

آموزش وزنی مدل های وابسته به متن در بازشناسی گفتار پیوسته

سید محمد احدی
استادیار

سید حسین شمس
کارشناسی ارشد

دانشکده مهندسی برق، دانشگاه صنعتی امیرکبیر

چکیده

در بازشناسی گفتار پیوسته مدل سازی وابسته به متن می تواند موجب افزایش دقت بازشناسی و در عین حال افزایش تعداد مدلها شود. این افزایش در تعداد مدلها می تواند مواجهه با کمبود داده های آموزشی را موجب گردد. روشهای مختلفی همانند گره زدن پارامترهای مدلها یا استفاده از تخمین MAP برای غلبه بر این مشکل به کار می روند. در این مقاله روشی جدید برای آموزش ارائه شده است که محاسن هر دو روش MAP و گره زدن را تماماً داراست و به نتایج بهتری نیز منجر می شود. مدلهای مورد استفاده از نوع سه آوانی می باشند.

کلمات کلیدی

مدلهای مارکوف پنهان با چگالی پیوسته، بازشناسی گفتار پیوسته فارسی، مدل سازی وابسته به متن، آموزش بام - ولش، گره زدن حالتها، آموزش MAP.

Weighted Training of Context Dependent Models for Continuous Speech Recognition

S.M. Ahadi
Assistant Professor

S.H. Shams
M.Sc. candidate

Electrical Engineering Department,
Amirkabir University of Technology

Abstract

Context dependent modelling is a well-known technique used to improve recognition rate in continuous speech recognition. However, it can also dramatically increase the overall system parameter count which usually leads to a grave problem known as sparse training data condition. Different techniques such as model parameter tying and MAP estimation are used to overcome this problem. In this paper, a new approach, namely weighted training, has been introduced which benefits from both MAP estimation and parameter tying strengths and also outperforms the widely used parameter tying approach. The triphone models are used in this implementation.

Keywords

Continuous density hidden Markov models, Continuous Persian speech recognition, Context dependent modelling, Baum-Welch training, State tying, MAP estimation.

HMM ها معمولاً برای مدل کردن سیستمهای غیر خطی متغیر با زمان مورد استفاده قرار می‌گیرند. استفاده از مدل‌های مارکوف پنهان در بازشناسی گفتار گسسته و پیوسته با دایره لغات بزرگ به نتایج قابل قبولی منجر شده است. انتظار می‌رود با استفاده از مدل‌سازی وابسته به متن^۱، به شرط آموزش مناسب، بتوان به نتایج بهتری دست پیدا نمود. مدل‌سازی وابسته به متن از آن جهت دارای اهمیت است که شرایط مختلف متنی را در مدل‌سازی آوای مورد نظر دخالت می‌دهد. به بیان دیگر، در این شیوه مدل‌سازی، اثر متن^۲ نیز در مدل دخالت می‌کند و بنابراین در شرایط متنی مختلف، مدل‌های مختلفی مورد استفاده قرار می‌گیرند. در عمل نشان داده شده است که مشخصات آکوستیک هر واج به شدت تابع واج‌های پیش و پس از آن می‌باشد. این اثر حتی می‌تواند تا چند واج قبل و بعد از آن واج را هم پوشش دهد [۱]. به همین دلیل برای غلبه بر این مشکل، بهترین راه حل مدل‌سازی آواها با توجه به شرایط متنی تشخیص داده شده است.

مدل‌سازی وابسته به متن، علیرغم مزیت بارزی که نسبت به مدل‌سازی بر اساس واج (ناوابسته به متن) دارد، از یک مشکل مهم رنج می‌برد و آن تعداد فوق‌العاده زیاد شرایط متنی است که هر واج می‌تواند در آن قرار گیرد. این امر منجر به افزایش فراوان تعداد مدل‌ها در مقایسه با شرایط عادی می‌گردد. این مسئله، علاوه بر افزایش حجم حافظه و پردازش در سیستم، باعث می‌شود که دستیابی به شرایط آموزش مقاوم^۳ بسیار مشکل گردد. به طور معمول، در اغلب سیستم‌های بازشناسی گفتار، حتی در مواردی که دادگان‌های بسیار بزرگی برای آموزش سیستم در اختیار می‌باشد، بسیاری از مدل‌های وابسته به متن از مسئله کمبود داده آموزشی رنج می‌برند. برای غلبه بر این مشکل چند روش می‌تواند مورد استفاده قرار گیرد:

الف- برای جلوگیری از آموزش نامعتبر، تنها مدل‌هایی که میزان مناسبی داده آموزشی داشته باشند آموزش داده شوند. عیب این روش در این است که از کلیه اطلاعات استفاده نشده است و روش بهینه‌ای نیست.

ب- پس از بدست آوردن مدل‌های اولیه، برای کم کردن پارامترهای سیستم از گره‌زدن استفاده شود. این کار می‌تواند موجب دستیابی به آموزش مقاوم شود. این شیوه و شیوه‌های مشابه در سیستم‌های بازشناسی زیادی پیاده سازی گردیده و به نتایج خوبی نیز منجر شده‌اند [۲][۳]. باید توجه داشت که هنگامی که خوشه‌بندی در این روش بر اساس داده صورت می‌پذیرد، چنانچه در ابتدا مدل‌ها بعلت کمی داده آموزشی بخوبی آموزش داده نشده باشند، الگوریتم خوشه‌بندی لزوماً باعث گره خوردن مدل‌هایی که در واقعیت باید به هم شباهت داشته باشند نمی‌شود.

