





دانشکده: برق و رباتیک

گروه: الکترونیک

پایان نامه دوره کارشناسی ارشد مهندسی برق - الکترونیک

کاربرد مدل آماری مخلوط لاپلاس در بهسازی گفتار

زینب محمدپوری

اساتید راهنما:

جناب آقای دکتر حسین مروی

جناب آقای دکتر امید رضا معروضی

استاد مشاور:

جناب آقای دکتر علی سلیمانی

زمستان ۱۳۸۹

تقدیم به

پدر و مادرم مهربانم

که هرچه دارم از وجود پرمهر آنهاست

و همسرم

که همواره حامی و مشوق من بودند.

تشکر و قدردانی

خداوند متعال را به جهت الطاف بی پایان و خاصه اش که هموار شامل حال من کرده بسیار شاکرم.

پس از آن بر خود لازم می دانم از زحمات و راهنمایی های اساتید گرانقدرم جناب آقای دکتر مروی، دکتر معروضی و دکتر سلیمانی تشکر نمایم. در ضمن از جناب آقای هادی ویسی به جهت راهنمایی ها و کمک هایشان و جناب آقای امیرحسین آفریدون هم به جهت زحماتشان برای این پایان نامه، کمال تشکر را دارم.

چکیده

بهسازی گفتار یا Speech enhancement بیانگر گروه بزرگی از روش‌هاست که با انجام پردازش‌هایی روی سیگنال‌های نویزی، نهایتاً منجر به بهبود کیفیت و قابلیت فهم گفتار می‌شود. از مهم‌ترین روش‌های بهسازی گفتار، روش‌های آماری هستند که به لحاظ کارایی بالاتر نسبت به سایر روش‌ها و ایجاد اعوجاج کمتر در سیگنال‌های نهایی بیشتر مورد توجه می‌باشد. در این روش یک مدل آماری برای گفتار و نویز فرض می‌شود و پردازش‌ها بر مبنای این مدل صورت می‌گیرد. این مدل آماری برای ضرایب فوریه یا پارامترهای دیگر هر قاب سیگنال گفتار یا نویز فرض می‌شود. در این پایان‌نامه، روش آماری جدیدی برای تخمین سیگنال تمیز از روی سیگنال نویزی در حضور نویز جمع‌شونده و مستقل از سیگنال تمیز، در حوزه فرکانس ارائه شده است. تا کنون از مدل‌های گوسی، لاپلاس، گاما و مخلوط گوسی، برای مدل کردن طیف سیگنال گفتار استفاده شده؛ اما در این پایان‌نامه، توزیع مخلوط لاپلاس برای مدل کردن بخش‌های حقیقی و موهومی طیف گفتار پیشنهاد شده است. پارامترهای این مدل (میانگین‌ها، واریانس‌ها و ضرایب) به صورت برون‌خطی و با استفاده از الگوریتم EM و سیگنال صحبت بدون نویز از پایگاه داده TIMIT محاسبه شده است. سپس رابطه تخمین گر MMSE یا کمترین میانگین مربعات خطا، با توجه به توزیع مخلوط لاپلاس به دست آورده شده است. واریانس نویز به صورت درون‌خطی و با استفاده از روش ردیابی کمینه‌ها به دست آمده است.

در نهایت برای ارزیابی و عملکرد روش پیشنهادی از سه معیار سیگنال به نویز قطعه‌ای، PESQ و Log Likelihood Ratio (LLR) استفاده شده؛ و روش پیشنهادی با روشهایی که مدل گوسی و لاپلاس را برای سیگنال تمیز فرض کرده‌اند، مقایسه شده است. بررسی نتایج بیانگر عملکرد قابل قبول روش پیشنهادی است.

کلمات کلیدی:

بهسازی گفتار، توزیع مخلوط لاپلاس، EM الگوریتم، تخمین گر MMSE، آمارگان کمینه

فهرست مقالات پذیرفته شده

- ۱- The Application Of Laplacian Mixture Model In Speech Enhancement: ICSAP CONFERENCE ۲۰۱۱, Singapore
- ۲- Speech Spectral Estimation Based on MMSE Estimator and Laplacian Mixture Model: ICECE CONFERENCE ۲۰۱۱, Dubai
- ۳- Speech Enhancement Using Mixture Of Laplacian Based MMSE Estimator: International Conference on Communication System, Zahedan, ۲۰۱۰

۱- فصل اول: مقدمه

۲	۱-۱- مقدمه ای بر بهسازی گفتار
۳	۲-۱- کاربردهای سیستم های بهسازی گفتار
۳	۳-۱- تاریخچه روش های بهسازی گفتار
۵	۴-۱- طبقه بندی روش های بهسازی گفتار
۶	۵-۱- اهداف پایان نامه
۷	۶-۱- ساختار پایان نامه

۲- فصل دوم: بهسازی گفتار

۹	۱-۲- مقدمه
۹	۲-۲- بهسازی گفتار
۱۱	۳-۲- طبقه بندی روش های بهسازی گفتار
۱۳	۴-۲- عوامل موثر در طراحی سیستم های بهسازی گفتار
۱۴	۵-۲- ارزیابی سیستم های گفتار
۱۵	۲-۵-۱- معیار ارزیابی سیگنال به نویز (SNR)
۱۸	۲-۵-۲- معیار ارزیابی PESQ
۱۹	۲-۵-۳- معیار ارزیابی LLR
۲۰	۶-۲- منابع و مشخصات نویزهای صوتی

۳- فصل سوم: مروری بر روش های بهسازی گفتار تک کاناله

۲۳	۱-۳- مقدمه
۲۳	۲-۳- تفریق طیفی
۲۶	۳-۳- زیر فضای سیگنال
۲۷	۴-۳- فیلتر کالمن
۲۹	۵-۳- روش های آماری

۴- فصل چهارم: بهسازی گفتار با روش های آماری

۳۲	۱-۴- مقدمه
۳۲	۲-۴- روش های آماری محض
۳۳	۳-۴- تخمین گر بیشترین شباهت
۳۵	۴-۴- تخمین گرهای بیز
۳۶	۵-۴- تخمین گر MMSE
۳۷	۱-۵-۴- فیلتر وینر
۳۹	۲-۵-۴- تخمین اندازه طیف (STSA-MMSE)
۴۲	۳-۵-۴- تخمین فاز با کمترین میانگین مربعات خطا
۴۳	۴-۵-۴- محاسبه سیگنال به نویز پیشین
۴۴	۱-۴-۵-۴- روش بیشترین شباهت
۴۵	۲-۴-۵-۴- روش تصمیم گرا
۴۶	۵-۵-۴- تخمین MMSE با توزیع های غیر گوسی
۴۸	۶-۴- تخمین گر بیشترین احتمال پسین MAP
۴۹	۷-۴- به کارگیری عدم قطعیت گفتار (SPU) در بهسازی گفتار

۵۰	۴-۷-۱- تخمین گر MMSE با SPU
۵۰	۴-۸- روش های تخمین نویز در بهسازی گفتار
۵۱	۴-۸-۱- روش میانگین گیری بازگشتی
۵۲	۴-۸-۲- روش مبتنی بر هیستوگرام
۵۴	۴-۸-۳- روش ردیابی کمینه ها
۵۷	۴-۹- روش های آماری مبتنی بر مدل

۵- فصل پنجم: روش پیشنهادی بر مبنای توزیع مخلوطی از لاپلاس ها

۶۱	۵-۱- مقدمه
۶۱	۵-۲- معرفی پایگاه داده
۶۳	۵-۳- تخمین طیف سیگنال تمیز
۶۳	۵-۳-۱- فرضیات
۶۴	۵-۳-۲- مدل کردن سیگنال تمیز با توزیع مخلوط لاپلاس
۶۵	۵-۳-۱- توزیع مخلوط لاپلاس (LMM)
۶۸	۵-۳-۳- تخمین طیف سیگنال تمیز با تخمین گر MMSE
۷۱	۵-۴- تخمین پارامترهای توزیع سیگنال تمیز با EM الگوریتم
۷۲	۵-۴-۱- بیشترین شباهت
۷۳	۵-۴-۲- EM الگوریتم
۷۵	۵-۴-۳- محاسبه پارامترهای توابع توزیع مخلوط
۷۸	۵-۴-۴- محاسبه پارامترهای LMM
۷۹	۵-۴-۵- پیاده سازی EM الگوریتم
۸۶	۵-۵- تخمین نویز
۸۷	۵-۶- ارزیابی روش پیشنهادی

۸۹	۱-۶-۵- معیار ارزیابی سیگنال به نویز قطعه ای
۹۳	۲-۶-۵- معیار ارزیابی LLR
۹۶	۳-۶-۵- معیار ارزیابی PESQ
۹۹	۷-۵- مقایسه روش پیشنهادی با روش های دیگر
۱۰۰	۱-۷-۵- معیار سیگنال به نویز قطعه ای
۱۰۲	۲-۷-۵- معیار LLR
۱۰۴	۳-۷-۵- معیار PESQ

۶- فصل ششم: نتیجه گیری و پیشنهادات

۱۰۸	۱-۶- مقدمه
۱۰۸	۲-۶- نتیجه گیری و جمع بندی
۱۱۱	۳-۶- پیشنهاد برای کارهای آینده
۱۱۲	فهرست منابع

فهرست شکل ها

صفحه

عنوان

- شکل (۱-۳): نمای کلی الگوریتم تفریق طیفی . ۲۵
- شکل (۲-۳) نمای کلی روش های آماری. ۳۰
- شکل (۱-۴): پریودوگرام $|Y(\lambda, k)|^2$ ، پریودوگرام هموار شده $P(\lambda, k)$ و تخمین نويز برای فرکانس $k=20$. ۵۶
- شکل (۱-۵): بلوک دیاگرام روش پیشنهادی. ۶۲
- شکل (۲-۵): توزیع لاپلاس به ازای واریانس و میانگین های متفاوت. ۶۵
- شکل (۳-۵): هیستوگرام بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز. ۶۶
- شکل (۴-۵): هیستوگرام بخش موهومی تبدیل فوریه زمان کوتاه سیگنال تمیز. ۶۷
- شکل (۵-۵): نقطه چین: تخمین توزیع بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز با یک لاپلاس، خط یکپارچه: تخمین توزیع بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز با مخلوطی از شش لاپلاس، خط چین: تخمین توزیع بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز با شش گوسی. ۶۸
- شکل (۶-۵): تخمین EM الگوریتم برای بخش حقیقی طیف سیگنال تمیز، (الف) پس از ۱۰ بار تکرار، (ب) پس از ۳۰ بار تکرار، (ج) پس از ۴۰ بار تکرار، (چ) پس از ۵۰ بار تکرار. ۸۴
- شکل (۷-۵): تخمین EM الگوریتم برای بخش حقیقی طیف سیگنال تمیز : (الف) به ازای یک لاپلاس، (ب) به ازای ده لاپلاس. ۸۵

شکل (۵-۸): تخمین EM الگوریتم برای بخش موهومی طیف سیگنال تمیز:

۸۶ (الف) به ازای یک لاپلاس، (ب) به ازای ده لاپلاس.

شکل (۵-۹): خط چین: پریدوگرام سیگنال نویزی که با نویز F۱۶ و

در سیگنال به نویز ۵ دسی بل، آغشته شده است. خط یکپارچه نازک:

۸۸ پریدوگرام هموار شده و خط یکپارچه کلفت: تخمین نویز.

شکل (۵-۱۰): نمودار سیگنال به نویز قطعه ای سیگنال بهسازی شده،

۹۱ پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز سفید.

شکل (۵-۱۱): نمودار سیگنال به نویز قطعه ای سیگنال بهسازی شده،

۹۱ پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز F۱۶.

شکل (۵-۱۲): نمودار سیگنال به نویز قطعه ای سیگنال بهسازی شده،

۹۲ پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز همهمه.

شکل (۵-۱۳): نمودار LLR برای سیگنال بهسازی شده، پس از اعمال

۹۴ روش پیشنهادی روی سیگنال آغشته به نویز سفید.

شکل (۵-۱۴): نمودار LLR برای سیگنال بهسازی شده، پس از اعمال

۹۵ روش پیشنهادی روی سیگنال آغشته به نویز F۱۶.

شکل (۵-۱۵): نمودار LLR برای سیگنال بهسازی شده، پس از اعمال

۹۵ روش پیشنهادی روی سیگنال آغشته به نویز همهمه.

شکل (۵-۱۶): نمودار PESQ برای سیگنال بهسازی شده، پس از اعمال

۹۷ روش پیشنهادی روی سیگنال آغشته به نویز سفید.

شکل (۵-۱۷): نمودار PESQ برای سیگنال بهسازی شده، پس از اعمال

۹۷ روش پیشنهادی روی سیگنال آغشته به نویز F۱۶.

شکل (۵-۱۸): نمودار PESQ برای سیگنال بهسازی شده، پس از اعمال

- ۹۸ روش پیشنهادی روی سیگنال آغشته به نویز همهمه.
شکل (۵-۱۹): مقایسه روش پیشنهادی با روش های
Lap-MMSE, Log-MMSE, MMSE بر حسب سیگنال
- ۱۰۱ به نویز قطعه ای و برای نویز F۱۶.
شکل (۵-۲۰): مقایسه روش پیشنهادی با روش های
Lap-MMSE, Log-MMSE, MMSE بر حسب سیگنال
- ۱۰۱ به نویز قطعه ای و برای نویز همهمه.
شکل (۵-۲۱): مقایسه روش پیشنهادی با روش های
Lap-MMSE, Log-MMSE, MMSE بر حسب
- ۱۰۲ LLR و برای نویز F۱۶.
شکل (۵-۲۲): مقایسه روش پیشنهادی با روش های
Lap-MMSE, Log-MMSE, MMSE بر حسب
- ۱۰۳ LLR و برای نویز همهمه.
شکل (۵-۲۳): مقایسه روش پیشنهادی با روش های
Lap-MMSE, Log-MMSE, MMSE بر حسب
- ۱۰۴ PESQ و برای نویز F۱۶.
شکل (۵-۲۴): مقایسه روش پیشنهادی با روش های
Lap-MMSE, Log-MMSE, MMSE بر حسب
- ۱۰۵ PESQ و برای نویز همهمه.

فهرست جدول ها

صفحه	عنوان
۱۶	جدول (۱-۲): مهم ترین روشهای ارزیابی سیستم های بهسازی گفتار.
۲۶	جدول (۱-۳): مروری بر روش تفریق طیفی.
۲۷	جدول (۲-۳): مروری بر روش زیر فضای سیگنال.
۳۰	جدول (۳-۳): خلاصه روش های آماری.
	جدول (۱-۵): پارامترهای LMM و I_{RL} برای بخش حقیقی تبدیل
	فوریه زمان کوتاه داده های TIMIT به ازای $N=3$ و بعد از ۱۰، ۳۰،
۸۱	۴۰، و ۵۰ تکرار EM الگوریتم.
	جدول (۲-۵): مقایسه سیگنال به نویزهای قطعه ای برای سه نویز همهمه،
۹۲	سفید و F1۶ در سه سیگنال به نویز ۰، ۵، ۱۰ و N های متفاوت.
	جدول (۳-۵): مقایسه LLR برای سه نویز همهمه، سفید و F1۶ در سه سیگنال
۹۶	به نویز ۰، ۵، ۱۰ و N های متفاوت.
	جدول (۴-۵): مقایسه PESQ برای سه نویز همهمه، سفید و F1۶ در سه سیگنال
۹۸	به نویز ۰، ۵، ۱۰ و N های متفاوت.
	جدول (۵-۵): مقایسه روش پیشنهادی با روش های
	Lap-MMSE, Log-MMSE, MMSE بر حسب سیگنال
۱۰۲	به نویز قطعه ای و برای نویز همهمه و F1۶.
	جدول (۶-۵): مقایسه روش پیشنهادی با روش های
	Lap-MMSE, Log-MMSE, MMSE بر حسب

۱۰۳

LLR و برای نویز همهمه و F۱۶.

جدول (۷-۵): مقایسه روش پیشنهادی با روش های

MMSE, Log-MMSE, Lap-MMSE بر حسب

۱۰۵

PESQ و برای نویز همهمه و F۱۶.

فصل اول

مقدمه

۱-۱- مقدمه ای بر بهسازی گفتار

همان طور که می دانیم مهمترین وسیله ارتباطی افراد از طریق صدا و گفتار می باشد. به همین جهت پردازش گفتار، به جهت کاربردهای فراوان در زندگی بشر، از سالیان قبل مد نظر دانشمندان و محققین بوده است. یکی از اساسی ترین بخش های پردازش گفتار، بهسازی گفتار و تلاش برای بهبود کیفیت عملکرد سیستم های ارتباط گفتاری می باشد.

شاید تا به حال برای شما پیش آمده باشد که از یک نوار کاست قدیمی که حاوی سخنرانی یا اطلاعات مورد نظر شماست، استفاده کرده باشید و احتمالاً به علت گذشت زمان یا مشکلات زمان ضبط، صداهای روی آن به سختی شنیده شود. قطعاً شما دوست دارید که به طریقی کیفیت صداهای ضبط شده روی نوار کاست را بالا ببرید. یا حالتی را در نظر بگیرید که در داخل خودرو با تلفن همراه صحبت می کنید. قطعاً صداهای مزاحم نیز مانند صدای باد، موتور ماشین و سرنشینان همراه با صدای شما وارد میکروفون می شوند.

این موارد تنها مثال های ساده ای از کاربردهای فراوان روش های کاهش یا حذف اثرات نویز بر سیگنال صوت است که به روش های بهسازی گفتار معروف هستند.

در قالب یک تعریف کلی، بهسازی گفتار عبارت است از تلاش برای بهبود کیفیت عملکرد سیستم های ارتباط گفتاری در مواردی که سیگنال گفتار تحت تاثیر نویز، انعکاس^۱ و سایر عوامل تخریبی واقع گردیده است.

نیاز به بهسازی گفتار از آنجا ناشی می شود که سیگنال گفتار معمولاً، یا از منبعی تولید می شود که در محیط نویزی قرار دارد، یا تحت تاثیر کانال انتشار و در اثر نویز یا انعکاس دچار تخریب می شود، یا در گیرنده به نویز آلوده می شود و یا ممکن است منشا نویز ترکیبی از سه حالت فوق باشد [۶۶].

^۱ - Echo

۱-۲- کاربردهای سیستم های بهسازی گفتار

سیستم های بهسازی گفتار کاربردهای وسیع و متفاوتی در حوزه های مختلف دارند که به طور کلی می توان آنها را به صورت زیر طبقه بندی کرد:

- تشخیص خودکار گفتار^۱ و تشخیص هویت گوینده^۲: در سیستم های تشخیص خودکار گفتار و هویت گوینده، عملکرد سیستم در محیط هایی با سیگنال به نویز پایین مناسب نمی باشد. ثابت شده که نرخ تشخیص با به کارگیری الگوریتم های بهسازی گفتار بهبود می یابد [۵۲].
- سیستم های تلفن: که در آنها گفتار اصلی به وسیله نویز زمینه یا نویز موجود در مسیر مخابره و نیز در اثر انعکاس صدای طرفین مکالمه خراب می شود.
- تلفن های عمومی: که در محیط های پر سر و صدا و شلوغ واقع اند.
- سیستم های مخابرات هوا به زمین: که در آنها نویز اتاق خلبان، پیام ارسال شده از سوی خلبان را خراب می کند.
- سمعک: که به عنوان یک تقویت کننده، هم سیگنال گفتار و هم نویز موجود در محیط را تقویت نموده و موجب ناراحتی کاربر می شود.
- حذف نویز در گفتارهای ضبط شده: که در آنها صدای ضبط شده اشخاص آغشته به نویز است و جهت بهبود کیفیت آنها، باید نویز را حذف کرد.
- کدینگ اطلاعات: که اثر نویز بر گفتار به شدت عملکرد آنها را تحت تاثیر قرار می دهد [۶۶].

۱-۳- تاریخچه روش های بهسازی گفتار

با توجه به کاربرد های مهم و فراوانی که برای روش های بهسازی گفتار بیان شد، و اهمیت آن، روش های مختلفی برای این کار پیشنهاد شده اند.

^۱ - Automatic Speech Recognition
^۲ - Speaker Identification

از مهمترین و قدیمی ترین روش های بهسازی گفتار می توان به فیلتر وینر^۱ [۶۰] و مشتقات آن اشاره کرد. در این روش ها، که برای اولین بار در سال ۱۹۴۹ وارد حوزه بهسازی گفتار گردید، سعی می شود که تخمین سیگنال تمیز بر اساس یک معیار بهینه سازی باشد. یکی دیگر از روش های قدیمی و بسیار معروف بهسازی گفتار، تفریق طیفی^۲ [۴] و [۶]، است که در سال ۱۹۷۵ پایه گذاری شد و در سال ۱۹۷۹ با پردازش در حوزه فرکانس مطرح شد. پس از آن بهبودهای زیادی بر روی این روش صورت گرفت. فیلتر کالمن^۳ هم در سال ۱۹۸۷ با مقاله Paliwal & Basu [۴۷] وارد این حوزه گردید.

روش های آماری محض^۴ به عنوان یکی از مهم ترین روش های بهسازی گفتار، با کارهای Lim & Malah & Ephraim [۲۲] و Oppenheim [۳۵] و McAullay & Malpass [۴۵] مقاله معروف [۲۲] و پیشنهاد شدند و بسیار مورد توجه محققین قرار گرفتند. روش های آماری مبتنی بر مدل^۵ نیز با کارهای Ephraim در [۲۴] و [۱۸] مطرح شدند و امروزه به عنوان یکی مهم ترین زمینه های علمی و کاربردی در بهسازی گفتار مورد توجه می باشند. می توان ادعا کرد که در کل روش های آماری چه از لحاظ کیفی و چه از لحاظ کمی نتایج بهتری نسبت به اغلب روشهای دیگر ارائه می دهند.

روش مهم دیگر بهسازی گفتار که مبتنی بر تئوری جبر خطی می باشد، زیر فضای سیگنال^۶ است که در ابتدا در [۱۴] و [۱۶] و [۲۵] و در سال ۱۹۹۱ ارائه گردید و سپس بهبودهای زیادی روی آن صورت پذیرفت. اخیراً روش های مبتنی بر کپسترال^۷، تبدیل موجک^۸ و تبدیل فوریه کسینوسی^۹ نیز وارد حوزه بهسازی گفتار شده اند.

علاوه بر این روش ها، روش های دیگری نیز برای بهسازی گفتار پیشنهاد شده اند که در آنها از دو یا

^۱ -Wiener Filter

^۲ - Spectral Subtraction

^۳ - Kalman Filter

^۴ - Pure Statistical Method

^۵ - Statistical Model-Based

^۶ - Signal Subspace

^۷ - Cepstral

^۸ - Wavelet

^۹ -Discrete Cosine Transform

چند میکروفن برای ضبط سیگنال نویزی استفاده می‌شود. از مهم ترین این روش ها می توان به روش فیلتر وقتی دوکاناله که در سال ۲۰۰۳ مطرح شد، اشاره کرد.

تا امروز مقالات و روش های بسیاری در زمینه بهسازی گفتار ارائه شده است، اما تخمین سیگنال تمیز از روی سیگنال نویزی هنوز یک مساله حل نشده در پردازش گفتار است و چالش های متعددی را در خود نگه داشته است. یکی از دلایل ساده برای اثبات این ادعا این است که علی رغم قدمت بیشتر بهسازی گفتار نسبت به برخی کاربردهای دیگر پردازش گفتار مثل بازشناسی گفتار، هنوز روش قطعی و غالبی در این زمینه وجود ندارد. از دلایل اصلی مشکل بودن حل این مساله، نامشخص بودن جواب دقیق برخی از سوالات کلیدی و ابهامات در این حوزه است. از جمله این سوالات عبارتند از:

- آیا می‌توان یک رابطه مشخص را برای معیارهای کیفیت و قابلیت فهم، به صورتی که انسان درک می کند، به دست آورد؟ پاسخ به این سوال، به نوبه خود منجر به یافتن معیار بهینه سازی معنی داری می شود که با ادراک انسان متناسب است. با داشتن روابط مشخص برای این دو، می توان یکی از آنها را بدون داشتن اثر منفی بر دیگری بهبود بخشید و می توان به این سوال نیز پاسخ داد که چه روش هایی برای بهسازی گفتار در کاربرد کمک شنیداری مناسب است؟ کدام یک برای کاربرد بازشناسی گفتار و کدام یک برای کدینگ؟
- چه پردازش هایی در سطوح مختلف سیستم شنوایی انسان انجام می‌گیرد و چگونه می‌توان آنها را مدل کرد؟ پاسخ این سوال نیز منجر به روشن شدن شبهات دیگری مانند نحوه ترکیب تکنیک های پردازش سیگنال با روش های ادراکی انسان می‌شود.

۱-۴- طبقه بندی روش های بهسازی گفتار

با توجه به تعدد روش های بهسازی گفتار، آنها را می‌توان به طرق مختلفی طبقه بندی کرد [۲۸] که در این فصل فقط آنها را نام می‌بریم. این طبقه بندی در فصل بعد توضیح داده خواهد شد.

۱- طبقه بندی بر اساس تعداد کانال ورودی

۲- طبقه بندی براساس حوزه فعالیت

۳- طبقه بندی بر اساس پارامتری و غیر پارامتری بودن

۱-۵- اهداف پایان نامه

در حال حاضر، رویکرد غالب در بهسازی گفتار تک کاناله، روش های آماری هستند که در کل کارایی بهتری در مقایسه با روش های دیگر دارند. این روش ها را که از دانش پیشین سیگنال تمیز و نویز بهره می گیرند، می توان در دو دسته آماری محض و آماری مبتنی بر مدل تقسیم بندی کرد. در روش های آماری مبتنی بر مدل، سعی می شود با استفاده از روش های مدل سازی، توزیع های گفتار تمیز و نویز به کمک داده های آموزشی به دست آید. سپس از این توزیع ها در تخمین سیگنال تمیز از روی سیگنال نویزی استفاده می شود. در مقابل، در روش های آماری محض، یک شکل توزیع خاص برای سیگنال و نویز فرض می شود و پارامترهای این توزیع ها با استفاده از سیگنال نویزی محاسبه می شوند. سپس با استفاده از این توزیع ها، روابط تخمین سیگنال تمیز از روی سیگنال نویزی به دست می آیند. در این پایان نامه از روش های آماری به منظور بهسازی استفاده شده است. اگرچه روش پیشنهادی را نمی توان به طور کامل متعلق به یکی از دو حوزه روش های آماری محض و آماری مبتنی بر مدل دانست، ولی شباهت آن به روش های آماری محض بیشتر می باشد. در این روش یک توزیع خاص برای سیگنال تمیز فرض می شود. ولی پارامترهای آن به جای استفاده از سیگنال نویزی، با کمک داده های آموزشی به دست می آیند. همچنین برای نویز نیز یک توزیع مشخص فرض می شود که پارامترهای آن از روی سیگنال نویزی محاسبه می شوند. هدف از انجام این پایان نامه و جهت گیری های اصلی آن را می توان در جملات زیر خلاصه کرد:

- انجام بهسازی گفتار به روش آماری محض، با به کارگیری توزیع مخلوط لاپلاس^۱ برای طیف

طیف

گفتار تمیز، به منظور تخمین سیگنال تمیز از روی سیگنال نویزی در حضور نویز جمع شونده

^۱-Laplacian Mixture Model

- بهبود روش های آماری محض مبتنی بر توزیع گوسی^۱ و لاپلاس^۲
- ارزیابی سیستم بهسازی با معیارهای سیگنال به نویز قطعه ای^۳، LLR^۴ و PESQ^۵

۱-۶- ساختار پایان نامه

این تحقیق شامل شش فصل می باشد: تعاریف اولیه و کلیاتی درباره بهسازی گفتار در فصل جاری مطرح شد. در فصل دوم مختصراً در مورد بهسازی گفتار، عوامل موثر در طراحی سیستم های بهسازی گفتار و چگونگی ارزیابی سیستم های بهسازی گفتار بحث می گردد. فصل سوم نیز مروری بر مهم ترین روش های بهسازی تک کاناله بهسازی گفتار از جمله تفریق طیفی، زیرفضای سیگنال، فیلتر کالمن و روش های آماری می باشد.

