

## کاربرد مدل‌های پنهانی مارکف برای درک مستقل از گوینده کلمات منفرد فارسی

بابک عظیمی سجادی\* - وحید طباطبائی\* - سید بهرام ظهیر اعظمی\*\* - کارولوکس\*\*\*

### چکیده

در این مقاله نحوه تحقیق یک سیستم درک صحبت به روش مدل پنهانی مارکف با مشاهده پیوسته تشریح می‌گردد. از این مدل جهت تشخیص تعداد محدودی از کلمات (اعداد «صفر» تا «ده») که به زبان فارسی بیان شده اند به صورت مستقل از گوینده استفاده شده است. در این روش برای هر کلمه یک مدل چپ به راست در نظر گرفته شده است. پارامترهای مدل با تخمین بیشترین شباهت از روی داده های یادگیری تعیین می‌شوند. این تخمین از روش  $k$  مقدار متوسط در هر قسمت به دست می‌آید. برای تعیین مقادیر اولیه پارامترهای مدل، روشی جدید بر مبنای چندی کردن برداری ارائه می‌شود که با استفاده از آن تمامی مراحل طراحی سیستم درک کننده صحبت به صورت خودکار انجام می‌شود. کارایی سیستم محقق شده  $97/00\%$  است.

---

\* فارغ التحصیلان کارشناسی ارشد دانشکده فنی دانشگاه تهران و کارشناسان مرکز تحقیقات الکترونیک

دانشگاه صنعتی شریف

\*\* کارشناس مرکز تحقیقات الکترونیک دانشگاه صنعتی شریف

\*\*\* دانشیار دانشکده فنی دانشگاه تهران

## مقدمه

ارتباط گفتاری بین انسان و ماشین نسبت به سایر موارد، مزایای عمده ای از قبیل سهولت کاربرد و عدم نیاز به آموزش‌های پیچیده دارد. از این رو توجه بسیاری از مراکز تحقیقاتی و صنعتی جهان به این مسأله جلب شده است. در ارتباط گفتاری، کاربر دستورات و یا اطلاعات را به وسیله گفتار به کامپیوتر انتقال می دهد، کامپیوتر هم با بیان کلمات و یا جملاتی با استفاده کننده ارتباط برقرار می نماید. برای ایجاد چنین امکانی کامپیوتر ابتدا باید صحبت گوینده را درک کند و سپس بتواند منظور خود را توسط کلام بیان نماید.

در این مقاله روشی برای ارتباط انسان با ماشین یا درک صحبت توسط کامپیوتر پیشنهاد می شود.<sup>۱</sup> برای طراحی یک سیستم درک صحبت ابتدا باید بتوان نحوه تولید صحبت را مدل کرد تولید صحبت در انسان عملی است هماهنگ که در آن اندامهای گفتاری دخالت دارند. در زمان ادای هر واج<sup>۲</sup> این اندامها در وضعیت خاصی قرار می گیرند به طوری که امواج صوتی حاصل از این عمل خواص معینی به دست می آورند. ولی تفاوت هایی در محل قرار گرفتن اندامها، محل وقوع هر واج در بین واجهای دیگر و تفاوت های فیزیکی و هندسی اندامهای گفتاری افراد مختلف خواص صوتی متفاوتی را در بین واجهای بیان شده ایجاد می کند [۲] این تفاوتها به شکل وسیعتری بین هجاها<sup>۳</sup>، واژه ها و جملات نیز پدید می آید ([۱ و ۲]).

مسأله اصلی در درک صحبت، استخراج خصوصیات مشترک در اشکال مختلف ادای هر واج (واژه) است. این خصوصیات باید دارای بیشترین شباهت در اداهای مختلف یک واج (واژه) و بیشترین اختلاف در ادای واجهای (واژه های) گوناگون باشد.

درک صحبت در دو سطح معنایی<sup>۴</sup> و ساختاری<sup>۵</sup> مطرح می شود. سطح ساختاری که در این مقاله مورد نظر است، از سه جنبه تعداد کلمات قابل تشخیص، محدودیت های گفتاری در بیان کلمات و

1 . Speech Recognition

۲ - حرف قابل گویش در یک زبان (Phoneme)

۳ - رشته ای پیوسته و بدون مکث از واجها (Syllable)

4 . Semantic

5 . Syntactic

تنوع استفاده کنندگان تقسیم بندی می شود. هر اندازه محدودیت‌های کمتری برای استفاده کننده قائل شویم به سیستم پیچیده تری نیاز خواهیم داشت. به طوری که هنوز سیستم قابل اعتمادی که بدون هیچ محدودیتی درک صحبت را انجام دهد ارائه نشده است. با وجود پیشرفت های قابل ملاحظه حتی طراحی سیستم هایی که محدودیت‌هایی برای استفاده کننده قائل می شوند، هنوز هم از مسایل در درست تحقیق به شمار می رود [۲].

سه روش عمده ای که در حال حاضر برای درک صحبت به کار می روند عبارتند از:

۱- مقایسه پارامترهای سیگنال صحبت با یک سری پارامترهای نمونه<sup>۱</sup> ([۳ تا ۵])

۲- استفاده از شبکه های عصبی ([۶ و ۷])

۳- تشخیص صحبت براساس خصوصیات آماری سیگنال صحبت ([۸ تا ۱۰])

در این مقاله تحقق روش مدل پنهانی مارکف<sup>۲</sup> به عنوان یک روش آماری موفق ارائه خواهد شد. با استفاده از این مدل یک سیستم درک صحبت مستقل از گوینده<sup>۳</sup> برای تشخیص اعداد "صفر" تا "ده" که به صورت تک کلمه ای<sup>۴</sup> به زبان فارسی ادا می شوند طراحی می شود. در بخش بعدی مدل پنهانی مارکف معرفی شده، مسائل و مشکلات آن توضیح داده می شود. سپس به نحوه استخراج مشخصه های سیگنال صحبت خواهیم پرداخت. آنگاه روش انتخاب مقادیر اولیه و مسائل عملی بررسی می شود و پس از آن نحوه آزمایش و نتایج آن ارائه خواهد شد. در انتها کارایی روش ارائه شده مورد ارزیابی قرار می گیرد و علیرغم وجود تفاوت اصولی در نتایج این آزمایشها با نتایج به دست آمده توسط دیگر محققان برای زبان انگلیسی، جهت تکمیل مطلب مقایسه ای صورت می گیرد.

