineering. – Kaunas: Technologija, 2005. – No. 2(58). – P. 61–64.

The paper deals with the use of dynamic programming for word endpoint detection in isolated word recognition. Endpoint detection is based on likelihood maximization. Expectation maximization approach is used to deal with the problem of unknown parameters. Speech signal and background noise energy is used as features for making decision. Performance of the proposed approach was evaluated using isolated Lithuanian words speech corpus. Performance evaluation based on recognition error rate showed that for low noise level endpoint detection based on dynamic programming slightly outperforms threshold based endpoint detection. For high background noise level preference of dynamic programming based algorithm is more noticeable. The main advantage of dynamic programming based approach is that this method does not need any threshold. This endpoint detection method is applied in development of isolated word recognition system. Ill. 1, bibl. 9 (in Lithuanian; summaries in Lithuanian, English and Russian).

Г. Тамулявичюс, А. Липейка. Об определении концов слова при распознавании изолированных слов // Электроника и электротехника. – Каунас: Технологія, 2004. – № 2(58). – P. 61–64.

Рассматривается применение динамического программирования для определения концов слова при распознавании изолированных слов. Определение концов слова базируется на максимизации функции правдоподобия. Для решения проблемы неизвестных параметров применяется метод максимизации математического ожидания. Принятие решения базируется на энергии речевого сигнала и фонового шума. Определение работоспособности предложенного метода производится с использованием базы данных изолированных слов литовского языка. Исследование работоспособности алгоритма показало, что для сигналов с низким уровнем шума определение концов слова с использованием динамического программирования дает несколько лучшие результаты по сравнению с методом, основанном на использовании порога. Преимущество предложенного метода выявляется при большом уровне фонового шума. Кроме того, предложенный метод не требует использования никакого заранее выбранного порога, метод автоматически приспосабливается к уровню фонового шума. Данный метод использован при создании системы распознавания изолированных слов. Ил. 1, библи. 9 (на литовском языке; рефераты на литовском, английском и русском яз.)

Literatūra

1. **Rabiner L., Juang B.-H.** Fundamentals of speech recognition. – Prentice Hall, 1993.
2. **Lamel L., Rabiner L., Rosenberg A., Wilpon J.** An improved endpoint detector for isolated word recognition // IEEE Trans. Acoustics, Speech, Signal Proc. – 1981. – Vol. c-29 (4). – P. 777–785.
3. **Young S., Kershaw D., Odell J., Ollason D., Valtchev V., Woodland P.** The HTK Book, Version 3.0. – Microsoft Corporation, 2000.
4. **Zhu J., Chen F.-I.** The analysis and application of a new endpoint detection method based on distance of autocorrelated similarity // Proceedings of the Eurospeech'99. – 1999. – P. 105–108.
5. **Wu G.-D., Lin Ch.-T.** Word boundary detection with Mel-scale frequency bank in noisy environment // IEEE Transactions on Speech and Audio Processing. – 2000. – Vol. c-8. – P. 541–555.
6. **Zelinski R., Class F.** A segmentation algorithm for connected word recognition based on estimation principles // IEEE Transactions on Acoustics, Speech, Signal Processing. – 1983. – Vol. c-31 (4). – P. 818–827.
7. **Lipeika A.** Optimal segmentation of random sequences // Informatica. – 2000. – Vol. c-11 (3). – P. 243–256.
8. **Duda R., Hart P., Stork D.** Pattern classification. – John Wiley & Sons, INC., 2001.
9. **Tamulevičius G., Lipeika A.** Žodžių atpažinimo sistemos kūrimas // Lietuvos matematikos rinkinys. – 2003. – 43 (spec. nr.). – P. 292–296.

Pateikta spaudai 2004 06 15