ج- آموزش MAP^۴: این شیوه آموزش که بر اساس تئوری پیز^۵ بنا گردیده، از یک مجموعه پارامتر پیشینه^۶ بهره می‌برد. این روش آموزش به گونه‌ای عمل می‌کند که در صورت وجود داده آموزشی کافی مدل به سمت مدل بدست آمده از آموزش درستمائی پیشینه (ML^۷) میل می‌کند و در صورت نبود داده آموزشی کافی به سمت مدل اولیه (پیشینه) تمایل دارد. استفاده از پارامترهای پیشینه موجب مقاوم شدن آموزش بویژه در شرایط داده آموزشی ناکافی می‌شود و از این رو MAP به عنوان شیوه‌ای موفق در آموزش در شرایط فوق محسوب شده و در روش‌های تطبیقی در بازشناسی گفتار نیز کاربرد فراوانی دارد [۴][۵]. مهم‌ترین مشکل در این روش بدست آوردن تخمین مناسب برای پارامترهای پیشینه است (رجوع به بخش ۵). در این تحقیق آموزش بگونه‌ای طرح ریزی شده که برای آموزش هر مدل از مدل‌های دیگر شبیه به آن مدل هم استفاده گردد (از این جهت شباهت به گره زدن دارد) ولی بگونه‌ای که اطلاعات مربوط به هر مدل در آموزش خود آن مدل مؤثرتر واقع می‌شود. در اینجا نیز در صورت وجود داده آموزشی کافی مدل به سمت مدل آموزش یافته به روش معمول (ML) میل می‌کند و در صورت نبود داده آموزشی کافی، شبیه به مدل ناوابسته به متن میشود (از این نظر شبیه به MAP است).

محاسن این روش بشرح زیر است:

الف - کلیه مدل‌ها آموزش می‌بینند.

ب - در صورت وجود داده آموزشی کافی اثر مدل‌های دیگر ناچیز می‌گردد.

ج - تفاوت بین مدل‌ها هر چند ناچیز در مدل نهایی ظاهر می‌شود.

د - آموزش مقاوم و دقیق خواهد بود.

می‌توان گفت که در این روش از محاسن روشهای MAP و گره زدن پارامترها تماماً استفاده می‌شود. پایگاه داده مورد استفاده در این تحقیق *فارس دات (FARSDAT)* میباشد که یک دادگان با دایره کلمات متوسط است. مقاله حاضر شامل نه بخش می‌باشد که بخش اول مقدمه و بخش دوم به مدلهای زیر لغوی می‌پردازند. بخش سوم به نحوه آموزش، بخش چهارم به نحوه گره زدن، بخش پنجم به روش MAP، بخش ششم به روش پیشنهادی و بخش هفتم به پیاده‌سازی اختصاص یافته‌اند. در بخش هشتم به نتایج پرداخته شده است. بخش نهم نیز نتیجه‌گیری و مسیر تحقیقات آینده را نشان می‌دهد.

۲- مدلهای زیر لغوی

انتخاب واحدهای پایه و مدل کردن آنها اولین قدم در طراحی سیستم بازشناسی می‌باشد. اگر عمل آموزش بطور کامل انجام گرفته باشد هر چه تنوع واحدها بیشتر شود نتیجه بازشناسی مطلوبتر خواهد بود. در عین حال افزایش تعداد مدل‌ها ممکن است به عدم آموزش مناسب منجر شود. در این تحقیق در مرحله اول مدلها بر اساس واجهای زبان فارسی انتخاب شدند و سپس به مدل وابسته به متن سه آوایی^۱ تعمیم یافتند.

۲-۱- مدلسازی وابسته به متن [۶]

منظور از مدلسازی وابسته به متن این است که برای هر واحد زیر لغوی با توجه به واحدهای زیر لغوی قبل یا بعد (شرایط متنی) مدل جداگانه‌ای در نظر گرفته شود. در زیر، گونه‌هایی از واحدهای پایه استفاده شده اعم از ناپسته یا وابسته به متن معرفی میشوند.

الف - مدلهای تک آوایی^۱: تنها یک مدل هر واج را در تمام محلهای وقوع مدل می‌کند. جدول شماره ۱ واجهای اصلی زبان فارسی و نمایش سمبلیک آنها را در این کاربرد نشان میدهد. مثال زیر به نحوه نمایش یک جمله بر این اساس می‌پردازد.

d@r baz '@st

_ d @ r + b a z + ' @ s t _

ب - مدلهای دو آوایی^۱: هر واج با توجه به واج سمت راست یا سمت چپ بطور جداگانه مدل می‌شود. برای سکوت و فاصله تنها یک مدل (بدون توجه به متن) در نظر می‌گیرند و در صورتی که واج کناری سکوت و یا فاصله باشد بصورت تک آوایی مدل می‌شود. به این ترتیب جمله بالا بصورت زیر نمایش داده می‌شود.

مدل‌سازی دو آوایی چپ:

_ d d<@ @<r + b b<a a<z + '<@ @<s s<t _

مدل سازی دو آوایی راست:

_ d>@ @>r r + b>a a>z z + '>@ @>s s>t t _

جدول (1) نمایش سمبلیک واحدهای استفاده شده بر اساس واحهای زبان فارسی با استفاده از کاراکترهای ASCII.

Phonetic Group	Phoneme	Example
Vowels	a	ketab
	@	s@br
	e	yek
	u	xub
	i	Sir
	o	sorx
	w	lwh
Liquids	r	rah
	/	lazem
Glides	y	yek
Nasals	m	morq
	n	niru
Plosives	b	b@d
	p	por
	t	tond
	d	dust
	q	morq
	k	kuh
	g	goft
,	e'lam	
Fricatives	f	fut
	s	sabun
	v	vorud
	x	xord
	z	ziba
	#	#ale
	\$	Sire
h	hale	
Affricates	j	javid
	c	c@medan
	+	فاصله بین کلمات
	-	سکوت

لازم به اشاره است که در نمایش قراردادی فوق، سمبل قرار گرفته در سمت راست علامت < و یا سمت چپ علامت > نشانگر واحد اصلی و سمبل دیگر نشانگر واحد متنی است.

ج - مدل های سه آوایی : هر واج با توجه به واج سمت راست و سمت چپ بطور جداگانه مدل می شود. همانند قبیل برای سکوت و فاصله بدون توجه به متن یک مدل در نظر می گیریم و در صورتی که واج کناری سکوت و یا فاصله باشد آنرا بصورت دو آوایی مدل می کنیم. به این ترتیب جمله بالا بصورت زیر نمایش داده می شود.