معرفی مدل پنهانی مارکف ([۱۰ تا ۱۲])

در این روش فرض می شود که مشخصه های سیگنال صحبت یک فرایند مارکف مرتبه اول با

1 . Dynamic Time Warping (DTW)

2 . Hidden Markov Model(HMM)

3 . Speaker Independent

4 . Isolated Word

مشاهده ناقص<sup>۱</sup> است. منظور از فرایند مارکف مرتبه اول فرایندی است که در آن:

$$p(x_n | x_{n-1}, x_{n-2}, \dots, x_0) = p(x_n | x_{n-1}) \quad (1)$$

به عبارت دیگر با داشتن حالت در لحظه  $n-1$ ، حالت در لحظه  $n$  از لحظات قبل از  $n-1$  مستقل است. مثالی خوب برای فرایند مارکف یک ماشین حالت<sup>۲</sup> است که گذار حالت در آن به صورت تصادفی با احتمال های مشخص انجام می شود. حال اگر در فرایند مارکف به جای حالت، یک متغیر تصادفی وابسته به حالت مشاهده شود، به آن فرایند مارکف با مشاهده ناقص می گویند. به مدل چنین فرایندی، مدل پنهانی مارکف گفته می شود. مدل پنهانی مارکف مدل خوبی برای بیان رفتار سیگنال صحبت است.

برای تشخیص صحبت به روش فوق به ازای هر یک از کلمات مقطع، یک مدل مارکف در نظر گرفته می شود. در مرحله یادگیری ابتدا پارامترهای مدلها تخمین زده شده، سپس در مرحله تشخیص از مدلهای تخمین زده شده استفاده می گردد و کلمه ای که مدل متناظر با آن دارای بیشترین احتمال باشد، انتخاب می شود.

پارامترهای مدل پنهانی مارکف در حالت مشاهده گسسته عبارتند از:

$$1 - N \text{ تعداد حالتها، حالتها ممکن } S_1, \dots, S_N$$

$$2 - M \text{ تعداد الفبای مشاهده، الفبای مشاهده } V_1, \dots, V_M$$

$$3 - A \text{ ماتریس گذار حالت،}$$

$$A = \{a_{ij}\}; a_{ij} = p[q_{t+1} = S_j | q_t = S_i] \quad (2)$$

در این رابطه  $1 \leq i, j \leq N$  و  $t$  مشخص کننده زمان است.

۴ - ماتریس احتمال صدور الفبای  $k$  در حالت  $j$ ،

$$B = \{b_{jk}\}; b_{jk} = p[V_k \text{ at } t | q_t = S_j] \quad ; \quad 1 \leq j \leq N \quad (3)$$

در این جا  $1 \leq k \leq M$  و  $t$  مشخص کننده زمان است.

۵ - توزیع آغازین حالت،

$$\pi = \{\pi_i\}; \pi_i = p[q_1 = s_i] \quad ; \quad 1 \leq i \leq N \quad (4)$$

سلسله<sup>۱</sup> مشاهدات را با  $O$  و سلسله حالات طی شده را با  $Q$  نمایش می دهیم:

$$O = o_1 o_2, \dots, o_T \quad (5)$$

$$Q = q_1 q_2, \dots, q_T \quad (6)$$

در این رابطه  $T$  کل زمان مشاهده است. در ضمن جهت اختصار، مدل پنهانی مارکف با  $\lambda$  نمایش داده می شود:

$$\lambda = (A, B, \pi) \quad (7)$$

در مراحل مختلف یادگیری و تشخیص صحبت توسط این مدل سه مسأله وجود دارد که عبارتند از:

مسئله اول: در مرحله تشخیص باید احتمال وقوع یک مشاهده را به شرط انتخاب هر مدل محاسبه کرد. به عبارت دیگر به ازای هر کلمه (در این جا ۱۱ کلمه) یک مدل  $\lambda_i$  وجود دارد و صحبت ادا شده به سلسله مشاهده  $O$  تبدیل شده است. حال، محاسبه  $p(O|\lambda_i), i = 0, \dots, 10$  مد نظر خواهد بود. پس از محاسبه، مقدار ماگزیمم مشخص شده و کلمه متناظر با مدل دارای بیشترین احتمال، به عنوان کلمه برگزیده معرفی می گردد.

در صورت استفاده از روش مستقیم در محاسبه  $P(O|\lambda_i)$  یعنی محاسبه

$$P(O|\lambda) = \sum_Q P(O|Q, \lambda).P(Q|\lambda) \quad (۸)$$

$$= \sum_{q_1, q_2, \dots, q_T} \prod_{q_1}^{all Q} b_{q_1}(o_1) a_{q_1 q_2} b_{q_2}(o_2) \dots a_{q_{T-1} q_T} b_{q_T}(o_T)$$

مرتبه محاسبات  $2T.N^T$  خواهد بود که برای  $T = 40$  و  $N = 6$  برابر  $10^{23}$  می شود. این میزان محاسبات از توان کامپیوترهای موجود خارج است. با استفاده از روش بازگشتی<sup>۱</sup> می توان در هر مرحله احتمال  $P(o_1, \dots, o_t, q_t = S_i|\lambda)$  را از روی اطلاعات یک مرحله قبل به صورت زیر محاسبه کرد:

$$P(o_1, \dots, o_t, q_t = S_i|\lambda) = \quad (۹)$$

$$\sum_{j=1}^N a_{ji} P(o_1, \dots, o_{t-1}, q_{t-1} = S_j|\lambda).b_i(o_t); i=1, \dots, N$$

و نهایتاً با داشتن احتمالات در زمان  $T$ ،  $P(O|\lambda)$  را محاسبه نمود.

$$P(O|\lambda) = \sum_{i=1}^N P(o_1, \dots, o_T, q_T = S_i|\lambda) \quad (10)$$

مرتبه محاسبات در این روش  $T.N^2$  است که برای مثال قبل ( $N = 6, T = 40$ ) برابر ۱۴۴۰ می‌شود. در نتیجه محاسبه  $P(O|\lambda)$  از این روش توسط کامپیوتر عملی خواهد بود.

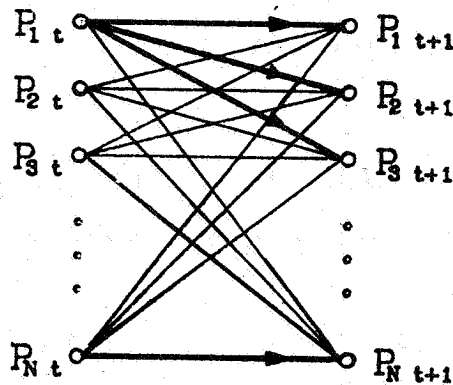
مسئله دوم: در مرحله یادگیری لازم است محتمل‌ترین مسیر حالت، با داشتن مشاهده  $O$  و مدل  $\lambda$  محاسبه شود. یعنی سلسله  $Q$  را به نحوی انتخاب کنیم که  $P(Q|O, \lambda)$  ماگزیمم شود. این کار معادل است با محاسبه

$$Q = \text{Arg Max}_{q_1, q_2, \dots, q_T} P[Q, O|\lambda] \quad (11)$$

بر اساس اصل بهینگی<sup>۱</sup> می‌دانیم اگر مسیر منتهی به هر حالت در زمان  $t$  بهترین مسیر ممکن (مسیر بهینه) باشد، با انتخاب یک سیاست بهینه در فاصله زمانی  $(t+1, T)$  در نهایت به یک سری جواب بهینه برای مسیرهای منتهی به هر حالت در زمان  $T$  خواهیم رسید [۱۳]. بهترین مسیر، مسیری است که تابع ارزش مربوط به آن حداکثر باشد. در این مسئله خاص تابع ارزش، احتمال هر مسیر بوده و هدف ماگزیمم کردن آن است. روش ویتربی<sup>۲</sup>، روشی تکراری برای یافتن مسیر بهینه است [۱۳]. شکل ۱ این روش را به طور خلاصه نشان می‌دهد.

