

ارائه درونیابی KNNGI^۱ و مقایسه آن با درونیابی FI^۲ در بازشناسی گفتار

ابوالقاسم صیادیان^{*}، کامبیز بدیع^{*}، محمد شهرام معین^{**} و نصرالله مقدم^{**}

دانشکده برق دانشگاه صنعتی امیرکبیر

مرکز تحقیقات مخابرات ایران

دانشکده فنی و مهندسی، دانشگاه تربیت مدرس

(دریافت مقاله: ۸۱/۹/۱۷ - دریافت نسخه نهایی: ۸۳/۳/۲۵)

چکیده - مدل سازی آماری HMM رویکردی پرکاربرد در سیستمهای بازشناسی گفتار پیوسته و گسسته است. توزیع احتمال بردارهای مشاهدات هر حالت پنهان مدل، به دو روش پیوسته^۳ یا گسسته^۴ تخمین زده می شوند. عملکرد توزیع احتمال پیوسته (با مدل سازی GMM^۵) بالاتر از عملکرد توزیع احتمال گسسته (با مدل سازی VQ^۶) است. ولی چنانچه بخواهیم از رویکرد HMM برای بازشناسی گفتار گسسته با دایره لغات وسیع استفاده کنیم، هزینه محاسباتی مرحله بازشناسی با افزایش تعداد لغات، به نحو چشمگیری افزایش می یابد. بدین لحاظ در بازشناسی گفتار گسسته با دایره لغات وسیع، از توزیع احتمال گسسته به منظور کاهش هزینه محاسباتی و امکان پیاده سازی بی درنگ^۷ استفاده می شود. برای جبران کاهش دقت و عملکرد مدل سازی DD-HMM، استفاده از درونیابی فازی FI مرسوم است. در این تحقیق روش درونیابی گوسی که دارای پشتوانه نظری قوی تر نسبت به FI است ارائه کرده ایم. کارایی دو روش درونیابی KNNGI و FI در بازشناسی ۱۵۰۰ کلمه فارسی مورد تحقیق و بررسی قرار دادیم. نتایج این تحقیق نشان می دهد که دقت و انعطاف پذیری درونیابی KNNGI بیشتر از روش FI است.

واژگان کلیدی: درونیابی گوسی، درونیابی فازی، مدل مارکف مخفی چگالی گسسته، بازشناسی تلفظ گسسته

Presentation of K Nearest Neighbor Gaussian Interpolation and Comparing it with Fuzzy Interpolation in Speech Recognition

A. Sayadiyan, K. Badi, M. Moin and N. Moghadam

Department of Electrical Engineering, AmirKabir University of Technology

Abstract: Hidden Markov Model is a popular statistical method that is used in continuous and discrete speech recognition. The probability density function of observation vectors in each state is estimated with discrete density or continuous density modeling. The performance (in correct word recognition rate) of continuous density is higher than discrete density HMM, but its computation complexity is very high, especially in very large discrete utterance recognition problems. For real time implementation of very large discrete utterance recognition, we must use discrete density HMM (DDHMM). To increase the performance of DDHMM, one usual solution is fuzzy interpolation. In this study, we present a new method named Gaussian interpolation. We implemented and compared the performance of two types of interpolation methods for 1500 Persian speech command words. Results show that precision and flexibility of Gaussian interpolation is better than those of the fuzzy interpolation.

Keywords: Gaussian Interpolation, Fuzzy Interpolation, Discrete Density HMM, Discrete utterance recognition

** - استادیار

* - دانشیار

q	متغیر نماینده حالت	a	متغیر نماینده احتمال وقوع حالات شرطی
t	متغیر نماینده توالی زمانی	b	متغیر نماینده احتمال وقوع هر بردار مرجع در هر حالت
V	متغیر نماینده نام بردار کتاب گد	c	نام یک کلاس اکوستیکی
W	متغیر نماینده وزن دادن	d	نام تابع فاصله اقلیدس
π	متغیر نماینده احتمال وقوع حالت در شروع	m	متغیر نماینده مقدار فازی بودن
α	متغیر نماینده توابع عضویت نمایش فازی	N	متغیر نماینده توزیع نرمال
σ	متغیر نماینده واریانس توزیع نرمال	o	بردار ویژگی مشاهده شده
		p	متغیر نماینده مقدار احتمال گوسی

۱- مقدمه

حالت (سه حالت برای تشخیص واجها و $N \geq 5$ حالت برای تشخیص کلمات و عبارات) مدل‌سازی می‌شوند. اطلاعات آماری مربوط به قواعد توالی زمانی تولید حالتها، توسط ماتریس A (با بعد $N \times N$) و بردار حالت اولیه π (با بعد $1 \times N$) بازنمایی می‌شوند. تغییرات آماری ویژگیهای طیفی در داخل هر حالت هر کلاس توسط یک pdf^{۱۶} پیوسته یا گسسته مدل‌سازی می‌شوند. روشهای تخمین (در فاز آموزش) و روشهای استفاده (در فاز بازشناسی) از ماتریس A و بردار π در هر دو روش CD-HMM و DD-HMM تقریباً یکسان است. وجه تمایز اصلی دو روش در نحوه مدل‌سازی pdf بردارهای طیفی در داخل حالات است. مدل‌سازی پیوسته pdf در سیستمهای بازشناسی پیوسته CSR^{۱۷} که تعداد کلاسهای اکوستیکی کم و محدود است. (متناظر با تعداد واجهای هر زبان) (حدود ۳۰ واج برای فارسی و حدود ۴۴ واج برای زبان انگلیسی)، بسیار مورد توجه است. ولی در سیستمهای بازشناسی گفتار گسسته با تعداد کلاسهای اکوستیکی بالا (حدوداً بالای ۱۰۰۰ لغت)، استفاده از pdf پیوسته به ویژه در فاز بازشناسی از نظر پیاده‌سازی بی‌درنگ مقرون به صرفه و عملی نیست. در این تحقیق تعداد کلاسهای اکوستیکی بالاتر از ۱۰۰۰ لغت و برابر ۱۵۰۰ کلمه است. در چنین مواردی از بازنمایی گسسته pdf (با استفاده از تکنیک VQ) یا از بازنمایی نیمه پیوسته^{۱۸} استفاده می‌کند. چون بازنمایی گسسته دارای بالاترین سرعت در هنگام