_ d>@ d<@>r @<r + b>a b<a>z. a<z + '>@ '<@>s @<s>t s<t _

توجه به این نکته لازم است که این گونه مدل سازی، سه آوایی داخل کلمه‌ای نامیده می شود. شیوه پیچیده تری به نام مدل سازی سه آوایی بین کلمه‌ای نیز وجود دارد که در مرزهای کلمات از مدل دو آوایی استفاده نمی کند بلکه به همراه فاصله، سکوت، آوای ابتدایی کلمه بعدی و یا آوای انتهایی کلمه قبلی یک مدل سه آوایی را بکار می گیرد. این شیوه موجب پیچیدگی فراوان در مدل سازی و بویژه در الگوریتم بازشناسی می گردد.

د - مدل‌های تعمیم یافته: اگر چه ممکن است به نظر برسد که هیچ محدودیتی بر اینکه این فرآیند (افزایش تعداد واحدهای کناری که در مدل‌سازی در نظر گرفته میشوند) تا کجا ادامه یابد وجود ندارد، در عمل عموماً مدل‌سازی به بیش از شرایط پنج آوایی^{۱۳} (تا دو آوا از هر طرف) گسترش پیدا نمی‌کند.

۳- آموزش

در مرحله آموزش ابتدا یک مدل اولیه مناسب برای واحدها بدست می‌آوریم (مدل تک آوایی ایجاد شده در [۷]). سپس مدل تک آوایی را به عنوان مدل اولیه برای سه آوایی‌هایی که هسته آنها همان آوا است در نظر می‌گیریم. آنگاه مدل‌های سه آوایی را آموزش می‌دهیم. بدیهی است که تعداد مدل‌های بدست آمده از این مرحله به مراتب بیشتر از تعداد مدل‌های اولیه (تک آوایی) می‌باشند.

۳-۱- الگوریتم Baum-Welch و تخمین پارامترها

اگر $\alpha(t, i)$ احتمال مشاهده t فریم اول باشد بطوری که در لحظه t در حالت i باشیم و $\beta(t, i)$ احتمال مشاهده $T-t-1$ فریم آخر باشد بطوری که در لحظه t در حالت i باشیم:

$$P(O(t) | \lambda) = \alpha(t, i) * \beta(t, i) \quad (1)$$

$$\alpha_{ij}^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{l=1}^{T_r-1} \alpha^{(q)r}(t, i) * a_{ij}^{(q)} * b_j^{(q)}(O_{l+1}^r) * \beta^{(q)r}(t+1, j)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{l=1}^{T_r} \alpha^{(q)r}(t, i) * \beta^{(q)r}(t, j)} \quad (2)$$

$$\mu_j^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{l=1}^{T_r} \alpha^{(q)r}(t, j) * \beta^{(q)r}(t, j) * O_j^r}{\sum_{r=1}^R \frac{1}{P_r} \sum_{l=1}^{T_r} \alpha^{(q)r}(t, j) * \beta^{(q)r}(t, j)} \quad (3)$$

$$\Sigma_j^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{l=1}^{T_r} \alpha^{(q)r}(t, j) * \beta^{(q)r}(t, j) * (O_j^r - \mu_j^{(q)})^2}{\sum_{r=1}^R \frac{1}{P_r} \sum_{l=1}^{T_r} \alpha^{(q)r}(t, j) * \beta^{(q)r}(t, j)} \quad (4)$$

که R تعداد مشاهدات هر مدل و T_r طول مشاهدات و q شماره مدل است. برای استفاده از این روابط در جملات بدون برچسب زمانی ابتدا مدل جمله از کنار هم قرار دادن تک تک مدل‌های زیر لغوی بدست می‌آید. به این مدل FSN^{۱۴} می‌گویند. سپس مقادیر جدید پارامترها با استفاده از روابط فوق تخمین زده می‌شود [۷][۱۸]. این الگوریتم باید بصورت مکرر روی داده‌های آموزشی اجرا گردد تا مقادیر مناسبی برای پارامترها بدست آید. برای ساده نویسی از جایگزین‌های زیر در معادلات فوق استفاده خواهیم کرد.

$$an_{ij}^{(q)} = \sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r-1} \alpha^{(q)r}(t,i) * a_{ij}^{(q)} * b_j^{(q)}(O_{t+1}^r) * \beta^{(q)r}(t+1,j) \quad (5)$$

$$ad_{ij}^{(q)} = \sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha^{(q)r}(t,i) * \beta^{(q)r}(t,j) \quad (6)$$

$$\mu n_j^{(q)} = \sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha^{(q)r}(t,j) * \beta^{(q)r}(t,j) * O_j^r \quad (7)$$

$$\mu d_j^{(q)} = \sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha^{(q)r}(t,j) * \beta^{(q)r}(t,j) \quad (8)$$

$$\Sigma n_j^{(q)} = \sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha^{(q)r}(t,j) * \beta^{(q)r}(t,j) * (O_j^r - \mu_j^{(q)})^2 \quad (9)$$

$$\Sigma d_j^{(q)} = \sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha^{(q)r}(t,j) * \beta^{(q)r}(t,j) \quad (10)$$

۴- گره زدن^{۱۵}

همانگونه که گفتیم هرچه تنوع واحدها و بنابراین تعداد مدلها افزایش یابد، با شرط آموزش مناسب، نتیجه بازشناسی بهتر خواهد بود. اما با افزایش مدلها و با داده های آموزشی محدود امکان آموزش نامناسب پیش می آید. بنابراین باید با تمهیداتی ضمن کاهش تعداد پارامترها، امکان آموزش صحیح را برای پارامترها فراهم کنیم.

یکی از این روشها گره زدن پارامترها است. در بدست آوردن مدلهای سه آوایی جدید حالتها برای رخدادهای شبیه ولی جدا از هم هر آوا با هم گره زده می شوند. شکل ۱ پارامترهای مختلفی، که در یک سیستم مبتنی بر HMM میتوانند با هم گره زده شوند، را نشان می دهد. دایرهها پارامترهایی را نشان میدهند که میتوانند گره زده شوند [۲].