در فصل چهارم نیز به بررسی روش های آماری، علی الخصوص روش آماری محض به سبب ارتباطی که با روش پیشنهادی دارد، پرداخته شده است. در این فصل تخمین گرهای بیشترین شباهت (ML)^۶ و روش بیز^۷ مرور شده و روش های بهسازی گفتار مرتبط با آنها معرفی می گردد. فصل پنجم هم به بیان روش پیشنهادی و نتایج آن می پردازد. این فصل شامل پنج بخش اساسی می باشد. این بخش ها عبارتند از: محاسبه روابط تخمین گر MMSE با فرض توزیع مخلوط لاپلاس برای طیف سیگنال گفتار تمیز، بسط EM^۸ الگوریتم برای محاسبه پارامترهای مخلوط لاپلاس، محاسبه واریانس نویز، ارزیابی نتایج روش پیشنهادی و مقایسه روش پیشنهادی با سه روش معروف بهسازی گفتار. فصل ششم نیز راجع به جمع بندی، نتیجه گیری و پیشنهادات برای کارهای آینده می باشد.

^۱ - Gaussian
^۲ - Laplacian
^۳ - Segmental Signal to Noise
^۴ - Log Likelihood Ratio
^۵ - Perceptual Evaluation of Speech Quality
^۶ - Maximum Likelihood
^۷ - Bayesian
^۸ - Expectation Maximization

فصل دوم

بهسازی گفتار

۲-۱- مقدمه

همان طور که قبلا بیان شد، بهسازی گفتار تلاش برای بهبود کیفیت عملکرد سیستم های ارتباط گفتاری در مواردی که سیگنال گفتار تحت تاثیر نویز، انعکاس و سایر عوامل تخریبی واقع گردیده است، می باشد. در این فصل به بررسی مفاهیم مرتبط با بهسازی گفتار پرداخته می شود و انواع طبقه بندی های موجود برای بهسازی گفتار بیان خواهد شد. همچنین مطالبی نیز راجع به روش های ارزیابی سیستم های بهسازی گفتار آورده شده است. منابع و مشخصات نویزهای صوتی نیز به عنوان یکی از عوامل مهم در شناخت ضرورت بهسازی گفتار، بررسی شده است.

۲-۲- بهسازی گفتار

بهسازی گفتار^۱ به روش هایی منتهی می شود که با انجام پردازش هایی روی سیگنال های نویزی، نهایتا منجر به بهبود کیفیت^۲ و قابلیت فهم^۳ گفتار می شود [۵۹].

کیفیت و قابلیت فهم دو معیار مهم ارزیابی یک سیگنال صوتی برای انسان می باشند. کیفیت یک معیار ذهنی^۴ ارزیابی سیگنال است که وضوح مطالب بیان شده و سطح نویز زمینه موجود در آن را اندازه می گیرد. می توان این معیار را با میزان خوشایند بودن سیگنال برای شنونده یا میزان تلاش شنونده برای درک مطلب بیان شده در سیگنال بیان کرد. در مقابل قابلیت فهم یک معیار عینی^۵ است که درصد کلماتی (که می توانند بی معنی باشند) که به درستی توسط شنونده تشخیص داده

^۱ - Speech enhancement

^۲ - Quality

^۳ - Intelligibility

^۴ - Subjective

^۵ - Objective

شده را بیان می کند. این دو معیار از هم مستقل هستند، یعنی یک سیگنال صوتی ممکن است کیفیت بالایی داشته باشد، در حالیکه قابلیت فهم آن پایین باشد یا بالعکس. به عنوان مثال صدای تولید شده توسط ماشین دارای کیفیت بسیار پایین است، اما قابلیت فهم بالایی دارد [۲۱].

بهبود کیفیت و افزایش قابلیت فهم سیگنال صوتی به صورت توأم، از اهداف بهسازی گفتار است. اما تقریباً تمامی الگوریتم های بهسازی تک کاناله فقط کیفیت را بهبود می بخشند، در حالیکه قابلیت فهم را نه تنها افزایش نمی دهند، بلکه به علت اعوجاجی^۱ که در سیگنال ایجاد می کنند، آن را کاهش نیز می دهند. این مشکل به علت تاثیر گذاری روش بهسازی بر بخش های نویزگونه انرژی پایین سیگنال گفتار مانند صامت ها^۲ به وجود می آید. زیرا این روش ها معمولاً صامت ها را نویز تشخیص داده و حذف می کنند. در حالیکه صامت ها نقش موثرتری نسبت به مصوت ها^۳ در فهم گفتار دارند [۲۰]. از این رو بیشتر شنونده ها اطلاعات بیشتری را از سیگنال نویزی در مقایسه با سیگنال بهسازی شده استخراج می کنند.

در برخی از کاربردها افزایش کیفیت سیگنال مورد توجه است. به ویژه زمانی که شنونده برای زمانی طولانی در معرض نویز شدید قرار دارد. مانند محیط کارخانه یا کابین خلبان هلی کوپتر، در مقابل در برخی کاربردها مانند سیستم های نظامی، قابلیت فهم سیگنال از کیفیت آن مهم تر است.

علی رغم اینکه روش های مختلفی برای این دو معیار وجود دارند، اما هنوز یک روش کمی دقیق که بتواند این دو معیار را بدون نیاز به شنونده انسانی محاسبه نماید، وجود ندارد. همین مساله یکی از مهم ترین دلایل مشکل بودن کار بهسازی گفتار است. چرا که تا به حال روش ریاضی دقیقی که با ادراک انسان سازگار باشد، برای این معیارها و در نتیجه محاسبه دقیق اعوجاج سیگنال، پیشنهاد نشده است. یکی از نکات کلیدی در طراحی یک سیستم بهسازی گفتار این است که روش بهسازی قادر به حذف نویز زمینه (افزایش کیفیت سیگنال) باشد، بدون اینکه سیگنال اصلی را دچار اعوجاج کند. این بدین معنی است که این سیستم ها در سیگنال بهسازی شده دو نوع اثر از خود به جای می

^۱ - Distortion
^۲ -consonant
^۳ -vowel

گذارند، یکی حذف نویز زمینه و دیگری اعوجاج سیگنال. در صورت دسترسی به معیاری دقیق برای فرموله کردن کیفیت و قابلیت فهم، می توان کیفیت را بدون ایجاد اعوجاج در سیگنال، افزایش داد.

۲-۳- طبقه بندی روش های بهسازی گفتار

در طول دهه های اخیر و با پیشرفت سخت افزاری و نرم افزاری روش های مختلف پردازش سیگنال گفتار، طیف وسیعی از رویکردها و روش ها برای بهسازی گفتار ارائه گردیده است. با توجه به تعدد روش های بهسازی گفتار، آنها را به روش های متفاوتی طبقه بندی می کنند.

یکی از روش های تقسیم بندی تکنیک های بهسازی گفتار بر پایه تعداد کانال (میکروفون) ورودی می باشد. بر این اساس دو خانواده تک کاناله^۱ و چند کاناله^۲ برای گروه بندی روش های مختلف موجود می باشد. در روش تک کاناله تنها یک میکروفون ورودی در دسترس بوده و کلیه اطلاعات لازم باید از همین سیگنال استخراج شود. اساس کار این خانواده از روش ها، مبتنی بر ایستادن^۳ بودن نویز می باشد. فرضی که در مواجهه با نویزهای غیرایستادن^۴ زیر سؤال رفته و عملکرد روش را دچار مشکل می نماید. در روش های چندکاناله، دو یا چند گیرنده (میکروفون) در ورودی سیستم بهسازی گفتار مورد استفاده واقع می شود. افزایش تعداد میکروفون ها و یا کانال های ورودی، قدرت روش را در پاکسازی سیگنال نویزی بالا می برد و کار بهسازی را ساده تر می کند، اما در مقابل بر هزینه و پیچیدگی پیاده سازی سیستم نیز می افزاید.

از روش های تک کاناله می توان به تفریق طیفی، فیلتر وینر، زیرفضای سیگنال، فیلتر کالمن، روشهای مبتنی بر تبدیل ویولت^۵، روشهای مبتنی بر کپسترال^۶ و روش های آماری اشاره کرد. از روش های چند کاناله نیز می توان فیلتر وقتی^۱ (بیشتر به صورت دوکاناله)، جداسازی سیگنال^۲ (قابل

قابل)

۱- Single Channel
۲- Multi Channel
۳-Stationary
۴-Non-Stationary
۵ - Wavlate
۶ - Cepstral

استفاده هم به صورت تک کاناله و هم به صورت چند کاناله) و روش های تشکیل پرتو^۳ را نام برد. روش های بهسازی گفتار را می توان براساس حوزه فعالیت نیز دسته بندی کرد که بر این اساس آنها را به دو دسته کلی حوزه زمان (مانند فیلترکالمن و فیلتروفقی) و حوزه فرکانس یا طیف^۴ (مانند تفریق طیفی و فیلتر وینر) تقسیم می کنند. هرچند برخی از روش ها مانند روش های مبتنی بر تبدیل موجک یا کپسترال به حوزه طیف نزدیکترند؛ اما نمی توان آنها را به طور مستقیم به یکی از دو دسته فوق نسبت داد.

بیشتر روش های بهسازی گفتار در حوزه ی فرکانس و روی طیف عمل می کنند. زیرا در حوزه طیف قابلیت جداسازی نویز و سیگنال بالاست و می توان روش های بهینه یا ابتکاری را به خوبی پیاده سازی کرد. همچنین، مستقل بودن مولفه های طیف امکان کار با هر کدام از آنها را به صورت مستقل فراهم می کند که این به نوبه خود منجر به ساده تر شدن روابط ریاضی و کاهش محاسبات می گردد. انتقال سیگنال به حوزه طیف در بیشتر مواقع به کمک تبدیل فوریه زمان کوتاه^۵ صورت می گیرد که منجر به تجزیه سیگنال به دو بخش دامنه^۶ و فاز^۷ یا دو بخش حقیقی و موهومی می شود. اغلب روش های بهسازی گفتار که در حوزه طیف عمل می کنند، فقط اندازه دامنه یا مجذور اندازه طیف را پردازش می کنند و از فاز سیگنال نویزی به عنوان فاز تخمینی سیگنال تخمین زده شده استفاده می کنند. یکی از دلایل اصلی این مساله حساسیت کمتر سیستم شنوایی انسان به فاز در مقایسه با دامنه است. برخی روش های دیگر نیز به تخمین بخش های حقیقی و موهومی سیگنال، به صورت جداگانه می پردازند. سپس این نتایج را با هم ترکیب می کنند تا به تخمین نهایی طیف سیگنال تمیز برسند. این روش ها اگر چه از روش های تخمین دامنه و مربع دامنه پیچیده تر به نظر می -

^۱ - Adaptive Filter

^۲ -Signal Separation

^۳ -Beam Forming

^۴ -Spectrum

^۵ -Short time Fourier Transform

^۶ - Amplitude

^۷ - Phase

رسند، اما نتایج آنها بر حسب سیگنال به نویز^۱ بهتر از این دو روش است [۹]. در اغلب روش های بهسازی گفتار، سیگنال نهایی، پس از اعمال تبدیل فوریه معکوس بر روی طیف سیگنال بهسازی شده، از روش هم پوشانی - جمع^۲ بدست می آید.

یکی دیگر از راه های دسته بندی کردن روشهای بهسازی گفتار تفکیک این روش ها براساس پارامتری یا غیر پارامتری بودن آنهاست. روش هایی مانند تفریق طیفی و زیر فضای سیگنال به دلیل کار بر روی خود سیگنال، روش های غیر پارامتری محسوب می شوند و روش های مبتنی برمدل مخفی مارکوف^۳ (HMM) و مدل مخلوط گوسی^۴ (GMM) به دلیل استفاده از مدل های آماری پارامتری برای سیگنال و نویز در ردیف روش های پارامتری قرار می گیرند [۲۸].

با وجود طبقه بندی های فوق، به دلیل تشابه زیاد روش های بهسازی گفتاری در ایده ها و پایه های بنیادی، تقسیم بندی آنها در قالب چند روش کاملاً مجزا، کار ساده ای نیست. به گونه ای که می توان گفت برخی از روش ها، حالت کلی یا تعمیم یافته برخی روش های دیگر بوده و یا این که برخی از روش ها صرف نظر از استفاده از تبدیل هایی متفاوت در نمایش سیگنال (مانند فوریه و موجک) دارای اساس مشترکی هستند.

۲-۴- عوامل موثر در طراحی سیستم های بهسازی گفتار

علی رغم وجود روش های مختلف در بهسازی گفتار، باید توجه داشت که طراحی یک سیستم بهسازی گفتار به عوامل مختلفی وابسته است. از جمله این عوامل می توان به موارد زیر اشاره کرد:

- کاربرد مورد نظر: در برخی کاربردها افزایش کیفیت مساله اصلی است مانند وسایل کمک شنوایی، ولی در برخی از آنها قابلیت فهم و درک (مانند کاربرد های نظامی و باز شناسی گفتار) مهم تر است.

^۱ SNR

^۲ - Overlap and Add

^۳ - Hidden Markov Model

^۴ -Gaussian Mixture Model

- نوع نویز و مشخصات آن: اینکه نویز ایستان یا غیر ایستان باشد، از جنس گفتار (مانند نویز همهمه)، جنس نویز زمینه (صدای ماشین، پنکه و باد) و یا از نوع پژواک اتاق و نویز کانال باشد، یا اینکه نویز دارای پهنای باند بالا یا کم باشد، در طراحی سیستم بهسازی گفتار موثر است.
- رابطه نویز و گفتار تمیز: اثر نویز برگفتار می تواند جمع شونده^۱ (مثل نویز سفید)، کانال-شونده^۲ (مثل نویز کانال تلفن یا نویز پژواک^۳) و یا ضرب شونده^۴ (نویز محوسازی) باشد.
- نویز و گفتار می توانند همبسته^۵ یا غیر همبسته^۶ باشند.
- تعداد میکروفون ها: هر چه تعداد میکروفون ها بیشتر باشد، بهسازی گفتار راحت تر است [۵۹].

۲-۵ - ارزیابی سیستم های بهسازی گفتار

مشابه سایر سیستم ها، برای ارزیابی کارایی روش های بهسازی گفتار و مقایسه تکنیک های مختلف با هم، نیازمند داشتن ابزارهایی هستیم که بتواند عملکرد یک سیستم بهسازی گفتار را به صورت کمی نشان دهند. هر چند یک انسان به عنوان شنونده، ممکن است معیارهای متفاوتی را برای سنجیدن کیفیت یک سیگنال در نظر داشته باشد، اما با توجه به محدودیت هایی که در فرموله کردن این معیارها وجود دارند، به روش های ساده تر و قابل اندازه گیری توجه می گردد.

ارزیابی یک سیستم بهسازی گفتار می تواند وابسته به کاربردی که از این سیستم مورد نظر است، باشد. مثلا اگر هدف از بهسازی گفتار، بهبود کارایی یک سیستم بازشناسی گفتار است، دقت یا صحت سیستم بازشناسی می تواند بهترین معیار ارزیابی برای روش های بهسازی گفتار باشد.

در یک نگاه کلی می توان تکنیک های ارزیابی کیفیت گفتار را به دو دسته ذهنی و عینی تقسیم بندی کرد.

معیارهای ذهنی، براساس نظر شنونده های انسانی بوده که به صورت شهودی رتبه ای را به سیگنال

^۱ - Additive
^۲ - Convolve
^۳ - Reverberation
^۴ - Multiply
^۵ - correlated
^۶ - Uncorrelated

گفتار در یک بازه ی امتیازی از قبل تعیین شده اختصاص می‌دهند. از معایب این معیارها این است که، نتایج یک آزمایش ذهنی تحت تاثیر عوامل زیادی از جمله محتوای آزمایش و علائق انجام دهنده آزمایش است. به علاوه این روش ها نیاز به دسترسی به تعداد افراد زیادی دارد که فرآیند ارزیابی را وقت گیر و پرهزینه می سازد.

در روش های عینی از تکنیک های محاسباتی جهت مقایسه سیگنال اصلی و سیگنال پردازش شده استفاده می‌کنند. اکثر روش های عینی اندازه گیری کمی کیفیت سیگنال را از روی مدل های سیستم شنوایی انسان انجام می‌دهند. امروزه به این روش ها در ارزیابی (علی رغم ارجحیت ارزیابی ذهنی) توجه بیشتری می‌شود [۳۷].

روش های عینی و ذهنی زیادی برای ارزیابی سیگنال گفتار ارایه شده است [۲۷] و [۳۷] که برخی از مهم ترین آنها در جدول (۱-۲) آورده شده است. از میان این روش ها محاسبه سیگنال به نویز، MOS^۱ و PESQ جهت ارزیابی کیفیت و روش DRT^۲ در ارزیابی قابلیت فهم پرکاربردتر می باشند. در ادامه این بخش معیارهای SNR، PESQ و LLR، با توجه به اینکه در فصل ۵ برای ارزیابی روش پیشنهادی استفاده شده اند، بررسی خواهند شد.

۲-۵-۱- معیار سیگنال به نویز (SNR)

محاسبه سیگنال به نویز از جمله قدیمی ترین و معمول ترین روش های ارزیابی عینی کیفیت سیگنال است که علاوه بر SNR کلی^۳، انواع مختلفی از آن مانند SNR قطعه ای و SNR قطعه ای وزن دار^۴ معرفی شده است. با فرض جمع شدن نویز $d(n)$ با سیگنال گفتار تمیز $s(n)$ ، می‌توان سیگنال-سیگنال-نال نویزی را به صورت $y(n) = s(n) + d(n)$ در نظر گرفت. اگر $\hat{s}(n)$ را سیگنال بهسازی شده در نظر بگیریم، آنگاه رابطه $e(n) = s(n) - \hat{s}(n)$ ، معادل خطای بین سیگنال اصلی و سیگنال تخمین زده شده برای نمونه n ام است. با استفاده از این فرضیات، SNR کلی برابر با نسبت انرژی

^۱ - Mean Opinion Score

^۲ - Diagnosis Rhythm Test

^۳ - Global SNR

^۴ - Frequency Weighted Segmental SNR

سیگنال اصلی به انرژی خطا تعریف می شود که رابطه آن بر حسب دسی بل به صورت زیر است:

$$SNR = 10 \cdot \log_{10} \frac{\sum_{n=1}^N s^*(n)}{\sum_{n=1}^N [s(n) - \hat{s}(n)]^2} \quad (1-2)$$

که در آن N تعداد کل نمونه های سیگنال است.

جدول (۱-۲) : مهم ترین روشهای ارزیابی سیستم های بهسازی گفتار [۵۹].

معیارهای عینی	
نوع ارزیابی	اسم روش
قابلیت فهم	Articulation Index (AI)
کیفیت	Copmposite Measure
کیفیت	Frequency Weighted Segmental SNR
کیفیت	Itakura Log-Likelihood Ratio(LLR)
کیفیت	Itakura Saito(IS)
کیفیت	Log Area Ratio (LAR)
کیفیت	Perceptual Evaluation of Speech Quality (PESQ)
کیفیت	Segmental SNR
کیفیت	Signal to Noise Ratio (SNR)
کیفیت	Weighted-Spectral Slope(WSS) (WSS)

معیارهای ذهنی	
نوع ارزیابی	اسم روش
قابلیت فهم	Diagnosis Rhythm Test (DRT)
قابلیت فهم	Hearing in Noise Test (HINT)
قابلیت فهم	Modified Rhythm Test (MRT)
قابلیت فهم	Speech Perception in Noise (SPIN)
قابلیت فهم	Speech Reception Test (SRT)
کیفیت	Diagnostic Accessibility Measure (DAM)
کیفیت	Isometric Absolute Judgment (IAJ)
کیفیت	ITU-T P.۸۳۵

کیفیت	Mean Opinion Score (MOS)
کیفیت	Paired Accessibility Rating (PAR)
کیفیت	Parametric Absolute Judgment (PAJ)
کیفیت	Quality Acceptance Rating Test (QUART)

بدیهی است که در این رابطه، به سیگنال تمیز نیاز داریم. بنابراین این رابطه بیشتر در شبیه سازی ها و برای زمانی که هر دو سیگنال تمیز و نویزی در دسترس هستند، کاربرد دارد. مهم ترین مزیت این

روش سادگی آن از نظر ریاضی است اما به علت عدم همبستگی با معیارهای ذهنی و یکسان فرض کردن کل سیگنال، در بسیاری از موارد تخمین مناسبی از کیفیت ارائه نمی دهد. می دانیم که در عمل، انرژی سیگنال گفتار با زمان تغییر می کند و بعضی از مولفه های سیگنال نقش بیشتری در تامین کیفیت ایفا می کنند. یکی از حالت هایی که می تواند به طور نادرست منجر به SNR بالا شود، وقتی است که سیگنال گفتار دارای قسمت های گفتاری متمرکزی باشد، در حالی که نویز بر روی قسمت های گفتار با انرژی کم اثر ادراکی شدیدی داشته باشد. از این رو اگر فرض کنیم که اعوجاج ناشی از نویز نوسانات انرژی کمی دارد، اندازه گیری می تواند قاب^۱ به قاب تغییر کند. این مسائل مبنای ارائه روش SNR قطعه ای است که SNR را در قاب های کوچک تر محاسبه می کنند و میانگین این مقادیر را به عنوان ارزیابی نهایی ارائه می کنند. رابطه این معیار به صورت زیر است که در آن M تعداد قاب ها و L طول هر قاب است.

$$SNR_{seg} = \frac{1}{M} \sum_{f=1}^M \log_{10} \left[\frac{\sum_{n=L_f}^{L_f+L-1} |s(n)|^2}{\sum_{n=L_f}^{L_f+L-1} [e(n) - \hat{e}(n)]^2} \right] \quad (2-2)$$

در برخی از موارد که قاب های سکوت نیز در سیگنال وجود دارند ممکن است این مقدار منفی شود. برای رفع این مشکل می توان بازه های سکوت را تشخیص داد و آنها را در محاسبه لحاظ نکرد. به علاوه، با در نظر گرفتن یک آستانه پایین و جایگزین کردن مقادیر کوچکتر از آن با مقدار آستانه، می توان از منفی شدن جلوگیری کرد. همچنین گوش انسان بین قاب هایی با مقادیر SNR بالاتر از ۳۵dB تمایز چندانی نمی گذارد، بنابراین یک آستانه بالا برای متعادل کردن SNR هایی که بسیار بالا هستند، به کار می رود. این دو آستانه که معمولاً ۰ dB و ۳۵ dB است، باعث می شوند تا جواب نهایی در بازه قابل قبولی باقی بمانند. توسعه رابطه SNR به حوزه فرکانس با نام SNR قطعه ای وزندار، به صورت زیر است که در آن K تعداد کل باندهای فرکانسی و W_k هر کدام از آنهاست. M

^۱ - frame

تعداد قاب ها، $F(f, k)$ و $\hat{F}(f, k)$ دامنه فیلتر بانک برای سیگنال تمیز و تخمینی در فرکانس k ام از قاب f است [۳۷].

$$SR_{FW-seg} = \frac{1}{M} \sum_{f=1}^M \frac{\sum_{k=1}^K W_k \log_2 \left[\frac{F^*(f, k)}{[F(f, k) - \hat{F}(f, k)]^2} \right]}{\sum_{k=1}^K W_k} \quad (3-2)$$

۲-۵-۲- معیار PESQ

برخی از معیارهای ارزیابی مورد استفاده در بهسازی گفتار از حوزه های دیگر پردازش گفتار، از جمله سنتز^۱ و کدینگ گفتار گرفته شده اند که روش PESQ یکی از آنهاست. در سال ۲۰۰۰ مسابقه ای توسط ITU^۲ برای معرفی یک معیار عینی جدید ترتیب داده شد. هدف اصلی از تعریف این معیار عینی جدید قابلیت اطمینان بالا در شبکه های مختلف و شرایط متفاوت بود و در این میان PESQ معرفی شد. به طور خلاصه، PESQ یک معیار عینی برای تخمین یا پیش بینی امتیاز یک آزمایش واقعی (معیار ذهنی MOS) است. در این آزمایش واقعی، کیفیت سیگنال از طریق گوش دادن افراد مختلف به سیگنال صوتی و امتیاز دادن به آن توسط افراد، سنجیده می شود. ضریب همبستگی بین این روش و روش های ذهنی ۹۳.۵٪ است. علی رغم این همبستگی بالا بر این نکته که این معیار نباید به جای روش های ذهنی به کار گرفته شود، تاکید شده است. (امتیاز این روش به وسیله پایگاه داده بزرگی از معیار های ذهنی کالیبره شده است.)

در این روش سیگنال بهسازی شده با سیگنال اصلی مقایسه می شود و نمره ای بین ۰.۵- تا ۴.۵ داده می شود. هر چقدر که این معیار بیشتر باشد به معنای کیفیت بالاتر سیگنال صوتی است. در اولین گام برای محاسبه امتیاز PESQ، ابتدا تاخیر بین دو سیگنال با مکانیزم هایی که در [ITU ۲۰۰۰] تشریح شده است، حذف می شود. سپس دو سیگنال با استفاده از یک مدل ادراکی با هم مقایسه می شوند. در این مقایسه سعی شده نمایشی که سیستم شنوایی انسان از سیگنال ایجاد می کند، مورد

^۱ - Speech Synthesis

^۲ - International Telecommunication Union

استفاده قرار گیرد که برای رسیدن به این هدف چند گام مانند تطبیق زمانی، تطبیق سطوح سیگنال ها با سطوح شنوایی، تبدیل زمان-فرکانس و سنجش شدت سیگنال در نظر گرفته شده است. به صورت خلاصه می توان گفت که الگوریتم PESQ به این ترتیب عمل می کند. ابتدا هر دو سیگنال (اصلی و اعوجاج یافته) به یک سطح استاندارد شنوایی نگاشته شده و با فیلتری میان گذر با فرکانس های قطع ۳۰۰ و ۳۰۰۰ هرتز، فیلتر می شوند. در مرحله بعد دو سیگنال از لحاظ زمانی همزمان شده و به این ترتیب تاخیرهای زمانی (که در کار کدینگ به خاطر عبور از شبکه معمول تر است) حذف می شود. پس از این مرحله، سیگنال ها توسط یک تبدیل شنوایی پردازش شده و میزان بلندی آنها به دست می آید. تفاوت بلندی دو سیگنال محاسبه شده و در طول زمان و فرکانس میانگین گیری می شوند تا تخمینی از معیار MOS به دست آید. از لحاظ پیچیدگی محاسباتی PESQ پیچیده ترین معیار عینی است [۳۷].

۲-۵-۳- معیار LLR

پس از بررسی دو معیار بالا به بررسی LLR (Log Likelihood Ratio) می پردازیم که به صورت رابطه (۲-۴) تعریف می شود:

$$LLR = \log \left(\frac{\sum_{n=1}^T R_{xx} a_n^2}{\sum_{n=1}^T R_{yy} a_n^2} \right) \quad (۲-۴)$$

که در آن a_n و a_n^2 به ترتیب ضرایب LPC^۱ سیگنال تمیز و سیگنال بهسازی شده و R_{xx} ماتریس همبستگی^۲ سیگنال تمیز است. معیار LLR برای همه پنجره های یک جمله حساب می شود. در واقع این مقدار برای هر قاب از سیگنال محاسبه شده و در نهایت از مقادیر به دست آمده برای قاب ها میانگین گرفته می شود [۳۱].

۲-۶- منابع و مشخصات نویزهای صوتی

^۱ - Linear Predictive Coding
^۲ - Autocorrelation

نویز یا صداهای مزاحم تقریباً در همه موارد همراه سیگنال صوت می‌باشند. زمانی که شما از خیابان عبور می‌کنید (صدای خودروها و ساخت و ساز و ...)، وقتی داخل خودرو نشسته اید و در حرکت هستید (صدای موتور ماشین، صدای باد و...) زمانی که به یک فروشگاه یا رستوران می‌روید (صدای همهمه دیگران)، وقتی که به دفتر محل کار خود می‌رسید (صدای کولر، زنگ تلفن و ...) و به طور خلاصه در همه حالات، شما با انواعی از نویزها، که هر کدام دارای ویژگی‌های خاصی هستند، سروکار دارید. به طور کلی نویزها را می‌توان به دو دسته ایستاد و غیر ایستاد تقسیم نمود. نویز ایستاد، نویزی است که مشخصات آماری آن با زمان تغییر نمی‌کند. مانند صدای کولر. در مقابل نویز غیر ایستاد، نویزی است که مشخصه‌های زمانی و طیفی آن متغیر باشد. مانند نویز فرودگاه که ترکیبی از نویز همهمه و سایر صداهاست.

اثر نویز برگفتار می‌تواند جمع شونده^۱ (مثل نویز سفید)، کانوالوشونده^۲ (مثل نویز کانال تلفن یا نویز پژواک) و یا ضرب شونده^۳ (مانند نویز محوسازی^۴) باشد. جنس نویز می‌تواند شبیه به گفتار باشد و با گفتار همبستگی داشته باشد، یا فقط یک نویز زمینه (مانند صدای کولر و پنکه) بوده و غیر همبسته با گفتار باشد.