رویکرد آماری HMM در کاربردهای عملی بازشناسی گفتار، رویکردی غالب نسبت به سایر روشها از جمله ANN^۱ و DTW^۲ است. از ویژگیهای بارز رویکرد HMM استفاده بهینه از اطلاعات و ویژگیهای طیفی و زبانی، در مدل‌سازی مؤثر و تشخیص صحیح وقایع اکوستیکی متفاوت و متمایز است. HMMها هم برای گفتار پیوسته^۱ و هم برای گفتار گسسته^{۱۱} مورد استفاده قرار می‌گیرند. در کاربردهای بازشناسی گفتار پیوسته، علاوه بر مدل‌سازی اکوستیکی^{۱۲} به مدل‌سازی زبانی^{۱۳} نیز نیازمندیم. با توجه به اینکه در زبان فارسی پایگاه داده گفتاری جامع و پوشا برای مدل‌سازی ATP^{۱۴} و همچنین پایگاه متنی جامع و پوشا برای مدل‌سازی SLM^{۱۵} موجود نیست، توجه اصلی این تحقیق به کاربرد HMMها در بازشناسی گفتار گسسته معطوف شده است. هدف سیستم بازشناسی گفتار گسسته عبارت است از، تعیین کلاس یک سیگنال اکوستیکی دریافتی (با فرض اینکه محدوده تقریبی ابتدا و انتهای آن با دقت مناسبی مشخص است): همچنین پارامترهای مدل آماری M کلاس اکوستیکی متمایز، از قبل (در فاز آموزش) مشخص است. هر کلاس اکوستیکی می‌تواند یک فرمان یا درخواست صوتی یک تا حداکثر پنج کلمه‌ای باشد (در این تحقیق). در مدل‌سازی HMM، به منظور استفاده بهینه از اطلاعات زمانی به منظور متمایزسازی کلاسها، هر کلاس اکوستیکی توسط N

قابل مطالعه است. در فاز بازشناسی با استفاده از الگوریتم برگشتی ^{19}FP ، نمره (یا احتمال) هر کلاس با دریافت هر دنباله از بردار مشاهدات اکوستیکی $O_1 O_2 \dots O_T$ به شرح زیر محاسبه می‌شود:

قدم اول: محاسبه در لحظه $t=1$

$$a_1(l, i) = \pi_i b_i(l, o_1) \quad 1 \leq i \leq N \quad (4)$$

قدم دوم: مرحله تکراری

$$a_{t+1}(l, j) = \left\{ \sum_{i=1}^N a_t(l, i) a_{ij}(l) \right\} b_j(l, o_{t+1})$$

$$1 \leq t \leq T-1 \quad \text{و} \quad 1 \leq j \leq N \quad (5)$$

قدم سوم: محاسبه نمره نهایی هر کلاس

$$P(l) = P(O_1, \dots, O_T / C_l) = \sum_{i=1}^N a_T(l, i) \quad (6)$$

$1 \leq t \leq T-1$ و $1 \leq j \leq N$
مراحل فوق برای L تا $l=1$ تکرار می‌شود. کلاسی که بیشترین $P(l)$ را تولید کند به عنوان کلاس برنده اعلام می‌شود. البته از الگوریتم ویتربی [5] نیز برای محاسبه نمره نهایی هر کلاس می‌توان استفاده کرد.

با توجه به محدود بودن رنج دینامیکی پارامترهای A ، π در مدل‌سازی آماری هر کلاس، بسیاری از محققان اعتقاد دارند که نقش اصلی در بالا بردن عملکرد سیستم بازشناسی به بهبود مدل‌سازی و تخمین دقیق پارامترهای ماتریس B نهفته است. مدل‌سازی پیوسته عناصر ماتریس B توسط مدل GMM به عنوان بهترین راه حل شناخته شده است. اما به دلایلی که در این تحقیق بیان شد، از بازنمایی گسسته می‌خواهیم استفاده کنیم.

۳- روشهای طراحی کتاب کد برای مدل‌سازی گسسته pdf بردارهای طیفی

در بخش قبل بیان کردیم که برای محاسبه عناصر ماتریس B به تعدادی الگوی مرجع به نام کتاب کد برای هر کلاس اکوستیکی نیازمندیم. فرض می‌کنیم TV_1, \dots, TV_S مجموعه بردارهای آموزشی متعلق به یک حالت از یک کلاس اکوستیکی باشد. با استفاده از مجموعه بردارهای آموزشی و دو الگوریتم کلاسیک ^{20}LBG یا ^{21}EM می‌توان کتاب کد طراحی کرد [7] و

پیاده‌سازی است (به علت امکان استفاده از ساختارهای مناسب برای جستجوی سریع در مرحله کدینگ (VQ)، مورد توجه و عنایت این تحقیق قرار گرفته است. بنابراین روش پایه در این تحقیق DDHMM است. عملکرد مدل‌سازی DDHMM از نظر دقت پایینتر از مدل‌سازی DDHMM است. روش کلاسیک برای افزایش نسبی دقت روش بازشناسی DDHMM، استفاده از درون‌یابی فازی (FI) است. این روش‌ها در بخش (4) مراجع [8] و [9] مورد بررسی قرار گرفته‌اند. در این تحقیق دو نوع درون‌یابی جدید به نامهای KNNGI (درون‌یابی گوسی K نزدیکترین همسایه) و KNNFI (درون‌یابی فازی K نزدیکترین همسایه) در بخش (5) ارائه کرده و عملکرد آنها را در بخش (5-3) و بخش (6) و (7) مورد ارزیابی قرار دادیم. نتایج ارزیابی نشان می‌دهند که روش جدید درون‌یابی KNNGI دارای بهترین عملکرد در میان روشها بازشناسی DDHMM است.