تقریباً کلیه پارامترهای HMM قابل گره زدن هستند. در فرآیند گره زدن، پارامترهای گره زده شده در قسمتهای مختلف یکسان در نظر گرفته می شوند. مزیت گره زدن پارامترها در HMM این است که فرمولهای اولیه بدون تغییر عمده قابل استفاده مجدد هستند و گره زدن همگرایی برنامه آموزش را نقض نمی کند. گره زدن باعث می شود پارامترهای مختلف بین مدلهای مختلف به اشتراک گذاشته شوند که در نتیجه با داده های آموزشی کم تخمین مقاومی از پارامترها بدست می آید [۹].

۴-۱- گره زدن حالتها

حالتها گزینه خوبی برای گره زدن هستند، زیرا در عین حال که جزئی از هر مدل هستند می توانند مشخص کننده جزء زمانی هر آوا هم باشند [۳].

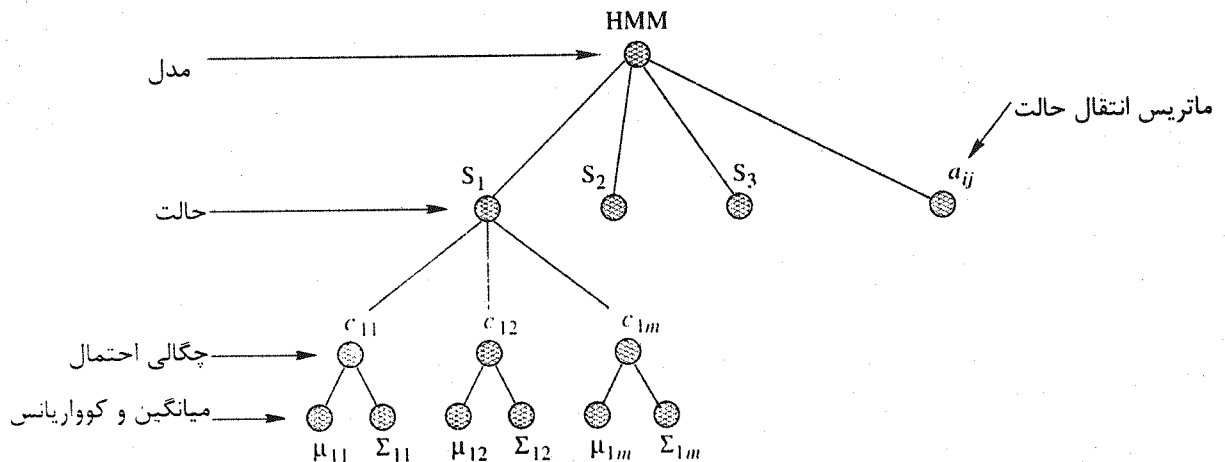
معادلات زیر چگونگی آموزش را در این حالت نشان می دهند.

$$a_{ij}^{(q)} = \frac{\sum_{n=1}^{N(q)} an_{ij}^{(M(q,n))} + an_{ij}^{(q)}}{\sum_{n=1}^{N(q)} ad_{ij}^{(M(q,n))} + ad_{ij}^{(q)}} \quad (11)$$

$$\mu_j^{(q)} = \frac{\sum_{n=1}^{N(q)} \mu_j^{(M(q,n))} + \mu_j^{(q)}}{\sum_{n=1}^{N(q)} \mu_j^{(M(q,n))} + \mu_j^{(q)}} \quad (12)$$

$$\Sigma_j^{(q)} = \frac{\sum_{n=1}^{N(q)} \Sigma_j^{(M(q,n))} + \Sigma_j^{(q)}}{\sum_{n=1}^{N(q)} \Sigma_j^{(M(q,n))} + \Sigma_j^{(q)}} \quad (13)$$

که $N(q)$ معرف تعداد حالاتی است که با حالت q ام خورده اند و $M(q,n)$ ، n امین حالتی است که با حالت q گره خورده است.



شکل (۱) ساختار سلسله مراتبی برای گره زدن پارامترها در HMM.

۴-۱- الگوریتم خوشه بندی

پیش از آنکه امکان گره زدن حالتها (یا بطور کلی پارامترها) در یک سیستم فراهم شود، باید مشخص نمود که کدام حالتها باید با یکدیگر گره بخورند. یک راه ممکن برای نیل به این هدف اجرای یک نوع خوشه بندی بر روی تمامی حالتهای سیستم است. در این شرایط معمولاً پارامترهای یک سیستم آموزش داده شده از نوع وابسته به متن مورد استفاده قرار می گیرند و خوشه بندی بر روی این پارامترها صورت می گیرد.

برای پیاده سازی عمل خوشه بندی می توان از چنین الگوریتمی استفاده نمود [۱۰].

۱- برای هر یک از حالتها سیستم یک دسته جدا در نظر می گیریم.

۲- دو دسته ای که کمترین فاصله را دارند پیدا می کنیم.

۳- اگر این فاصله از مقدار از قبل تعیین شده‌ای کمتر بود این دو دسته را ادغام می‌کنیم و به قسمت ۲ برمی‌گردیم.

۴- پایان.

حالت‌هایی را که در پایان مراحل فوق در یک دسته قرار می‌گیرند با هم گروه می‌زنیم.