مسیر انتخاب شده در گذار از زمان  $t$  به  $t+1$  برای هر حالت  $i$ ، بر اساس تابع ارزش در زمان  $t$  از رابطه زیر محاسبه می‌شود:

$$P_i(t+1) = \text{Max}_j (P_j(t) \cdot a_{ji}) \cdot b_i(o_{t+1}) \quad (12)$$



شکل ۱ - نمایش انتخاب مسیرهای بهینه با روش ویتربی، تا لحظه  $t+1$  با استفاده از مسیرهای بهینه در لحظه  $t$

به غیر از مسیری که رابطه بالا را برای هر حالت ماگزیمم می کند، بقیه مسیرها حذف می شوند. (مسیرهایی که در شکل با خطوط نازکتر رسم شده، حذف گردیده اند.) بنابراین در هر گذار زمان فقط  $N$  مسیر حفظ خواهد شد، در خاتمه در زمان  $T$ ،  $N$  مسیر بهینه که هر کدام به یکی از حالت های  $S_1, \dots, S_N$  ختم شده اند باقی می مانند و مسیر بهینه مسیری است که احتمال آن بیشترین مقدار را داشته باشد.

مسأله سوم: در مرحله یادگیری لازم است مدلی که بیشترین شباهت را به مشاهده  $O$  دارد، تخمین زده شود. برای حل این مسأله روشهای متفاوتی وجود، که هیچکدام رسیدن به جواب بهینه را تضمین نمی کند. روشی که برای طراحی سیستم حاضر در نظر گرفته شده عبارت است از تخمین درست‌نمایی بیشینه<sup>۱</sup> مدل  $\lambda$  به فرض آنکه مسیر  $Q$  طوری انتخاب شده باشد که  $P(O, Q | \lambda)$  حداکثر شود ([۱۰ تا ۱۲]):

$$\lambda = \underset{\lambda}{\text{ArgMax}} \{ \underset{Q}{\text{Max}} P[Q, O | \lambda] \} \quad (۱۳)$$



جهت حل این مسأله ابتدا از یک روش مناسب حدس اولیه ای برای مدل  $\lambda$  زده می شود، سپس با توجه به سلسله های مشاهدات با یک روش تکراری تخمین بهتری ( $\bar{\lambda}$ ) از روی مقدار قبلی  $\lambda$ ، به دست می آید [۱۴]. این عمل تا رسیدن به  $\lambda$  ی بهینه تکرار می شود. در به دست آوردن روابط تکراری تخمین  $\lambda$ ، مشاهده می شود که تخمین ضرائب ماتریس B مستقل از تخمین ضرائب ماتریس A است. این امر باعث می شود که روش به دست آمده تا اندازه ای ساده شود. به طور خلاصه این روش که "روش k مقدار متوسط برای هر قسمت" نامیده می شود، به صورت زیر است [۱۴]:

۱- به دست آوردن مسیر حالت  $Q^*$  برای مشاهده O با ماگزیمم کردن تابع  $P(O, Q|\lambda)$  (حل مساله دوم).

۲- بدست آوردن مجموعه های  $T_{ij}; 1 \leq i, j \leq N$ :

$$T_{ij} = \{t : q_{t-1}^* = i, q_t^* = j\} \quad (14)$$

۳- به دست آوردن  $a_{ij}$  از روی مجموعه های  $T_{ij}$  به صورت زیر:

$$\bar{a}_{ij} = \frac{\|T_{ij}\|}{\sum_{k=1}^N \|T_{ik}\|}; 1 \leq i, j \leq N \quad (15)$$

که در آن  $\|\cdot\|$  برای نمایش تعداد اعضای یک مجموعه بکار می رود.

تعبیر رابطه بالا بسیار ساده است، در واقع تعداد دفعات گذار از حالت i به j محاسبه شده و به تعداد دفعات حضور در حالت i تقسیم می شود.

۴- مجموعه های  $O_i; i = 1, \dots, N$  به صورت زیر بدست می آید:

$$O_i = \{o = o_t ; q_t = S_i\} \quad (16)$$

مجموعه  $O_i$  شامل کلیه مشاهداتی است که در حالت  $S_i$  اتفاق افتاده اند.  
 ۵- تخمین ماتریس  $B$  از روی مجموعه های  $O_i$  به دست می آید. اگر مشاهدات به صورت گسسته باشند،  $b_{jk}$  معادل است با تکرار  $V_k$  در مجموعه  $O_i$ ، تقسیم بر کل تعداد اعضای مجموعه  $O_i$ .  
 تاکنون فرض بر این بود که مشاهدات، کمیات گسسته ای هستند. در حالی که می توان مشاهدات را به صورت بردارهایی پیوسته نیز در نظر گرفت. در این صورت سلسله مشاهدات  $O = o_1 o_2 \dots o_T$  است که  $o_i$  یک بردار پیوسته است. در نتیجه به جای تخمین تابع توزیع احتمال گسسته، (ماتریس  $B$ ) باید تابع چگالی احتمال پیوسته بردار  $o$  برای هر حالت تخمین زده شود. به منظور تسهیل مسأله تخمین تابع چگالی احتمال برای هر حالت، مجموع  $M$  تابع توزیع نرمال فرض می شود.

$$b_j(o) = \sum_{m=1}^M c_{jm} \cdot N(o, \mu_{jm}, \sigma_{jm}) \quad (17)$$

در رابطه ۱۷،  $M$  تعداد خوشه ها (ترکیب ها) و  $c_{jm}$  ضریب تأثیر خوشه  $m$  در حالت  $j$  است که به صورت  $P(m|j)$  احتمال شرطی خوشه  $m$  در حالت  $j$  نیز در نظر گرفته می شود.  $\mu_{jm}$  بردار متوسط تابع توزیع خوشه  $m$  در حالت  $j$ ،  $\sigma_{jm}$  ماتریس کواریانس خوشه  $m$  در حالت  $j$ ، و  $N$  نمایانگر تابع توزیع نرمال برداری است. تابع توزیع پیوسته  $b_j(o)$  بدین علت به صورت ترکیبی از چند تابع توزیع نرمال در نظر گرفته شده است که معمولاً توابع توزیع واقعی برای مشاهده  $o$  پیچیده تر از آن است که با یک تابع توزیع نرمال قابل توصیف باشد. برای نمونه ممکن است چندین قله در شکل تابع توزیع واقعی وجود داشته باشد که بطور وضوح با یک تابع توزیع نرمال قابل تقریب زدن نخواهد بود.