۲- بررسی مدل‌سازی DD-HMM در بازشناسی گفتار گسسته

فرض می‌کنیم تعداد L کلاس اکوستیکی (کلمه یا عبارت) داریم. هر کلاس اکوستیکی توسط N حالت مدل‌سازی زمانی می‌شود. اطلاعات زمانی ماندن یا گذر از حالتی به حالتی دیگر برای هر کلاس به صورت زیر توصیف می‌شود:

$$A_l = \{a_{ij}(l)\} \quad \text{that} \quad a_{ij}(l) = P\{q_{t+1}(l) = j / q_t(l) = i\} \quad (1)$$

$$l=1 \text{ تا } L \quad 1 \leq i, j \leq N$$

$$\pi_l = \{\pi_i(l)\} \quad \text{that} \quad \pi_i(l) = P\{q_{t=1}(l) = i\} \quad (2)$$

$$l=1 \text{ تا } L \quad 1 \leq i \leq N$$

فرض می‌کنیم بردارهای مشاهدات هر حالت توسط M بردار مرجع کد می‌شوند، در این صورت توزیع احتمال گسسته حالات به صورت زیر توصیف می‌شود:

$$B_l = \{b_i(l, k)\} \quad \text{that} \quad b_i(l, k) = P\{o_t = v_k(l) / q_t(l) = i\} \quad (3)$$

$$l=1 \text{ تا } L \quad 1 \leq i \leq N \quad 1 \leq k \leq M$$

روشهای تخمین پارامترهای فوق‌الذکر از روی نمونه‌های آموزشی (در فاز آموزش) در مراجع متعدد از جمله [1] و [5]

چندی کردن گسسته ویژگیهای طیفی به تعداد معینی سطوح، نویز چندی شده کارایی سیستم را از نظر دقت بازشناسی به ویژه برای تلفظها و گویندگانی که در مرحله آموزش شرکت نداشته اند، پایین می آورد. کاهش دقت سیستم بازشناسی چگالی گسسته نسبت به چگالی پیوسته امری مسلم و ثابت شده برای محققان است. تلاش دانشمندان برای افزایش کارایی مدل DD-HMM منجر به ابداع روشهای FVQ، SCD شد. این روشها در حد قابل قبولی نقطه ضعف DD-HMM را از نظر دقت بازشناسی برطرف کرده اند. در روش SCD، ابتدا یک کتاب کد همانند حالت چگالی گسسته طراحی می شود. سپس با استفاده از مجموعه بردارهای آموزشی نسبت داده شده به هر کلمه کد، یک چگالی پیوسته با تعداد مخلوطهای محدودی (نوعاً یک تا شش) توسط مدل سازی GMM طراحی می شود. استفاده از مدل سازی SCD موجب افزایش دقت سیستم در حد مدل سازی پیوسته CD شده است. [۲ و ۴]. اگر چه استفاده از روش SCD بسیار مفید است، ولی چون در عمل مجبوریم تعداد کلمات کد کتاب کد اصلی (DD) را کم انتخاب کنیم، دیگر قادر نیستیم از کارایی مدل DD به صورت کامل (برای ایجاد ساختارهای جستجوی سریع) استفاده کنیم. در نتیجه از هدف اصلی خود در به کارگیری مدل سازی گسسته pdf دور خواهیم شد. در روش FVQ علاوه بر استفاده از کلمه کد برنده (در مرحله بازشناسی) از تعداد I کد با کمترین فاصله نسبت به بردار مشاهده O_t نیز برای بالا بردن دقت سیستم به صورت زیر استفاده می کنند [۸ و ۹]. فرض می کنیم $V_{i1}, V_{i2}, \dots, V_{iI}$ کدهای برنده با کمترین فاصله نسبت به بردار مشاهده O_t باشند یعنی داریم:

$$d(v_{i1}, o_t) \leq d(v_{i2}, o_t) \leq \dots \leq d(v_{iI}, o_t) \quad (۸)$$

$$1 \leq j \leq M$$

$$j \neq i_1, \dots, i_I$$

در روش کلاسیک DD، مستقل از فواصل O_t نسبت به سایر کدها، نمره نهایی سیستم محاسبه می شود. در روش FVQ، مقدار $b_j(O_t)$ به صورت زیر درون یابی می شود:

در مدل DD کلاسیک: $b_j(O_t) = b_j(V_{ij})$

در مدل FVQ:

۱۰ - ۱۳]. فرض می کنیم تعداد کلمات کد هر کتاب (یا تعداد بردارهای مرجع) برابر $M \ll S$ باشد. همچنین فرض می کنیم ویژگیهای طیفی ناهمبسته باشند (در اغلب کاربردهای عملی از ویژگی MFCC و مشتقات مرتبه اول و دوم آن استفاده می شود که تقریباً شرایط مطرح شده را دارا هستند). در این صورت خروجی هر دو الگوریتم LBG و EM شامل M بردار میانگین V_i ($i=1$ تا M)، بردار واریانس VAR_i ($i=1$ تا M) (اگر ویژگیها ناهمبسته نباشند به جای بردار واریانس ماتریس کوواریانس تولید خواهد شد)، و تعداد M اسکالر وزنی W_i (M تا $i=1$) (با محدودیت $\sum_{i=1}^M w_i = 1$) خواهند بود. w_i ها در الگوریتم EM به عنوان ضرایب مخلوط و در الگوریتم LBG به عنوان فرکانس نسبی برنده شدن هر کد ذکر می شوند. در مدل سازی کلاسیک DD-HMM صرفاً از بردارهای متوسط V_i ($i=1$ تا M) در مرحله کدینگ و بازشناسی استفاده می شود. یعنی با دریافت هر بردار O_t آن را به یکی از M بردار V_i طبق منطق زیر نسبت می دهند:

$$\text{Code}(O_t) = V_k \quad \text{if} \quad k = \text{argmin} [d(V_1, O_t), d(V_2, O_t), \dots, d(V_M, O_t)]$$