برای یافتن فاصله دسته‌هایی که دارای چندین عضو هستند فاصله تمام اعضاء هر دسته را با تمام عضوهای دسته دیگر محاسبه می‌کنیم و سپس از آنها میانگین می‌گیریم. در این الگوریتم از رابطه دیورژانس، به شرح زیر برای بدست آوردن فاصله بین دسته‌ها استفاده کرده‌ایم [۱۱].

$$d(S_1, S_2) = \left(\frac{1}{N} \sum_{i=1}^N \frac{(\mu_{1i} - \mu_{2i})^2}{\sigma_{1i} \sigma_{2i}} \right)^{\frac{1}{2}} \quad (14)$$

۵ - تخمین MAP

تخمین MAP که گاه تحت عنوان روش بیزین^{۱۷} در آموزش HMMها نیز نامیده می‌شود، بر اساس تخمین بیشینه پسین^{۱۸} عمل می‌کند. در این روش، در تخمین پارامترهای HMMها علاوه بر مقادیر مشاهدات از مقادیر پیشینه نیز کمک گرفته می‌شود. به این ترتیب، بر خلاف تخمین درست‌نمایی بیشینه (ML) که در شرایط داده‌های آموزشی پراکنده^{۱۹} به تخمینهای غیر قابل اعتماد منجر می‌شود، در این شرایط می‌توان تخمین‌هایی نه چندان دور از واقعیت بدست آورد. البته رسیدن به این هدف مشروط بر انتخاب پارامترهای پیشینه مناسب است.

فرمولهای آموزش MAP برای یک HMM دارای چگالی مشاهدات تک گوسین بصورت زیر است [۹].

$$a_{ij}^{(q)} = \frac{(\eta_{ij} - 1) + a n_{ij}^{(q)}}{\sum_{k=1}^n (\eta_{ik} - 1) + a d_{ij}^{(q)}} \quad (15)$$

$$\mu_j^{(q)} = \frac{\tau_j \hat{\mu}_j^{(q)} + \mu n_j^{(q)}}{\tau_j + \mu d_j^{(q)}} \quad (16)$$

$$\Sigma_j^{(q)} = \frac{u_j^{(q)} + \tau_j (\hat{\mu}_j^{(q)} - \mu_j^{(q)})^2 + \Sigma n_j^{(q)}}{(\gamma_j - p) + \Sigma d_j^{(q)}} \quad (17)$$

که در این روابط p بعد بردار میانگین (و ماتریس کواریانس) و مقادیر η ، τ ، γ ، $\hat{\mu}_j^{(q)}$ و $u_j^{(q)}$ پارامترهای توزیع پیشینه می‌باشند که در این میان τ و $\hat{\mu}_j^{(q)}$ پارامترهای توزیع گوسین چند متغیره پیشینه، γ و $u_j^{(q)}$ پارامترهای توزیع و یشارت پیشینه^{۲۰} و η پارامتر توزیع دیریشله^{۲۱} پیشینه می‌باشند [۹].

آنچه در تخمین به شیوه MAP عمدتاً مشکل‌ساز است، تخمین مناسب پارامترهای پیشینه می‌باشد. از آنجا که مزیت اصلی تخمین MAP بر تخمین ML، استفاده از پارامترهای پیشینه است که در شرایط کمبود داده آموزشی از تخمین نامناسب پارامترهای مدل جلوگیری می‌نمایند، لزوم داشتن پارامترهای پیشینه مناسب انکار ناپذیر است. لازم به تذکر است که علیرغم تلاشی که در این زمینه صورت گرفته، یک راه تحلیلی مناسب برای حل این مسأله یافت نشده است [۱۲].

۶- آموزش وزنی

تلاش در آموزش مدل‌های وابسته به متن در اینجا بر این نکته متمرکز گشته که تا حد امکان از شباهت و نزدیکی مدل‌ها در

حین آموزش و برای جبران مشکل کمبود داده‌ها استفاده نمود. به همین منظور روابط آموزش (باز تخمین) پارامترهای CDHMM به شرح زیر تغییر داده شده‌اند.

$$\alpha_{ij}^{(q)} = \frac{\varphi \sum_{n=1}^{N(q)} a n_{ij}^{(M(q,n))} + a n_{ij}^{(q)}}{\varphi \sum_{n=1}^{N(q)} a d_{ij}^{(M(q,n))} + a d_{ij}^{(q)}} \quad (18)$$

$$\mu_j^{(q)} = \frac{\varphi \sum_{n=1}^{N(q)} \mu n_j^{(M(q,n))} + \mu n_j^{(q)}}{\varphi \sum_{n=1}^{N(q)} \mu d_j^{(M(q,n))} + \mu d_j^{(q)}} \quad (19)$$

$$\Sigma n_j^{(q)} = \frac{\varphi \sum_{n=1}^{N(q)} \Sigma n_j^{(M(q,n))} + \Sigma n_j^{(q)}}{\varphi \sum_{n=1}^{N(q)} \Sigma d_j^{(M(q,n))} + \Sigma d_j^{(q)}} \quad (20)$$

که $N(q)$ معرف تعداد حالاتی است که با حالت q ام شبیه در نظر گرفته شده‌اند و $M(q,n)$ ، n امین حالتی است که شبیه به حالت q است. φ پارامتری است بین ۰ و ۱ که میزان اثر مدل‌های دیگر را در مدل q تعیین می‌کند.

بطوری که مشاهده می‌شود، روابط فوق به هر دو روش گره زدن و MAP شباهتهایی دارند. در واقع اگر در این روابط φ برابر یک قرار داده شود به فرمولهای گره زدن می‌رسیم، و هر چه این پارامتر کوچکتر شود نتیجه به سمت مدل سه آوایی گره نخورده نزدیک می‌شود. این روش مشکل عمده MAP یعنی تعداد زیادی پارامتر آزاد که تاثیر عمده ای بر نتیجه آموزش داشته و باید به درستی تخمین زده شوند را ندارد و تنها یک پارامتر است که باید تعیین شود. در حالی که محاسن روش MAP از جهت کمک گرفتن از پارامترهای مدل های شبیه (که در آنجا مقادیر پیشینه نامیده می شوند) را داراست. در عین حال، بر خلاف شیوه گره زدن پارامترها، تفاوت بین مدلها را، حتی در پایین ترین سطوح، هم حفظ می‌کند.