برای برقراری رابطه  $\int_{-\infty}^{\infty} b_j(\Omega) d\Omega = 1$  باید شرط  $\sum_{m=1}^M c_{jm} \geq 0$  و  $\sum_{m=1}^M c_{jm} = 1$  برقرار باشد.

تخمین زدن  $\mu_{jk}^{\tau+1}$  و  $\sigma_{jk}^{\tau+1}$  با فرض موجود بودن این مقادیر در زمان  $\tau$  و اینکه رابطه

$$b_j(\Omega) = \sum_{K=1}^M c_{jk}^{\tau+1} N(\Omega, \mu_{jk}^{\tau+1}, \sigma_{jk}^{\tau+1}) \quad 1 \leq j \leq 6 \quad (18)$$

تخمینی از توابع چگالی احتمال اعضای مجموعه  $A_j$  است، صورت می‌گیرد. مجموعه  $A_j$  به صورت زیر به دست می‌آید [۱۴].

$$A_j = \{O_{tk} : q_{tk} = S_j, 1 \leq k \leq L\} \quad (19)$$

$L$  کل تکرار کلمه مورد نظر است.

در روابط بالا نشاندهنده حالت و  $M$  تعداد توابع توزیع نرمال به ازای هر حالت است. روش تخمین مقادیر جدید  $\mu_{jk}^{\tau+1}$  و  $\sigma_{jk}^{\tau+1}$  از روی مجموعه های  $A_j, j = 1, \dots, 6$  مشابه به تخمین در روش  $k$  مقدار متوسط در حالت مشاهدات گسسته است. در این روش ابتدا احتمال شرطی هر مشاهده  $O_t \in A_j$ ، نسبت به  $M$  خوشه محاسبه می‌شود، بیشترین احتمال معین می‌کند که  $O_t$  به کدام خوشه مربوط می‌شود. از این طریق اعضای مجموعه  $A_j$  نیز به  $M$  دسته (در اینجا  $M = 5$ ) تقسیم می‌شوند. با استفاده از این تقسیم بندی تخمین پارامترهای مورد نظر به صورت زیر به دست می‌آید:

$$c_{jk}^{\tau+1} = \frac{\text{تعداد اعضای خوشه } k \text{ در تکرار } \tau}{\text{کل تعداد اعضای مجموعه } A_j \text{ در تکرار } \tau} \quad (20)$$

$$\mu_{jk}^{\tau+1} = \text{متوسط داده های موجود در خوشه } k \text{ در حالت } j \text{ در تکرار } \tau \quad (21)$$



## استخراج مشخصه های سیگنال صحبت

به منظور استخراج مشخصه های سیگنال صحبت ابتدا باید محدوده کلمات ادا شده، از محدوده سکوت جدا شود. عمل تفکیک صحبت از سکوت با استفاده از مقدار آستانه انرژی و پاره ای فرضیات انجام می گیرد. به این ترتیب که ابتدا در حالت سکوت حداکثر انرژی محاسبه شده و مقادیر آستانه از روی آن محاسبه می گردد. آنگاه یک ماشین حالت چهار وضعیتیتی ابتدا و انتهای هر کلمه را تشخیص می دهد ([۱۵ و ۱۶]).

پس از تشخیص ابتدا و انتهای کلمه، مشخصه های سیگنال صحبت در طی زمان ادای کلمه محاسبه می شوند. یک دسته از مشخصات مفید، ضرایب کپسترال سیگنال صحبت هستند ([۱۷ تا ۱۹]). انرژی، ضرایب کپسترال و ضرایب دلتا کپسترال پس از انجام پیش پردازش بر روی سیگنال محاسبه می شوند (شکل ۳). دلیل استفاده از ضرائب دلتا کپسترال، وجود اطلاعات مربوط به تغییرات بین قطعه های مجاور در این دسته ضرائب است [۱۷].

جزئیات مراحل پیش پردازش و پردازش به قرار زیر است:

\* فیلتر آنالوگ پایین گذر

\* مبدل A/D ۱۲ بیتی، با نرخ نمونه برداری ۸۰۰۰ هرتز

\* پیش تاکید<sup>۱</sup> با یک صفر در  $z = 0.9$ ,  $H(z) = 1 - 0.9z^{-1}$

\* تقسیم بندی به قطعه<sup>۲</sup> های ۳۷/۵ میلی ثانیه ای (این قطعه ها هر بار ۱۲/۵ میلی ثانیه به جلو برده می شوند. به عبارت دیگر هر قطعه شامل ۳۰۰ نمونه و میزان جابجایی ۱۰۰ نمونه است و هر دو قطعه مجاور در ۲۰۰ نمونه مشترک هستند).

\* پنجره همینگ<sup>۳</sup> [۲۰]

\* آنالیز LPC<sup>۴</sup> با مرتبه هشت که از روی آن پارامترهای فیلتر تمام قطب محاسبه می شوند [۲۱].

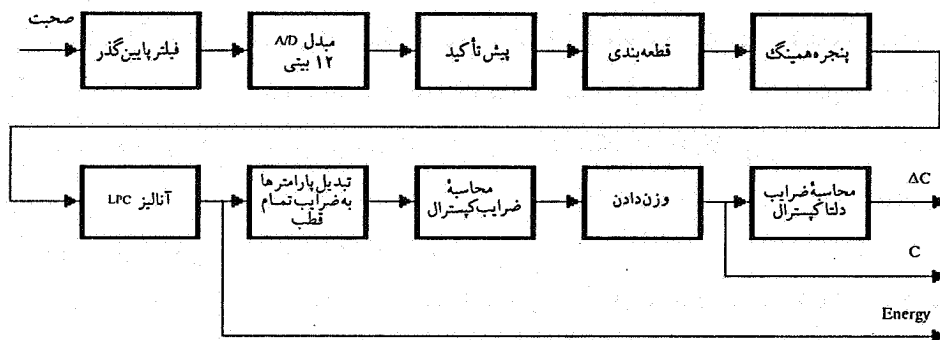
\* محاسبه دوازده پارامتر اول کپسترال از روی پارامترهای فیلتر تمام قطب [۲۲].