$d(X, Y)$ تابع فاصله اقلیدسی بین دو بردار Y, X است. بار اصلی هزینه محاسباتی مرحله کدینگ تعیین M تابع فاصله اقلیدسی $d(O_t, V_i)$ ($i=1$ تا M) است. نکته برجسته در استفاده از روش بازشناسی با چگالی گسسته آن است که، با استفاده از ساختارهای جستجوی سریع می توان فرایند تعیین کد برنده MDWC^{۳۳} را به جای M بار محاسبه به $M \log_2 k$ بار محاسبه تابع فاصله کاهش داد (k نوعاً بین ۲ تا ۴ قرار دارد). البته کارایی ساختارهای جستجوی سریع موقعی قابل ملاحظه و قابل استفاده است که M باندازه کافی بزرگ باشد. این مزایا وقتی قابل توجه است که بخواهیم سیستم بازشناسی با دایره لغات وسیع و به صورت بی درنگ طراحی کنیم.

۴- استفاده از بازنمایی SCD^{۳۳} و FVQ^{۲۴} برای جبران کاهش دقت مدل سازی DD-HMM

استفاده از چگالی گسسته اگرچه دارای مزایای سادگی و سرعت پیاده سازی به ویژه در فاز بازشناسی است، اما به دلیل

مناسب در فاز آموزش به طور سریع تعیین می‌شود. فاصله بردار مشاهده O_t را با $V_i, V_{i1}, \dots, V_{ik}$ به شرح زیر تعیین می‌کنیم.

$$d(O_t, V_i), d(O_t, V_{i1}), \dots, d(O_t, V_{ik})$$

با استفاده از فواصل فوق توابع عضویت را به شرح زیر به دست می‌آوریم:

$$\alpha_j(O_t, V_i) = \frac{[d(O_t, V_{ij})]^{\frac{1}{m}}}{\sum_{k=0}^{K+1} [d(O_t, V_{ik})]^{\frac{1}{m}}} \quad (12)$$

$$j=0, 1, \dots, k$$

$$V_{i0} \triangleq V_i$$

تعداد توابع عضویت برابر $k+1$ است. با استفاده از مقادیر $k+1$ تابع عضویت $\alpha_0, \dots, \alpha_k$ مقدار تابع احتمال $b_j(O_t)$ به شرح زیر درونیابی می‌شود:

$$b_j(O_t) = \sum_{l=0}^k \alpha_l \cdot b_j(V_{il}) \quad (13)$$

$$V_{i0} \triangleq V_i$$

ملاحظه می‌شود که درونیابی صرفاً با استفاده از فاصله بردار مشاهده O_t از کد برنده و k نزدیکترین همسایه آن انجام می‌پذیرد. یعنی تعیین $b_j(O_t)$ به فواصل O_t از سایر کلمات کد بستگی ندارد. بنابراین در این روش نیز می‌توان از ساختارهای جستجوی سریع استفاده کرد.

۵-۲- ارائه درون یابی KNNGI

بردارهای مرجع V_1, \dots, V_M (یا کلمات کد) نوعاً توسط الگوریتم LBG یا EM طراحی می‌شوند. در پایان اجرای الگوریتم LBG (یا EM)، علاوه بر بردارهای مرجع، بردارهای واریانس متناظر نیز محاسبه می‌شوند (با فرض قطری بودن ماتریس کوواریانس بردارها). در روش درون یابی KNNFI از اطلاعات واریانس هر کلمه کد هیچ گونه استفاده‌ای نمی‌شود. اما می‌دانیم که واریانس هر بردار مرجع، نماینده چگونگی پراکندگی بردارهای آموزشی است که آن بردار مرجع نماینده آنهاست. سؤال اساسی برای محققان این نوشتار این بود که چگونه از این اطلاعات اضافی برای بهبود عملکرد فرایند درونیابی استفاده کنیم. با توجه به کارایی مدل‌های GMM، ایده

$$b_j(O_t) = \sum_{i=1}^I \alpha_i \cdot b_j(O_t = V_i) \quad (9)$$

که در آن V_i ها، کدهای برنده با قواعد ذکر شده‌اند.

ضرایب وزنی درونیابی α_i ها (I تا 1) به صورت فازی به شرح زیر محاسبه می‌شوند:

$$\alpha_i = \frac{(d_i)^{\frac{1}{m}}}{\sum_{k=1}^I (d_k)^{\frac{1}{m}}} \quad \text{that} \quad \sum_{k=1}^I \alpha_i = 1 \quad (10)$$

که مقدار فازی بودن^{۲۵} سیستم (عددی بین ۱ تا ۲) و α_i ها توابع عضویت در بازسازی $b_j(O_t)$ هستند [۹]. کارایی FVQ به مراتب بیشتر از VQ (روش DD کلاسیک) است. همچنان که ملاحظه می‌شود، در این روش نیز مجبوریم کلیه توابع فاصله بین بردار مشاهده O_t و کلیه اعضای کتاب کد را به دست آورده و سپس آنها را مرتب کرده و کلمه کد برنده با تابع فاصله‌های متناظر را برای محاسبه α_i مورد استفاده دهیم. بنابراین نمی‌توانیم از ساختارهای جستجوی سریع برای اهداف این تحقیق استفاده کنیم.

۵- ارائه درون یابی KNNFI, KNNGI

در روش KNNFI و KNNGI ساختارهایی به شرح زیر بر روی اعضای کتاب کد (در مرحله آموزش) ایجاد می‌کنیم. فرض می‌کنیم (V_1, \dots, V_M) اعضای کتاب کد یک کلاس اکوستیکی باشند. برای هر بردار مرجع، تعداد k نزدیکترین همسایه بر مبنای کمترین فاصله به شرح زیر به دست می‌آوریم. فرض می‌کنیم $\hat{I}_1, \hat{I}_2, \dots, \hat{I}_k$ تعداد k نزدیکترین همسایه کلمه کد V_i باشند یعنی داریم:

$$d(V_i, V_{i1}) \leq d(V_i, V_{i2}) \leq \dots \leq d(V_i, V_{ik}) \quad (11)$$

فرایند تعیین k نزدیکترین همسایه را برای کلیه کلمات کد انجام می‌دهیم. ایندکس k نزدیکترین همسایه هر کلمه کد به همراه سایر اطلاعات کتاب کد را در فاز آموزش به دست آورده و ذخیره می‌کنیم. توجه می‌شود که فرایند تعیین ایندکس k نزدیکترین همسایه هر کلمه کد، صرفاً یکبار در مرحله آموزش محاسبه و ذخیره می‌شود.