مزیت دیگر این روش نسبت به MAP این است که عباراتی که در صورت یا مخرج با هم جمع می‌شوند ممنوع هستند و در ضمن نوعی بیانگر میزان نمونه‌های آموزشی استفاده شده در تخمین هستند. برای هر مدل، هر چه که تعداد نمونه‌های آموزشی بیشتر باشد، مقدار بدست آمده از روابط (۵) تا (۱۰) بیشتر خواهد بود. تفاوت دیگر این روش با روش گره زدن بر اساس داده در این است که مستقیماً از مدل تک آوایی به مدل سه آوایی می‌رسیم و نیازی به دو مرحله آموزش، یکی پیش از مرحله خوشه بندی و گره زدن و دیگری بعد از آن نمی‌باشد. به این ترتیب، روش آموزش وزنی ضمن ارائه مزایای روش گره زدن، جایگزین مناسبی نیز برای روش MAP می‌باشد.

۷- پیاده سازی

سیستم بازشناسی موجود با استفاده از داده‌های فارس دات آموزش داده شده و آزمایش گردیده‌اند. این پایگاه داده در شکل اصلی خود متشکل از ۶۰۰۰ ادای جمله است که از ۳۰۰ گوینده بدست آمده است. هر گوینده بطور تصادفی ۲۰ جمله از ۴۰۰ جمله موجود را خوانده است که ۲ جمله بین تمام گوینده‌ها مشترک است. با توجه به اینکه گوینده‌ها از نقاط مختلف کشور و با لهجه‌های مختلف انتخاب شده‌اند، جملات مربوط به گوینده‌های تهرانی استخراج گردیده که پس از بررسی‌ها و

انتخاب نهایی حدود ۲۷۰۰ جمله بدست آمده است. این جملات به دو دسته ۱۸۰۰ و ۹۰۰ تایی تقسیم شده و از دسته اول برای آموزش و از دسته دوم برای تست استفاده گردید [۱۳].

پایگاه داده فارس دات حدود ۱۲۰۰ کلمه مختلف را شامل میشود و در آن حدود ۲۴۰۰ سه آوایی مجزا وجود دارد. با توجه به اینکه تعداد آواهای پایه در زبان فارسی، در این تحقیق، ۳۰ در نظر گرفته شده است، با یک محاسبه ساده، تعداد سه آوایی‌ها باید به ۲۷۰۰۰ بالغ گردد. در این شرایط، تنها با در نظر گرفتن گروه‌های همخوان و واکه و نیز ساختار هجائی زبان فارسی، از ۸ ترکیب ممکن سه گانه واکه و همخوان، تنها ۴ ترکیب CVC، CCV، VCV و VCC در زبان فارسی ممکن است. به این ترتیب، با توجه به تعداد ۷ واکه و ۲۳ همخوان موجود، تعداد سه آوایی‌های حاصل به ۱۲۲۳۶ کاهش خواهد یافت. البته در عمل، بسیاری از این ترکیب‌های سه آوایی ممکن نمی‌باشند [۱۴] و لذا تعداد واقعی سه آوایی‌ها در زبان فارسی به کسری از عدد فوق محدود می‌گردد. با این حال، و با توجه به اینکه مطالعه زیادی در زبان فارسی بر روی سه آوایی‌ها صورت نگرفته، آمار دقیقی در این خصوص در منابع مختلف یافت نشده و احتمال می‌رود که تعداد واقعی سه آوایی‌ها در زبان فارسی به بیش از دو برابر تعداد موجود در دادگان فارس دات هم برسد.

محدود بودن حجم دادگان فارس دات و پراکنده بودن سه آوایی‌ها در آن (از نظر تعداد دفعات وقوع) به خوبی نشان‌دهنده این واقعیت است که آموزش مستقیم آنها مواجه با مشکل داده پراکنده بوده و نتیجه مطلوبی نخواهد داشت. به عنوان مثال، از میان بیش از ۲۴۰۰ سه آوایی بدست آمده از فارس دات، ۲۳۶ سه آوایی کمتر از سه نمونه آموزشی و ۸۱۱ سه آوایی کمتر از پنج نمونه آموزشی داشته‌اند. آمار کاملتری از تعداد دفعات وقوع سه آوایی‌ها در متن آموزشی در [۱۵] ارائه گردیده است. همین امر به خوبی بستر لازم را برای پیاده سازی الگوریتمی که در شرایط داده پراکنده بتواند با استفاده از اطلاعات سایر مدل‌ها اقدام به آموزش مقاوم مدل‌های با داده آموزشی کم نماید فراهم می‌کند. بنابراین، الگوریتم آموزش وزنی، الگوریتمی کاملاً مناسب برای آموزش سه آوایی‌ها در شرایط حاضر می‌باشد.

هدف از آموزش مدل‌ها استفاده از آنها برای بازشناسی می‌باشد. لذا طراحی و پیاده سازی یک الگوریتم بازشناسی قوی و در عین حال سریع اهمیت فراوانی دارد. روشی که جهت بازشناسی بکار گرفته شده است موسوم به الگوریتم انتقال نشانه^{۲۳} می‌باشد [۱۷]. مدل‌های زبان مورد استفاده برای آزمایش سیستم طراحی شده، زوج کلمه^{۲۳} و بدون گرامر^{۲۴} می‌باشند.

۸ - نتایج

در پیاده‌سازی از یک توزیع گوسین برای هر حالت استفاده شده است. بدیهی است که استفاده از مخلوط گوسین موجب بهبود نتیجه بازشناسی خواهد شد. علت عدم استفاده از آن در اینجا به اشباع رسیدن سیستم بازشناسی بویژه با مدل زبان زوج کلمه می‌باشد. برای آموزش وزنی، سه آوایی‌های با هسته مشترک را شبیه به هم گرفته‌ایم و ضریب ϕ برابر ۰/۰۰۵ قرار داده شده است. همچنین در بدست آوردن مدل سه آوایی گره زده نشده، برای جلوگیری از آموزش نامناسب، مدل‌هایی که کمتر از سه نمونه آموزشی داشته‌اند دست نخورده باقی مانده‌اند. مدل‌های سه آوایی گره زده شده بر اساس الگوریتم بحث شده در بخش ۴-۱ بدست آمده‌اند.