1. Preemphasis

2. Frame

3. Hamming Window

4. Linear Predictive Coding



شکل ۳ - دیاگرام بلوکی استخراج مشخصه های صحبت

\* وزن دادن<sup>۱</sup> پارامترهای کپسترال محاسبه شده توسط رابطه زیر [۱۹]:

$$\hat{c}_1(m) = [1 + Q \sin(\pi \cdot m / Q) / 2] \cdot c_1(m) ; 1 \leq m \leq Q \quad Q = 12 \quad (24)$$

(در این رابطه  $c_1(m)$  ضریب کپسترال  $m$  در زمان  $I$  است.)

\* محاسبه ضرایب دلتا کپسترال<sup>۲</sup> که توسط رگرسیون خطی مرتبه دو از روی پارامترهای کپسترال قطعه های مجاور آن به صورت زیر محاسبه می شوند [۱۷]:

$$\Delta \hat{c}_1(m) = \left[ \sum_{k=-2}^2 k \hat{c}_1(m-k) \right] \times 0.375 ; 1 \leq m \leq Q \quad Q = 12 \quad (25)$$

(ضریب  $0.375$  برای متعادل نمودن واریانس این ضرایب با ضرایب کپسترال بکار می رود).  
بردار مشاهده مربوط به هر قطعه شامل ضرایب کپسترال، ضرایب دلتا کپسترال و لگاریتم انرژی است.

## تحقق مدل پنهانی مارکف

همانطور که قبلاً توضیح داده شد، یادگیری به "روش  $k$  مقدار متوسط برای هر قسمت" که یک روش جستجوی محلی<sup>۱</sup> است انجام می‌گیرد. با توجه به این موضوع که مسأله دارای چندین ماگزیمم محلی است، الگوریتم به ماگزیممی همگرا می‌شود که شرایط اولیه در بستر جذب آن قرار داشته‌اند. در نتیجه انتخاب شرایط اولیه در کارآیی نهایی سیستم اثر زیادی خواهد داشت. البته مسلم است که با توجه به روش بکار رفته هیچگاه نمی‌توان مطمئن بود که نقطه بهتری برای همگرا شدن وجود نداشته است. از میان پارامترهای مختلف مدل  $\lambda$  اهمیت بردارهای متوسط تابع چگالی نرمال،  $\mu, \Sigma$ ، بیشتر است. در مرحله اول مقادیر اولیه پارامترهای  $\mu, \Sigma$  به صورت تصادفی انتخاب شدند. در این حال به علت وجود خطای ته ریز<sup>۲</sup> اصولاً الگوریتم با مشکل روبرو شد و به جوابی نرسید. دلیل این امر مکان نامناسب نقاط آغازین است. توجه کنید  $\mu, \Sigma$  یک بردار ۲۵ بعدی است. انتخاب تصادفی مقدار آغازین  $\mu, \Sigma$  در واقع انتخاب تصادفی یک نقطه در فضای ۲۵ بعدی است. لذا امکان آن که هیچ یک از پنج بردار متوسط مربوط به پنج خوشه در نزدیکی بردارهای مشاهده نباشند زیاد است. از طرف دیگر تابع چگالی احتمال مربوط به هر خوشه بصورت نرمال است. به دلیل وجود فرم نمائی در تابع چگالی نرمال با زیاد شدن فاصله بردار متوسط از بردار مشاهدات احتمال بردار مشاهده به شدت کاهش می‌یابد و در عمل خطای ته ریز پیش می‌آید. برای اجتناب از این امر باید در نزدیکی هر یک از بردارهای مشاهده حداقل یکی از بردارهای متوسط قرار داشته باشد.

برای اینکه جواب نهایی روش محقق شده ماگزیمم نسبتاً بزرگی باشد و همچنین برای اینکه مشکل ته ریز باعث بروز خطا نشود، روشی برای تعیین مقادیر اولیه  $\mu, \Sigma$  ابداع شده است که در ادامه به توضیح آن خواهیم پرداخت ([۱۵] و [۲۳]).

ایده اصلی این روش این است که هر حالت را می‌توان به یک واج (واجگونه) درون کلمه مورد نظر نسبت داد. این واجها اغلب پشت سر هم بیان می‌شوند (به غیر از مواردی که واجها خورده می‌شوند و بیان نمی‌شوند). بردارهای مشاهده مربوط به هر واج در یک محدوده کوچک در فضای

۲۵ بعدی قرار دارند. اگر در این محدوده های نسبتاً کوچک بردار متوسط محاسبه شود، شانس نسبت داده شدن بردارهای مشاهده هرواج به ناحیه مربوطه زیاد شده و علاوه بر اینکه این موضوع به درستی عمل تشخیص کمک می کند، جلوی کوچک شدن بیش از اندازه احتمالات محاسبه شده را نیز می گیرد. با تقسیم بندی کلمات ادا شده به تکه های مجزا، که هر تکه یک واج (واجگونه) را شامل می شود، می توان پایگاه داده ها را به  $N$  دسته ( $N$  تعداد حالات است) تقسیم کرد که هر دسته به یک حالت مربوط می شود. این تقسیم بندی می تواند خیلی دقیق و یا با تقریب (با در نظر گرفتن  $N$  قسمت برابر در هر کلمه) صورت پذیرد.

برای محاسبه مقادیر  $\mu$  از روی این  $N$  دسته در حالت مشاهده پیوسته کافی است با استفاده از یک روش چندی کردن برداری [۲۴] با  $M$  بردار نمونه، داده های هر دسته را به  $M$  خوشه تقسیم کرده و متوسط هر خوشه را به عنوان مقادیر اولیه  $\mu$  انتخاب نمود. با این روش این تضمین بوجود می آید که بردارهای متوسط  $\mu_{jk}$   $j=1, \dots, N$  و  $k=1, \dots, M$  در محدوده بردارهای مشاهده قرار بگیرند. همچنین به کمک این روش مرحله یادگیری در مدل پنهانی مارکف به صورت کاملاً خودکار انجام می گیرد و عوض کردن مجموعه کلمات و مجموعه صداها وقتی از شخص نخواهد گرفت. جهت انجام چندی کردن برداری از روش LBG<sup>۱</sup> استفاده شد.