با توجه به تنوع موجود در شکل و مشخصات نویزها، واضح است که قبل از هر کاری نیاز به شناسایی دقیق نویزها از نظر مشخصات زمانی و طیفی و درک تفاوت بین آنها هستیم. از ویژگی‌های مهمی که انواع نویز را از هم متمایز می‌کند شکل طیف و توزیع انرژی آنها درحوزه فرکانس است. به عنوان مثال انرژی نویز ناشی از باد بیشتر در فرکانس‌های پایین حدود ۵۰۰ هرتز متمرکز است. اما توزیع انرژی نویز رستوران، پهنای باند فرکانسی وسیع تری را پوشش می‌دهد.

یکی دیگر از عواملی که در طراحی سیستم بهسازی گفتار به آن نیاز داریم، شناسایی میزان SNR در کاربردها و شرایط مختلف است.

^۱ - additive
^۲ - convolutive
^۳ - Multiplicative
^۴ - Fading Noise

یک مطالعه جامع در هشت محیط مختلف نشان داده SNR در محیط های کم نویز مثل بیمارستان و کلاس بین ۵ تا ۱۵ دسی بل و برای محیط های دارای نویز زیاد مثل قطار حتی به ۰ dB هم می رسد. از این رو می توان گفت فاصله SNR موثر برای عملکرد سیستم های بهسازی گفتار ۱۵-۰ dB است. همچنین فاصله بین فرستنده (گوینده) و گیرنده (شنونده) یکی دیگر از عوامل موثر در کاهش نسبت سیگنال به نویز است. مثلا در کلاس درس کیفیت صدای معلم برای افراد ردیف جلو بالاتر از افراد ردیف های آخر است [۵۹].

فصل سوم

مروری بر روش های بهسازی
گفتارتک کاناله

۳-۱- مقدمه

بهسازی گفتار، به جهت اهمیت و کاربردهای مهم و فراوانی که دارد، یکی از شاخه های مورد توجه پردازش گفتار می باشد و تا کنون روش های متعددی برای آن پیشنهاد شده است. در این فصل مرور مختصری بر مهم ترین روش های بهسازی گفتار تک کاناله خواهیم داشت. روش های زیادی برای این مساله ارائه شده اند که می توان آنها را در دو گروه اصلی روش های آماری و غیرآماري تقسیم کرد. علی رغم وجود ماهیت تصادفی و آماری سیگنال، منظور از روش های غیرآماري، روش هایی مانند تفریق طیفی و زیر فضای سیگنال است که در آنها هیچ فرض خاصی در مورد شکل توزیع سیگنال یا طیف آن و اطلاعات پیشین مرتبط انجام نمی شود. در مقابل در روش های آماری فرض می شود که سیگنال یا طیف آن از توزیع خاصی پیروی می کنند. روش های آماری که خود به دو گونه آماری محض و مبتنی بر مدل تقسیم می شوند، در این فصل به طور خلاصه بررسی می شوند. در این فصل روش های تفریق طیفی، زیرفضای سیگنال، فیلتر کالمن و روش آماری، به عنوان مهم ترین و معروف ترین روش های بهسازی گفتار، بررسی می شود.

۳-۲- تفریق طیفی

تفریق طیفی از قدیمی ترین و رایج ترین الگوریتم های بهسازی گفتار است که به خاطر سادگی و پایین بودن محاسبات آن هنوز هم مورد توجه محققین قرار دارد. ایده این روش بسیار واضح و ساده است. اگر فرض شود که سیگنال گفتار تمیز $s(n)$ در اثر جمع شدن با نویز مستقل $d(n)$ به سیگنال نویزی $y(n)$ تبدیل شود (۳-۱)، آنگاه با اعمال تبدیل فوریه زمان کوتاه، رابطه (۳-۱)، به صورت رابطه (۳-۲) خواهد شد.

$$y(n) = s(n) + d(n) \quad (1-3)$$

$$Y(\omega_k) = S(\omega_k) + D(\omega_k) \Rightarrow Y_k e^{j\theta_y(\omega_k)} = S_k e^{j\theta_s(\omega_k)} + D_k e^{j\theta_d(\omega_k)} \quad (2-3)$$

که در آن $\omega_k = 2k\pi/N$ $k=0,1,2,\dots,N$ و N طول هر قاب سیگنال گفتار است. در این رابطه Y_k ، S_k و D_k به ترتیب دامنه طیف سیگنال نویزی، تمیز و نویز هستند و $\theta_y(k)$ ، $\theta_s(k)$ و $\theta_d(k)$ به ترتیب فاز سیگنال های نویزی، تمیز و نویز هستند.

در این روش برای تخمین دامنه طیف سیگنال تمیز از روی سیگنال نویزی می توان به سادگی دامنه طیف نویز را از دامنه طیف سیگنال نویزی پیدا کرد. یعنی:

$$\hat{S}_k = Y_k + D_k \quad (3-3)$$

نهایتاً در تخمین نهایی می توان از فاز سیگنال نویزی به عنوان فاز سیگنال تمیز استفاده کرد.

$$\hat{S}(\omega_k) = \hat{S}_k e^{j\theta_s(k)} = [Y_k + D_k] e^{j\theta_y(k)} \quad (4-3)$$

همانگونه که از رابطه ی (۴-۳) مشخص است برای پیاده سازی به تخمینی از دامنه ی طیف نویز نیاز داریم. روشهای مختلفی برای این کار وجود دارند که در تفریق طیفی معمولاً از VAD^۱ یا آشکار کننده ی گفتار از غیرگفتار استفاده می شود. در این روش، فرض اساسی ایستادن بودن موضعی نویز در فاصله میان دو سکوت ضروری است. زیرا فرض می شود مشخصات نویز در فاصله دو سکوت ثابت بوده و تخمین طیف نویز در این فاصله به روز نمی شود.

با وجود سادگی روش تفریق طیفی عمل تفریق باید با دقت انجام شود تا باعث ایجاد اعوجاج در

سیگنال تخمینی نشود. در شرایطی به علت تخمین بیش از حد نویز ممکن است مقداری از اطلاعات سیگنال گفتار حذف شود. از طرف دیگر با تخمین کم نویز سیگنال بهسازی شده دارای نویز خواهد بود. یکی از مشکلاتی که معمولاً در عمل در تفریق به وجود می آید منفی شدن برخی از مولفه های طیف تخمینی است که این معمولاً به علت تخمین نامناسب طیف نویز ایجاد می شود. معمول ترین راه حل رفع این مشکل استفاده از یکسوسازهای نیم موج است که در آن مانند رابطه زیر مقادیر منفی با صفر یا یک مقدار مثبت کوچک جایگزین می گردد.

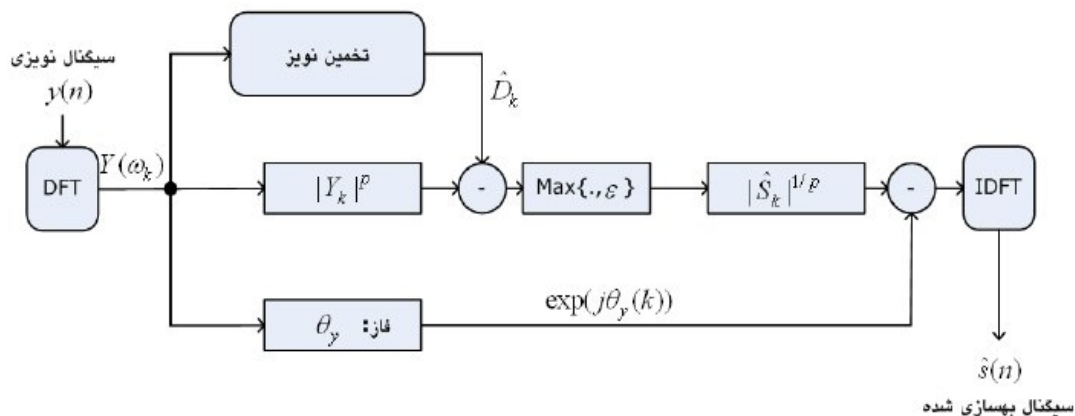
$$\text{اگر: } |Y(\omega)| > |D(\omega)|$$

^۱ - Voice Activity Detection

$$|\hat{S}(\omega)| = |Y(\omega)| - |\hat{D}(\omega)| \quad \text{آنگاه:}$$

$$|\hat{S}_i(\omega)| = * \quad \text{در غیر این صورت:}$$

شکل (۱-۳) نمای کلی الگوریتم تفریق طیفی را نشان می دهد. همانگونه که اشاره شد روش های خانواده تفریق طیفی دارای مزایای سادگی و محاسبات پایین هستند. اما در مقابل دارای کاستی هایی نیز هستند. عمده ترین عیب این روش ایجاد اعوجاج در سیگنال نهایی و تولید نویزی موسوم به نویز موزیکال^۱ یا موسیقی است [۶] و [۴]. بسیاری از بهبودهای پیشنهاد شده روی این خانواده از روش ها برای حذف این اعوجاج ها و نویزهای ناخواسته بوده است. دلیل اصلی ایجاد این نویزها را می توان در منفی شدن طیف تخمینی که خود نتیجه ی تخمین نادرست طیف نویز است جستجو کرد. روش یکسوسازی نیم موج به علت غیر خطی بودن باعث ایجاد قله های کوچک و جدا از هم که توزیع آنها نیز تصادفی است در طیف سیگنال و در فرکانس های متفاوت می شود. این صداهای اضافی همان نویزهای موسیقی هستند که اثر آنها بر شنونده گاه از نویز اولیه هم آزردهنده تر است. یکی دیگر از کاستی های این روش که در مقابل نویز موسیقی از اهمیت کمتری برخوردار



شکل (۱-۳): نمای کلی الگوریتم تفریق طیفی [۵۹].

جدول (۱-۳): مروری بر روش تفریق طیفی.

^۱ - Musical Noise

فرضیات اساسی	نویز ایستان یا ایستان زمان کوتاه، مستقل و جمع شونده با گفتار گفتار ایستان زمان کوتاه
مزایا	ساده و دارای محاسبات پایین
معایب	ایجاد نویز زمینه (نویز موسیقی) در سیگنال بهسازی شده

است، استفاده از فاز سیگنال نویزی به عنوان فاز سیگنال بهسازی شده است. این جایگزینی در برخی از موارد می تواند باعث تغییرات زیاد و مداوم در کیفیت سیگنال شود. علت این جایگزینی کم اثر بودن فاز سیگنال نویزی بر روی کیفیت سیگنال گفتار و عدم حساسیت زیاد سیستم شنوایی انسان به آن است [۴۸]. اما فقط در صورت بالا بودن SNR سیگنال گفتار (بالاتر از ۸dB) می توان فاز سیگنال نویزی را بدون از دست دادن کیفیت به جای فاز سیگنال تمیز استفاده کرد [۵۸]. جدول (۳-۱) اطلاعات کلی راجع به روش تفریق طیفی ارائه می دهد.

۳-۳- زیر فضای سیگنال

اغلب روش های بهسازی براساس تئوری پردازش سیگنال و تخمین آماری هستند. اما روش های زیر فضای سیگنال مبتنی بر تئوری جبرخطی هستند. در واقع این روش بر این اصل استوار است که می توان سیگنال نویزی را به زیر فضاهایی تجزیه کرد که در آنها بخش های سیگنال تمیز و نویز متمایز می شوند و با صفر کردن بخش های نویزی در این زیر فضا می توان تخمینی از سیگنال تمیز بدست آورد. تجزیه سیگنال نویزی به دو زیر فضای سیگنال تمیز و نویز با استفاده از روش های شناخته شده تجزیه ماتریس های متعامد در جبرخطی مثل تجزیه مقادیر تکین^۱ (SVD) یا تجزیه براساس بردارهای ویژه^۲ و مقادیر ویژه^۳ امکان پذیر است. روش تجزیه بر اساس مقادیر ویژه و

^۱ - Singular Value Decomposition

^۲ - Eigen Vector

^۳ - Eigen Value

بردارهای ویژه به اسم تبدیل کارهونن- لاف^۱، تبدیل هتلینگ^۲ و تحلیل اجزای اصلی^۳ نیز معروف است [۵۹]. جزئیات بیشتر درباره این روش در [۵۹] آمده است. مزایا و معایب این روش به طور خلاصه در جدول (۲-۳) آورده شده است.

جدول (۲-۳): مروری بر روش زیر فضای سیگنال .

فرضیات اساسی	نویز و گفتار مستقل بوده و اطلاعات آماری درجه دوم آنها در دسترس است
مزایا	ساده و دارای محاسبات پایین
معایب	عدم کارایی بالا

۳-۴- فیلتر کالمن

فیلتر کالمن را می توان از دسته روش های پارامتری به حساب آورد، که در این بخش به طور خلاصه مرور می گردد. در این روش با استفاده از مدل AR^۴، اطلاعات سیگنال گفتار استخراج شده و مورد استفاده قرار می گیرد.

اولین استفاده از فیلتر کالمن برای بهسازی گفتار را می توان در مقاله paliwal در سال ۱۹۸۷ یافت [۴۷] که در آن فیلتر برای حذف نویز سفیدی که با گفتار جمع شده بود به کار رفت. پس از آن نیز روش های متفاوتی جهت بهبود و تعمیم فیلتر کالمن براساس ایده اصلی آن ارائه شد، که در اینجا فقط به بیان ایده اصلی اکتفا می کنیم.

فیلتر کالمن در واقع روشی برای کاهش مربعات خطا به صورت تطبیقی است که با استفاده از یک روش بازگشتی برای جداسازی سیگنال و نویز مورد استفاده قرار می گیرد. این الگوریتم از معادلات فضای حالت^۵ برای مدل کردن فرآیند تولید سیگنال گفتار و از معادلات مشاهده برای مدل کردن

^۱ - Karhunen- Loeve Transform(KLT)

^۲ - Hotelling Transform

^۳ -Principal Component Analysis (PCA)

^۴ - Auto Regressive

^۵ - State Space

سیگنال آغشته به نویز استفاده می کند. برای این کار ابتدا یک مدل AR از مرتبه P به صورت رابطه (۵-۳) برای بردار گفتار تمیز در نظر می گیریم که در آن $\{\alpha_i, i = 1, 2, \dots, P\}$ ضرایب پیش بینی خطی و $w(n)$ نویز سفید گوسی با میانگین صفر و واریانس σ_{ww}^2 است.

$$s(n) = \sum_{i=1}^P \alpha_i s(n-i) + w(n) \quad (5-3)$$

اگر فرض کنیم:

$$\mathbf{S}(n) \triangleq [s(n), s(n-1) \dots s(n-p+1)]^T \quad (6-3)$$

رابطه (۵-۳) را می توان به صورت زیر نوشت :

$$\mathbf{S}(n) = \mathbf{F}\mathbf{S}(n-1) + \mathbf{g}w(n) \quad (7-3)$$

که در آن:

$$\mathbf{F} = \begin{bmatrix} \alpha_1 & \alpha_2 & \dots & \alpha_{p-1} & \alpha_p \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix}$$

$$\mathbf{g} = [0, 0, \dots, 0, 1]^T \quad (8-3)$$

همانگونه که قبلا بیان شد، داریم:

$$\mathbf{y}(n) = \mathbf{s}(n) + \mathbf{d}(n) \quad (9-3)$$

که در آن $\mathbf{y}(n)$ مشاهدات یا همان سیگنال نویزی است. رابطه فوق را می توان به صورت زیر نوشت:

$$\mathbf{y}(n) = \mathbf{h}^T \mathbf{S}(n) + \mathbf{d}(n) \quad (10-3)$$

که نمایش فضای حالت سیگنال نویزی است و در آن $\mathbf{h}=\mathbf{g}$ می باشد. در این معادلات $\mathbf{S}(n)$ بردار $1 \times p$ بعدی سیگنال تمیز، معادل بردار حالت در لحظه n می باشد. \mathbf{F} ماتریس $p \times p$ بعدی است که آن را ماتریس انتقال می نامند و واسط بین حالت n و $n-1$ است. همچنین $\mathbf{d}(n)$ نیز بردار $1 \times p$ بعدی نویز گوسی با میانگین صفر است [۴۷].

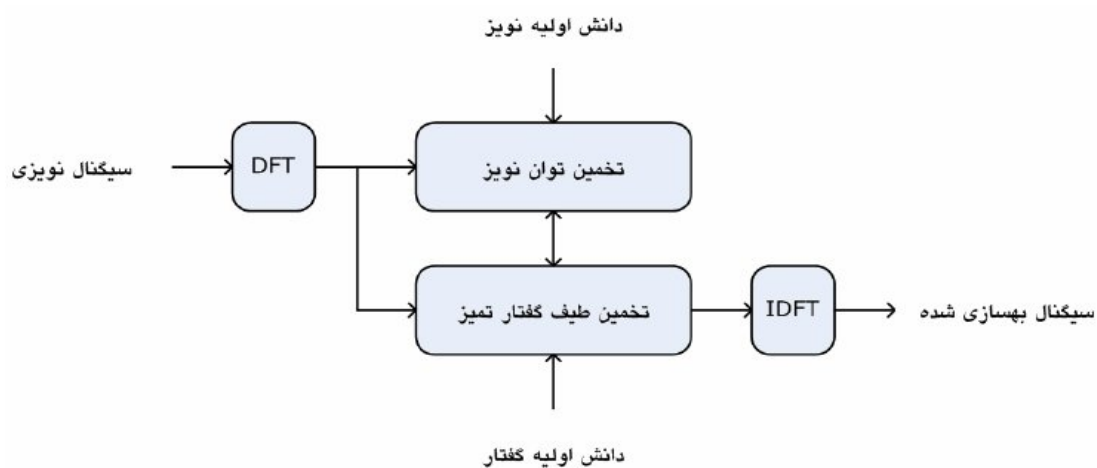
۳-۵- روش های آماری

روش های آماری به تکنیک هایی گفته می شود که در آنها یک مدل آماری برای سیگنال گفتار تمیز و نویز در نظر گرفته می شود و پردازش براساس این مدل ها صورت می گیرد. این روش به دو دسته روش های آماری محض و روش های مبتنی بر مدل تقسیم می شود: روش های آماری محض که در آن یک مدل آماری خاص (مثلا توزیع گاوسی) برای ضرایب فوریه یا پارامترهای دیگر هر قاب سیگنال گفتار و نویز فرض می شود و سپس روابط بهسازی گفتار متناسب با این مدل فرض شده به دست آید. این روش ها منجر به روابط و فیلترهایی می شوند که در سطح قاب عمل کرده و اطلاعات سیگنال گفتار و نویز را در همین سطح مدل سازی می کند. از این رو در برخی منابع به این روش ها، روش های مبتنی بر مدل کوتاه مدت یا روش های آماری بدون آموزش گفته می شود.

دسته دیگر از روش های آماری روش های مبتنی بر مدل یا روش های long term یا روش آموزش یافته هستند. در این روش ها به کمک داده های آموزشی مدل یا مدل هایی ایجاد می شود که خصوصیات مورد نظر را (که معمولا تابع توزیع احتمال و پارامترهای آن است) در خود دارد. در این روش با استفاده از داده های آموزشی، خصوصیات بلند مدت و مشترک بین قاب های مختلف فرا گرفته می شود.

بهسازی گفتار با روش های آماری محض توسط [۳۵] Lim and Oppenheim و McAullay and [۴۳] Malpass و مقاله معروف [۲۱] Ephraim and Malah شروع شده و امروزه به عنوان یکی از مهم ترین زمینه های علمی و کاربردی در بهسازی گفتار مورد توجه است. در خانواده روش های آماری مبتنی بر مدل و به طور خاص HMM نیز می توان به کارهای Ephraim [۱۷] و [۲۳] و [۱۹] اشاره کرد که بعد از آن استفاده از این رویکرد در بهسازی گفتار گسترش زیادی داشته است. شکل (۲-۳) نمای کلی روش های آماری را نمایش می دهد. همانطور که از شکل نیز مشخص است، در روش های آماری، علاوه بر توزیع های آماری که بخش جدانشدنی الگوریتم هستند، نیاز به یک

معیار اعوجاج یا تخمین گر نیز وجود دارد که بر اساس کمینه کردن اعوجاج یا بهینه کردن رابطه مربوط به تخمین گر، روابط تخمین سیگنال تمیز به دست می آید. از معروفترین معیارهای مورد استفاده در روش های آماری، می توان به تخمین گرهای کمترین میانگین مربعات خطا (MMSE)، بیشترین شباهت (ML) و بیشترین احتمال پسین (MAP) اشاره کرد که روابط مربوط به آنها در فصل بعد خواهد آمد.



شکل (۲-۳) نمای کلی روش های آماری [۵۹].

امروزه به جرات می توان ادعا کرد که روش های آماری چه از لحاظ کیفی و چه از لحاظ کمی نتایج بهتری نسبت به اغلب روش های دیگر ارائه می دهند. جدول (۳-۳) خلاصه روش های آماری را بیان می کند. در فصل بعد توضیحات بیشتری درباره این روش بیان خواهد شد.

جدول (۳-۳): خلاصه روش های آماری.

فرضیات اساسی	نویز و گفتار مستقل وجود دانش قبلی از سیگنال تمیز و نویز مثلا نوع توزیع آنها
مزایا	کارایی بالاتر نسبت به سایر روش ها ایجاد نویز موسیقی کم اعوجاج کمتر در سیگنال
معایب	محاسبات و پیچیدگی زیاد

فصل چهارم

بهبودی گفتار با روش های آماری

۴-۱- مقدمه

از مدل های آماری در کاربردهای مختلف پردازش گفتار مانند بازشناسی گفتار، بهسازی گفتار و کدینگ به صورت موفقیت آمیزی استفاده شده است. رویکرد آماری در حل مساله بهسازی گفتار را می توان جزو روش های پارامتری به حساب آورد که در آن مدل های آماری، که معمولاً توابع چگالی احتمال هستند، مجموعه ای از پارامترها را تشکیل می دهند. با استخراج توابع چگالی احتمال، رفتار تصادفی سیگنال مدل شده و می توان در کنار معیارهای اعوجاج، سیستم های بهینه برای بهسازی گفتار را به دست آورد. راه حل رایج برای مدل سازی آماری این است که سیگنال گفتار یک فرآیند تصادفی فرض شود که به صورت محلی و در سطح قاب های ۲۰ تا ۴۰ میلی ثانیه ای ایستاد هستند. بر این اساس گفتار باید قاب به قاب پردازش شود. این روش ها به دو دسته روش های آماری محض و آماری مبتنی بر مدل تقسیم می شوند که در این فصل بررسی خواهند شد. همچنین تخمین گرهای بیشترین شباهت و روش بیز نیز مرور شده و روش های بهسازی مرتبط با آنها معرفی گردیده است.

۴-۲- روش های آماری محض

همان طور که قبلاً گفته شد، در روش های آماری محض یک مدل آماری خاص (مثلاً توزیع گوسی) برای ضرایب فوریه یا پارامترهای هر قاب سیگنال گفتار و نویز فرض می شود و سپس روابط بهسازی گفتار متناسب با این مدل فرض شده به دست می آید. هرچند می توان مدل های آماری را در حوزه زمان و در سطح نمونه های سیگنال نیز بدست آورد، اما این مدل ها معمولاً در حوزه تبدیل فوریه فرض می شوند. در این روش ها ابتدا توزیع های آماری خاصی مثلاً گوسی برای سیگنال بدون نویز و نویز در نظر گرفته می شود. سپس با استفاده از یک معیار تخمین گر، روابط تخمین سیگنال تمیز

بدست می آید.

قبلا گفته شد که یکی از اجزای روش های بهسازی گفتار آماری، تخمین گرها هستند. مهم ترین تخمین گرهای مورد استفاده در این حوزه، تخمین گر بیشترین شباهت، کمترین میانگین مربعات خطا و بیشترین احتمال پسین هستند که به طور مفصل توضیح داده خواهند شد.

مهم ترین روش های آماری محض عبارت اند از:

- فیلتر وینر با فرض گوسی بودن سیگنال تمیز و نویز و تخمین گر MMSE^۱ [۲] و [۶۰]
- تخمین دامنه طیف با MMSE [۲۲]
- تخمین لگاریتم دامنه طیف و توان دوم و P دامنه طیف سیگنال تمیز بر اساس تخمین MMSE [۲۳] و [۱] و [۶۳]
- تخمین دامنه ی طیف با تخمین گر ML [۴۵]
- تخمین طیف سیگنال تمیز به روش MAP^۲ [۶۲] و [۳۸] و [۶۱]

۴-۳- تخمین گر بیشترین شباهت

بیشترین شباهت یکی از پرکاربردترین روش های تخمین آماری است که در کاربردهای مختلفی برای تخمین پارامترهای مورد نظر به کار می رود. در این روش فرض می شود که L نمونه داده از $y(n)$ داریم که به پارامتر نامعلوم δ وابسته اند و تابع چگالی احتمال y که با $p(y; \delta)$ نشان داده می شود در دسترس باشد. هدف این تخمین یافتن مقدار پارامتر δ به گونه ای است، که تابع چگالی احتمال y بیشینه شود:

$$\delta_{ML} = \operatorname{argmax} p(y; \delta) \quad (1-4)$$

ابتدا فرض می شود متغیر نامعلوم معین^۳ و غیر تصادفی است. این مساله تخمین بیشترین شباهت را از تخمین بیزی که در آن فرض می شود پارامتر نامعلوم یک متغیر تصادفی^۱ است، متفاوت می کند.

^۱ - Minimum Mean Square Error

^۲ - Maximum A Posteriori

^۳ - Deterministic

برای حل مساله فوق ابتدا یک شکل توزیع خاص مثلا گوسی برای تابع چگالی احتمال که به آن تابع شباهت^۲ نیز گفته می شود، فرض می شود و از آن نسبت به پارامتر نامعلوم مشتق گرفته مساوی صفر قرار داده می شود. می توان برای ساده تر شدن محاسبات از لگاریتم تابع شباهت^۳ استفاده کرد. مشابه قبل، ابتدا فرض می شود که سیگنال گفتار تمیز $s(n)$ در اثر جمع شدن با نویز $d(n)$ سیگنال نویزی $y(n)$ را طبق رابطه (۲-۴) تولید کرده است. نمایش این رابطه در حوزه فرکانس، به صورت رابطه (۳-۴) می باشد که در آن $\omega_k = 2k\pi/N$ $k=0,1,2,\dots,N$ و N طول هر قاب گفتار است. در این رابطه D_k و S_k ، Y_k به ترتیب دامنه طیف سیگنال نویزی، تمیز و نویز و $\theta_{d(k)}$ و $\theta_{s(k)}$ ، $\theta_{y(k)}$ به ترتیب فاز سیگنال های نویزی، تمیز و نویز هستند.