۵-۱- ارائه درون یابی KNNFI

فرض می‌کنیم O_t یک بردار مشاهده در فاز بازشناسی و V_i نزدیکترین بردار مرجع به آن باشد (که با ایجاد ساختارهای

مجموع هزینه محاسباتی متناسب با 50×1024 عمل جمع خواهد شد. البته در تخمین هزینه محاسباتی فوق الذکر، از وارد کردن حجم محاسبات مربوط به تحلیل هر فریم و مقایسه پویا در محاسبه نمرات کلاسهای اکوستیکی، تشخیص ابتدا و انتهای کلمه و ...، به علت مشترک بودن فرایندهای فوق الذکر در هر دو مدل DDHMM و CDHMM صرف نظر شده است. در محاسبات تحلیلی فوق، مشاهده می شود که هزینه محاسباتی DDHMM حدود نصف CDHMM است. این نسبت کاهش محاسباتی، ارزش توجه و به کارگیری مدل DDDHMM در مقابل CDHMM را ندارد (با توجه به مقدار افت دقت بازشناسی مدل DDHMM در مقابل مدل CDHMM). هزینه محاسباتی تخمین زده شده برای مدل DDHMM با فرض استفاده از ساختار جستجوی کامل است. در سیستمهای کاربردی از ساختارهای جستجوی سریع، مانند ساختار باینری و یا ساختار تلفیقی باینری و جستجوی تصادفی k نزدیکترین همسایه استفاده می کنند. در حالت دوم، تعداد عملیات برای محاسبه نمره یک فریم متناسب با $2p \cdot (\log_2 1024)$ می شود. (k بین ۲ تا ۴ بوده و در این تحقیق $k=3$ حاصل شده است). بنابراین در مدل سازی DDHMM با ساختار تلفیقی جستجوی سریع، هزینه محاسباتی برای تعیین نمره یک فریم متناسب با $1/5 \times 1024$ عمل جمع می شود. مقایسه هزینه محاسباتی مدل DDHMM با ساختار تلفیقی جستجوی سریع و مدل CDHMM، علت توجه محققان به مدل سازی DDHMM را به خوبی آشکار می کند (کاهش هزینه محاسباتی به نسبت $1/2$).

دو روش درون یابی KNNFI و KNNGI که برای جبران کاهش دقت بازشناسی مدل DDHMM در طی این تحقیق ارائه و پیاده سازی شده اند، از ساختار تلفیقی جستجوی سریع مدل DDHMM (برای کاهش بار محاسباتی) استفاده می کنند. بنابراین هزینه محاسباتی آنها تقریباً معادل روش DDHMM با ساختار تلفیقی جستجوی سریع است. در روشهای ارائه شده (درفاز بازشناسی)، ابتدا بردار مرجع برنده (از میان 1024 بردار مرجع) را به روش جستجوی سریع تعیین می کنیم. آن گاه

زیر را مورد توجه قرار دادیم:

فرض می کنیم v_i کد برنده با مشاهده بردار O_t باشد، همچنین v_{i_1}, \dots, v_{i_k} تعداد k نزدیکترین همسایه بردار مرجع v_i باشد. ضرایب وزنی درون یابی را به شرح زیر پیشنهاد و مورد استفاده قرار دادیم:

$$p_j \triangleq N(o_t, v_{ij}, \sigma_{ij}) \quad (14)$$

$N(-)$ نماینده تابع گوسی با بردار میانگین v_{ij} و بردار واریانس σ_{ij} است.

$$\alpha_j \triangleq \frac{p_j}{\sum_{l=0}^k p_l} \quad \text{that} \quad j=0 \text{ تا } k \quad (15)$$

ملاحظه می شود که در درون یابی KNNGI، نیز صرفاً از فاصله گوسی بردار مشاهده O_t و کد برنده O_t و k نزدیکترین همسایه آن استفاده شده است. بنابراین استفاده از درون یابی KNNGI به ساختار جستجوی سریع خدشه ای وارد نمی کند.

۳-۵ بررسی کارایی محاسباتی KNNFI و KNNGI؛

با توجه به اینکه روشهای ارائه شده در این تحقیق برای سیستمهای بازشناسی با دایره لغات وسیع طراحی شده اند، در نتیجه پیاده سازی هر دو مدل بازشناسی (CDHMM و DDHMM) به صورت گره زدن مخلوطها (TM)^{۲۸} انجام پذیرفته است. در مدل CDHMM از ۵۱۲ مخلوط گوسی و برای مدل DDHMM از $1024 \times VQ$ برای بازنمایی pdf حالات تمامی کلمات مربوط به تمامی گویندگان استفاده کردیم. برای محاسبه مقدار درست نمائی هر فریم دریافتی در مدل DDHMM، به حداقل $512 \times 3p$ عمل (ضرب و جمع و تقسیم) و ۵۱۲ محاسبه تابع اکسپانسیل و ۵۱۲ عمل ضرب نیازمندیم. (p بعد بردار ویژگی است که در این تحقیق برابر ۲۵ است). فرض می کنیم هزینه محاسباتی هر عمل ضرب دو برابر عمل جمع و هر عمل تقسیم چهار برابر عمل جمع و محاسبه تابع اکسپانسیل ده برابر عمل جمع باشد. بنابراین مجموع تعداد عملیات محاسباتی برای تعیین نمره یک فریم دریافتی متناسب با 93×1024 عمل جمع می شود. برای حالت DDHMM،