از مسائل عملی که در تحقق مدل مورد نظر پیش می آید، خطای ته ریز است. برای مثال در حل مسأله اول به روش بازگشتی مقادیر  $p(\Omega_1, \dots, \Omega_t, q_t = S_i | \lambda)$  با زیاد شدن زمان بسیار کوچک می شوند، به طوری که از حدود دقت کامپیوتر خارج می گردند. جهت اجتناب از این مشکل، در هر مرحله  $p(\Omega_1, \dots, \Omega_t, q_t = S_i | \lambda)$  را در عدد مناسبی ضرب می کنیم به طوری که حاصل در محدوده دقت کامپیوتر باقی بماند. اگر عدد نرمالیزاسیون ( $c_t$ ) به صورت زیر انتخاب شود:

$$c_t = \frac{1}{\sum_{i=1}^N P(\Omega_1, \dots, \Omega_t, q_t = S_i | \lambda)} ; t = 1, \dots, T \quad (26)$$



$p(O|\lambda)$  نیز نهایتاً توسط رابطه زیر بدست می آید:

$$P(O|\lambda) = \frac{1}{\prod_{t=1}^T c_t} \quad (27)$$

و یا

$$\text{Log } P(O|\lambda) = -\sum_{t=1}^T \text{Log } c_t \quad (28)$$

مشکل ته ریز در محاسبه تابع چگالی احتمال پیوسته  $(b_j(o_t))$  نیز وجود دارد. برای حل این مشکل ابتدا در هر حالت احتمال خوشه ای که حداکثر احتمال را دارد در نظر گرفته می شود. سپس حداقل این مقادیر در بین حالات مختلف، بدست آمده و تمام احتمالات محاسبه شده به این مقدار نرمالیزه می شود. در این صورت می توان امیدوار بود که مقادیر به دست آمده در محدوده دقت کامپیوتر باقی بمانند. به سادگی می توان نشان داد که با این نرمالیزاسیون روش حل مساله دوم (یافتن مسیر حالت بهینه) تغییری نمی کند. حل مساله اول (محاسبه  $p(O|\lambda)$ ) نیز با انجام عمل عکس نرمالیزاسیون<sup>۱</sup> تصحیح می شود.

### روش آزمایش

جهت تخمین پارامترهای مدلها و آزمایش سیستم محقق شده از سه مجموعه سیگنال صحبت به نامهای UV، YW، ZX استفاده شد. هر کدام از این مجموعه ها شامل صدای ۲۵ زن و ۲۵ مرد است که هر یک از افراد اعداد صفر تا ده را دوبار بیان کرده اند. ضبط صدا در محیط آزمایشگاه و توسط میکروفون و تقویت کننده صورت گرفته است. عمل ضبط به صورت ۱۲ بیتی بر روی کامپیتر مجهز به A/D صورت پذیرفته است. از دو مجموعه ZX و YW برای تخمین پارامترهای مدل و از مجموعه UV برای آزمون کارایی سیستم استفاده شد. نتایج این آزمایش در جدول ۱ آمده است. همان گونه که در این جدول مشاهده می شود، میزان خطای کل سیستم ۳/۰۰٪ برآورد شده است.

1 . Denormalization

لازم به تذکر است که در آزمایش صدای UV نیز عمل شناسایی و ضبط صدا در دو مرحله مختلف صورت گرفته است.

استفاده از مدل پنهانی مارکف به صورت چپ به راست با این فرض همراه است که هر یک از حالات مدل با یکی از ساختارهای صحبت موجود در کلمه (واج) متناظر است. البته این تناظر در کلیه موارد به طور دقیق به چشم نمی خورد. به هر حال شکل ۴ تقسیم بندی یک سیگنال نمونه به حالت‌های تشکیل دهنده اش توسط این مدل را نشان می دهد. همان طوری که در شکل ۴ ملاحظه می شود هر یک از حالتها به یکی از واج ها (یا واجگونه ها) نسبت پیدا می کند.

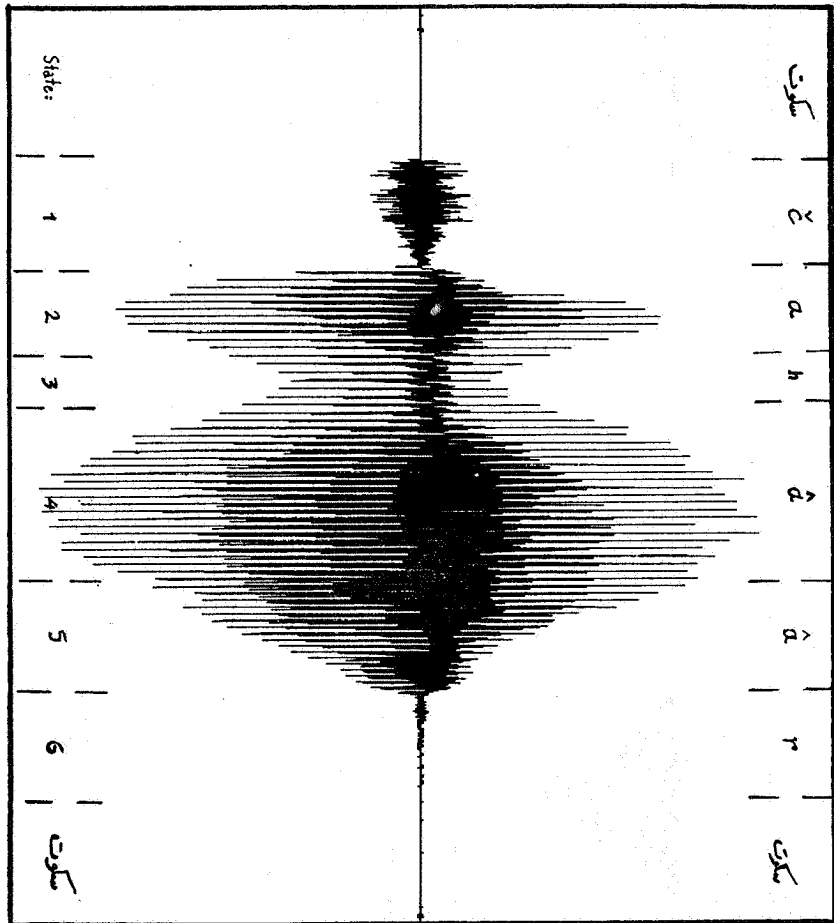
یکی از مسائل مهم در مدل پنهانی مارکف پیوسته تخمین توابع چگالی احتمال مربوط به هر حالت است. همان طوری که گفته شد برای هر حالت تابع چگالی احتمال به صورت مجموع پنج تابع چگالی نرمال در نظر گرفته شده است. استفاده از چند تابع نرمال می تواند به تخمین توابع پیچیده تر کمک کند. در شکل ۵ نحوه تخمین چند نمودارستونی<sup>۱</sup> توسط این روش نشان داده شده است. برای نمونه در شکل ۴-۵ وجود دو ماگزیم در نمودار ستونی نیز به خوبی تخمین زده شده است.