تخمین گر بیشترین شباهت اولین بار توسط McAulay & Malpass در بهسازی گفتار به کار گرفته شد [۴۵]. آنها فرض کردند مشاهدات، سیگنال نویزی باشد. همچنین، دامنه و فاز سیگنال تمیز، پارامترهای نامعلوم و معین بوده و ضرایب فوریه نویز، گوسی با میانگین صفر و واریانس (در هر دو بخش حقیقی و موهومی) $\sigma_d(k)/2$ باشند. بر اساس این فرضیات، تابع توزیع ضرایب فوریه سیگنال نویزی $Y(\omega_k)$ ، نیز گوسی، با میانگین $s_k e^{j\theta_s(k)}$ و واریانس $\sigma_d(k)$ خواهد بود. بنابراین رابطه (۴-۴) را برای تابع احتمال مورد نظر می توان نوشت.

$$y(n) = s(n) + d(n) \quad (2-4)$$

$$Y(\omega_k) = S(\omega_k) + D(\omega_k) \Rightarrow Y_k e^{j\theta_y(k)} = S_k e^{j\theta_s(k)} + D_k e^{j\theta_d(k)} \quad (3-4)$$

$$p\left(Y(\omega_k); s_k, \theta_s(k)\right) = \frac{1}{\pi\sigma_d(k)} \exp\left[-\frac{|Y(\omega_k) - s_k e^{j\theta_s(k)}|^2}{\sigma_d(k)}\right]$$

$$= \frac{1}{\pi\sigma_d(k)} \exp\left[-\frac{Y_k^2 - 2S_k \operatorname{Re}\{Y(\omega_k) e^{-j\theta_s(k)}\} + S_k^2}{\sigma_d(k)}\right] \quad (4-4)$$

بیشینه کردن این تابع که از دو متغیر نامعلوم دامنه و فاز تشکیل شده است، کار پیچیده ای است. از این رو ابتدا تابع احتمال فوق را با فرض یکنواخت^۴ برای فاز و استفاده از رابطه انتگرالی (۵-۴)

^۱ -Random

^۲ -Likelihood

^۳ -Likelihood Function

^۴ - Uniform Distribution

تبدیل به تابعی با یک متغیر نامعلوم به صورت رابطه (۴-۶) می نمایند.

$$p_m(Y(\omega_k); s_k) = \int_{-\pi}^{\pi} p(Y(\omega_k); s_k, \theta_s(k)) p(\theta_s(k)) d\theta_s(k) \quad (۵-۴)$$

$$p_m(Y(\omega_k); s_k) = \frac{1}{\pi \sigma_d(k)} \exp\left[-\frac{Y_k^* + X_k^*}{\sigma_d(k)}\right] \int_{-\pi}^{\pi} \exp\left[\frac{Y_k^* \operatorname{Re}\{Y(\omega_k) e^{-j\theta_s(k)}\}}{\sigma_d(k)}\right] d\theta_s(k) \quad (۶-۴)$$

در این رابطه $\operatorname{Re}(\cdot)$ بیانگر بخش حقیقی طیف است و انتگرال فوق تابع بسط تغییر یافته^۱ نوع اول است که می توان آن را با تابع نمایی تقریب زد. حال از این تابع که فقط بر حسب متغیر s_k است، مشتق گرفته و برابر صفر قرار داده می شود تا تخمین بیشترین شباهت دامنه به صورت رابطه (۴-۷) به دست آید. سپس با قرار دادن مقدار فاز سیگنال نویزی به جای فاز سیگنال تمیز؛ تخمین طیف سیگنال تمیز، مطابق رابطه (۴-۸) به دست می آید.

$$\hat{s}_k = \frac{1}{\pi} \left[Y_k + \sqrt{Y_k^* - \sigma_d(k)} \right] \quad (۷-۴)$$

$$\hat{S}(\omega_k) = \hat{s}_k e^{j\theta_y(k)} = \hat{s}_k \frac{Y(\omega_k)}{Y_k} = \left[\frac{1}{\pi} + \frac{1}{\pi} \sqrt{\frac{Y_k^* - \lambda_d(k)}{Y_k^*}} \right] Y(\omega_k)$$

$$= \left[\frac{1}{\pi} + \frac{1}{\pi} \sqrt{\frac{Y_k^* - 1}{Y_k^*}} \right] Y(\omega_k) = G_{ML}(Y_k) Y(\omega_k) \quad (۸-۴)$$

در این روابط $Y_k = Y_k^* / \sigma_d(k)$ سیگنال به نویز پسین^۲ و $G_{ML}(Y_k)$ را اصطلاحاً تابع بهره^۳ تخمین ML می نامند که تابعی از سیگنال به نویز پسین است. سیگنال به نویز پسین Y_k معادل SNR مشاهده شده جز k ام طیف است، زمانی که سیگنال نویز به سیگنال افزوده شده باشد.

۴-۴- تخمین گر بیز

در تخمین گر بیشترین شباهت فرض شد که پارامتر نامعلوم مورد نظر δ ، یک متغیر معین است. در حالیکه در تخمین گر بیز، فرض بر این است که پارامتر نامعلوم یک متغیر تصادفی است. از آنجا که در فرآیند تخمین از قانون بیز استفاده می شود، آن را تخمین گر بیز می نامند. در این روش،

^۱ - Modified Bessel Function
^۲ - A Posteriori SNR
^۳ - Gain Function

همچنین فرض می شود که یک دانش اولیه (پیشین)^۱ مانند توزیع متغیر نامعلوم، در دسترس است. با توجه به توانایی این تخمین گر در استفاده از این دانش، در کاربردهایی که تخمین مناسبی از توزیع وجود داشته باشد، این روش کارایی بهتری نسبت به روش بیشترین شباهت دارد. این تخمین گر شامل روشهای کمترین مربعات خطا (MMSE) و بیشترین احتمال پسین (MAP) است [۳]. علاوه بر روش کمترین مربعات خطا چند روش برگرفته از آن مانند Log-MMSE یا MMSE با توان های بیش از یک نیز وجود دارند.

۴-۵- تخمین گر MMSE

فرض کنید که \hat{X} تخمینی از متغیر X باشد که این تخمین بر اساس متغیردیگری مانند Y که مقدار آن مشخص است (مشاهدات)، به دست آمده باشد. میانگین مربعات خطا (MSE)، همانگونه که از اسمش پیداست، برای متغیر X و تخمین آن \hat{X} ، به صورت رابطه (۴-۹) تعریف می شود. در این رابطه $p(Y, X)$ توزیع توأم X و Y و $p(X|Y)$ تابع چگالی احتمال X بعد از مشاهده Y است.

$$\begin{aligned} e &= E\{|\hat{X} - X|^2\} = \iint (\hat{X} - X)^2 p(Y, X) dY dx \\ &= \int \left(\int (\hat{X} - X)^2 p(X|Y) dx \right) p(Y) dY \end{aligned} \quad (۴-۹)$$

هدف تخمین گر کم ترین میانگین مربعات خطا (MMSE) کمینه کردن رابطه (۴-۹) است. برای انجام این کار کافیسست که عبارت انتگرال داخلی کمینه شود. با مشتق گیری از این رابطه نسبت به \hat{X} و برابر صفر قرار دادن آن، رابطه (۴-۱۰) به عنوان تخمین MMSE نتیجه می شود [۴۸].

$$\min_{\hat{X}} e = E\{|\hat{X} - X|^2\} \Rightarrow \hat{X} = E(X|Y) \quad (۴-۱۰)$$

در بهسازی گفتار متغیر شناخته شده (مشاهدات)، طیف گفتار نویزی $Y = [Y(\omega_1), \dots, Y(\omega_{N-1})]$ و پارامتر نامعلوم همان طیف گفتار تمیز $s(\omega_k)$ است. بنابراین تخمین بهینه اندازه ی طیف و تخمین بهینه طیف سیگنال تمیز با MMSE به ترتیب به صورت روابط (۴-۱۱) و (۴-۱۲) خواهد بود.

^۱ - A Priori

^۲ - Joint Distribution

$$\hat{s}_k = E[s_k|Y] = E[s_k|Y(\omega_0), Y(\omega_1), \dots, Y(\omega_{N-1})] \quad (11-4)$$

$$\hat{s}(\omega_k) = E[s_k|Y] = E[s(\omega_k)|Y(\omega_0), Y(\omega_1), \dots, Y(\omega_{N-1})] \quad (12-4)$$

۴-۵-۱- فیلتر وینر

فیلتر وینر یکی از روش های شناخته شده در زمینه فیلتر کردن آماری سیگنال است و در زمینه های مختلفی مانند پردازش سیگنال دیجیتال، پیش بینی خطی و ... کاربرد دارد. این روش که ابتدا توسط نوربرت وینر (۱۹۴۹) فرموله شده است، به حل مساله زیر می پردازد.

فرض کنید که یک سیگنال ورودی مانند $s(n)$ از یک سیستم خطی نامتغیر با زمان عبور کرده و سیگنال خروجی $y(n)$ را تولید کند. هدف این است که فیلتر خطی نامتغیر با زمان h به گونه ای طراحی شود که معکوس سیستم فوق را شبیه سازی کند و با اعمال آن روی $y(n)$ ، سیگنال $\hat{s}(n)$

را که تخمینی از $s(n)$ است، نتیجه دهد. یک راه حل ممکن برای این مساله این است که خطای بین $s(n)$ و $\hat{s}(n)$ یعنی $e(n) = s(n) - \hat{s}(n)$ را در نظر گرفته شده و سعی کنیم آن را کمینه نماییم.

این روش ابتدا توسط وینر برای حل این مساله در حالت پیوسته، مورد استفاده قرار گرفته است.

فیلتر می تواند در حوزه زمان (ضرایب فیلتر) و یا حوزه فرکانس به دست آید. در هر دو حوزه میانگین مربعات خطا محاسبه شده و با مشتق گیری از آن و برابر صفر قرار دادن، کمینه می شود.

یکی از کاربرد های فیلتر وینر، در بهسازی گفتار است. در این کاربرد سیگنال تمیز $s(n)$ با عبور از

یک سیستم توسط نویز جمع شونده $d(n)$ به سیگنال نویزی $y(n)$ تبدیل می شود. فیلتر وینر h در

حوزه زمان، به صورت $\hat{s}(n) = h \cdot y$ روی سیگنال نویزی اعمال شده و تخمین سیگنال تمیز را نتیجه

می دهد. در اینجا $h = [h_0, h_1, \dots, h_{M-1}]$ بردار ضرایب فیلتر و $y = [y_0, y_1, \dots, y_{M-1}]$ بردار

سیگنال نویزی برای M نمونه پشت سر هم است. برای استخراج ضرایب h_k این فیلتر، در حوزه

زمان از $J = E\{e^2(n)\}$ نسبت به h_k مشتق گرفته و برابر صفر قرار داده می شود، یعنی:

$$\frac{\partial J}{\partial h_k} = \frac{\partial E\{e^2(n)\}}{\partial h_k} = \frac{\partial E\{[s(n) - h \cdot y]^2\}}{\partial h_k}$$

$$= \frac{\partial (E\{s^*(n)\} - s^*(n) - \nabla h \cdot E\{y \cdot s(n)\} + h \cdot E\{y \cdot y^T\} h^T)}{\partial h_k}$$

$$= -\nabla E\{e(n)y(n-k)\} = 0, k = 0, 1, 2, \dots, M-1 \quad (13-4)$$

در این رابطه $E\{y \cdot s(n)\} = E\{[y_0, y_1, \dots, y_{M-1}] \cdot s(n)\}$ که یک بردار $M \times 1$ است، همبستگی بین سیگنال تمیز و سیگنال نویزی را نشان می دهد و با r_{ys} نشان داده می شود. همچنین $E\{y \cdot y^T\}$ یک ماتریس $M \times M$ است که ماتریس خودهمبستگی سیگنال نویزی است و با R_{yy} نشان داده می شود. رابطه فوق شرط لازم و کافی برای کمینه شدن تابع هزینه $J = E\{e^*(n)\}$ است. به صورت مشابه، مشتق گیری از شکل برداری- ماتریسی رابطه فوق منجر به دستگاه معادلات زیر (14-4) می گردد که شامل M معادله و M مجهول (همان ضرایب فیلتر) است. حل این معادلات که به معادلات Wiener&Hopf معروف اند، ضرایب بهینه فیلتر را نتیجه می دهد.

$$\frac{\partial J}{\partial h} = -\nabla r_{ys} + \nabla h \cdot R_{rr} = 0 \Rightarrow r_{ys} = \nabla h \cdot R_{rr} \Rightarrow h = R_{yy}^{-1} r_{ys}$$

$$\sum_{k=0}^M h_k r_{yy}(m-k) = r_{yd}(-m), m = 0, 1, \dots, M-1 \quad (14-4)$$

فیلتر وینر را می توان در حوزه فرکانس نیز استخراج کرد. در بهسازی گفتار استفاده از این حوزه رایج تر بوده و عمده مطالعات تکمیلی نیز در این حوزه بوده است. می دانیم که عبور سیگنال $y(n)$ از فیلتر خطی نامتغیر با زمان $h(n)$ منجر به خروجی $\hat{s}(n) = h(n) * y(n)$ می شود و در حوزه فرکانس داریم: $\hat{S}(\omega_k) = H(\omega_k)Y(\omega_k)$ و $E(\omega_k) = S(\omega_k) - \hat{S}(\omega_k)$. در این حالت کمینه کردن میانگین مربعات خطا منجر به رابطه (15-4) می شود.

$$\frac{\partial (E\{|S(\omega_k) - H(\omega_k)Y(\omega_k)|^2\})}{\partial H(\omega_k)} = H(\omega_k)p_{yy}(\omega_k) - p_{ys}(\omega_k) = 0$$

$$\Rightarrow H(\omega_k) = \frac{p_{ys}(\omega_k)}{p_{yy}(\omega_k)} \quad (15-4)$$

که در آن:

$$p_{yy}(\omega_k) = E\{|Y(\omega_k)|^2\} = E\{[S(\omega_k) + D(\omega_k)][S(\omega_k) + D(\omega_k)]^*\}$$

$$\begin{aligned}
&= E\{S(\omega_k) + S^*(\omega_k)\} + E\{S(\omega_k) + D^*(\omega_k)\} \\
&+ E\{S^*(\omega_k) + D(\omega_k)\} + E\{D(\omega_k) + D^*(\omega_k)\} \\
&= p_{ss}(\omega_k) + p_{dd}(\omega_k)
\end{aligned} \tag{۱۶-۴}$$

9

$$\begin{aligned}
p_{ys}(\omega_k) &= E\{[Y(\omega_k)S^*(\omega_k)]\} = E\{[S(\omega_k) + D(\omega_k)]S^*(\omega_k)\} \\
&= E\{S(\omega_k)S^*(\omega_k)\} + E\{S^*(\omega_k)D(\omega_k)\} = p_{ss}(\omega_k)
\end{aligned} \tag{۱۷-۴}$$

همچنین با فرض مستقل بودن نویز و سیگنال تمیز داریم:

$$E\{S(\omega_k)D^*(\omega_k)\} = E\{S^*(\omega_k)D(\omega_k)\} = 0 \tag{۱۸-۴}$$

در نتیجه فیلتر وینر نهایی به صورت زیر خواهد بود که در آن ξ_k را سیگنال به نویز پیشین^۱ برای فرکانس K ام می گویند که به صورت رابطه (۴-۲۰) تعریف می شود. سیگنال به نویز پیشین ξ_k ، مقدار SNR واقعی جز k ام طیف است که بعداً روش های محاسبه آن توضیح داده خواهد شد.

$$H(\omega_k) = \frac{p_{ys}(\omega_k)}{p_{yy}(\omega_k)} = \frac{p_{ss}(\omega_k)}{p_{ss}(\omega_k) + p_{dd}(\omega_k)} = \frac{\xi_k}{\xi_k + 1} \tag{۱۹-۴}$$

$$\xi_k = p_{ss}(\omega_k) / p_{nn}(\omega_k) = E\{|S(\omega_k)|^2\} / E\{|D(\omega_k)|^2\} \tag{۲۰-۴}$$

۴-۵-۲- تخمین اندازه طیف (STSA-MMSE)

در بخش قبل، روش وینر به عنوان یک تخمین گر خطی بهینه برای تخمین طیف مختلط سیگنال مرور شد. در این بخش یک تخمین بهینه برای دامنه طیف براساس معیار MMSE به دست می آید که به آن تخمین کمترین میانگین مربعات خطای دامنه طیف زمان کوتاه (STSA-MMSE) نیز می گویند. در این روش فرض می شود که ضرایب فوریه ی مختلط گفتار دارای توزیع گوسی با میانگین صفر و واریانس متغیر هستند. به علاوه فرض شده این ضرایب مستقل و در نتیجه ناهمبسته هستند. فرض گوسی بودن ضرایب براساس قضیه حد مرکزی^۳ است. این قضیه بیان می کند که توزیع

^۱ -A Priori SNR

^۲ -Minimum Mean Square Error Estimation of Short-Time Spectral Amplitude

^۳ - Central Limit Theorem

مجموع N متغیر تصادفی، وقتی که N بزرگ باشد به سمت توزیع گوسی میل می کند و هر ضریب

فوریه در یک قاب با طول N ، برابر مجموع وزندار N متغیر تصادفی آن قاب است. پس توزیع آن

گوسی فرض می شود. اما این فرض هنگامی صحیح است که طول قاب زیاد باشد.

همانگونه که قبلا بیان شد رابطه تخمین گر MMSE برای تخمین دامنه به صورت رابطه (۴-۱۱)

می باشد. با توجه به مستقل بودن ضرایب فوریه سیگنال تمیز و نویز این رابطه به صورت رابطه (۴-۴)

(۲۱) در می آید.

$$\hat{S}_k = E(S_k | Y(\omega_1) Y(\omega_2) \dots Y(\omega_{N-1})) = E[S_k | Y(\omega_k)] \quad (21-4)$$

گام اول در محاسبه امید ریاضی بیان شده در این رابطه، تعیین تابع توزیع پسین S_k یا همان

$p(S_k | Y(\omega_k))$ است. برای این کار قانون بیز به کار گرفته می شود که منجر به رابطه (۴-۲۲) می

شود.

$$p(S_k | Y(\omega_k)) = \frac{p(Y(\omega_k) | S_k) p(S_k)}{\int_{-\infty}^{\infty} p(Y(\omega_k) | S_k) p(S_k) dS_k} \quad (22-4)$$

با استفاده از رابطه (۴-۲۲)، امید ریاضی بیان شده به صورت رابطه (۴-۲۳) محاسبه می شود.

$$\begin{aligned} \hat{S}_k &= E(S_k | Y(\omega_k)) = \int_{-\infty}^{\infty} S_k p(S_k | Y(\omega_k)) dS_k \\ &= \frac{\int_{-\infty}^{\infty} S_k p(Y(\omega_k) | S_k) p(S_k) dS_k}{\int_{-\infty}^{\infty} p(Y(\omega_k) | S_k) p(S_k) dS_k} = \frac{\int_{-\infty}^{\infty} S_k p(Y(\omega_k) | S_k) p(S_k) dS_k}{\int_{-\infty}^{\infty} p(Y(\omega_k) | S_k) p(S_k) dS_k} \end{aligned} \quad (23-4)$$

از طرفی می توان نوشت:

$$p(Y(\omega_k) | S_k) p(S_k) = \int_{-\pi}^{\pi} p(Y(\omega_k) | S_k, \theta_s(k)) p(S_k, \theta_s(k)) d\theta_s(k) \quad (24-4)$$

که در آن $\theta_s(k)$ متغیر تصادفی، معادل فاز $S(\omega_k)$ است. بنابراین داریم:

$$\hat{S}_k = \frac{\int_{-\infty}^{\infty} \int_{-\pi}^{\pi} S_k p(Y(\omega_k) | S_k, \theta_s(k)) p(S_k, \theta_s(k)) d\theta_s(k) dS_k}{\int_{-\infty}^{\infty} \int_{-\pi}^{\pi} p(Y(\omega_k) | S_k, \theta_s(k)) p(S_k, \theta_s(k)) d\theta_s(k) dS_k} \quad (25-4)$$

حال برای حل این رابطه باید $p(S_k, \theta_s(k))$ و $p(Y(\omega_k) | S_k, \theta_s(k))$ محاسبه شوند. به همین

منظور از مدل آماری گوسی برای سیگنال و نویز استفاده می شود. با توجه به این مدل می توان

گفت که $Y(\omega_k)$ نیز دارای توزیع گوسی است، زیرا از جمع دو متغیر گوسی با میانگین صفر تشکیل

شده است. به همین ترتیب می توان نتیجه گرفت که $p(Y(\omega_k)|S_{k'}, \theta_s(k))$ نیز گوسی است.

بنابراین داریم:

$$p(Y(\omega_k)|S_{k'}, \theta_s(k)) = \frac{1}{\pi\sigma_d(k)} \exp\left[-\frac{|Y(\omega_k) - S(\omega_k)|^2}{\sigma_d(k)}\right] \quad (26-4)$$

که در آن $\sigma_d(k)$ بیانگر واریانس نویز در مولفه k ام طیف است. از طرفی اندازه دامنه و فاز یک متغیر تصادفی گوسی مستقل از هم بوده [49] و می توان نوشت:

$$p(S_{k'}, \theta_s(k)) = p(S_{k'})p(\theta_s(k)) \quad (27-4)$$

به علاوه، اندازه یک متغیر تصادفی گوسی دارای توزیع رایلی^۱ و توزیع فاز آن در بازه $(-\pi, \pi)$ یکنواخت است [49]. بنابراین:

$$p(S_{k'}, \theta_s(k)) = \frac{S_k}{\pi\sigma_d(k)} \exp\left[-\frac{S_k^2}{\sigma_s(k)}\right] \quad (28-4)$$

که در آن $\sigma_s(k)$ بیانگر واریانس سیگنال تمیز در مولفه k ام طیف است. با جایگزینی روابط (28-4) و (26-4) در رابطه (25-4)، در نهایت تخمین اندازه طیف به صورت زیر به دست می آید [22]:

$$\begin{aligned} \hat{S}_k &= \frac{\sqrt{v_k}}{Y_k} \Gamma(1, \sigma) \Phi(-*, \sigma, 1; -v_k) Y_k \\ &= \frac{\sqrt{\pi}}{\gamma} \frac{\sqrt{v_k}}{Y_k} \exp\left(-\frac{v_k}{\gamma}\right) \left[(1 + v_k) I_0\left(\frac{v_k}{\gamma}\right) + v_k \left(\frac{v_k}{\gamma}\right) \right] Y_k = G(\xi_k, \gamma_k) Y_k \end{aligned} \quad (29-4)$$

که در آن:

$$\Gamma(a) = \int_0^{\infty} t^{a-1} e^{-t} dt \quad (30-4)$$

$$\Phi(a, b; x) = 1 + \frac{a}{b} \frac{x}{\gamma} + \frac{a(a+1)}{b(b+1)} \frac{x^2}{\gamma^2} + \frac{a(a+1)(a+2)}{b(b+1)(b+2)} \frac{x^3}{\gamma^3} + \dots \quad (31-4)$$

$$\sigma_k = \frac{\sigma_s(k)\sigma_d(k)}{\sigma_s(k) + \sigma_d(k)} = \frac{\sigma_s(k)}{1 + \xi_k} \quad (32-4)$$

$$v_k = \frac{\xi_k}{\xi_k + 1} Y_k \quad (33-4)$$

$$\xi_k = \frac{\sigma_s(k)}{\sigma_d(k)} \quad (34-4)$$

$$Y_k = \frac{Y_k^2}{\theta_d(k)} \quad (35-4)$$

^۱ - Rayleigh

و $G(\xi_k, \gamma_k)$ تابع بهره تخمین است که به دو معیار سیگنال به نویز پسین و پیشین وابسته است. تحلیل تابع بهره که وابسته به دو پارامتر سیگنال به نویز پیشین و پسین است، در درک بیشتر این تخمین گر و ارزیابی کارایی آن نقش موثری دارد. وابستگی تابع بهره MMSE به این دو پارامتر یکی از عوامل موثر در کاهش نویز موزیکال این روش در مقایسه با سایر روش ها و از جمله فیلتر وینر است. نکات مهم این تخمین عبارت اند از: نقش سیگنال به نویز پیشین در حذف نویز پرنرنگتر است. همچنین سیگنال به نویزهای بالا MMSE دامنه، حساسیت زیادی به تغییر سیگنال به نویز پیشین به میزان اندک ندارد. MMSE به فرو-تخمین سیگنال به نویز پیشین بیشتر از فرا تخمین آن حساسیت دارد [۲۲]. مقایسه عینی انجام شده در [۳۰] نیز کارایی بهتر تخمین دامنه با MMSE را در مقایسه با سایر روش ها (تفریق طیفی، زیر فضای سیگنال و فیلتر وینر) نشان می دهد.

۴-۵-۳- تخمین فاز با کمترین میانگین مربعات خطا

همانگونه که می دانیم برای تخمین سیگنال تمیز از روی سیگنال نویزی، علاوه بر تخمین دامنه به تخمین فاز نیز نیاز است. با داشتن فاز سیگنال تمیز، نهایتاً طیف کامل سیگنال تخمینی به صورت $\hat{S}(\omega_k) = \hat{S}_k B_k \exp(j\theta_y(k))$ به دست می آید و از روی آن سیگنال در حوزه طیف تخمین زده شده و سپس در حوزه زمان محاسبه می شود. برای تخمین فاز، Ephraim & Malah در مقاله معروف خود [۲۲] دو راه حل پیشنهاد دادند که در ادامه مرور خواهند شد.

در روش اول مشابه مراحل که برای دامنه انجام شد، تخمین MMSE برای بخش موهومی $\hat{S}(\omega_k)$ که به صورت $\exp(j\theta_x(k)) = E\{\exp(j\theta_x(k)) | Y(\omega_k)\}$ است، انجام می شود. با توجه به مدل فرض شده، تخمین بخش نمایی طیف به صورت رابطه (۴-۳۶) نتیجه می شود.

$$\begin{aligned} \exp(j\theta_x(k)) &= \frac{\sqrt{\pi}}{\gamma} \sqrt{v_k} \exp\left(-\frac{v_k}{\gamma}\right) \left[I_0\left(\frac{v_k}{\gamma}\right) + I_1\left(\frac{v_k}{\gamma}\right) \right] \exp(j\theta_y(k)) \\ &\triangleq B_k \exp(j\theta_y(k)) \end{aligned} \quad (۴-۳۶)$$

حال با ترکیب این تخمین برای فاز و تخمین دامنه که در بخش قبل به دست آمد، رابطه تخمین

طیف سیگنال تمیز به صورت $S(\omega_k) = S_k B_k \exp(j\theta_y(k))$ به دست می آید.

با توجه به اینکه تخمین دامنه و فاز به صورت مستقل از هم به دست آمده اند، تخمین رابطه فوق با دامنه جدید $S_k B_k$ الزاما بهینه نخواهد بود. چرا که $|B_k|$ برابر یک نیست و در نتیجه این مساله بهینه بودن S_k را تحت تاثیر قرار می دهد. از این رو راه حل دومی نیز که در همان مقاله ارائه شده، یک تخمین گر MMSE است که در آن یک شرط به مساله اضافه شده است. در این حالت جواب با این محدودیت که قدرمطلق قسمت موهومی برابر یک باشد، به دست آورده شده است، یعنی صورت مساله به شکل رابطه زیر درآمده است:

$$\min_{\exp(j\hat{\theta}_s(k))} E \left\{ \left| \exp(j\hat{\theta}_s(k)) - \exp(j\theta_s(k)) \right|^2 \right\} \text{ و}$$

$$\left| \exp(j\hat{\theta}_s(k)) \right| = 1 \quad (37-4)$$

حل این مساله با کمک گرفتن از ضرایب لاگرانژ منجر به پاسخ زیر می شود.

$$\exp(j\hat{\theta}_s(k)) = \exp(j\theta_s(k)) \quad (38-4)$$

این رابطه بیان می کند که فاز بهینه برای سیگنال از نظر MMSE همان فاز سیگنال نویزی است که استفاده از آن در بیشتر سیستم های بهسازی گفتار به جای فاز سیگنال تمیز، امری رایج است.

۴-۵-۴- محاسبه سیگنال به نویز پیشین

همانگونه که در روابط قبلی آورده شد، تخمین بهینه MMSE برای طیف سیگنال منجر به استفاده از مقادیر واریانس نویز و واریانس سیگنال تمیز جهت محاسبه سیگنال به نویز پیشین و پسین است. در عمل که فقط به سیگنال نویزی دسترسی داریم، می توان با فرض ایستادن بودن نویز، واریانس آن را از روی بخش های غیرگفتار سیگنال با استفاده از VAD یا تخمین گره های نویز محاسبه کرد. اما تخمین واریانس سیگنال تمیز جهت محاسبه سیگنال به نویز پیشین، با توجه به در دسترس نبودن سیگنال تمیز، کار آسانی نیست. برای تخمین این نسبت روش های زیادی ارائه شده است که عمده آنها مبتنی بر دو روش ارائه شده در [۲۲] هستند. این دو روش یکی مبتنی بر روش بیشترین

شبهات و دیگری براساس رابطه این پارامتر با سیگنال به نویز پسین با نام تصمیم گرا^۱ ارائه شده است که اغلب کارهای بعدی توسعه و بهبود یکی از این دو روش است.

۴-۵-۴-۱- روش بیشترین شبهات

در این رویکرد فرض می شود واریانس سیگنال تمیز یک پارامتر نامعلوم و معین است و سعی می شود با استفاده از روش بیشترین شبهات تخمین زده شود. برای این کار فرض می شود که تخمین $\sigma_s^2(k)$ در قاب m ام به کمک L نمونه قبلی مشاهده شده از اندازه دامنه سیگنال نویزی یعنی $\mathbf{Y}_k(m) \triangleq \{Y_k(m), Y_k(m-1), Y_k(m), \dots, Y_k(m-L+1)\}$ انجام شود.