جدول ۲ نشان دهنده نتایج بدست آمده از آزمایش‌های بازشناسی گفتار می‌باشد. نتایج اعلام شده برای آزمایش‌های بدون گرامر با استفاده از حدود ۱۲۰ جمله آزمایشی بدست آمده‌اند و مبین دقت بازشناخت کلمه^{۲۵} می‌باشند. این جملات از ابتدای مجموعه آزمایشی و به ترتیب انتخاب و مورد استفاده قرار گرفته‌اند. علت استفاده از تنها ۱۲۰ جمله برای آزمایش در این مرحله، بکارگرفته نشدن مدل زبان می‌باشد که باعث طولانی شدن فرایند آموزش (زمان بازشناسی) می‌گردد. بنابراین به کار گرفتن ۹۰۰ جمله در آزمایش‌ها مشکلات عملی به همراه دارد. بررسی‌های انجام شده نشان دهنده این واقعیت می‌باشند که تفاوت چندانی بین نتایج بدست آمده از ۱۲۰ و ۹۰۰ جمله در این حالت وجود ندارد و نتایج بدست آمده کاملاً قابل اعتماد می‌باشند.

جدول (۲) نتایج بازشناسی

مدل	تعداد حالت‌ها	مدل زبان	دقت بازشناسی
تک آوایی	۹۶	بدون گرامر	٪۴۶/۲
سه آوایی بدون گره زدن	۷۲۷۸		٪۶۰/۱
سه آوایی گره خورده	۱۴۰۰		٪۶۷/۱۸
سه آوایی با آموزش وزنی	۷۲۷۸		٪۷۲/۹
تک آوایی	۹۶	زوج کلمه	٪۹۳/۱
سه آوایی بدون گره زدن	۷۲۷۸		٪۷۲/۱۰
سه آوایی گره خورده	۱۴۰۰		٪۹۴/۱
سه آوایی با آموزش وزنی	۷۲۷۸		٪۹۶/۵

نتایج بخش بدون گرامر بخوبی حاکی از برتری شیوه آموزش پیشنهادی بر روش معمول آموزش یعنی درستمائی بیشینه است. در حالیکه مدل‌های سه‌آوایی، پیش از گره زدن، بوضوح نسبت به مدل‌های تک‌آوایی بهتر عمل نموده‌اند، گره‌زدن حالت در آنها موجب بهبودی بیش از پیش نتایج گردیده است. با این همه، شیوه پیشنهادی آموزش توانسته در این شرایط حتی به نتایج برتری نیز دست یابد. اعمال مدل زبان زوج کلمه، طبیعتاً موجب بهبودی در نتایج بدست آمده از کلیه آزمایش‌های فوق می‌گردد. با این همه، این بهبودی در سیستم مبتنی بر مدل‌های تک‌آوایی بیش از حد انتظار است. علت این امر پایین بودن ضریب "سرگشتگی" ^{۲۶} در دادگان تصور می‌شود. این مطلب می‌تواند دقت بازشناسی را بشدت افزایش دهد. با این همه، در استفاده از مدل‌های سه‌آوایی گره نخورده، با توجه به تعداد زیاد مدل‌ها و اینکه بسیاری از آنها بخوبی آموزش ندیده‌اند، احتمال خطا در بازشناخت بعضی کلمات زیادتر است و محدودیت انتخاب کلمات بعدی در بسیاری موارد از اصلاح نتیجه جلوگیری می‌نماید. با این حال، در این شرایط هم روش پیشنهادی آموزش به بهترین نتایج دست یافته است.

ذکر یک نکته نیز در اینجا لازم است و آن اینکه ستون تعداد حالت‌ها در جدول ۲ نشان‌دهنده کل تعداد حالت‌های موجود در سیستم است. به عنوان مثال، برای مدل‌های تک آوایی، ۳۲ مدل سه حالت (۹۶ حالت) و برای سه آوایی‌های بدون گره زدن، ۲۴۲۶ مدل سه حالت (۷۲۷۸ حالت) وجود دارند. همین تعداد مدل و حالت در سه آوایی‌های با آموزش وزنی نیز دیده می‌شود. در مورد سه آوایی‌های گره خورده، با توجه به فرایند گره زدن حالت‌ها، ۷۲۷۸ حالت منطقی به ۱۴۰۰ حالت فیزیکی کاهش یافته است. لازم به اشاره است که تعداد حالت‌های فیزیکی بدست آمده کاملاً تابع سطح آستانه مورد استفاده در الگوریتم خوشه‌بندی بخش ۴-۱-۱ است.

جدول (۳) نتایج بازشناسی به ازاء ضرایب مختلف تضعیف در آموزش وزنی.

ϕ	۱	۰/۱	۰/۰۵	۰/۰۲	۰/۰۱	۰/۰۰۵	۰/۰۰۱
NG	٪۴۶/۲۴	٪۶۵/۲۶	٪۶۸/۱۹	٪۷۱/۱۱	٪۷۱/۶۶	٪۷۲/۹۴	٪۶۸/۳۷
WP	٪۹۴/۱۶	٪۹۶/۷۳	٪۹۶/۹۳	٪۹۷/۴۵	٪۹۷/۳۲	٪۹۶/۵۵	٪۹۱/۳۱

جدول (۳) نتایج بازشناسی را به ازاء ضرایب تضعیف مختلف نشان می‌دهد. همانگونه که دیده می‌شود با کوچک شدن ضریب تضعیف ابتدا دقت سیستم افزایش می‌یابد، اما بعد همانطور که انتظار داریم رو به کاهش می‌نهد. در هر دو حالت بدون گرامر و استفاده از گرامر زوج کلمه با این وضعیت مواجه هستیم ولی میزان پارامتر ϕ برای دستیابی به نرخ بازشناسی بیشینه در این دو حالت متفاوت است.