فوق‌الذکر را نشان می‌دهد. در آزمون مربوط به نمونه‌های آموزشی، $\frac{1}{2}$ گویندگان مرد و زن برای آموزش و $\frac{1}{4}$ گویندگان باقیمانده برای آزمایش مورد استفاده قرار گرفته‌اند. فرایند فوق پنج مرتبه با تغییر گویندگان آموزشی (به صورت تصادفی) انجام پذیرفته است. نتیجه اعلام شده متوسط و انحراف معیار این آزمون است. با مروری به نتایج چهار آزمون انجام پذیرفته به جمع‌بندی زیر می‌رسیم:

۱- دقت بازشناسی بر روی نمونه‌های آزمایشی همیشه کمتر از نمونه‌های آموزشی است. اما قدرت سیستم طراحی شده می‌بایستی به گونه‌ای باشد که افت دقت در حد قابل قبول باشد (حدود یک تا حداکثر دو درصد) افت دقت بر روی نمونه‌های آزمایشی در مدل CD-HMM قابل قبول بوده ولی برای مدل DD-HMM (بدون درون‌یابی) غیرقابل قبول است.

۲- قضاوت نهایی برای بررسی عملکرد یک سیستم بازشناسی، بررسی نتایج دقت بازشناسی بر روی نمونه‌های آزمایشی است. بدین لحاظ جدول (۲) نتایج بازشناسی دو روش درون‌یابی KNNFI و KNNGI را بر روی داده‌های آزمایشی نشان می‌دهد. ملاحظه می‌شود، که دقت بازشناسی دو روش درون‌یابی ارائه شده در مقایسه با مدل DD-HMM به شدت بهبود یافته و به دقت سیستم CD-HMM نزدیک شده است.

۳- دقت بازشناسی درون‌یابی KNNGI بر روی نمونه‌های آزمایشی به طور متوسط حدود ۰/۷۹ درصد بهتر از دقت بازشناسی درون‌یابی KNNFI است.

۴- دقت بازشناسی DD-HMM با درون‌یابی KNNGI در حدود دقت بازشناسی CD-HMM (با حدود ۰/۵ درصد افت دقت) است بنابراین با اطمینان کافی می‌توان از آن برای سیستم‌های کاربردی استفاده کرد.

در آزمون دیگری، برنامه‌ای به زبان C (در محیط VC++6) برای بازشناسی گفتار ۱۵۰۰ کلمه فارسی به دو روش CD-HMM و DD-HMM (با درون‌یابی KNNGI) نوشته شده است. سیستم بازشناسی به روش DD-HMM (با درون‌یابی KNNGI) به طور متوسط حدود ۰/۷۲ ثانیه پس از دریافت

فاصله فازی یا فاصله گوسی بردار فریم دریافتی و k بردار مرجع نزدیکترین همسایه بردار مرجع برنده را محاسبه می‌کنیم. ایندکس k نزدیکترین همسایه هر بردار مرجع در فاز آموزش تعیین می‌شود. بنابراین در فاز بازشناسی به بار محاسباتی اضافی نیازی نداریم. مقدار بهینه k (تعداد همسایه‌ها) بین ۵ تا ۶ گزارش می‌شود. در این تحقیق پس از آزمایش‌های تجربی فراوان k را برابر ۶ در نظر گرفتیم.

۶- نتایج پیاده‌سازی

مدلهای درون‌یابی ارائه شده در بازشناسی ۱۵۰۰ کلمه فارسی ناوابسته به گوینده^{۲۹} مورد استفاده و ارزیابی قرار گرفتند. پایگاه داده گفتاری شامل ۳۰ مرد و ۲۰ زن بوده است. از هر گوینده حدود ۶۰ تلفظ برای هر کلمه ضبط کردیم. (ده تلفظ طی ۶ روز متفاوت). فرکانس نمونه‌برداری ۸ kHz و نمونه‌ها به صورت ۱۶ بیتی ضبط شده‌اند. داده‌های گفتاری توسط میکروفون برای کیفیت خوب در محیط عادی ضبط شده‌اند. تعداد حالات هر کلمه در مدل‌سازی HMM را برابر پنج (N=5) قرار دادیم. ویژگی‌های مورد استفاده ۱۲ ضریب MFCC^{۳۰} و مشتق مرتبه اول آن و مشتق لگاریتم گین یعنی مجموعاً ۲۵ ویژگی بوده‌اند. روش تحلیل برای استخراج ویژگی‌های طیفی روش PLP^{۳۱} [۵] بوده است. ویژگی‌های فوق به ازای هر ۶/۲۵ میلی ثانیه یکبار استخراج شده‌اند. برای هموار کردن کانتر ویژگی‌ها، از یک فیلتر غیرخطی (۷ و ۳) MM^{۳۲} استفاده کردیم. برای تخمین پارامترهای آماری مدل HMM، از الگوریتم بام و لشل [۵ و ۶] استفاده کردیم. تعداد اعضاء کتاب کد برای DDHMM برابر ۱۰۲۴ بردار و برای CDHMM برابر ۵۱۲ مخلوط در نظر گرفته شدند. چندین آزمون به شرح زیر انجام پذیرفته است:

۱- پیاده‌سازی مدل CD-HMM

۲- پیاده‌سازی مدل DD-HMM بدون درون‌یابی

۳- پیاده‌سازی مدل DD-HMM با درون‌یابی KNNFI

۴- پیاده‌سازی مدل DD-HMM با درون‌یابی KNNGI

جداول (۱) و (۲) نتایج بازشناسی برای چهار نوع مدل‌سازی

جدول ۱- نتایج دقت بازشناسی نا وابسته به گوینده مدل CDHMM و DDHMM

نوع مدل	نمونه‌های آموزشی	نمونه‌های آزمایشی
CDHMM	۹۸/۲۱ (۱/۰۵)	۹۷/۶۵ (۱/۸۵)
DDHMM	۹۹/۱۸ (۰/۵۱)	۸۹/۲۷ (۲/۸۲)

جدول ۲- نتایج دقت بازشناسی ناوابسته به گوینده مدل DD-HMM با درون‌یابی KNNFI و KNNGI