بفرض استقلال میان L نمونه و گوسی بودن توزیع، تابع شبهات زیر به دست می آید.

$$p(\mathbf{Y}_k(m); \sigma_s^2(k), \sigma_n^2(k)) = \prod_{i=1}^{L-1} \frac{1}{\pi[\sigma_s^2(k) + \sigma_n^2(k)]} \exp\left(-\frac{Y_k^2(m-i)}{\sigma_s^2(k) + \sigma_n^2(k)}\right) \quad (39-4)$$

که مشتق گیری از آن نسبت به $\sigma_s^2(k)$ و یافتن نقطه بیشینه تابع فوق منجر به تخمین زیر می شود.

$$\hat{\sigma}_s^2(k, m) = \begin{cases} \frac{1}{L} \sum_{i=1}^{L-1} Y_k^2(m-i) - \lambda_d(k, m) & \text{if non - negativ} \\ \cdot & \text{else} \end{cases}$$

$$= \max\left(\frac{1}{L} \sum_{i=1}^{L-1} Y_k^2(m-i) - \lambda_d(k, m), \cdot\right) \quad (40-4)$$

باتوجه به رابطه (۴۰-۴) سیگنال به نویز پیشین به صورت زیر محاسبه می شود:

$$\xi_k(m) = \max\left(\frac{1}{L} \sum_{i=1}^{L-1} \gamma_k^2(m-i) - \alpha, \cdot\right) \quad (41-4)$$

که در این رابطه همان سیگنال به نویز پسین برای قاب m ام است و عملگر \max باعث می شود که مقدار تخمین همیشه غیرمنفی باشد. در عمل میانگین L نمونه از سیگنال به نویز پسین در قاب m را به صورت رابطه (۴۲-۴) محاسبه می کنند که در آن پارامترهای $\alpha \leq 1$ و $\beta \geq 1$ به ترتیب، ضرایب ثابت هموارسازی و همبستگی هستند.

$$\bar{Y}_k(m) = \alpha \bar{Y}_k(m-1) + (1-\alpha) \frac{Y_k(m)}{\beta} \quad (42-4)$$

استفاده از این رابطه منجر به تغییر رابطه $\xi_k(m)$ به صورت $\xi_k(m) = \max(\bar{Y}_k(m) - 1, 0)$ می شود.

۴-۵-۴-۲- روش تصمیم گرا

پراکاردترین روش تخمین سیگنال به نویز پیشین روش تصمیم گراست که اولین بار در [۲۲] ارائه گردید. در این روش از رابطه بین سیگنال به نویز پیشین و پسین جهت تخمین این پارامتر استفاده می شود. با توجه به تعریف سیگنال به نویز پیشین در رابطه (۴-۳۴) داریم:

$$\begin{aligned} \xi_k(m) &= \frac{\sigma_s(k, m)}{\sigma_d(k, m)} = \frac{E\{Y_k^s(m)\}}{\sigma_d(k, m)} = \frac{E\{Y_k^s(m) - D_k^s(m)\}}{\sigma_d(k, m)} \\ &= \frac{E\{Y_k^s(m)\}}{\sigma_d(k, m)} - E\{Y_k(m)\} - 1 \end{aligned} \quad (43-4)$$

از رابطه فوق می توان رابطه (۴-۴۴) را نتیجه گرفت و بر اساس آن رابطه نهایی تخمین ξ_k را به صورت (۴-۴۵) نوشت که در آن $0 < a < 1$. ضریب وزن جایگزین شده به جای $\frac{1}{\beta}$ در رابطه (۴-۴۴) است و $S_k^s(m-1)$ تخمین اندازه به دست آمده برای قاب قبلی است.

$$\xi_k(m) = E \left\{ \frac{1}{\beta} \frac{S_k^s(m)}{\sigma_d(k, m)} + \frac{1}{\beta} [Y_k(m) - 1] \right\} \quad (44-4)$$

این رابطه را تخمین تصمیم گرا می نامند چون بر اساس تخمین اندازه دامنه در قاب قبلی، کار تخمین را انجام می دهد. این رابطه بیان می کند که تخمین سیگنال به نویز پیشین برابر با جمع وزندار سیگنال به نویز پیشین قاب قبل و سیگنال به نویز پسین قاب جاری است. برای مقداردهی اولیه به رابطه (۴-۴۵) در قاب اول که $m=0$ است، توصیه می شود که از رابطه $\xi_k(m) = a + (1-a) \max\{\gamma_k(0) - 1, 0\}$ استفاده شود که در آن مقدار $a = 0.98$ منجر به

بهترین نتیجه می شود [۲۲].

$$\xi_k(m) = a \frac{\hat{s}_k^2(m-1)}{\sigma_a^2(km-1)} + (1-a) \max\{\gamma_k(m) - 1, 0\} \quad (4-45)$$

نتایج مقایسه روش MMSE که در آن از روش تصمیم گرا ($a = 0.98$) جهت تخمین سیگنال به نویز پیشین استفاده شده، با روش های تفریق طیفی و فیلتر وینر در حضور نویز سفید نشان می دهد که در صورت استفاده از این روش سیگنال بهسازی شده دارای نویز موزیکال نخواهد بود. در مقابل اگر از روش بیشترین شباهت جهت تخمین سیگنال به نویز پیشین استفاده شود، سیگنال خروجی در روش MMSE دارای نویز موزیکال خواهد بود. همچنین در صورت استفاده از روش بیشترین شباهت جهت تخمین سیگنال به نویز پیشین در روش های وینر و MMSE، این دو روش تقریباً مشابه عمل می کنند. در صورتی که اگر از روش تصمیم گرا در روش وینر استفاده شود، سیگنال خروجی اعوجاج بیشتری نسبت به MMSE دارد.

نتایج نشان می دهد که کم کردن مقدار a منجر به کاهش اعوجاج در سیگنال بهسازی شده اما افزایش نویز زمینه باقیمانده در سیگنال می شود که بیانگر نقش این پارامتر در کنترل میزان اعوجاج و نویز زمینه است [59]. به علاوه مقایسه روش MMSE و روش بیشترین شباهت نشان می دهد که روش بیشترین شباهت نویز موزیکال بیشتری را ایجاد می کند.

۴-۵-۵- تخمین MMSE با توزیع های غیر گوسی

یکی از فرض های اساسی برای بدست آوردن تخمین های MMSE، گوسی بودن طیف سیگنال بود که بیان شد این فرض براساس قضیه حد مرکزی است. از طرفی این فرض هنگامی صحیح است که طول قاب زیاد باشد. ولی در کاربردهای پردازشی طول هر قاب کوتاه (بین ۲۰ تا ۴۰ میلی ثانیه) است. بر همین اساس، می توان انتظار داشت استفاده از توزیع های غیرگوسی منجر به نتایج بهتری شوند. این نکته که ابتدا در [50] مطرح شد، اخیراً بسیار مورد توجه محققین قرار گرفته است. در کارهای مختلفی از توزیع های گاما [7]، رایلی [3]، لاپلاس [43] و مخلوط گوسی [26] به جای گوسی در مدل کردن بخش های حقیقی و موهومی ضرایب فوریه سیگنال استفاده شده است. همچنین در برخی از مقالات، برای نویز هم از توزیع های غیر گوسی استفاده شده است.

در این بخش یک تخمین گر MMSE، با فرض گاما بودن ضرایب سری فوریه سیگنال تمیز و گوسی بودن ضرایب فوریه نویز به دست آورده شده است. ابتدا فرض می شود که $Y_R(k) = \text{Re}\{Y(\omega_k)\}$ قسمت حقیقی طیف سیگنال نویزی و $Y_I(k) = \text{Im}\{Y(\omega_k)\}$ قسمت حقیقی طیف سیگنال نویزی باشد، یعنی $Y = Y_R + jY_I$ (برای سادگی بیشتر از اندیس k صرف نظر می شود). بر همین اساس طیف مختلط سیگنال تمیز نیز به صورت $S = S_R + jS_I$ خواهد بود. با فرض گاما بودن توزیع بخش های حقیقی و موهومی طیف سیگنال تمیز، خواهیم داشت:

$$p(S_R) = \frac{\sqrt{1.5}}{\sqrt{\pi}\sigma_S} |S_R|^{-1.5} \exp\left(-\frac{\sqrt{1.5}|S_R|}{\sqrt{\sigma_S}}\right) \quad (46-4)$$

$$p(S_I) = \frac{\sqrt{1.5}}{\sqrt{\pi}\sigma_S} |S_I|^{-1.5} \exp\left(-\frac{\sqrt{1.5}|S_I|}{\sqrt{\sigma_S}}\right) \quad (47-4)$$

که در این روابط $\sigma_S/2$ برابر واریانس بخش های حقیقی و موهومی طیف سیگنال تمیز است. برای به دست آوردن تخمین MMSE طیف سیگنال تمیز از روی سیگنال نویزی، فرض می شود که بخش های حقیقی و موهومی سیگنال مستقل هستند. با این فرض، تخمین هر کدام از این دو بخش به صورت مجزا به دست آمده و سپس با هم ترکیب می شوند. یعنی:

$$E[s(\omega_k)|Y(\omega_k)] = E[S|Y] = E[S_R|Y_R] + jE[S_I|Y_I] \quad (48-4)$$

حال فرض می شود که نویز دارای توزیع گوسی باشد. در این صورت برای بخش حقیقی طیف داریم

$$\tilde{S}_R = E[S_R|Y_R] = \frac{\int_{-\infty}^{+\infty} S_R p(Y_R|S_R) p(S_R) dS_R}{p(Y_R)} \quad (49-4)$$

برای محاسبه این رابطه باید $p(Y_R|S_R)$ ، $p(Y_R)$ و $p(S_R)$ (که همان توزیع گاما است) محاسبه شوند. با توجه به اینکه بخش حقیقی طیف سیگنال نویزی برابر مجموع بخش های حقیقی طیف سیگنال تمیز و نویز است، برای $p(Y_R|S_R)$ داریم:

$$p(Y_R|S_R) = p_D(Y_R - S_R) = \frac{1}{\sqrt{\pi}\sigma_d} \exp\left(-\frac{(Y_R - S_R)^2}{\sigma_d}\right) \quad (50-4)$$

این رابطه با توجه به گوسی بودن توزیع طیف نویز به دست آمد. محاسبه $p(Y_R)$ نیز با استفاده از

رابطه (4-51) قابل انجام است.

$$p(Y_R) = \int_{-\infty}^{+\infty} p(S_R|Y_R) p(S_R) dS_R \quad (51-4)$$

بنابراین تخمین بخش حقیقی طیف به صورت رابطه (52-4) به دست خواهد آمد.

$$\begin{aligned} \hat{S}_R &= \frac{\sqrt[4]{1.0}}{\sqrt{\pi} \sqrt{\sigma_d \sigma_s} p(Y_R)} \int_{-\infty}^{+\infty} S_R |S_R|^{-1.0} \exp\left(-\frac{(Y_R - S_R)^2}{\sqrt{\sigma_d}} - \frac{\sqrt[4]{1.0} |S_R|}{\sqrt{\sigma_s}}\right) dS_R \\ &= \frac{\sqrt{\sigma_d Z_r}}{\sqrt[4]{2} \sqrt{Z_1}} \end{aligned} \quad (52-4)$$

که در آن :

$$C_p(z) = \frac{\exp\left(-\frac{z^2}{\gamma}\right)}{\Gamma(-p)} \int_0^{\infty} x^{-p-1} \exp\left(-xz - \frac{x^2}{\gamma}\right) dx \quad (53-4)$$

$$G_1 = \frac{\sqrt{1.0 \sigma_d}}{\sqrt[4]{2} \sqrt{\sigma_s}} + \frac{Y_R}{\sqrt{\sigma_d}} = \frac{\sqrt{1.0}}{\sqrt[4]{2} \sqrt{\xi}} + \frac{Y_R}{\sqrt{\sigma_d}} \quad (54-4)$$

$$G_r = \frac{\sqrt{1.0 \sigma_d}}{\sqrt[4]{2} \sqrt{\sigma_s}} + \frac{Y_R}{\sqrt{\sigma_d}} = \frac{\sqrt{1.0}}{\sqrt[4]{2} \sqrt{\xi}} + \frac{Y_R}{\sqrt{\sigma_d}} \quad (55-4)$$

$$Z_1 = \exp\left(\frac{G_1^2}{\gamma}\right) C_{-1.0}(\sqrt{\gamma} G_1) - \exp\left(\frac{G_r^2}{\gamma}\right) C_{-1.0}(\sqrt{\gamma} G_r) \quad (56-4)$$

$$Z_r = \exp\left(\frac{G_r^2}{\gamma}\right) C_{-1.0}(\sqrt{\gamma} G_r) - \exp\left(\frac{G_1^2}{\gamma}\right) C_{-1.0}(\sqrt{\gamma} G_1) \quad (57-4)$$

به صورت مشابهی می توان بخش موهومی را به دست آورد که منجر به تخمینی مشابه تخمین

بخش حقیقی می شود. ترکیب این دو تخمین به صورت $\hat{S}_1 = \hat{S}_R + j\hat{S}_I$ منجر به تخمین نهایی

طیف که یک تخمین مختلط است، می شود. مقایسه این تخمین و فیلتر وینر بیانگر عملکرد مشابه

آنها از نظر اعوجاج سیگنال در SNRهای بالا و متفاوت بودن عملکرد آنها در SNRهای پایین است.

4-6- تخمین گر بیشترین احتمال پسین MAP

علاوه بر تخمین گرهای ML و MMSE، تخمین گر MAP یکی از روش های تخمین در بهسازی

گفتار است [28] و [61] و [62]. در MAP هدف بیشینه کردن احتمال پسین $p(\mathbf{s}_{k_i} | Y(\omega_k))$ است.

تخمین MAP را می توان به عنوان یک راه حل جایگزین برای MMSE در مواردی که تخمین

بیشینه احتمال پسین از محاسبه میانگین آن ساده تر است به کاربرد. یکی دیگر از مزایای تخمین

MAP این است که می توان تخمین دامنه و فاز را به صورت همزمان انجام داد.

در ادامه تخمین MAP به صورت همزمان برای دامنه و فاز آورده خواهد شد. تخمین همزمان دامنه و فاز به روش MAP منجر به رابطه (۵۸-۴) می شود که در بخش آخر آن، با توجه به مستقل بودن $p(Y(\omega_k))$ از s_k و θ_s ، از آن صرفنظر شده است. با فرض گوسی بودن توزیع ها رابطه (۵۹-۴) نتیجه می شود.

$$(\hat{s}_k, \hat{\theta}_s) = \underset{s_k, \theta_s}{\operatorname{argmax}} p(s_k, \theta_s | Y(\omega_k)) = \underset{s_k, \theta_s}{\operatorname{argmax}} \frac{p(Y(\omega_k) | s_k, \theta_s) p(s_k, \theta_s)}{p(Y(\omega_k))} = \underset{s_k, \theta_s}{\operatorname{argmax}} P(Y(\omega_k) | s_k, \theta_s) p(s_k, \theta_s) \quad (58-4)$$

$$p(Y(\omega_k) | s_k, \theta_s) p(s_k, \theta_s) = \frac{s_k}{\pi^2 \sigma_s(k) \sigma_d(k)} \exp\left(-\frac{|Y(\omega_k) - s_k e^{j\theta_s}|^2}{\sigma_s(k)} - \frac{s_k^2}{\sigma_s(k)}\right) \quad (59-4)$$

با مشتق گیری از این رابطه (یا از لگاریتم آن) نسبت به θ_s و برابر صفر قرار دادن آن، تخمین $\hat{\theta}_s = \theta_s$ نتیجه می شود که همان تخمین نتیجه شده از MMSE است و همچنین مشتق گیری از رابطه فوق نسبت به s_k و یافتن نقطه بیشینه آن، رابطه (۶۰-۴) برای تخمین دامنه نتیجه می شود:

$$\hat{s}_R = \frac{s_k + \sqrt{s_k^2 + \tau(1 + \xi_k)s_k}}{\tau(1 + \xi_k)} Y_k = G_{MAP} Y_k \quad (60-4)$$

۷-۴- به کارگیری عدم قطعیت گفتار (SPU)^۱ در بهسازی گفتار

در روش های بیان شده در بخش های قبلی، یکی از فرض های اصلی همه روش ها این است که گفتار در همه زمان ها در سیگنال مورد بررسی وجود دارد. اما در عمل، در بسیاری از بخش های سیگنال گفتار، به علت وجود بخش هایی همچون سکوت و توقف گفتار وجود ندارد و یا حداقل به میزان کمی وجود دارد. به علاوه، حتی در بخش های گفتاری سیگنال نیز ممکن است در برخی از باندهای فرکانسی گفتار وجود نداشته باشد. در بهسازی گفتار می توان برای هر باند فرکانسی دو حالت حضور و عدم حضور گفتار را متصور بود و براساس آن کار بهسازی گفتار را انجام داد. در اینجا از مساله فوق به این صورت که گفتار ممکن است در بعضی از زمان ها و در برخی از فرکانس ها وجود نداشته باشد، استفاده شده است. برای فرموله کردن بیان فوق، یک مدل دو حالتی برای

^۱ - Speech Present Uncertainty

سیگنال گفتار فرض می شود که در آن هر حالت بیانگر یکی از دو فرضیه حضور یا عدم حضور گفتار است. به عبارتی داریم:

$$H_1^k: \text{Speech absent } |Y(\omega_k)| = |D(\omega_k)| \quad (۴-۶۱)$$

$$H_0^k: \text{Speech present } |Y(\omega_k)| = |S(\omega_k) + D(\omega_k)| \quad (۴-۶۲)$$

در این رابطه H_1^k بیانگر فرضیه تهی^۱ است که به معنی عدم حضور گفتار در فرکانس k ام است و H_0^k بیانگر فرضیه وجود گفتار در سیگنال نویزی در فرکانس k ام است. حال برای استفاده از این مدل دودویی در روش های تخمین بیان شده، بایستی از جمع وزن دار دو فرضیه استفاده کرد [۲۲].

۴-۷-۱- تخمین گر MMSE با SPU

برای روشن تر شدن موضوع به عنوان مثال تخمین گر MMSE با SPU به طور خیلی خلاصه توضیح داده شده است. برای تخمین گر MMSE که $\hat{s}_k = E\{S_k | Y_k\}$ است، داریم:

$$\hat{s}_k = p(H_1^k | Y_k) E\{S_k | Y_k, H_1^k\} + p(H_0^k | Y_k) E\{S_k | Y_k, H_0^k\} \quad (۴-۶۳)$$

در این رابطه $E\{S_k | Y_k, H_1^k\}$ بیانگر متوسط دامنه سیگنال گفتار تمیز به شرط داشتن دامنه طیف و عدم حضور گفتار است که می توان آن را صفر فرض کرد. با این فرض تخمین گر جدید MMSE به صورت رابطه (۴-۶۴) خواهد بود:

$$\hat{s}_k = p(H_1^k | Y_k) E\{S_k | Y_k, H_1^k\} \quad (۴-۶۴)$$

۴-۸- روش های تخمین نویز در بهسازی گفتار

یکی از نیازهای کلیه روش های بهسازی گفتار تک کاناله، تخمین نویز یا مشخصات آن است که کارایی روش بهسازی گفتار را به شدت تحت تاثیر قرار می دهد. در روش های بیان شده تا به حال فرض شد که تخمینی از نویز در دسترس است. در این قسمت روش های تخمین نویز بررسی

^۱ - Null Hypothesis

خواهند شد.

ساده ترین روش تخمین نویز VAD است. در این روش با استفاده از انرژی، نرخ عبور از صفر^۱، متناوب بودن و . . . برای هر قاب، در مورد گفتار یا غیر گفتار بودن آن تصمیم گیری می شود. تصمیم گیری اغلب براساس مقایسه این مقادیر با مقادیر آستانه است. روش های مختلفی برای VAD تا به حال پیشنهاد شده اند که مرور کاملی بر آنها در [۳۳] آورده شده است. برخی از روش ها مبتنی بر انرژی و نرخ عبور از صفر هستند [۳۲] و برخی دیگر بر اساس ضرایب کپسترال [۲۸]، معیار فاصله LPC ایتاکورا [۵۱] و متناوب بودن گفتار [۵۷] عمل می کنند. در این روش مشخصات نویز در بخش های سکوت و غیرگفتاری سیگنال به روز می شود. این روش در SNRهای بالا و برای نویزهای ایستاد مناسب است. به علاوه برخی از آنها نیاز به برخی پارامترها مانند مقادیر آستانه و برخی تنظیمات دیگر دارند که ممکن است از محیطی به محیط دیگر تغییر کند. سیستم های بهسازی گفتار برای داشتن عملکرد مطلوب در کاربردهای واقعی باید از توانایی تخمین مشخصات نویز در شرایط غیر ایستاد و دنبال کردن این مشخصات برای همه زمان ها حتی، در شرایطی که گفتار وجود داشته باشد، برخوردار باشند. از این رو روش های تخمین نویز به عنوان یکی از بخش های مهم در سیستم های بهسازی گفتار مطرح شده اند.

در ادامه سه روش از مهم ترین روش های تخمین نویز توضیح داده می شود. در همه این روش ها فرض می شود که توان نویز نسبت به گفتار تغییرات کمتری دارد و ایستاد تر است.

۴-۸-۱- روش میانگین گیری بازگشتی^۲

مشاهده شده است که علاوه بر بخش های سکوت در سیگنال گفتار نویزی، بخش های دیگری نیز وجود دارد که در زمان حضور گفتار نیز انرژی آنها بسیار کم و نزدیک به انرژی نویز است. از این بخش ها می توان به بخش های فرکانس پایین (به ویژه زیر ۲ کیلو هرتز) و بخش های فرکانس

^۱ - Zero Crossing

^۲ - time-recursive averaging

بالا) به ویژه بالای ۴ کیلوهرتز) در مصوت‌ها اشاره کرد. به طور کلی می‌توان گفت برای هر نوع نویزی، باند‌های فرکانسی خاصی از سیگنال که احتمال حضور گفتار در آنها کمتر است و یا دارای SNR پایین هستند، برای به روز کردن مشخصات نویز مناسب‌تر هستند. براساس این مشاهده روش میانگین‌گیری بازگشتی، طیف نویز را به صورت میانگین وزندار مقادیر قبلی تخمین نویز و مقدار جاری طیف سیگنال نویزی، طبق رابطه (۴-۶۵) به روز می‌کند. وزن این رابطه متغیر با فرکانس و زمان بوده و مقدار آن برای هر باند فرکانسی به SNR آن باند و یا به احتمال حضور گفتار در آن باند مرتبط است.

$$\hat{\sigma}_d^2(f, \omega_k) = \alpha(f, \omega_k) \hat{\sigma}_d^2(f-1, \omega_k) + (1 - \alpha(f, \omega_k)) |Y(f, \omega_k)|^2 \quad (۴-۶۵)$$

در این رابطه $\hat{\sigma}_d^2(f, \omega_k)$ توان نویز تخمین زده شده در فرکانس ω_k و قاب f ، $|Y(f, \omega_k)|^2$ مجذور اندازه طیف سیگنال نویزی و $\alpha(f, \omega_k)$ ضریب هموارسازی وابسته به زمان و فرکانس است. تمامی روش‌های میانگین‌گیری بازگشتی از رابطه فوق پیروی می‌کنند و تفاوت آنها با یکدیگر در نحوه محاسبه ضریب $\alpha(f, \omega_k)$ است. برخی از محققان مقدار این ضریب را بر اساس تخمین SNR هر باند فرکانسی محاسبه کرده [۳۶] (روش‌های مبتنی بر SNR) و برخی دیگر مقدار ضریب هموارسازی را ثابت گرفته و در عوض تنها در برخی شرایط مقدار $\hat{\sigma}_d^2(f, \omega_k)$ را به روز کرده‌اند. (روش‌های مبتنی بر تصمیم‌گیری) [۲۹] و [۵۳]. همچنین تعدادی دیگر از محققین مقدار آن را بر حسب احتمال حضور/عدم حضور گفتار در باند‌های فرکانسی حساب کرده‌اند [۱۱] و [۳۹].

۴-۸-۲- روش مبتنی بر هیستوگرام^۱

هیستوگرام باند‌های فرکانسی سیگنال نشان می‌دهند که بیشترین تکرار (مقدار بیشینه هیستوگرام) در یک باند فرکانسی خاص، متناسب با سطح نویز در آن باند است. در مواردی هیستوگرام انرژی طیف دارای دو قله است که یکی از آنها کم انرژی و مربوط به بخش‌های سکوت و قسمت‌های کم انرژی سیگنال نویزی است و دومی انرژی بالاست و مربوط به بخش‌های صدادر سیگنال است. در

^۱ - Histogram-Based

بیشتر موارد قله کم انرژی بالاتر از قله با انرژی بالاست (تکرار بیشتر در هیستوگرام)، هر چند که این مساله به نوع نویز و طول سیگنال مورد بحث نیز وابسته است. این مشاهده اساس روش های مبتنی بر هیستوگرام است که در بیشتر آنها تخمین نویز از روی هیستوگرام مقادیر قبلی طیف انجام می شود. بدین معنی که برای هر قاب در هر فرکانس، هیستوگرام توان طیف از روی چند صد میلی ثانیه قبلی سیگنال (مثلا با استفاده از یک پنجره ۴۰۰ میلی ثانیه ای) محاسبه می شود و بر اساس بیشترین تکرار موجود در هیستوگرام، نویز تخمین زده می شود. همچنین گاهی برای هموارسازی طیف نویز تخمین زده شده از یک رابطه بازگشتی استفاده می شود. با توجه به موارد فوق می توان الگوریتم تخمین نویز به روش هیستوگرام را به صورت زیر خلاصه کرد:

- برای هر قاب f از گفتار نویزی مراحل زیر را انجام دهید:

- توان طیف سیگنال نویزی، $|Y(f, \omega_k)|^2$ ، را محاسبه کنید.

- توان طیف سیگنال نویزی را به صورت زیر هموار کنید.

$$X(f, \omega_k) = \alpha X(f, \omega_k) + (1 - \alpha) |Y(f, \omega_k)|^2 \quad (۴-۶۶)$$

- هیستوگرام D نمونه از توان طیف های قاب های قبلی، یعنی $\{X(f-1, \omega_k), X(f-2, \omega_k), \dots, X(f-1, \omega_k)\}$

را محاسبه نمایید.

- در هیستوگرام از R بین (Bin) استفاده شود.

- با فرض اینکه $c = [c_1, c_2, \dots, c_R]$ تعداد تکرار هر بین در هیستوگرام بوده و

$s = [s_1, s_2, \dots, s_R]$ مراکز متناسب آنها باشد، رابطه زیر را حساب کنید:

$$c_{max} = \arg \max_{1 < i < R} c_i \quad (۴-۶۷)$$

- حال با استفاده از c_{max} ، تخمین نویز $H_{max}(f, \omega_k)$ را به صورت $S_{max} = H_{max}(f, \omega_k)$

به دست آورید که s_{max} مرکز بین متناسب با c_{max} در هیستوگرام است.

- توان طیف تخمینی نویز را با رابطه زیر هموار کنید.

$$\hat{\sigma}_d^2(f, \omega_k) = \beta \hat{\sigma}_d^2(f-1, \omega_k) + (1 - \beta) H_{max}(f, \omega_k) \quad (۴-۶۸)$$

در این الگوریتم $\hat{\sigma}_d^2(f, \omega_k)$ توان طیف تخمین زده شده برای نویز و α و β به ترتیب ضرایب هموار سازی توان طیف سیگنال نویزی و توان طیف نویز هستند. نتایج بررسی این روش نشان می دهد که اگر هموار سازی روی توان طیف سیگنال نویزی اعمال نشود، کارایی این روش افت کرده و دچار فرو تخمین می شود. طول پنجره مورد استفاده جهت محاسبه هیستوگرام اهمیت زیادی دارد. هرچه طول پنجره کمتر باشد، تخمین نویز دچار فرا تخمین می شود، چرا که در این حالت رخ دادن یک بخش دارای انرژی بالا در گفتار ممکن است قله هیستوگرام را از نویز به گفتار عوض کند. این مشکل به ویژه در فرکانس های پایین که انرژی سیگنال بالاست، برجسته تر است. در فرکانس های بالا که انرژی کمتری در آنها متمرکز است، استفاده از پنجره با طول کوچک چندان مطرح نیست. برای حل مشکل فرا تخمین می توان به دو روش عمل کرد. روش اول این است که D بزرگ باشد. اما بزرگ گرفتن D باعث می شود که این روش نتواند با سرعت کافی تغییرات موجود در نویز را دنبال کند. راه حل دوم این است که از گذاشتن قاب با مقدار توان بالا در محاسبه هیستوگرام ممانعت شود. یعنی قاب های حاوی بخش گفتاری سیگنال فیلتر شده و کنار گذاشته شوند. برای تحقق این امر می توان از سیگنال به نویز پسین استفاده کرد و تنها از قاب هایی که مقدار این پارامتر برای آنها از مقداری کمتر است، در محاسبه هیستوگرام استفاده شود [۵۳].