نتیجه‌ای که در پایان از آزمایش مدل محقق شده روی مجموعه های S1 (داده‌های یادگیری) و S2 (داده‌های آزمایش) گرفته شد، در جدول ۱ آمده است:

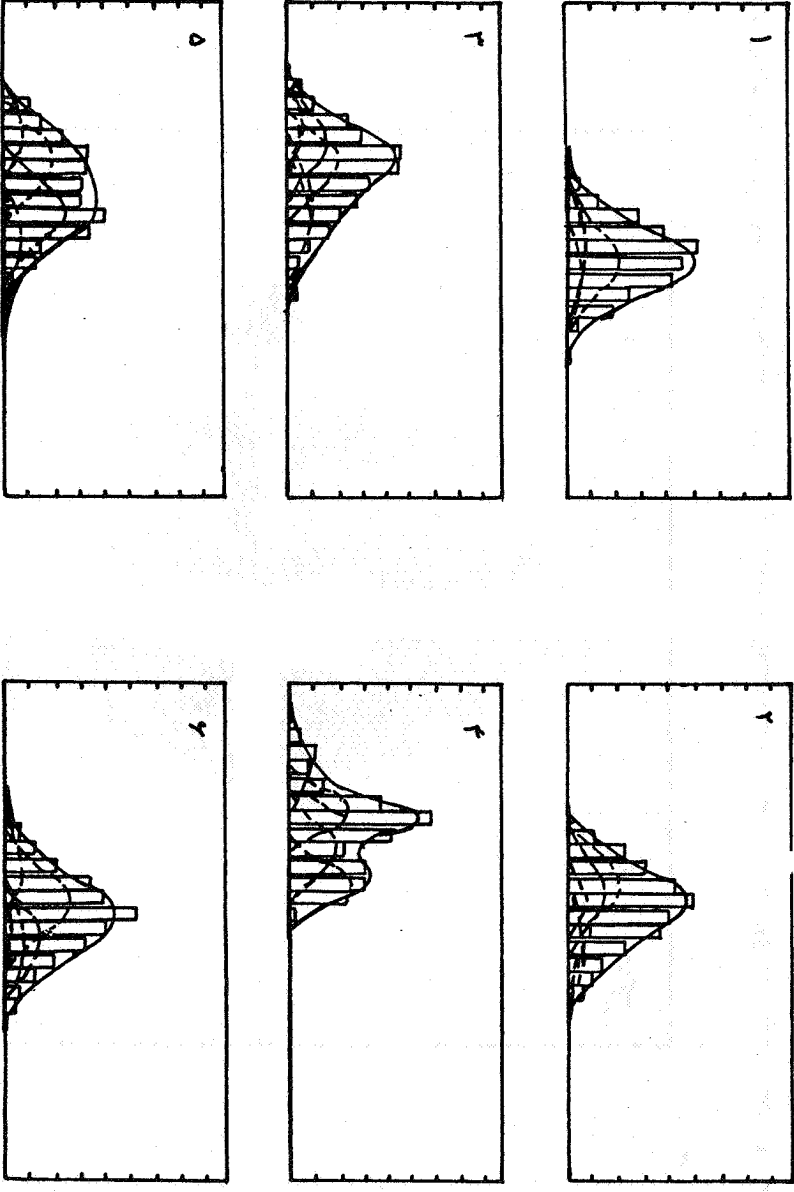
همان طور که مشاهده می شود بیشترین خطا مربوط به کلمات ("sefr/صفر" و "se/سه") است و این به خاطر شباهت این دو کلمه به یکدیگر است. خصوصاً که بسیاری از افراد واج آخر کلمه "صفر" را به درستی بیان نمی کنند. همین طور دو کلمه ("do/دو" و "noh/نه") نیز گاهی با هم اشتباه می شوند و کلمات ("haft/هفت" با "hašt/هشت") نیز دیگر موارد عمده خطا را تشکیل می دهند.

### نتیجه گیری

در این مقاله تحقق یک سیستم درک صحبت به روش مدل پنهانی مارکف ارائه شد. برای انتخاب مقادیر اولیه از چندی کردن برداری داده های یادگیری تقسیم شده بین حالتها در هر کلمه استفاده



شکل ۴ - نحوه تقسیم بندی کلام به حالت های مختلف (متناظر با واحدها) توسط HMM؛ نمونه مورد استفاده کلمه چهار می باشد.



شکل ۵ - هیستوگرام و تابع چگالی تخمینی برای یک متغیر (CS) در کلمه "چهار"، در حالت‌های ششگانه. هیستوگرام نمایش دهنده توزیع تجربی، منحنی پررنگ تخمین تابع چگالی احتمال و منحنی های کم رنگ حاصل ضرب تابع چگالی احتمال هر خوشه در احتمال آن خوشه می باشند.

جدول ۱- میزان خطا برای کلمات و مجموعه های مختلف

مجموعه کلمه	ZX آموزش			YW آموزش			UV آزمون		
	زن	مرد	زن و مرد	زن	مرد	زن و مرد	زن	مرد	زن و مرد
۰	۳/۵۰	۶/۵۰	۹/۱۰۰	۴/۵۰	۲/۵۰	۶/۱۰۰	۰	۱/۵۰	۱/۱۰۰
۱	۰	۰	۰	۰	۰	۰	۰	۰	۰
۲	۰	۰	۰	۰	۰	۰	۱/۵۰	۶/۵۰	۷/۱۰۰
۳	۲/۵۰	۰	۲/۱۰۰	۲/۵۰	۱/۵۰	۳/۱۰۰	۸/۵۰	۴/۵۰	۱۲/۱۰۰
۴	۰	۰	۰	۰	۰	۰	۰	۰	۰
۵	۰	۰	۰	۰	۰	۰	۰	۰	۰
۶	۰	۰	۰	۰	۰	۰	۱/۵۰	۰	۱/۱۰۰
۷	۱/۵۰	۰	۱/۱۰۰	۰	۱/۵۰	۱/۱۰۰	۰	۳/۵۰	۳/۱۰۰
۸	۰	۰	۰	۰	۰	۰	۰	۰	۰
۹	۰	۰	۰	۲/۵۰	۲/۵۰	۴/۱۰۰	۰	۵/۵۰	۵/۱۰۰
۱۰	۰	۰	۰	۰	۰	۰	۲/۵۰	۲/۵۰	۴/۱۰۰
۰-۱۰	٪۱/۰۹	٪۱/۰۹	٪۱/۰۹	٪۱/۴۵	٪۱/۰۹	٪۱/۲۷	٪۲/۱۸	٪۳/۸۱	٪۳/۰