۹- نتیجه گیری و پیشنهادات

بررسی‌های انجام شده نشان دهنده این واقعیتند که استفاده از مدل‌های وابسته به متن نتایج بهتری از مدل‌های ناپسته به متن دارند. با این همه در مدلسازی وابسته به متن تعداد پارامترها می‌تواند بشدت افزایش یابد. بنابراین، برای رسیدن به نتیجه بهتر باید تمهیداتی برای مقاوم کردن آموزش اندیشید.

روش آموزش وزنی پیشنهادی، که بر اساس استفاده از پارامترهای مدل‌های هم خانواده عمل می‌نماید، هم در حالت بدون گرامر و هم در شرایط اعمال مدل زبان (در اینجا مدل زوج کلمه) به نتایج بهتری نسبت به روش معمول آموزش دست یافته است. این نتایج در مقایسه با دقت حاصل از یک سیستم با حالت‌های گره زده شده (روش معمول مورد استفاده امروزه در سیستم‌های بازشناسی گفتار در سطح جهان) نیز حکایت از برتری نسبی روش آموزش وزنی دارند. پیشنهادات زیر برای ادامه کار مفید به نظر می‌رسند:

- الف - استفاده از توزیع گوسین با تعداد عناصر مخلوط بزرگتر.
- ب - استفاده از پایگاه داده کاملتر.
- ج - استفاده از مدل‌های زبان پیچیده‌تر (مشروط به رعایت سد ب).
- د - استفاده از مدل‌های زیر لغوی وابسته به متن بین کلمه‌ای.
- ه - تحقیق بیشتر در مورد اینکه چه مدل‌هایی را شبیه به هم در نظر بگیریم. برای این کار می‌توان همانند گره زدن عمل کرد، ولی در پایان از روابط (۱۸) تا (۲۰) برای آموزش استفاده کرد.
- و - گره زدن پارامترها پس از آموزش وزنی و استفاده از این روش برای آموزش مجدد.
- ز - تحقیق در مورد بدست آوردن پارامتر φ بصورت تحلیلی. احتمالاً این کار همانند روش MAP قابل انجام است.
- ح - استفاده از آموزش وزنی در کاربردهائی نظیر تطبیق گوینده.

زیر نویس‌ها

- 1- Context Dependent Modeling
- 2- Context
- 3- Robust
- 4- Maximum a Posteriori
- 5- Bayes
- 6- Prior
- 7- Maximum Likelihood
- 8- Triphones
- 9- Monophones
- 10- Biphones
- 11- Word-internal triphones
- 12- Cross-word triphone modeling
- 13- Quinphones
- 14- Maximum Likelihood
- 15- Finite State Network
- 16- Tying
- 17- Clu16- stering
- 18- Bayesian
- 19- Maximum a Posteriori
- 20- Sparse Training Data
- 21- Wishart distribution
- 22- Dirichlet distribution
- 23- Token passing
- 24- Word-pair
- 25- No-gram
- 26- Word recognition accuracy
- 27- Perplexity

- [1] P.C. Woodland, C.J. Leggetter, J.J. Odell, V. Valtchev and S.J. Young, "The Development of the 1994 HTK Large Vocabulary Speech Recognition System", In Proc. ARPA Spoken Language Technology Workshop, Bartoncreek, 1995.
- [2] S.J. Young, "The general use of Tying in Phoneme-Based HMM Speech Recognisers", Proc. ICASSP-93, Vol. 1, pp. 569-572, San Francisco.
- [3] S.J. Young, J.J. Odell and P.C. Woodland, "Tree-Based State Tying for High Accuracy Acoustic Modelling", In ARPA Human Language Technology Workshop, pp. 307-312, Plainsboro, March 1994.
- [4] J-L. Gauvain, and C-H. Lee, "Bayesian Learning of Gaussian Mixture Densities for Hidden Markov Models", Speech Communication, vol. 11, 1992, pp.205-213.
- [5] J. Ming, P. O'Boyle, M. Owens and F.J. Smith, "A Bayesian Approach for Building Triphone Models for Continuous Speech Recognition", IEEE Trans. Speech and Audio Proc., vol. 7, No. 6, Nov. 1999, pp. 678-684.
- [6] J.J. Odell, "The Use of Context in Large Vocabulary Speech Recognition", PhD Thesis, Cambridge University Engineering Department, 1995.
- [۷] م. پرویزی، س. م. احدی، "بازشناسی گفتار فارسی پیوسته با دایره کلمات متوسط"، مجموعه مقالات کنفرانس مهندسی برق ایران، ۱۳۸۰.
- [8] L.R. Rabiner and B.H. Juang, Fundamentals of Speech Recognition, New Jersey, Prentice-Hall, 1993.
- [9] J-L. Gauvain and C-H. Lee, "Maximum a Posteriori Estimation for Multivariate Gaussian Mixture Observations of Markov Chains", IEEE Transactions on Speech and Audio Processing, Vol. 2, No. 2, pp. 291-298, April 1994.
- [10] S.M. Ahadi, "Reduced Context Sensitivity in Persian Speech Recognition via Syllable Modeling", In Proc. SST-2000, pp. 492-497, Canberra.
- [11] HTK: Hidden Markov Model Toolkit Ver. 1.4, Reference Manual, Cambridge University Engineering Department, p. 76, 1992.
- [12] Q. Huo and C. Chan, "Bayesian Adaptive Learning of the Parameters of Hidden Markov Model for Speech Recognition", Technical Report, Department of Computer Science, University of Hong Kong, September 1992.
- [13] S.M.Ahadi, "Recognition of Continuous Persian Speech Using a Medium-Sized Vocabulary Speech Corpus", in Proc. EUROSPEECH'99, Budapest.
- [۱۴] ی. ثمره، آواشناسی زبان فارسی، ویرایش دوم، مرکز نشر دانشگاهی، تهران ۱۳۸۰.
- [۱۵] س. ح. شمس، "مدل سازی وابسته به متن در بازشناسی گفتار پیوسته فارسی"، پایان نامه کارشناسی ارشد، دانشکده مهندسی برق، دانشگاه صنعتی امیرکبیر، ۱۳۸۰.