۴-۸-۳- روش ردیابی کمینه ها^۱

توان سیگنال نویزی حتی در بخش های گفتاری، در برخی از باندهای فرکانسی به سطح توان نویز کاهش می یابد. از این رو می توان با دنبال کردن مقادیر کمینه طیف سیگنال نویزی در تمام باندهای فرکانسی تخمینی از سطح نویز به دست آورد. این ایده مبنای روش ردیابی کمینه ها برای تخمین نویز است که خود به دو دسته تقسیم می شود:

۱- آمارگان کمینه^۲ است که ابتدا توسط مارتین در ۱۹۹۴ ارائه گردید [۴۰] و سپس بهبودهایی روی

^۱ - Minimal Tracking

^۲ - Minimum Statistic

آن داده شد [۴۱]. در این روش کمینه توان طیف در یک پنجره محدود (به طول قسمت آنالیز) دنبال می شود. در روش آمارگان کمینه بعد از محاسبه پریودوگرام^۱ سیگنال نویزی، توان این سیگنال با توجه به میزان نوسان بالای پریودوگرام، به صورت رابطه (۴-۶۹) محاسبه می شود. در این رابطه $0 \leq \alpha \leq 1$ یک ضریب هموار سازی است. با ردیابی کمینه $p(f, \omega_k)$ بر روی یک پنجره با طول محدود، تخمینی از توان نویز به دست می آید. طول پنجره مورد نظر در [۴۰] بین ۰.۸ تا ۱.۴ ثانیه پیشنهاد شده است. شکل (۴-۱) پریودوگرام، پریودوگرام هموار شده و توان تخمینی نویز را برای یک سیگنال نویزی و برای فرکانس $k=25$ نشان می دهد.

$$p(f, \omega_k) = \alpha p(f - 1, \omega_k) + (1 - \alpha) |Y(f, \omega_k)|^2 \quad (۴-۶۹)$$

یافتن کمینه مقدار $p(f, \omega_k)$ بر روی D مقدار پشت سر هم این تابع، به صورت رابطه زیر می باشد.

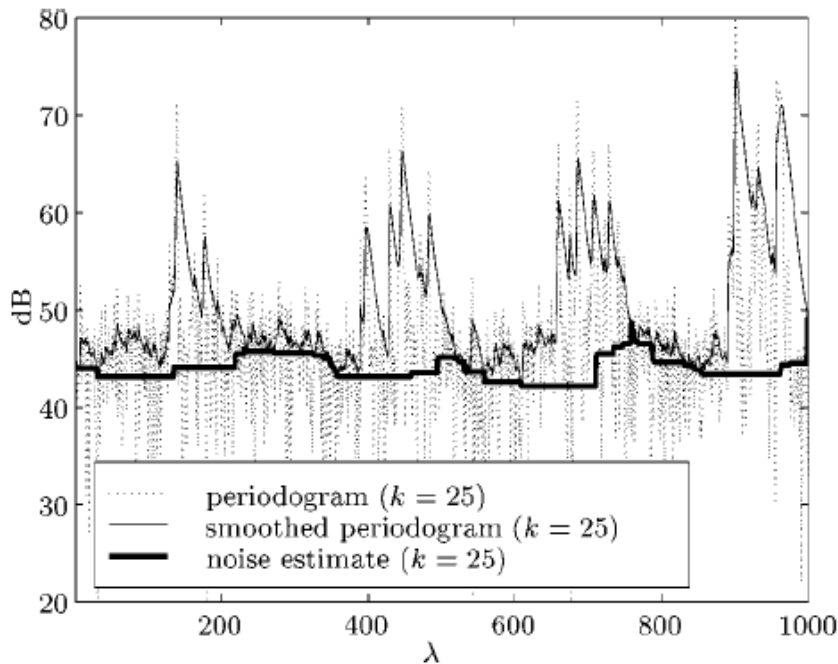
$$p_{min}(f, \omega_k) = \min\{p(f, \omega_k), p(f - 1, \omega_k), \dots, p(f - D + 1, \omega_k)\} \quad (۴-۷۰)$$

تخمین نویز با این روش و مقایسه آن با طیف واقعی نویزی نشان می دهد که مقدار نویز به سمت مقادیر پایین تر بایاس شده و در برخی موارد کمتر از مقدار واقعی تخمین زده می شود. برای به دست آوردن مقدار بایاس و حذف اثر فروتخمین نویز می توان از این واقعیت استفاده کرد که در توزیع های احتمال مرتبط به تخمین پریودوگرام، کمینه مجموعه ای از متغیرهای تصادفی از مقدار میانگین این متغیرها کمتر است و در نتیجه تخمین نویز با روش فوق منجر به بایاس شدن آن به سمت مقادیر پایین می شود. بنابراین برای جبران اثر بایاس شدن می توان مقدار میانگین کمینه ها را بدست آورد و با ضرب معکوس آن در تخمین توزیع توان نویز به صورت زیر این اثر را از بین برد.

$$\hat{p}_d^y(f, \omega_k) = B_{min}(f, \omega_k) \cdot p_{min}(f, \omega_k) \quad (۴-۷۱)$$

که $B_{min}(f, \omega_k)$ ضریب جبران سازی بایاس است که برابر با معکوس میانگین $p_{min}(f, \omega_k)$ است. همچنین بررسی مقادیر مختلف α بیانگر این واقعیت است که مقادیر بالا برای α (مثلاً $\alpha = 0.9$) منجر به هموارتر شدن پریودوگرام و حذف فرورفتگی های مختلف موجود آن می شود. این مساله به

^۱ - Periodogram



شکل (۴-۱): پریودوگرام $|Y(\lambda, k)|^2$ ، پریودوگرام هموار شده $P(\lambda, k)$ و تخمین نویز برای فرکانس $k=25$ [۴۱].

نوبه خود منجر به وسیع تر شدن قله های پریودوگرام و نزدیک شدن تخمین نویز به این قله ها شده و تخمین نویز را دچار تخمین رو به بالا می کند. در واقع برای دنبال کردن نویز بهتر آن است که برای بخش های حاوی گفتار، مقدار این پارامتر نزدیک به صفر و در نتیجه از $p(f, \omega_k) \approx |Y(f, \omega_k)|^2$ که همان حالت هموار نشده پریودوگرام است، استفاده شود. پس با به کارگیری ضریب هموارسازی α به صورت متغیر با زمان و فرکانس، تخمین نویز با روش آمارگان کمیته بهبود می یابد. برای تغییر ضریب ثابت هموارسازی α به $\alpha(f, \omega_k)$ ، از کمیته کردن خطای بین طیف تخمینی و طیف واقعی به رابطه (۴-۷۲) می رسیم.

$$\hat{\alpha}(f, \omega_k) = \frac{\alpha_{max} + \alpha_c}{1 + \left| \frac{p(f-1, \omega_k) / \hat{p}_a^2(f-1, \omega_k) - 1 \right|^T} \quad (۷۲-۴)$$

که در آن:

$$\alpha_c(f, \omega_k) = \frac{1}{1 + \left[\frac{\sum_{i=1}^{M-1} p(f-1, \omega_k) / \sum_{i=1}^{M-1} |Y(f-1, \omega_k)|^2 - 1 \right]} \quad (۷۳-۴)$$

و α_{max} که در [۴۱] برابر با ۰.۹۶ قرار داده شده است برای محدود کردن سقف ضریب هموارسازی

است. جزئیات بیشتر درباره این روش در [۴۱] آمده است.

۲- در این روش مینیمم به صورت پیوسته و بدون نیاز به پنجره دنبال می شود.^۱ یکی از معایب روش قبل عدم توانایی دنبال کردن تغییرات سریع طیف نویز است. به همین دلیل روشی متفاوت در [۳۷] پیشنهاد شد. در این روش تخمین نویز به صورت پیوسته در هر فرکانس، با یک رابطه هموارسازی غیرخطی، به روز می شود. در این روش نیز برای دنبال کردن مینیمم توان نویزی از پرلودوگرام هموار شده استفاده می شود.

$$P(f, \omega_k) = \alpha p(f - \nu, \omega_k) + (1 - \alpha) |Y(f, \omega_k)|^2 \quad (74-4)$$

که در آن α فاکتور هموارسازی و در بازه $0.7 \leq \alpha \leq 0.9$ می باشد.

رابطه غیر خطی که برای تخمین نویز بر پایه دنبال کردن مینیمم طیف توان سیگنال نویزی، در هر فرکانس استفاده می شود، به صورت زیر است:

$$\begin{aligned} & \text{If } p_{\min}(f - \nu, \omega_k) < p(f, \omega_k) \\ & p_{\min}(f, \omega_k) = \lambda p_{\min}(f - \nu, \omega_k) + \frac{1-\lambda}{1-\beta} (p(f, \omega_k) - \beta p(f - \nu, \omega_k)) \\ & \text{Else} \\ & \quad p_{\min}(f, \omega_k) = p(f, \omega_k) \\ & \text{End} \end{aligned} \quad (75-4)$$

که در آن $p_{\min}(f, \omega_k)$ تخمین نویز است. مزیت این روش نسبت به روش قبل پیچیدگی محاسباتی کمتر است.

۹-۴- روش های آماری مبتنی بر مدل

در این روش ها به کمک داده های آموزشی مدل یا مدل هایی ایجاد می شود که خصوصیات مورد نظر را (که معمولاً تابع توزیع احتمال پارامترهاست) در بردارد. در این روش با استفاده از داده های آموزشی، خصوصیات بلندمدت و مشترک بین قاب های مختلف آموزش داده می شود. استفاده از رویکرد مدل سازی آماری در پردازش گفتار به ویژه بازشناسی گفتار بسیار رایج و موفق بوده است. به کارگیری این مفهوم در بهسازی گفتار را می توان به کار [۴۵] نسبت داد که در آن فرض شده که هر قاب سیگنال در یکی از دو حالت سکوت و غیر سکوت قرار می گیرد. استفاده از

^۱ - Continuous Spectral Minimum Tracking

چند حالت مختلف (کلاس) برای سیگنال در بهسازی گفتار قبلا در [۱۷] ارائه شده بود که در آن گفتار به پنج کلاس 'stop'، 'fricative'، 'vowel'، 'glids' و 'nasal' تقسیم شده و سیگنال نویز به یکی از این کلاس ها نسبت داده می شود. سپس برای هر کدام از حالات از تخمین گر خاص آن کلاس استفاده شده است. امروزه ایده حالات و مدل های فوق به صورت کلی تری تعمیم داده شده است که از روش های پرکاربرد آنها می توان به مدل مخفی مارکوف HMM، مدل مخلوط گوسی GMM و چندی سازی بردار^۱ VQ اشاره کرد.

مدل سازی در این روش را می توان شامل دو فاز آموزش و آزمون دانست. در فاز آموزش با استفاده از روش های HMM، GMM یا VQ مدل هایی برای گفتار و نویز بر اساس داده های آموزشی به دست می آید. در این فاز پارامترهای مدل ها به کمک روش های آموزش متناسب از روی داده های آموزشی استخراج می شوند. این پارامترها در واقع همان توزیع های گفتار و نویز هستند. فاز استفاده در این روش ها شامل به کارگیری توزیع های به دست آمده در فاز آموزش برای تخمین سیگنال تمیز از روی سیگنال نویزی است. برای این کار از تکنیک های تخمین که رایج ترین آنها MMSE و MAP است، استفاده می شود. همچنین برای سادگی بیشتر در آموزش و استفاده از مدل ها، مدل سازی روی روش های پارامتری کوتاه مدت مانند مدل AR انجام می گیرد.

می توان HMM را حالت کلی تر روش های GMM و VQ در نظر گرفت. روش GMM را می توان یک HMM یک حالتی در نظر گرفت که تنها تفاوت آن با HMM حذف فرض مارکوفی بودن سیگنال است. استفاده از مدل مارکوف چند حالتی برای مدل کردن گفتار مناسب تر است چرا که به نظر می رسد گفتار ذاتا دارای چند حالت است و این حالات مختلف ساختار غیر ایستانی گفتار و وابستگی زمانی و فرکانسی بین قاب های مختلف سیگنال را بهتر نمایش می دهند. از روش VQ نیز هر چند به صورت مستقل برای کاربردهای ذکر شده استفاده می شود، اما در بیشتر موارد به عنوان روشی برای مقداردهی اولیه آموزش استفاده شده است.

^۱ - Vector Quantisation

فصل پنجم

روش پیشنهادی بر مبنای توزیع
مخلوطی از لاپلاس ها

۵-۱- مقدمه

پس از توضیحاتی که در فصل‌های قبل درباره اصول روش‌های بهسازی گفتار و معرفی روش‌های متداول آماری بیان شد، در این فصل روش پیشنهادی مورد بررسی قرار خواهد گرفت. این روش که در حوزه روش‌های آماری جای می‌گیرد، مبتنی بر توزیع مخلوط لاپلاس است. روش پیشنهادی را می‌توان به سه بخش کلی زیر تقسیم کرد:

(۱) به دست آوردن رابطه تخمین گر MMSE با فرض توزیع مخلوط لاپلاس برای سیگنال تمیز.

(۲) الگوریتم امید ریاضی بیشینه (EM) برای تخمین پارامترهای توزیع مخلوط لاپلاس.

(۳) تخمین پارامترهای نویز

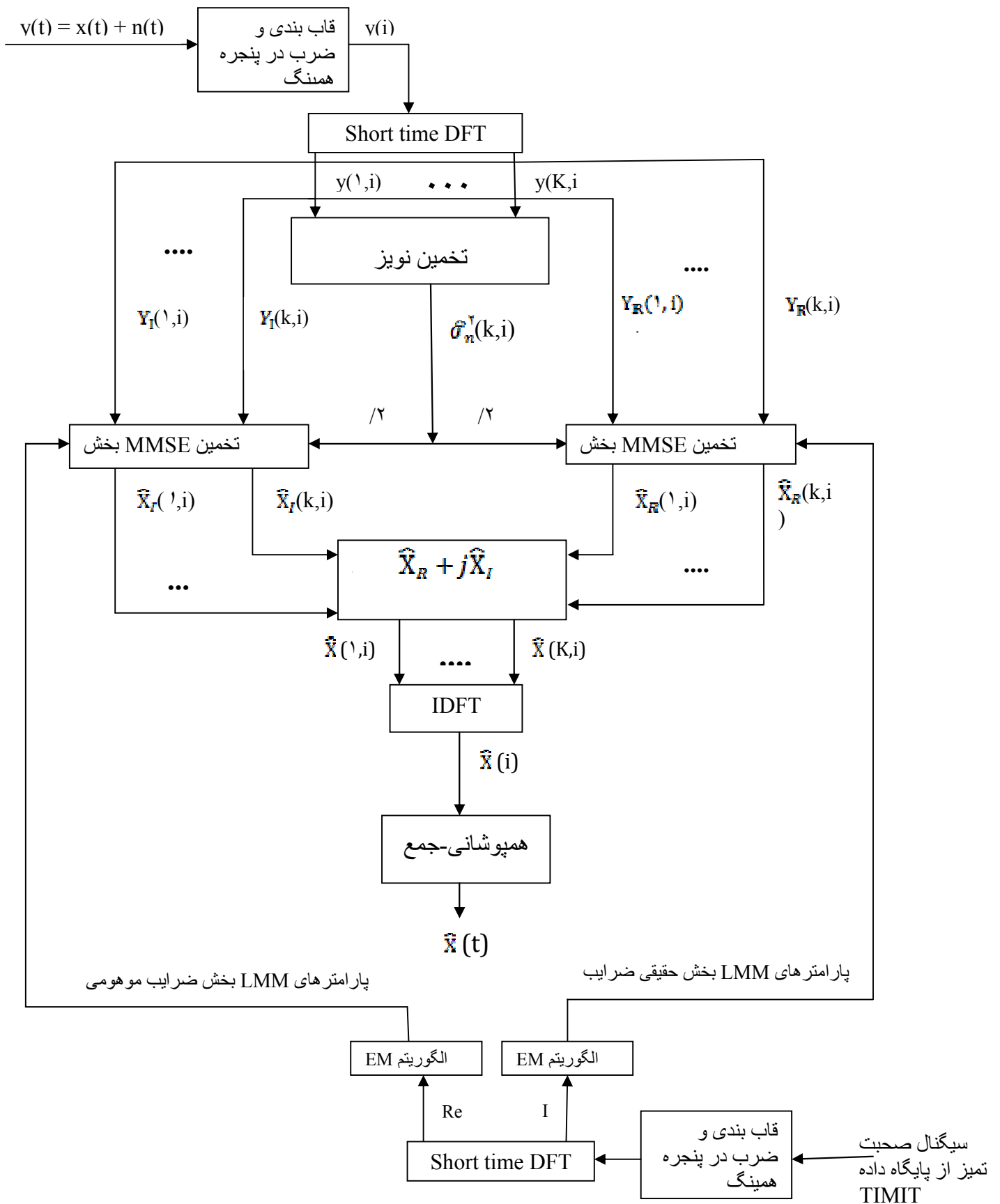
بلوک دیاگرام روش پیشنهادی در شکل (۵-۱) آورده شده است. کلیه بخش‌های این بلوک در ادامه فصل و بعد از بخش معرفی پایگاه داده، توضیح داده خواهند شد.

در انتهای فصل نیز به بررسی نتایج ارزیابی روش پیشنهادی و مقایسه آن با سه روش دیگر بهسازی گفتار براساس معیارهای ارزیابی بیان شده، پرداخته می‌شود.

۵-۲- معرفی پایگاه داده

در این پایان‌نامه از پایگاه داده TIMIT^۱ استفاده شده است. علت استفاده از این پایگاه داده معتبر و پر استفاده بودن آن است. در واقع در بیشتر مقالات معتبر در زمینه بهسازی گفتار از این پایگاه داده استفاده شده است. این پایگاه داده مجموعه‌ای از جملات است، که برای مطالعات در حوزه گفتار، علی‌الخصوص تشخیص خودکار گفتار ایجاد شده است. TIMIT شامل ۶۳۰۰ جمله می‌باشد. این ۶۳۰۰ جمله در واقع ده جمله است، که به وسیله ۶۳۰ گوینده از هشت منطقه از ایالات متحده

^۱ - Texas Instruments and Massachusetts Institute of Technology



شکل (۵-۱): بلوک دیاگرام روش پیشنهادی .

آمریکا، با لهجه های متفاوت، ادا شده است. این افراد، کسانی بودند که از بچگی در آن منطقه زندگی کرده اند و کاملاً بر لهجه آن منطقه مسلط می باشند.

جملات این پایگاه داده به دو دسته آموزشی و تست تقسیم شده اند که بر اساس نام و جنسیت گوینده، نامگذاری شده اند. پایگاه داده TIMIT نتیجه تلاش مشترک MIT^۱ و TI^۲ است [۶۵].

۵-۳- تخمین طیف سیگنال تمیز

در فصل گذشته روش های بهسازی آماری توضیح داده شدند. روش پیشنهادی در این پایان نامه در حوزه روش های آماری است که در آن از تخمین گر MMSE استفاده شده است. همچنین قبلاً تخمین طیف سیگنال تمیز هم با فرض گوسی بودن مولفه های طیف سیگنال تمیز و سیگنال نویز (فیلتر وینر) و هم با فرض غیر گوسی بودن مولفه های طیف سیگنال تمیز (گاما) و گوسی بودن طیف سیگنال نویز بررسی شد. در روش پیشنهادی نیز طیف سیگنال تمیز با فرض توزیع مخلوطی از لاپلاس ها برای سیگنال تمیز و توزیع گوسی برای نویز، به دست خواهد آمد.

۵-۳-۱- فرضیات

قبل از هر چیز لازم است فرضیات اساسی به کار رفته در روش پیشنهادی مطرح شود. در روش پیشنهادی فرض می شود که مولفه های طیف سیگنال تمیز مستقل از هم بوده [۴۲] و تابع چگالی احتمال آنها مخلوطی از لاپلاس ها با میانگین غیر صفر می باشد. فرض دیگر این است که نویز از سیگنال تمیز مستقل بوده [۴۲] و مولفه های طیف آن از توزیع گوسی با میانگین صفر پیروی می کنند. فرض بعدی که به طور عمده ای مسئله تخمین را ساده تر می کند، مستقل بودن بخش حقیقی و موهومی طیف سیگنال است [۴۴].

برای به دست آوردن روابط تخمین گر MMSE، مساله به صورت زیر مطرح می شود. ابتدا فرض می

شود که سیگنال گفتار تمیز $s(n)$ در اثر جمع شدن با نویز $d(n)$ سیگنال نویزی $y(n)$ را طبق

^۱ - Massachusetts Institute of Technology

^۲ - Texas Instruments

رابطه (۱-۵) تولید کرده است.

$$y(i)=s(i)+n(i) \quad (1-5)$$

سپس سیگنال ها به وسیله تبدیل فوریه زمان کوتاه به حوزه فرکانس برده می شود، یعنی سیگنال به قاب هایی که با هم همپوشانی^۱ دارند تقسیم شده و سپس قاب ها در یک پنجره مثل همینگ^۲ یا هنینگ^۳ ضرب شده و از آنها تبدیل فوریه گرفته می شود. علت این که سیگنال به قاب های کوچک تقسیم می شود، این است که سیگنال صوت عملاً یک سیگنال غیر ایستان است، ولی در این قاب های کوچک ایستان فرض می شود. در این پایان نامه طول قاب ها ۳۲ میلی ثانیه و با همپوشانی ۵۰٪ انتخاب شده و از پنجره همینگ برای ضرب در قاب های سیگنال، استفاده می شود [۲۲] و [۳۷]. پنجره همینگ و طول قاب ۳۲ میلی ثانیه مواردی است که در معتبرترین منابع و مراجع بهسازی گفتار پیشنهاد شده است. رابطه (۱-۵) به صورت زیر در حوزه فرکانس بازنویسی می شود:

$$Y(m, k) = S(m, k) + N(m, k) \quad (2-5)$$

که در آن m شماره قاب و k اندیس متناظر با فرکانس $\frac{2k\pi}{L}$ می باشد. (L تعداد نمونه های یک قاب است). برای راحتی نمایش تبدیل فوریه زمان کوتاه سیگنال نویزی، سیگنال تمیز و نویز را در قاب m ام و فرکانس k ام، به ترتیب با Y, S, N و نمایش می دهیم. $Y_R, Y_I, S_R, S_I, N_R, N_I$ به ترتیب نمایش دهنده بخش حقیقی و موهومی سیگنال نویزی، سیگنال تمیز و نویز هستند.

۵-۳-۲- مدل کردن سیگنال تمیز با توزیع مخلوط لاپلاس

در این بخش، برای واضح تر شدن مطالب، ابتدا توضیح لاپلاس بررسی شده است.

توزیع لاپلاس به صورت زیر نمایش داده می شود [۶۴]:

$$L(x.c.m) = ce^{-\tau c|x-m|} \quad (3-5)$$

^۱ - overlapp
^۲ - Hamming
^۳ - Hanning

که در آن $\frac{1}{c}$ بیانگر واریانس و m بیانگر میانگین توزیع می باشد. شکل (۵-۲) توزیع لاپلاس را به

ازای میانگین و واریانس های متفاوت نشان می دهد.

همانطور که قبلا ذکر شد، تابع چگالی احتمال سیگنال صحبت در حوزه زمان بهتر است با توزیع گاما یا لاپلاس به جای توزیع گوسی مدل شود. در حوزه تبدیل فوریه زمان کوتاه (سایز پنجره کمتر از ۱۰۰ میلی ثانیه) نیز توزیع های گاما و لاپلاس برای مدل کردن بخش حقیقی و موهومی تابع چگالی احتمال ضرایب طیف سیگنال صحبت مناسب تر می باشند [۴۳].

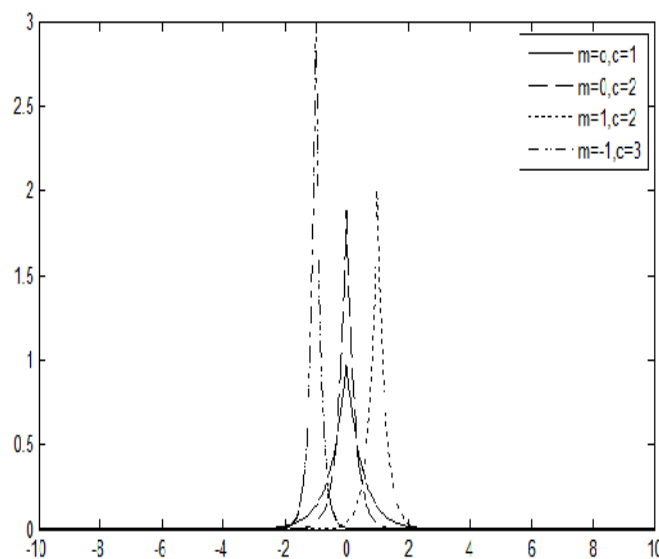
۵-۳-۲-۱- توزیع مخلوط لاپلاس (LMM)

توزیع مخلوطی از لاپلاس ها (LMM) به صورت زیر تعریف می شود [۴۶]:

$$p(x) = \sum_{i=1}^N \alpha_i L(x; c_i, m_i) = \sum_{i=1}^N \alpha_i c_i e^{-c_i |x - m_i|} \quad (4-5)$$

که در آن α_i ، m_i و c_i به ترتیب ضریب، میانگین، و واریانس هر کدام از لاپلاس ها و N تعداد لاپلا-س هاست و داریم: $\sum_{i=1}^N \alpha_i = 1$. متداول ترین روش برای تخمین این پارامترها EM الگوریتم است که در بخش بعدی توضیح داده خواهد شد.

در شکل های (۵-۳) و (۵-۴) برای نمونه هیستوگرام^۱ بخش حقیقی و موهومی تبدیل فوریه زمان

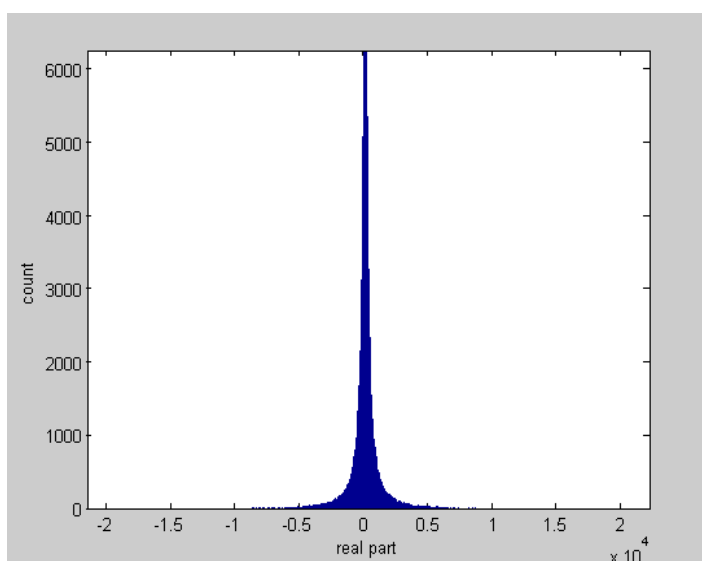


^۱ - Histogram

شکل (۲-۵): توزیع لاپلاس به ازای واریانس و میانگین های متفاوت .

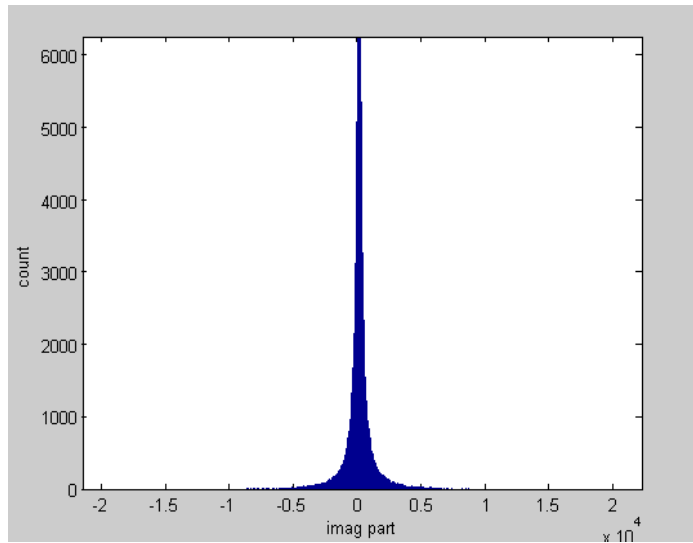
کوتاه یک میلیون داده بدون نویز که از پایگاه داده TIMIT استخراج شده اند، رسم شده است. همانطور که مشخص است، هیستوگرام ها شباهت زیادی به توزیع لاپلاس دارند. ولی با بررسی دقیق تر نمودارها و مقادیر آنها، و مقایسه آنها با نمودار توزیع لاپلاس که در شکل (۲-۵) آورده شد، مشخص می شود که توزیع داده ها را نمی توان تنها با یک لاپلاس با میانگین صفر تقریب زد. زیرا هیستوگرام ها حول مبدا متقارن نیستند، ولی توزیع لاپلاس با میانگین صفر حول مبدا متقارن است. از طرفی اگر از یک لاپلاس با میانگین غیر صفر هم استفاده شود، دنباله های^۱ توزیع داده ها به درستی مدل نمی شود. همچنین مجموع چند لاپلاس با میانگین صفر نیز باز حول مبدا متقارن بوده و برای تقریب مناسب نیست. به همین جهت به نظر می رسد، مناسب ترین گزینه برای مدل کردن طیف سیگنال تمیز توزیع مخلوطی از چند لاپلاس با میانگین های غیر صفر باشد.