نوع درون‌یابی	نمونه‌های آموزشی	نمونه‌های آزمایشی
KNNGI	۹۸/۳۱ (۰/۹۷)	۹۷/۱۶ (۱/۹۲)
KNNFI	۹۸/۴۸ (۰/۹۱)	۹۶/۳۷ (۲/۱۲)

قدرت تعمیم پذیری مدل DDHMM، بازنمایی گسسته فضای pdf توسط تعداد محدودی بردار (اعضای کتاب کد) است. در روش CDHMM به علت پیوسته بودن بازنمایی فضای pdf، بردارهایی که در فاز آموزش دیده نشده‌اند نیز دارای مقادیر احتمال مناسبی برای همان کلاس‌اند. این توانمندی به علت گسسته شدن فضای pdf در مدل DDHMM در حد مناسبی برقرار نیست. البته مدل DDHMM نسبت به مدل CDHMM دارای دو برتری اساسی است؛ ۱- قدرت متمایز سازی بیشتر بر روی نمونه‌های آموزشی ۲- سرعت اجرای بسیار سریعتر. تنها مشکل آن پایین بودن قدرت تعمیم پذیری است. (یعنی پایین بودن دقت بازشناسی بر روی نمونه‌های آزمایشی که قبلاً در آموزش شرکت نداشته‌اند). روش درون‌یابی فازی FI که قبلاً توسط محققان دیگر [۸ و ۹] ارائه شده است، تا حدودی نقطه ضعف مدل DDHMM را از نظر قدرت تعمیم پذیری برطرف کرده است. روش درون‌یابی فازی کلاسیک، از ساختارهای جستجوی سریع استفاده نمی‌کند. زیرا به ازای هر بردار تست ورودی، فاصله فازی آن با تمام بردارهای مرجع محاسبه شده و سپس تعداد ۵ تا ۶ بردار مرجع که فاصله فازی آن از بقیه اعضای کتاب کد کمتر است، برای درون‌یابی نهایی انتخاب می‌شوند. بنابراین درون‌یابی فازی کلاسیک، قابلیت پیاده‌سازی سریع مدل DDHMM را کم رنگ می‌کند. در روش KNNFI که در این تحقیق ارائه کردیم، ضمن استفاده از خواص درون‌یابی فازی برای افزایش قابلیت تعمیم پذیری

کلمه مورد نظر نتایج بازشناسی را اعلام می‌کند سیستم بازشناسی به روش CD-HMM به طور متوسط حدود دو دقیقه پس از دریافت کلمه مورد نظر نتایج بازشناسی را اعلام می‌دارد. برنامه‌ها بر روی رایانه پنتیوم ۴ (۲/۴ GHZ) با حافظه ۵۱۲ mbyte اجرا شدند. بنابراین از نظر اجرای بی درنگ، به نظر می‌رسد که روش CD-HMM قابلیت کاربرد چندانی برای دایره لغات وسیع ندارد. ولی روش DD-HMM (با درون‌یابی KNNGI) دارای قابلیت کاربرد در سیستمهای کاربردی با اجرای بی درنگ است.

۷- بررسی کارایی درون‌یابی KNNFI و KNNGI

مدل CDHMM، توزیع آماری بردارهای ویژگی داخل هر حالت هر کلاس اکوستیکی را توسط یک pdf پیوسته و پارامتریک، به نام GMM بازنمایی می‌کند. تعداد مخلوطهای توزیع پارامتریک GMM هر حالت، بستگی به نوع گفتار و تعداد گوینده‌ها، متغیر است. توزیع پیوسته و پارامتریک GMM، توازن مناسبی بین قدرت متمایزسازی و قدرت تعمیم پذیری کلاسهای اکوستیکی ایجاد می‌کند. با افزایش تعداد بردارهای کتاب کد، قادریم قدرت متمایزسازی مدل DDHMM را بیشتر از مدل CDHMM کنیم. این توانمندی مدل DDHMM در اغلب تحقیقات محققان دیگر نیز گزارش شده است. متأسفانه در اغلب کاربردهای عملی، به قدرت متمایزسازی و قدرت تعمیم پذیری به طور همزمان نیازمندیم. علت اصلی پایین بودن

۸- جمع‌بندی و نتیجه‌گیری

هدف این تحقیق بازشناسی گفتار گسسته (با دایره لغات وسیع) با امکان پیاده‌سازی بی‌درنگ (با تأخیری در حدود نیم تا حداکثر یک ثانیه) بر روی رایانه‌های PC قابل دسترس بوده است. با توجه به این که مدل‌سازی HMM روش غالب و موفق در سیستم‌های بازشناسی گفتار پیوسته و گسسته است، از این روش برای کاربرد ذکر شده استفاده کردیم. دو رویکرد چگالی پیوسته CD و چگالی گسسته DD، برای بازنمایی pdf بردارهای طیفی مورد استفاده قرار می‌گیرند. رویکرد اول CD دارای دقت بازشناسی مناسب است. متأسفانه به علت حجم بالای محاسبات (برای دایره لغات وسیع) امکان استفاده از آن در بازشناسی بی‌درنگ با رایانه‌های PC موجود و قابل دسترس وجود ندارد. رویکرد دوم DD اگرچه دارای دقت پایبتری (نسبت به CD) است ولی پیاده‌سازی بی‌درنگ آن بر روی pc‌های قابل دسترس امکانپذیر است. روش کلاسیک برای افزایش دقت رویکرد دوم (DD) (با حفظ ساختارهای جستجوی سریع) استفاده از درون‌یابی فازی FI است. در این تحقیق روش درون‌یابی جدیدی به نام KNNGI ارائه و مورد استفاده قرار دادیم. درون‌یابی گوسی، ضمن بالا بردن دقت بازشناسی (در حد رویکرد اول CD) به راحتی بر روی رایانه‌های موجود به صورت بی‌درنگ قابل پیاده‌سازی است. پیاده‌سازی دو مدل درون‌یابی برای بازشناسی ناوابسته به گوینده ۱۵۰۰ کلمه مجزای فارسی نشان می‌دهد که عملکرد دقت روش درون‌یابی گوسی بیشتر از درون‌یابی فازی و در حد رویکرد اول یعنی CDHMM است.