مقادیر اولیه از چندی کردن برداری داده های یادگیری تقسیم شده بین حالتها در هر کلمه استفاده گردید. با آزمایش سیستم محقق شده خطای سیستم ۳/۰۰ درصد برآورد شد. مقایسه کارایی سیستم با سایر سیستمهای درک صحبت به دلیل عدم دسترسی به اطلاعات سیستم های مشابه برای زبان فارسی میسر نشد. مقایسه این سیستم با سیستم هایی که برای زبانهای بیگانه طراحی شده است نیز خالی از اشکال نیست. با اینحال اگر با تسامح بخواهیم مقایسه ای انجام بدهیم، دیده می شود که در سیستمی که مشابه سیستم حاضر در زبان انگلیسی می باشد و برای تشخیص اعداد "صفر" تا "نه" به کار می رود [۱۰] و تنها مرحله مقدار دهی اولیه پارامترهای آن با سیستم حاضر متفاوت است (در واقع در مورد روش مقدار دهی اولیه آن سیستم اطلاعی در دست نیست)، مشاهده می شود که

کارآیی سیستم فارسی حدود یک درصد کمتر است، همان طور که مشاهده می شود در مجموعه کلمات فارسی زوج های ("دو"، "نه") و ("صفر"، "سه") همچنین کلمات ("هفت"، "هشت"، "پنج") شباهت زیادی دارند. چنین شباهتی برای ارقام "صفر" تا "نه" به زبان انگلیسی فقط در زوج های (Five, Nine) و (Six, Three) دیده می شود که درجه شباهت نیز از موارد مشابه در زبان فارسی کمتر است. چنانچه اگر یکی از موارد شباهت در فارسی مثلاً "صفر" به کلی حذف شود، میزان خطا به  $1/8$  درصد می رسد که در حدود درصد خطای سیستم انگلیسی ( $1/3$  و  $1/8$  درصد برای دو مجموعه آزمایشی مختلف) خواهد بود.

#### قدردانی

نویسندگان مقاله لازم می دانند که از مسئولین محترم مرکز تحقیقات الکترونیک دانشگاه صنعتی شریف بویژه آقایان دکتر محمود تیبانی و مهندس همایون برهانی و همچنین از معاونت پژوهشی دانشگاه صنعتی شریف و مسئولین محترم امور پژوهشی دانشگاه تهران به جهت در اختیار گذاردن شرایط و امکانات انجام این پروژه تشکر و قدردانی به عمل آورند.

### مراجع

- ۱- ی. ثمره، آواشناسی زبان فارسی، آواها و ساخت آوایی هجا، مرکز نشر دانشگاهی، ۱۳۶۴.
2. Levinson, E. and Roe, D.B. , "A Perspective on Speech Recognition", *IEEE Comm. Mag.*, Vol. 28, No. 1, Jan. 1990.
3. Slutsker, G., " Non - Linear Method of Analysis of Speech Signal ", *Trudy N.I.I.R.*, 1968.
4. Venlichko, V. and Zagoruyko, N., "Automatic Recognition of 200 words", *Int. J. Man - Machine Studies*, 2, pp. 223-234, 1970.
5. Vintsyuk, T.K., "Speech Recognition by Dynamic Programming", *Kybernetika*, 4, 1, pp. 81-88, 1968.
۶. و. طباطبایی، ب. عظیمی سجادی، ب. ظهیر اعظمی و ک. لوکس، بررسی و مقایسه دو روش تشخیص کلمات منفرد فارسی، کارنامه اولین کنفرانس آموزش، پژوهش و کاربرد کامپیوتر در ایران، تهران، بهمن ۱۳۷۱.
7. Waibel, A. et al, "Phoneme Recognition Using Time Delay Neural Networks", *IEEE Trans. on ASSP*, Vol. ASSP-37, No. 3, March 1989.
8. Avarbuch, A. et al, "An IBM - PC Based Large - Vocabulary Isolated Utterance Speech Recognizer", in *Proceedings IEEE ICASSP-86*, Tokyo, Japan, April 1986.
9. Lec, K.F., *Large - Vocabulary Speaker Independent Continuous Speech Recognition: The SPHINX System*, Ph.D. Dessertation, Computer Science Department, Carnegie Mellon University, 1988.
10. Rabiner, L.R., "A Tutorial on Hidden Markov Models and Selected

- Applications in Speech Recognition" ,*IEEE Proc.*, Vol. 77, No. 2, Feb. 1989.
11. Baum, L.E. and Petrie, T., "Statistical Inference for Probabilistic Functions of Finite State Markov Chains", *Ann. Math. Stat.*, Vol.37,1966.
  12. Baum, L.E., Petrie, T. , Soules, G. and Weiss, N., "A Maximization Technique Occurring in the Statistical Analysis of Probabilistic Functions of Markov Chains", *Ann. Math. Stat.* ,Vol. 41, No. 1, 1970.
  13. Forney, G.D , "The Viterbi Algorithm" , *Proc. of the IEEE*, Vol. 61, No.3 ,Mar. 1973.
  14. Juang, B.H. and Rabiner, L.R., "The Segmented K - Means Algorithm for Estimating Parameters of Hidden Markov Models", *IEEE Trans. on ASSP*, Vol. ASSP-38, No. 9, Sep. 1990.
  ۱۵. و. طباطبایی ، استخراج مشخصه ها و پیش پردازش سیگنال صحبت به منظور درک کلمات ، پایان نامه کارشناسی ارشد، گروه برق دانشکده فنی دانشگاه تهران، ۱۳۷۱.
  16. Rabiner, L.R.and Sambur, M.R., "An Algorithm for Determining the End Point of Isolated Utterances", *B.S.T.J.*, Vol 54, No. 2, Feb. 1975.
  17. Furui, S., "Speaker Independent Isolated Word Recognition Using Dynamic Features of Speech Spectrum", *IEEE Trans. on ASSP*, Vol. ASSP-34, No.1, Feb. 1986.
  18. Furui, S., *Digital Speech Processing, Synthesis and Recognition*, Copyright 1989 by Marcel Decker, Inc.
  19. Juang, B.H., Rabiner, L.R. and Wilpon, J.G., "On the Use of Bandpass Liftering in Speech Recognition", *IEEE Trans. on ASSP*, Vol. ASSP-35, No. 7, July 1987.
  20. Oppenheim, A.V. and Schafer, R.W., *Discrete - Time Signal Processing* , Prentice Hall, INC. Englewood Cliffs, NJ.,1989.



21. Parsons, T.W., *Voice and Speech Processing*, McGraw-Hill, 1986.
22. Viswanathan, R. and Makhoul, J., "Quantization Properties of Transmission Parameters in Linear Predictive Systems", *IEEE Trans. on ASSP*, Vol. ASSP-23, June 1975.
۲۳. ب. عظیمی سجادی، درک صحبت مستقل از گوینده برای کلمات منفرد فارسی با تعداد محدود، پایان نامه کارشناسی ارشد، گروه برق دانشکده فنی دانشگاه تهران، ۱۳۷۱.
24. Makhoul, J., Roucos S. and Gish, H., "Vector Quantization in Speech Coding", *IEEE Proc.*, Vol. 73, No. 11, Nov. 1985.