در شکل (۵-۵) هیستوگرام بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز به ترتیب با یک لاپلاس با میانگین غیر صفر، مخلوطی از شش لاپلاس با میانگین غیر صفر و مخلوطی از شش گوسی با میانگین غیر صفر مدل شده است. (تصویر برای وضوح بیشتر در حالت بزرگنمایی است) با توجه به



^۱ - tail

شکل (۳-۵): هیستوگرام بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز .



شکل (۴-۵): هیستوگرام بخش موهومی تبدیل فوریه زمان کوتاه سیگنال تمیز .

شکل مشاهده می شود که توزیع مخلوطی از لاپلاس ها تقریب بهتری برای هیستوگرام داده ها، نسبت به سایرین می باشد. لذا در این پایان نامه، توزیع مخلوطی از لاپلاس ها با میانگین غیر صفر برای مدل کردن بخش حقیقی و موهومی پیشنهاد می شود.

لازم به ذکر است که، پیش از این از توزیع مخلوط لاپلاس در پردازش تصویر و حوزه های دیگر پردازش صوت، مثل تشخیص گفتار، استفاده شده است. ولی در زمینه بهسازی گفتار، تحقیقی در این زمینه مشاهده نشده است.

با توجه به توضیحات قبل، پیشنهاد می شود که مطابق روابط (۵-۵) و (۶-۵) بخش حقیقی و موهومی طیف سیگنال تمیز با LMM مدل شوند. یعنی:

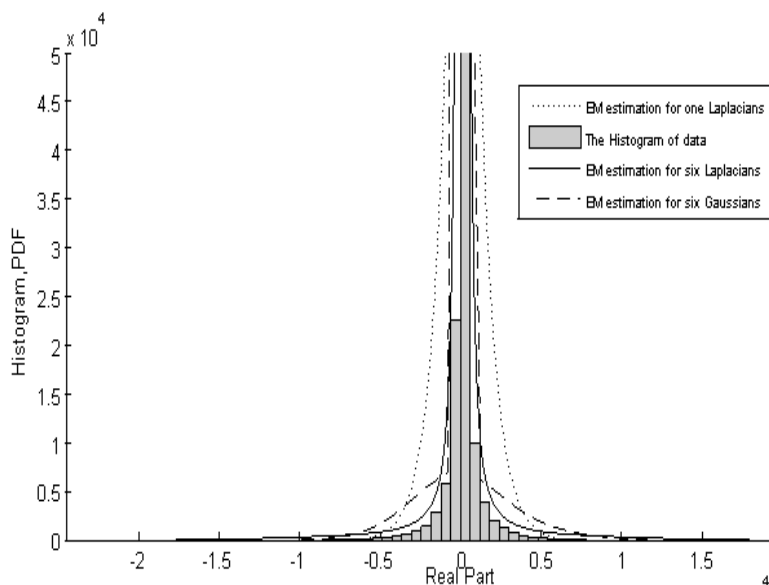
$$p(S_R) = \sum_{i=1}^N \alpha_i c_i e^{-\gamma c_i |S_R - m_i|} \quad (۵-۵)$$

$$p(S_I) = \sum_{j=1}^N \alpha_j c_j e^{-\gamma c_j |S_I - m_j|} \quad (۶-۵)$$

همچنین برای نویز، توزیع گوسی با میانگین صفر به صورت روابط (۷-۵) و (۸-۵) پیشنهاد می شود.

$$P(N_R) = \frac{1}{\sqrt{\pi} \sigma_n} \exp\left(\frac{-(N_R)^2}{\sigma_n^2}\right) \quad (۷-۵)$$

$$P(N_I) = \frac{1}{\sqrt{\pi} \sigma_n} \exp\left(\frac{-(N_I)^2}{\sigma_n^2}\right) \quad (۸-۵)$$



شکل (۵-۵): نقطه چین: تخمین توزیع بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز با یک لاپلاس، خط یکپارچه : تخمین توزیع بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز با مخلوطی از شش لاپلاس، خط چین: تخمین توزیع بخش حقیقی تبدیل فوریه زمان کوتاه سیگنال تمیز با مخلوطی از شش گوسی.

۵-۳-۳- تخمین طیف سیگنال تمیز با تخمین گر MMSE

در فصل گذشته روابط تخمین گر MMSE برای تخمین طیف سیگنال تمیز، با فرض توزیع گاما برای مولفه های طیف، به دست آمد. تخمین سیگنال تمیز از روی سیگنال نویزی با تخمین گر MMSE همانطور که قبلا توضیح داده شد، با فرض مستقل بودن بخش حقیقی و موهومی طیف سیگنال [۴۳]، به صورت زیر خواهد بود:

$$E\{S|Y\} = E\{S_R|Y_R\} + jE\{S_I|Y_I\} \quad (۹-۵)$$

تخمین هر کدام از بخش های حقیقی و موهومی به صورت جداگانه و ترکیب آنها، تخمین نهایی طیف را خواهد داد. تخمین MMSE بخش حقیقی طیف به صورت رابطه (۱۰-۵) می باشد.

$$\bar{S}_R = E\{S_R|Y_R\} = \int_{-\infty}^{+\infty} S_R P(S_R|Y_R) dS_R = \frac{\int_{-\infty}^{+\infty} S_R P(Y_R|S_R) P(S_R) dS_R}{P(Y_R)} \quad (۱۰-۵)$$

که برای محاسبه آن، ابتدا باید $P(Y_R|S_R)$ و $P(Y_R)$ محاسبه شوند. $P(S_R)$ نیز توزیع مخلوط لاپلاس است که پارامترهای آن باید محاسبه شوند. با توجه به اینکه بخش حقیقی طیف سیگنال نویزی برابر

مجموع بخش های حقیقی طیف سیگنال تمیز و نویز است و با فرض گوسی بودن توزیع نویز

$P(Y_R|S_R)$ به صورت زیر در نظر گرفته می شود:

$$P(Y_R|S_R) = \frac{1}{\sqrt{\pi}\sigma_n} \exp\left(\frac{-(Y_R - S_R)^2}{\sigma_n^2}\right) \quad (11-5)$$

با جایگذاری رابطه (11-5) به جای $P(Y_R|S_R)$ و $P(S_R)$ برای $P(S_R)$ در رابطه (10-5)، این رابطه به

صورت رابطه (12-5) درخواهد آمد.

$$E(S_R|Y_R) = \frac{1}{\sqrt{\pi}\sigma_n P(Y_R)} \int_{-\infty}^{+\infty} S_R \sum_{i=1}^N \alpha_i c_i \exp\left(\frac{-(Y_R - S_R)^2}{\sigma_n^2}\right) \exp(-\gamma c_i |S_R - m_i|) ds_R \quad (12-5)$$

برای محاسبه این انتگرال، از نتایج انتگرال های حل شده در [43] استفاده شده است. این انتگرال

ها با تغییر متغیر و اندکی تغییرات و دستکاری به صورت انتگرال (12-5) در می آیند. سپس از

جواب های این انتگرال ها برای حل رابطه (12-5) استفاده می شود. روابط (13-5) و (14-5) نشان

دهنده این انتگرال هاست.

$$\int_{-\infty}^{+\infty} S_R \exp\left(\frac{-(Y_R - S_R)^2}{\sigma_n^2}\right) \exp(-\gamma c |S_R|) ds_R = \frac{\sqrt{\pi}\sigma_n^2 \exp(c^2 \sigma_n^2)}{\gamma} [L_{R\gamma}(Y_R) \exp(\gamma c Y_R) \operatorname{erfc}(L_{R\gamma}(Y_R)) - L_{R\gamma}(Y_R) \exp(-\gamma c Y_R) \operatorname{erfc}(L_{R\gamma}(Y_R))] \quad (13-5)$$

$$\int_{-\infty}^{+\infty} \exp\left(\frac{-(Y_R - S_R)^2}{\sigma_n^2}\right) \exp(-\gamma c |S_R|) ds_R = \frac{\sqrt{\pi}\sigma_n \exp(c^2 \sigma_n^2)}{\gamma} [\exp(\gamma c Y_R) \operatorname{erfc}(L_{R\gamma}(Y_R) + \exp(-\gamma c Y_R) \operatorname{erfc}(L_{R\gamma}(Y_R)))] \quad (14-5)$$

دراین روابط erfc متمم تابع خطا (complementary error function) است که از رابطه زیر محاسبه

می شود:

$$\operatorname{erf}(x) = \frac{\gamma}{\sqrt{\pi}} \int_x^{\infty} e^{-t^2} dt = 1 - \operatorname{erfc}(x) \quad (15-5)$$

و $L_{R\gamma}(Y_R)$ و $L_{R\gamma}(Y_R)$ به صورت زیر تعریف می شوند:

$$L_{R\gamma}(Y_R) = c\sigma_n + \frac{Y_R}{\sigma_n} \quad (16-5)$$

$$L_{R^*}(Y_R) = c\sigma_n - \frac{Y_R}{\sigma_n} \quad (17-5)$$

با عوض کردن جای انتگرال و سیگما در رابطه (5-12) و تغییر متغیر در آن به صورت زیر :

$$S_R - m_i = t_i \Rightarrow S_R = t_i + m_i \quad , \quad ds_R = dt_i \quad (18-5)$$

و با استفاده از روابط (5-13) تا (5-17)، رابطه (5-12) به صورت رابطه (5-19) در می آید.

$$\begin{aligned} E(S_R|Y_R) = & \frac{1}{\sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2)} \sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2) \{ \sigma_n [L_{R^*}(Y_R - m_i) \exp(\tau c_i (Y_R - m_i)) \\ & \operatorname{erfc}(L_{R^*}(Y_R - m_i)) - L_{R^*}(Y_R - m_i) \exp(-\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i))] \\ & + m_i [\exp(\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i)) + \exp(-\tau c_i (Y_R - m_i)) \\ & \operatorname{erfc}(L_{R^*}(Y_R - m_i))] \} \end{aligned} \quad (19-5)$$

از رابطه (5-19) مشخص می شود که، $P(Y_R)$ نیز باید محاسبه شود. از آنجایی که $P(Y_R)$ برابر با

رابطه $\int_{-\infty}^{+\infty} P(S_R) P(Y_R|S_R) dS_R$ است، با توجه به روابط (5-13) تا (5-17) و تغییر متغیر قبل،

رابطه (5-20) برای محاسبه آن به دست می آید.

$$\begin{aligned} P(Y_R) = & \frac{1}{\sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2)} \sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2) [\exp(\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i)) + \\ & \exp(-\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i))] \end{aligned} \quad (20-5)$$

با جایگذاری $P(Y_R)$ در رابطه (5-17) تخمین نهایی بخش حقیقی طیف به صورت رابطه (5-21) به

دست خواهد آمد.

$$\bar{S}_R = E(S_R|Y_R) = \quad (21-5)$$

$$\frac{\sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2) \left[\sigma_n (L_{R^*}(Y_R - m_i) \exp(\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i)) - L_{R^*}(Y_R - m_i) \exp(-\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i))) + m_i (\exp(\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i)) + \exp(-\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i))) \right]}{\sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2) [\exp(\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i)) + \exp(-\tau c_i (Y_R - m_i)) \operatorname{erfc}(L_{R^*}(Y_R - m_i))]}$$

برای بخش موهومی طیف سیگنال هم به طریق مشابه عمل شده و رابطه (5-22) به دست می آید.

$$\hat{S}_I = E(S_I|Y_I) = \quad (22-5)$$

$$\frac{\sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2) \left[\sigma_n (L_{R^*}(Y_I - m_i) \exp(\tau c_i (Y_I - m_i)) \operatorname{erfc}(L_{R^*}(Y_I - m_i)) - L_{R^*}(Y_I - m_i) \exp(-\tau c_i (Y_I - m_i)) \operatorname{erfc}(L_{R^*}(Y_I - m_i))) + m_i (\exp(\tau c_i (Y_I - m_i)) \operatorname{erfc}(L_{R^*}(Y_I - m_i)) + \exp(-\tau c_i (Y_I - m_i)) \operatorname{erfc}(L_{R^*}(Y_I - m_i))) \right]}{\sum_{i=1}^N \alpha_i c_i \exp(c_i^2 \sigma_n^2) [\exp(\tau c_i (Y_I - m_i)) \operatorname{erfc}(L_{R^*}(Y_I - m_i)) + \exp(-\tau c_i (Y_I - m_i)) \operatorname{erfc}(L_{R^*}(Y_I - m_i))]}$$

دو رابطه (۲۱-۵) و (۲۲-۵) مطابق رابطه (۵-۹) با هم ترکیب شده و تخمین نهایی طیف را خواهند داد. سپس با تبدیل معکوس فوریه زمان کوتاه گرفتن از این تخمین، سیگنال بهسازی شده قاب بندی شده (ماتریسی که سطرها یا ستون های آن قاب های سیگنال هستند) در حوزه زمان به دست خواهد آمد. با اعمال روش همپوشانی-جمع بر روی این ماتریس به سیگنال اصلی خواهیم رسید. در روش همپوشانی - جمع، همپوشانی بین قاب ها حذف کرده و آنها را با هم ادغام می کنند.

۵-۴- تخمین پارامترهای توزیع سیگنال تمیز با EM الگوریتم

در قسمت قبلی گفته شد که برای تخمین سیگنال تمیز از روی سیگنال نویزی با روش MMSE، توزیع سیگنال تمیز مورد نیاز است. توزیع سیگنال تمیز نیز که LMM فرض شد، دارای پارامترهای نامعلومی مانند میانگین، واریانس و ضریب است که باید به دست آورده شوند. برای به دست آوردن این پارامترها از EM الگوریتم استفاده شده که در این بخش به معرفی آن پرداخته شده است. در واقع ابتدا، به وسیله داده های آموزشی مناسب و با EM الگوریتم پارامترهای LMM به دست می آید و سپس این پارامترها در رابطه های (۲۱-۵) و (۲۲-۵) جایگذاری می شود.

به طور کلی EM (امید ریاضی بیشینه) الگوریتم در مسائلی به کار می رود که می خواهیم مجموعه ای از پارامترهای θ را تخمین بزنیم که مبتنی بر یک توزیع احتمالی اند و تنها بخشی از داده هایی که توسط این توزیع احتمالی تولید شده اند را داریم نه کل آنها را. ارائه کامل این الگوریتم، در مقاله اولیه Dempster; Larid and Rubin (۱۹۷۷) [۳۴] و در کتاب های Little and Rubin (۱۹۸۷) و McLachlan and Krishnan (۱۹۹۷) یافت می شود. در این بخش به تشریح روش بیشترین شباهت و نحوه ارتباط آن با EM الگوریتم برای مسئله تخمین پارامترهای نامعلوم می پردازیم. سپس به تشریح تخمین پارامترهای تابع چگالی احتمال مخلوطی از لاپلاس ها خواهیم پرداخت.

۵-۴-۱- بیشترین شباهت

درباره این روش در فصل های قبل توضیح مختصری داده شد که در اینجا برای بررسی نحوه ارتباط آن با EM الگوریتم دوباره تکرار می شود. فرض می شود N داده با توزیع فرضی مشخصی با پارامترهای θ که نامعلوم و معین است، موجود باشند و هدف به دست آوردن پارامترهای توزیع فرضی مشخص باشد. همچنین فرض می شود که بردار داده ها مستقل و به طور یکسان توزیع شده^۱ (i. i. d) با توزیع p باشند. در نتیجه تابع چگالی داده ها به صورت زیر می باشد:

$$p(\mathbf{x}|\theta) = \prod_{i=1}^N p(X_i|\theta) = \mathcal{L}(\theta|\mathbf{x}) \quad (۲۳-۵)$$

که به آن تابع شباهت نیز می گویند. حال θ به گونه ای به دست می آید که تابع شباهت را طبق رابطه (۲۴-۵) بیشینه کند.

$$\theta^* = \operatorname{argmax} \mathcal{L}(\theta|\mathbf{x}) \quad (۲۴-۵)$$

معمولا به علت سادگی بیشتر از لگاریتم تابع شباهت در محاسبات استفاده می شود. درجه پیچیدگی مساله تخمین بستگی به فرم $p(\mathbf{x}|\theta)$ دارد. مثلا اگر این تابع یک گوسی و پارامترهای مجهول میانگین و واریانس آن باشد، به راحتی از لگاریتم تابع شباهت نسبت به پارامترهای مجهول مشتق گرفته و مساوی صفر قرار داده می شوند و پارامترها با حل این معادله به دست می آیند. ولی در اغلب مسایل به دست آوردن پارامترها با مساوی صفر قرار دادن مشتق، غیر ممکن است و باید از تکنیک های دیگری استفاده شود [۵].

۵-۴-۲- EM الگوریتم

EM الگوریتم یک روش کلی برای به دست آوردن تخمین بیشترین شباهت پارامترهای توزیعی است که داده های مربوط به آن توزیع کامل نیستند. یکی از کاربردهای مهم دیگر این الگوریتم در مواقعی است که بهینه کردن تابع شباهت یا لگاریتم آن، کار مشکلی باشد، اما بتوان شکل این تابع را

^۱ - independent and identically distributed.

با فرض وجود داده های دیگری ساده تر کرد. به این داده ها، داده های گمشده^۱ یا مشاهده نشده گفته می شود. فرض می شود، کل داده ها اجتماع داده های مشاهده شده^۲ (که به آنها، داده های ناقص^۳ نیز گفته می شود) و داده های گمشده باشند. داده های مشاهده شده با X نمایش داده می شوند و فرض می شود که این داده ها از توزیع خاصی پیروی می کنند. داده های گمشده فرضی نیز با Y نشان داده شده و کل داده ها هم با $Z=(X,Y)$ نمایش داده می شوند. تابع توزیع توأم X و Y که به آن تابع توزیع داده های کامل^۴ نیز گفته می شود، به صورت زیر تعریف می شود:

$$p(z|\theta) = p(x,y|\theta) = (y|x,\theta)p(x|\theta) \quad (25-5)$$

حال تابع شباهت داده های کامل به صورت $L(\theta|z) = L(\theta|x,y) = p(x,y|\theta)$ و تابع شباهت داده های ناقص (مشاهده شده) به صورت $L(\theta|x) = p(x|\theta)$ تعریف می شوند. در این روش سعی می شود داده های مشاهده نشده به صورتی تعریف شوند که شکل تابع شباهت داده های کامل ساده تر از تابع شباهت داده های ناقص باشد. این نکته قابل توجه است که، تابع شباهت داده های کامل یک متغیر تصادفی است. زیرا داده های ناقص و پارامترهای نامعلوم، معین و داده های مشاهده نشده ماهیتی تصادفی دارند. EM الگوریتم ابتدا، امید ریاضی لگاریتم تابع شباهت داده های کامل را که به صورت رابطه (۲۶-۵) تعریف می شود، محاسبه می کند. این مرحله، مرحله E (محاسبه امید ریاضی) الگوریتم نامیده می شود.

$$Q(\theta, \theta^{(i-1)}) = E[\log P(x,y|\theta) | x, \theta^{(i-1)}] \quad (26-5)$$

که در آن $\theta^{(i-1)}$ پارامترهای جاری یا پارامترهایی هستند که از مرحله قبل تخمین محاسبه شده اند و برای محاسبه امید ریاضی و پارامترهای جدید استفاده می شوند، و θ پارامترهایی است که با بهینه شدن Q به دست می آیند. در واقع مقدار این امید ریاضی^۵ با استفاده از داده های ناقص و

^۱ - missing value
^۲ - observed
^۳ - incomplete
^۴ - complete
^۵ - Expectation

پارامترهای جاری محاسبه می شود. سمت راست رابطه (۵-۲۶) را می توان به صورت رابطه (۵-۲۷) بازنویسی کرد:

$$E[\log P(x, y|\theta) | x, \theta^{(i-1)}] = \int_{y \in Y} \log P(x, y|\theta) f(y|x, \theta^{(i-1)}) dy \quad (۵-۲۷)$$

که در آن $f(y|x, \theta^{(i-1)})$ تابع توزیع مرزی داده های مشاهده نشده می باشد که هم به پارامترهای جاری و هم به داده های مشاهده شده بستگی دارد. در بهترین حالت این توزیع، تابع ساده ای از پارامترهای جاری و شاید داده های مشاهده شد می باشد و در بدترین حالت ممکن است فراهم کردن آن کار بسیار پیچیده ای باشد. □ هم نشان دهنده فضایی است که y می تواند مقادیر خود را روی آن اختیار کند. سمت راست رابطه (۵-۲۷) یک تابع معین است که می توان از آن مشتق گرفت. در مرحله بعدی EM الگوریتم، $Q(\theta, \theta^{(i-1)})$ نسبت به θ بهینه می شود. یعنی از تابع $Q(\theta, \theta^{(i-1)})$ نسبت به θ مشتق گرفته می شود و این مشتق مساوی با صفر قرار داده می شود. به این مرحله M یا مرحله بیشینه کردن^۱ می گویند.

خلاصه EM الگوریتم را می توان بدین صورت بیان کرد: EM الگوریتم یک روش بیشینه کردن لگاریتم تابع شباهت داده های ناقص، با استفاده از ماکزیمم کردن امید ریاضی لگاریتم تابع شباهت داده های کامل، به صورت تکراری می باشد. الگوریتم با یک مقداردهی اولیه θ^E در فضای θ آغاز شده و امید ریاضی $L(\theta | x, y)$ با کمک داده های مشاهده شده و θ^E به عنوان پارامترهای توزیع داده های مشاهده شده، محاسبه می شود. این مرحله، مرحله E (محاسبه امید ریاضی) الگوریتم نامیده می شود. در مرحله بعدی، امید ریاضی محاسبه شده نسبت به θ بیشینه شده و θ جدید محاسبه می گردد. به این مرحله M یا مرحله بیشینه کردن می گویند. سپس با θ جدید امید ریاضی لگاریتم تابع شباهت دوباره محاسبه می شود. این دو مرحله اینقدر تکرار می شوند که مقدار پارامترها در دو تکرار متوالی برابر شوند. اما در عمل به خاطر کم شدن حجم محاسبات، الگوریتم را

^۱ - maximization

وقتی مقدار پارامترهای حساب شده در دو مرحله به مقدار کوچکی اختلاف داشته باشند، متوقف می کنند.

در EM الگوریتم:

- در هر تکرار EM الگوریتم، امید ریاضی لگاریتم تابع شباهت داده های کامل و لگاریتم تابع شباهت داده های ناقص حتما افزایش می یابد.
- EM الگوریتم، به جواب تخمین روش بیشترین شباهت برای بیشینه کردن تابع شباهت داده های ناقص همگرا می شود.
- در EM الگوریتم همگرایی به یک ماکزیمم محلی تضمین شده است.
- استخراج پارامترهای مدل از تابع شباهت داده های کامل به مراتب راحت از استخراج آنها از تابع شباهت داده های ناقص می باشد [۳۴].

۵-۴-۳- محاسبه پارامترهای توابع توزیع مخلوط

به دست آوردن پارامترهای مخلوطی از چگالی ها یکی از مهمترین و گسترده ترین کاربرد های EM الگوریتم در شناسایی الگو می باشد. در این قسمت به دست آوردن پارامترهای توابع توزیع مخلوط را برای حالت کلی به دست می آوریم و سپس آن را برای چگالی مخلوطی از لاپلاس ها بسط خواهیم داد.

مدل احتمالاتی زیر را در نظر بگیرید:

$$p(\mathbf{X}|\theta) = \sum_{i=1}^M \alpha_i p_i(\mathbf{X}|\theta_i) \quad (28-5)$$

که در آن p_i ها توابع چگالی احتمالی هستند که با ضرایب α_i با هم جمع شده اند و

$(\alpha_1, \dots, \alpha_M, \theta_1, \dots, \theta_M)$ پارامترهایی هستند که با توجه به قید $\sum_{i=1}^M \alpha_i = 1$ محاسبه خواهند شد.

لگاریتم تابع شباهت داده های ناقص، به صورت زیر تعریف می شود:

$$\log(\mathcal{L}(\theta|\mathbf{x})) = \log \prod_{i=1}^N p(\mathbf{X}_i|\theta) = \sum_{i=1}^N \log(\sum_{j=1}^M \alpha_j p_i(\mathbf{X}_i|\theta_j)) \quad (29-5)$$

که بهینه کردن آن مشکل می باشد، زیرا مجموعی از چندین لگاریتم است. اگر داده های مشاهده نشده را به صورت $y = \{y_i\}_{i=1}^N$ فرض کنیم، که مقادیر آن بیان کننده این است که کدام داده به وسیله کدام جز^۱ از چگالی احتمال مخلوط تولید می شود، تابع شباهت تا حد زیادی ساده خواهد شد. در واقع برای هر i داریم: $y_i \in \{1, \dots, M\}$ و اگر k امین نمونه با k امین جز مخلوط تولید شد، خواهیم داشت: $y_i = k$ در نتیجه رابطه (۲۹-۵) به صورت زیر در خواهد آمد:

$$\log(\mathcal{L}(\theta|x,y)) = \log(P(x,y|\theta)) = \sum_{i=1}^N \log(P(x_i|y_i)P(y)) = \sum_{i=1}^N \log(\alpha_{y_i} P_{y_i}(x_i|\theta_{y_i})) \quad (30-5)$$

که شکل خاصی از توابع چگالی مخلوط را می دهد و با انواع روش ها می توان آن را بهینه کرد ولی مشکل اصلی این است که مقدار y نامشخص است و اگر فرض شود یک بردار تصادفی است، روابط ساده تر خواهند شد. در ابتدا باید توزیع داده های مشاهده نشده بیان شوند. فرض می شود $\theta^g = (\alpha_1^g, \dots, \alpha_M^g, \theta_1^g, \dots, \theta_M^g)$ پارامترهای مناسب برای تابع شباهت $\mathcal{L}(\theta^g|x,y)$ است. با داشتن این پارامترها به راحتی می توان $P_j(x_i|\theta_j^g)$ را برای هر i و j حساب کرد. اگر α_j احتمال پیشین جز j ام مخلوط لاپلاس تعریف شود، با استفاده از قاعده بیزین می توان تابع چگالی احتمال داده های مشاهده نشده را به صورت رابطه (۳۱-۵) بیان کرد.

$$P(y_i|\theta^g, x_i) = \frac{\alpha_{y_i}^g P_{y_i}(x_i|\theta_{y_i}^g)}{p(x_i|\theta^g)} = \frac{\alpha_{y_i}^g P_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g P_k(x_i|\theta_k^g)} \quad (31-5)$$

و

$$p(y|x, \theta^g) = \prod_{j=1}^N p(y_j|x_j, \theta^g) \quad (32-5)$$

در حالی که $y = (y_1, \dots, y_N)$ تعدادی از داده های مشاهده نشده است که به طور مستقل بیرون کشیده شده اند.