مدل DDHMM، به سرعت پیاده‌سازی مدل DDHMM نیز خدشه‌ای وارد نمی‌کند. زیرا در روش KNNFI، از ساختارهای جستجوی سریع می‌توان بر راحتی استفاده کرد. بنابراین درون‌یابی KNNFI در مقابل روش FI کلاسیک، دارای قدرت پیاده‌سازی سریعتر است (که یکی از اهداف اصلی این تحقیق بوده است).

درون‌یابی KNNGI که نیز در این تحقیق ارائه و پیاده‌سازی شده است، ضمن حفظ قابلیت پیاده‌سازی سریعتر (همانند درون‌یابی KNNFI) مقدار مناسبی نیز به دقت بازشناسی سیستم افزوده است. این حقیقت در نتایج مندرج در جدول (۲) به خوبی آشکار است. همان طوری که از جدول (۲) مشاهده می‌شود، دقت درون‌یابی KNNGI حدود ۰/۷۹ درصد بیشتر از درون‌یابی KNNFI شده است. این افزایش دقت معادل کاهش ۲۲ درصدی خطای سیستم است. دلیل تحلیلی کاهش خطا، استفاده از پارامتر نرمالیزه‌کننده واریانس و تابع تحلیلی گوسی می‌تواند باشد. همان طور که می‌دانیم، واریانس نماینده پراکندگی نمونه‌ها حول مقدار میانگین است. در درون‌یابی فازی، مستقل از واریانس خوشه، فاصله بردار دریافتی با مرکز خوشه محاسبه شده و از یک فاصله اقلیدسی نیز برای محاسبه فاصله فازی استفاده می‌شود. با مقایسه روابط (۱۲) و (۱۴) می‌توان گفت که علت دوم برتری درون‌یابی KNNGI، استفاده از تابع فاصله گوسی به جای تابع فاصله اقلیدسی است. بنابراین در مجموع می‌توان گفت که مدل DDHMM به همراه درون‌یابی KNNGI جایگزین مناسبی برای دسترسی به دقت بالا و پیاده‌سازی سریع سیستم‌های بازشناسی با دایره لغات وسیع است.

واژه‌نامه

- | | | |
|--|----------------------------------|-------------------------------------|
| 1. K nearest neighbor Gaussian interpolation | 9. dynamic time warpping | 17. continous speech recognition |
| 2. fuzzy interpolation | 10. continous speech | 18. semi continous |
| 3. continous density HMM | 11. discrete (utterance) speech | 19. forward procedure |
| 4. discrete density HMM | 12. acoustic modeling | 20. Linde, Bazo, Gray |
| 5. Gaussian mixture model | 13. language modeling | 21. expectation maximization |
| 6. vector quantization | 14. acoustic to phonetic | 22. minimum distance winer codeword |
| 7. real time | 15. statistical LM | 23. semi continous density |
| 8. artifitial neural net | 16. probability density function | 24. fuzzy vector quantization |

- | | | |
|--|----------------------------------|------------------------|
| 25. fuzziness | 29. speaker dependent | 32. mean median filter |
| 26. K nearest neighbor fuzzy interpolation | 30. mel frequency cepstrum | |
| 27. KNN Gaussian interpolation | coeffitient | |
| 28. TIED mixture | 31. perceptual linear prediction | |

مراجع

1. Rabiner, L., and Juang, B.H., *Fundamentals of Speech Recognition*, Prentice Hall Inc, 1993.
2. Huang, X. D., "Phoneme Classification Using Semi Continuous Hidden Markov Models," *IEEE Transactions on signal processing*, Vol. 40, No. 5, pp.45 – 56, May 1992.
3. Zhang, Y., Togneri, R., and Alder, M., "Phoneme Based Vector Quantization in a Discrete HMM Speech Recognizer," *IEEE Transaction on Speech and Audio Processing*, Vol. 5, No.1, January 1997.
4. Furui, S., and Sondhi, M. M, *Advances in Speech Signal Proceesing*, Marcel Dekker, Inc., 1992.
5. Becchetti, C., Ricotti, L. P., *Speech Recognition, Theory and C Implementation*, John wiley & Sons, Ltd, 1999.
6. Jurafsky, D., Martin, J. H., *Speech and Language Processing*, Prentice Hall. Inc., 2000.
7. Linde, Y., Buzo, A., Gray, R. M., "An Algorithm for Vector Quantizer design," *IEEE Transaction on Communication*, Vol. Com 28, pp. 84-95, 1980.
8. Tseng, H. P., Sabin, M., Lee, E., "Fuzzy Vector Quantization applied to Hidden Markov Modeling," *In proc, ICASSP-87*, 1987.
9. Schwartz, R., Kimball, O., kubala, F., Feng, M., Chow, Y., Barry, C., Makhoul, J., "Robust Smothing Methods for Discrete Hidden Markiov Models," *In proc, ICASSP-89*, 1989.
10. Dempster, A. P., Laird, N. M., Rubin, D. B., "Maximum Likelihood from Incomplete Data Via the LM Algorithm," *J. Roy. Stat. Soc. Ser. B(methodological)*, Vol. 39, pp. 1-38, 1977.
11. Render, R. A., Walke, H. F. "Mixture Densities, Maximum Likelihood, and the EM Algorihm," *SIAM Rer.*, Vol. 26, pp. 195-239, 1984.
12. Bailey, T. L., Elkan, C., "Fitting a Mixture Model by Expectation Maximization to Discove Motifs in Biopolymers," *proc. of Biology*, pp. 28-30, AAAI PRESS, (1994).
13. Gauvain, J. L., Lee, C. H., "Maximum a Posteriori Estimation Gaussian Mixture Observations of Mardov Chain," *IEEE Trans. On speech and Audio processing*, No. 2, pp. 291-298, Apr. 1994.