حال می توان معادله (۲۶-۵) را به صورت زیر باز نویسی کرد:

$$Q(\theta, \theta^g) =$$

^۱ - Component

$$\begin{aligned}
\sum_{y \in Y} \log(\mathcal{L}(\theta|x, y)) p(y|x, \theta^{\otimes}) &= \sum_{y \in Y} \sum_{i=1}^N \log(\alpha_{y_i} p_{y_i}(x_i|\theta_{y_i})) \prod_{j=1}^N p(y_j|x_j, \theta^{\otimes}) \\
&= \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{i=1}^N \sum_{v=1}^M \delta_{v, y_i} \log(\alpha_i p_i(x_i|\theta_i)) \prod_{j=1}^N p(y_j|x_j, \theta^{\otimes}) = \\
\sum_{i=1}^M \sum_{v=1}^M \log(\alpha_i p_i(x_i|\theta_i)) \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \delta_{i, y_1} \prod_{j=1}^M p(y_j|x_j, \theta^{\otimes}) \quad (33-5)
\end{aligned}$$

و برای $l \in 1, \dots, M$ داریم:

$$\begin{aligned}
&\sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \delta_{l, y_1} \prod_{j=1}^M p(y_j|x_j, \theta^{\otimes}) \\
&= \left(\sum_{y_1=1}^M \dots \sum_{y_{l-1}=1}^M \sum_{y_{l+1}=1}^M \dots \sum_{y_N=1}^M \prod_{j=1, j \neq l}^M p(y_j|x_j, \theta^{\otimes}) \right) p(l|x_l, \theta^{\otimes}) \\
&= \prod_{j=1, j \neq l}^M \left(\sum_{y_j=1}^M p(y_j|x_j, \theta^{\otimes}) \right) p(l|x_l, \theta^{\otimes}) = p(l|x_l, \theta^{\otimes}) \quad (34-5)
\end{aligned}$$

و چون $\sum_{i=1}^M p(i|x_l, \theta^{\otimes}) = 1$ ، با استفاده از رابطه (34-5)، رابطه (33-5) را می توان به صورت رابطه (35-5) بازنویسی کرد. برای بیشینه کردن این عبارت، باید عبارات شامل α_1 و θ_1 به طور مستقل بیشینه شوند.

$$\begin{aligned}
Q(\theta, \theta^{\otimes}) &= \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l p_l(x_i|\theta_l)) p(l|x_i, \theta^{\otimes}) \\
&= \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l|x_i, \theta^{\otimes}) + \sum_{l=1}^M \sum_{i=1}^N \log(p_l(x_i|\theta_l)) p(l|x_i, \theta^{\otimes}) \quad (35-5)
\end{aligned}$$

برای به دست آوردن α_1 تحت قید $\sum_l \alpha_l = 1$ ، با کمک گرفتن از ضرایب لاگرانژ، معادله زیر به دست خواهد آمد:

$$\begin{aligned}
\frac{\partial}{\partial \alpha_l} \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l|x_i, \theta^{\otimes}) + \lambda \sum_l \alpha_l - 1 &= 0 \\
\sum_{i=1}^N p(l|x_i, \theta^{\otimes}) + \lambda &= 0 \Rightarrow \quad (36-5)
\end{aligned}$$

که از حل معادله قبل جواب زیر برای محاسبه ضرایب لاپلاس ها به دست خواهد آمد [5]:

$$\alpha_l = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \theta^{\otimes}) \quad (37-5)$$

۵-۴-۴- محاسبه پارامترهای LMM

در این قسمت با توجه به روابط (۳۵-۵) و (۳۷-۵) که در بخش قبل، برای یک توزیع دلخواه به دست آورده شد، تلاش می شود که پارامترهای توزیع LMM که در رابطه (۳۸-۵) آورده شده است، به دست آید.

$$p(\theta) = \sum_{i=1}^N \alpha_i L(\theta, c_i, \theta_i) = \sum_{i=1}^N \alpha_i c_i e^{-\gamma c_i |\theta_n - \theta_i|} \quad (۳۸-۵)$$

اگر تعداد لاپلاس ها N و تعداد کل داده ها T باشد، رابطه (۳۷-۵) برای محاسبه α_i ها یا همان ضرایب هر جز از مخلوط لاپلاس به صورت رابطه (۳۹-۵) به دست خواهد آمد که در آن $p(i|\theta_n)$ یا احتمال پیشین جز i ام به صورت رابطه (۴۰-۵) تعریف می شود.

$$\alpha_i = \frac{1}{T} \sum_{n=1}^T p(i|\theta_n) \quad (۳۹-۵)$$

$$p(i|\theta_n) = \frac{\alpha_i c_i e^{-\gamma c_i |\theta_n - \theta_i|}}{\sum_{i=1}^N \alpha_i c_i e^{-\gamma c_i |\theta_n - \theta_i|}} \quad (۴۰-۵)$$

پس از آن رابطه (۳۵-۵) برای به دست آوردن θ_i ها و c_i ها به صورت رابطه (۴۱-۵) تشکیل می شود و چون بخش اول رابطه مذکور به این دو پارامتر وابسته نبوده و در مشتق گیری نسبت به این دو پارامتر نقشی ندارد، فقط بخش دوم آن آورده شده است.

$$J(c_i, \theta_i) = \sum_{n=1}^T \sum_{i=1}^N (\log c_i - \gamma c_i |\theta_n - \theta_i|) p(i|\theta_n) \quad (۴۱-۵)$$

با مشتق گرفتن از رابطه فوق نسبت به c_i و مساوی صفر قرار دادن آن، رابطه (۴۲-۵) برای محاسبه c_i ، برای هر جز از مخلوط لاپلاس به دست خواهد آمد. (واریانس هر جز $\frac{1}{c_i^2}$ می باشد).

$$c_i = \frac{\sum_{n=1}^T p(i|\theta_n)}{\gamma \sum_{n=1}^T |\theta_n - \theta_i| p(i|\theta_n)} \quad (۴۲-$$

۵)

همچنین با مشتق گرفتن از رابطه (۴۱-۵) نسبت به θ_i و مساوی صفر قرار دادن آن، رابطه (۴۳-۵) برای محاسبه میانگین هر جز به دست خواهد آمد. مشاهده می شود که در این رابطه به خاطر وجود

تابع علامت (sgn)، به دست آوردن مقدار دقیق برای θ_i ممکن نیست. اما می توان مقدار جدید آن را با کمک تخمین قبلی $|\theta_n - \theta_i|$ به دست آورد. در نتیجه رابطه نهایی (۴۴-۵) برای تقریب میانگین هر جز به دست آمد [۴۶].

$$\frac{\partial J}{\partial \theta_i} = \sum_{n=1}^T (-r c_i \frac{\partial}{\partial \theta_i} |\theta_n - \theta_i|) p(i|\theta_n) =$$

$$\sum_{n=1}^T -r c_i \operatorname{sgn}(\theta_n - \theta_i) p(i|\theta_n) = 0 \quad \Rightarrow$$

$$\sum_{n=1}^T \operatorname{sgn}(\theta_n - \theta_i) p(i|\theta_n) = \sum_{n=1}^T \frac{\theta_n - \theta_i}{|\theta_n - \theta_i|} p(i|\theta_n) = 0 \quad (۴۳-۵)$$

$$\sum_{n=1}^T \frac{\theta_n}{|\theta_n - \theta_i|} p(i|\theta_n) = \sum_{n=1}^T \frac{\theta_i}{|\theta_n - \theta_i|} p(i|\theta_n)$$

$$\Rightarrow \theta_i = \frac{\sum_{n=1}^T \frac{\theta_n}{|\theta_n - \theta_i|} p(i|\theta_n)}{\sum_{n=1}^T \frac{1}{|\theta_n - \theta_i|} p(i|\theta_n)} \quad (۴۴-۵)$$

۵-۴-۵- پیاده سازی EM الگوریتم

همانطور که قبلا گفته شد، برای مدل کردن سیگنال بخش حقیقی و موهومی طیف سیگنال تمیز توزیع LMM پیشنهاد شد. همچنین گفته شد که برای محاسبه پارامترهای LMM با استفاده از داده های آموزشی از EM الگوریتم استفاده می شود.

برای پیاده سازی EM الگوریتم، از داده های آموزشی برای تخمین پارامترهای توزیع مورد نظر استفاده می شود. هر چه تعداد این داده ها بیشتر باشد، مدل دقیق تر و بهتری ایجاد می شود. در این پایان نامه از حدود ده دقیقه سیگنال تمیز از پایگاه داده TIMIT (۲۰۰ جمله از ۱۰۰ مرد و ۱۰۰ زن) استفاده شده است. این داده ها با نرخ دو کاهش داده شده است. سپس از داده ها، تبدیل فوریه زمان کوتاه گرفته و بخش حقیقی و موهومی طیف از کل قاب ها به عنوان داده آموزشی به طور جداگانه به الگوریتم داده شده تا پارامترهای LMM برای بخش ها حقیقی و موهومی به طور جداگانه محاسبه شود. در واقع EM الگوریتم دو بار اجرا می شود. یکبار با داده های بخش حقیقی و یکبار هم با داده های بخش موهومی و در نهایت دو LMM به دست خواهد آمد.

مقدار دهی اولیه EM الگوریتم برای به دست آوردن پارامترهای LMM را بدین طریق انجام شده است:

مقدار اولیه α_i ها

(ضرایب لاپلاس ها)، برابر $1/N$ در نظر گرفته شده است.

مقدار اولیه θ_i و c_i

را هم به صورت تصادفی انتخاب شده است. بدین صورت که از تابع rand

نرم افزار

MATLAB استفاده شده و خروجی آن برای θ_i در میانگین کل داده ها و برای c_i در واریانس داده ها ضرب شده است. به ازای هر N و برای پارامترهای LMM دو مقدار دهی به صورت تصادفی به طریقی که توضیح داده شد، انجام شده است و سپس با توجه به معیاری که در رابطه (۴۵-۵) بیان شده، [۴۴] جواب بهتر انتخاب شده است.

$$I_{KL} = \sum_{\mathbf{x}} P_H(\mathbf{x}) \log \left(\frac{P_H(\mathbf{x})}{P(\mathbf{x})} \right) \quad (45-5)$$

در این رابطه $P_H(\mathbf{x})$ هیستوگرام داده ها و $P(\mathbf{x})$ مقدار توزیع مخلوط لاپلاس به دست آمده با شرایط اولیه مذکور، به ازای X است. هر چه این معیار کمتر باشد، بیانگر این است که توزیع مخلوط لاپلاس به دست آمده، با خطای کمتری بر توزیع داده ها برازنده شده است.

همانطور که قبلاً گفته شد، EM الگوریتم یک الگوریتم تکراری است که شرط توقف آن می تواند تعداد تکرار خاص یا عدم تغییر محسوس پارامترها باشد.

در این پایان نامه، شرط توقف تعداد تکرار ۴۵ در نظر گرفته شده است. این مقدار با توجه به عدم تغییرات محسوس در مقادیر پارامترهای LMM، در آزمایشات متعددی که با N های متفاوت انجام شد، به دست آمد.

به عنوان مثال، جدول (۵-۱)، نشان دهنده پارامترهای LMM و I_{KL} ، برای بخش حقیقی تبدیل فوریه زمان کوتاه داده های TIMIT به ازای $N=3$ و بعد از ۱۰، ۳۰، ۴۰ و ۵۰ تکرار EM الگوریتم

است. K بیانگر تعداد تکرار الگوریتم می باشد. قبلا گفته شد که کمتر بودن مقدار I_{KL} بدین معنی است که توزیع به دست آمده با تقریب بهتری PDF^۱ داده ها را نمایش می دهد. با توجه به جدول، کم شدن مقدار I_{KL} در هر تکرار بیانگر این است که در هر مرحله تکرار الگوریتم توزیع به دست آمده

جدول (۵-۱): پارامترهای LMM و I_{KL} برای بخش حقیقی تبدیل فوریه زمان کوتاه داده های TIMIT به ازای $N=3$ و بعد از ۱۰، ۳۰، ۴۰ و ۵۰ تکرار EM الگوریتم.

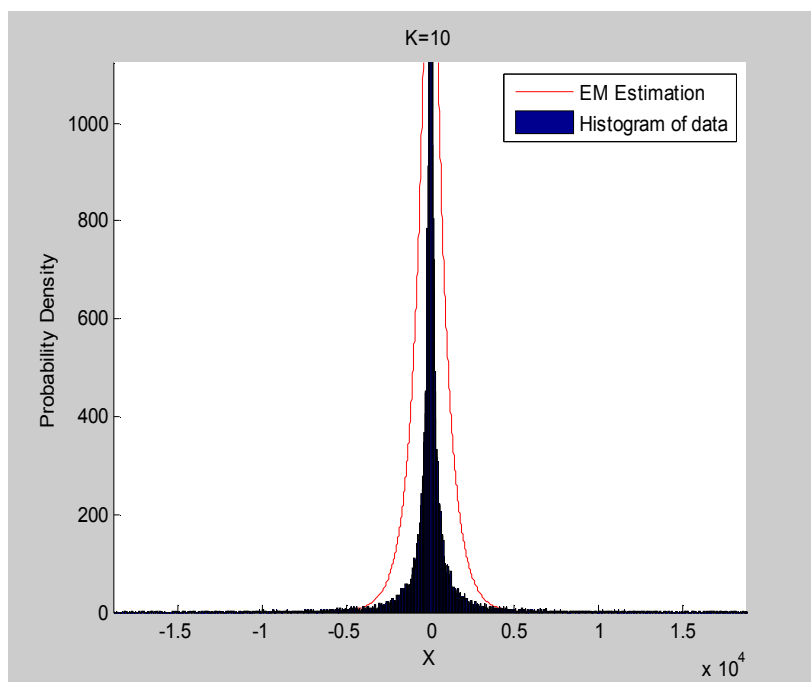
K	α_i	c_i	θ_i	I_{KL}
۱۰	۰.۰۸۲ ۰.۴۶۲۱ ۰.۴۵۵۹	۰.۰۰۰۱ ۰.۰۱۱۹ ۰.۰۰۱	۴۵.۴۴۶۹ ۱.۲۸۲۳ ۸.۸۷۵۹	۳۳۶۱۴۰
۳۰	۰.۰۸۳۷۰ ۰.۴۲۳۱ ۰.۴۹۳۲	۰.۰۰۰۱ ۰.۰۱۴۱ ۰.۰۰۱	۱۵۱.۳۲۱ ۰.۰۶۷۷ ۳.۳۵۱۶	۳۳۵۹۵۰
۴۰	۰.۰۸۵۴ ۰.۴۱۸۷ ۰.۴۹۵۹	۰.۰۰۰۱ ۰.۰۱۴۳ ۰.۰۰۱۱	۱۶۴.۲۸۷۷ -۰.۲۳۲ ۲.۹۸۸۱	۳۳۵۹۶۰
۵۰	۰.۰۸۶۲ ۰.۴۱۶۸ ۰.۴۹۷	۰.۰۰۰۱ ۰.۰۱۴۴ ۰.۰۰۱۱	۱۶۴.۸۱۱۲ -۰.۲۴۳۶ ۲.۹۱۵۷	۳۳۵۹۶۰

بهبتر بر PDF داده ها برازنده می شود و عدم تغییر محسوس آن در دو تکرار، نشان دهنده این است که توزیع های به دست آمده در دو تکرار زیاد متفاوت نیستند. همچنین با توجه به جدول، اختلاف همه پارامترهای LMM برای $K=40$ و $K=50$ بسیار کم است. (این اختلاف با افزایش N کمتر هم می شود) در حالیکه اختلاف این پارامترها برای $K=10$ و $K=30$ بیشتر است. نکته دیگر این است که با افزایش K زمان اجرای برنامه و پیچیدگی محاسباتی بیشتر می شود، به همین خاطر میزان تکرارها ۴۵ در نظر گرفته شده است.

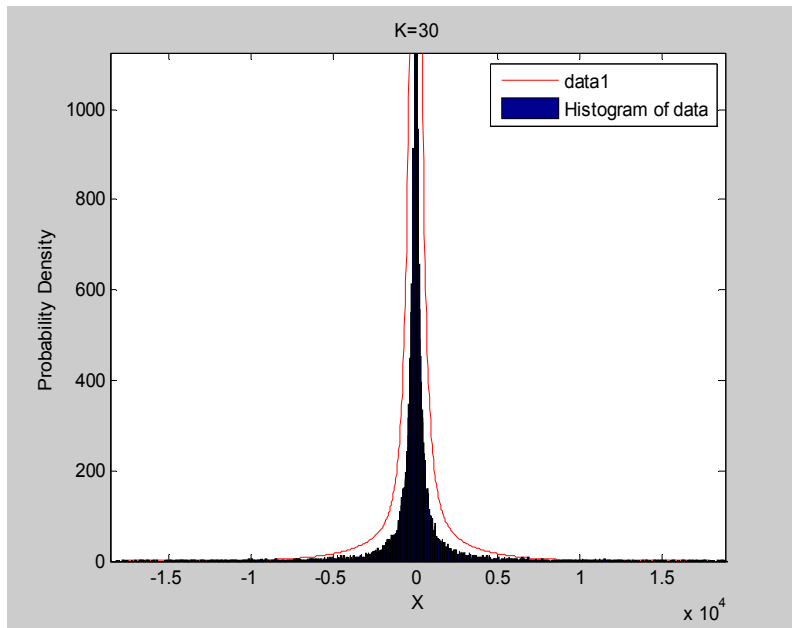
شکل های (۵-۶) نیز نتیجه EM الگوریتم را برای مثال قبل نمایش می دهند. همان طور که از نمودارها مشخص است، بین تخمین EM الگوریتم برای $K=10$ و $K=30$ و $K=40$ تفاوت نسبتا زیادی وجود دارد، در حالی که اختلاف برای $K=40$ و $K=50$ بسیار کم است.

^۱ - Probability Density Function

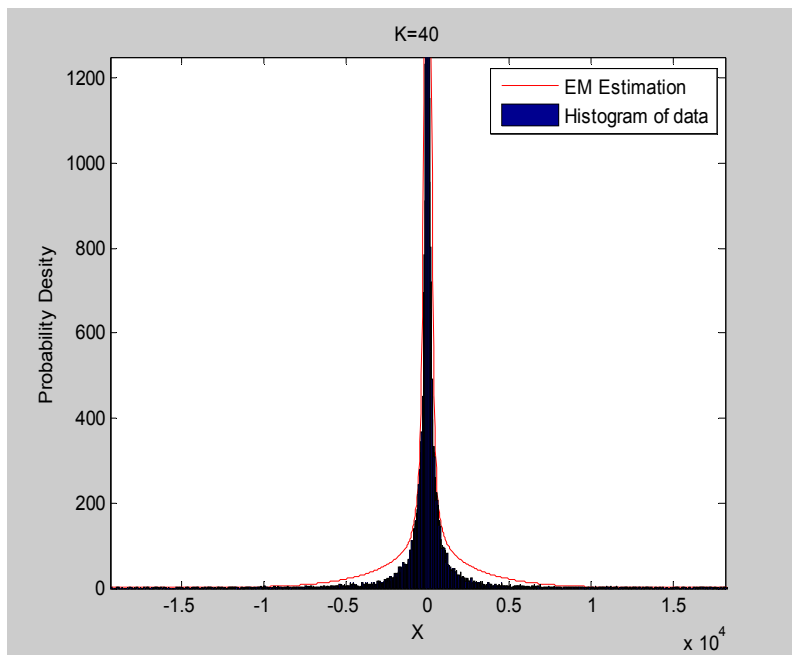
شکل های (۷-۵) و (۸-۵) نیز تخمین EM الگوریتم را برای بخش حقیقی و موهومی طیف سیگنال تمیز، به ازای N های متفاوت نشان می دهند. با توجه به نمودارها، هر چه N بزرگتر باشد، تخمین EM الگوریتم به هیستوگرام داده ها نزدیکتر است. برای اینکه با یک معیار دقیق و عددی نشان دهیم که با افزایش N منحنی تخمین زده شده بهتر بر PDF داده ها برازنده می شود، از معیار I_{KL} که قبلاً معرفی شد، استفاده نمودیم. مشاهده می شود که با افزایش N این مقدار کم می شود. مثلاً برای $N=1$ مقدار این معیار برابر 355140 و برای $N=10$ این معیار برابر با 335520 است. قبلاً ذکر شد که



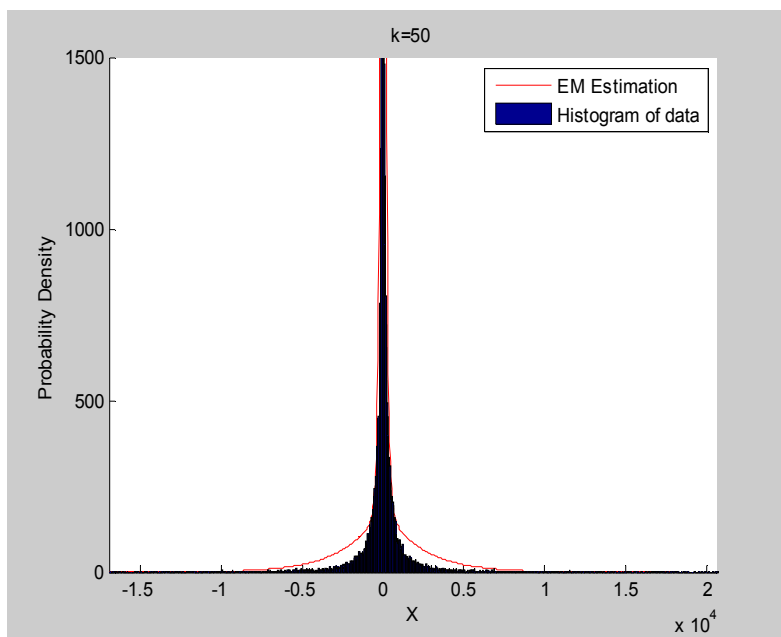
(الف)



(ب)

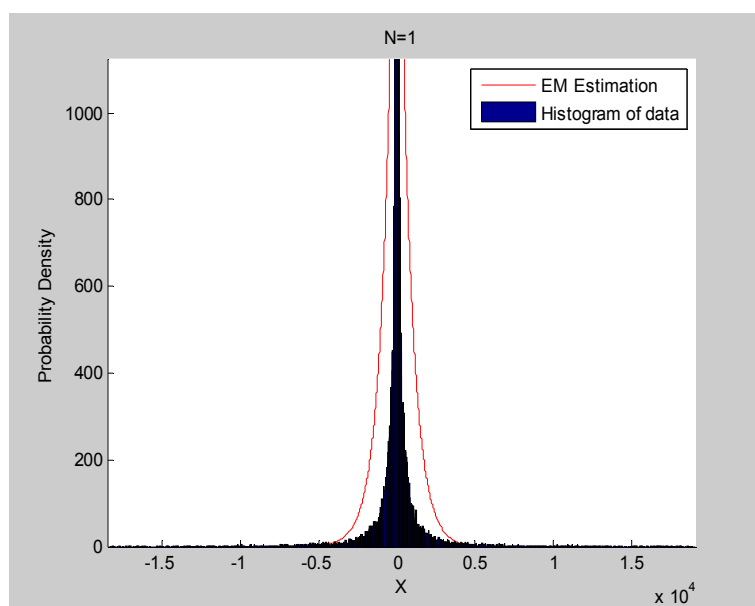


(ج)

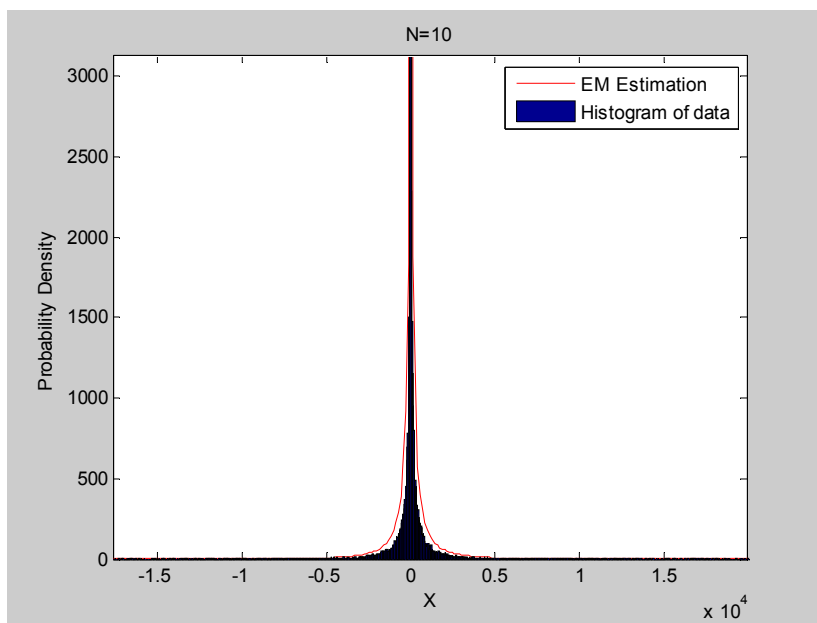


(چ)

شکل (۵-۶): تخمین EM الگوریتم برای بخش حقیقی طیف سیگنال تمیز با $N=3$ ، (الف) پس از ۱۰ بار تکرار، (ب) پس از ۳۰ بار تکرار، (ج) پس از ۴۰ بار تکرار، (چ) پس از ۵۰ بار تکرار.

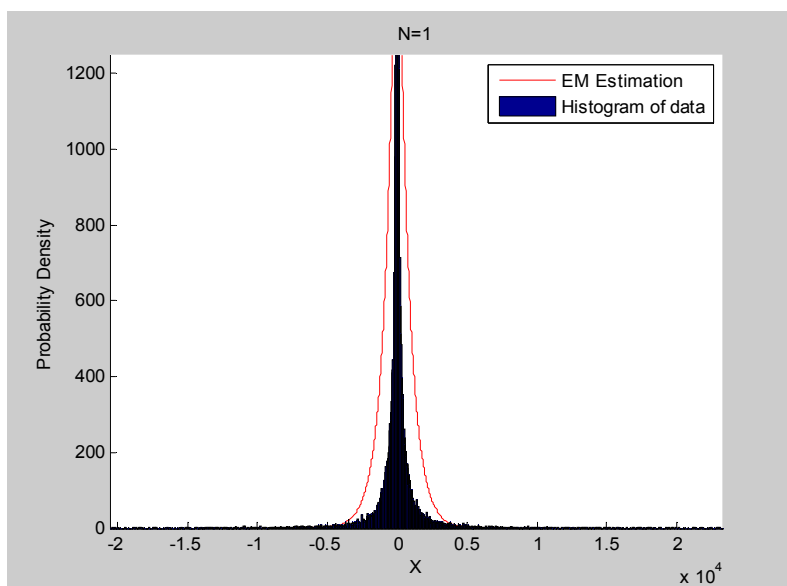


(الف)

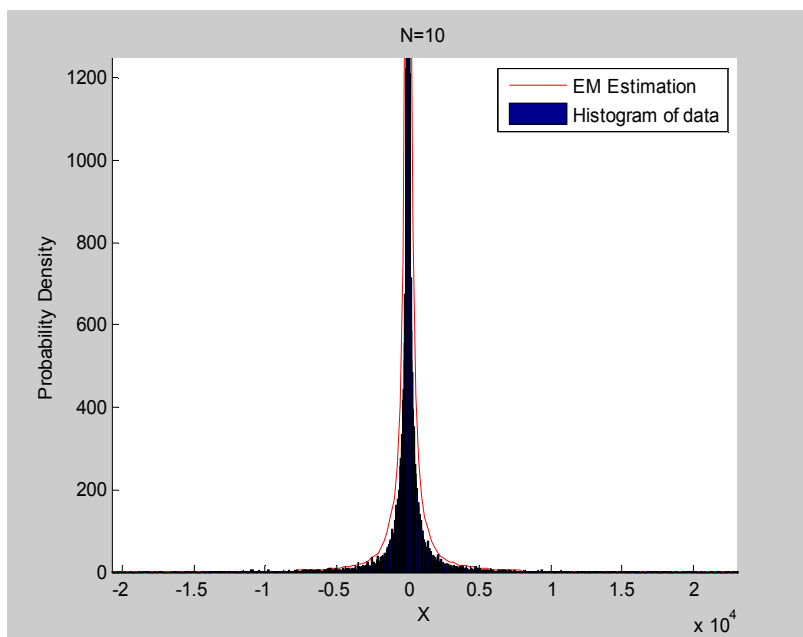


(ب)

شکل (۷-۵): تخمین EM الگوریتم برای بخش حقیقی طیف سیگنال تمیز: (الف) به ازای یک لاپلاس، (ب) به ازای ده لاپلاس.



(الف)



(ب)

شکل (۵-۸): تخمین EM الگوریتم برای بخش موهومی طیف سیگنال تمیز : (الف) به ازای یک لاپلاس، (ب) به ازای ده لاپلاس .

کمتر بودن این معیار به معنای بهتر بودن تقریب است که نشان می دهد $N=10$ نسبت به $N=1$ تقریب مناسب تری از سیگنال ارائه می دهد. نکته دیگر این است که با افزایش N تغییرات این معیار کم می شود. یعنی اختلاف این معیار به ازای دو N متفاوت کم و کمتر می شود. این نشان می دهد که با افزایش N اختلاف بین توزیع های به دست آمده از EM الگوریتم کم می شود. به عنوان مثال مقدار این معیار برای $N=1$ برابر 355140 ، برای $N=10$ برابر 335520 و برای $N=20$ برابر 32470 است. واضح است اختلاف این معیار برای دو حالت $N=10$ و $N=20$ کمتر از اختلاف این معیار برای دو حالت $N=1$ و $N=10$ می باشد.

۵-۵- تخمین نویز

همان طور که در رابطه تخمین MMSE مشاهده شد، برای به دست آوردن تخمینی از سیگنال تمیز، علاوه بر پارامترهای LMM، به واریانس نویز نیز نیاز است. با توجه به اینکه فقط سیگنال

نویزی و وارپانس آن در دسترس است، لازم است که وارپانس نویز به طریقی از آن استخراج گردد. روش های تخمین نویز در فصل های قبلی توضیح داده شد. در این روش ها فرض می شود که نویز غیرایستاد است. ولی تغییرات توان آن به آرامی و کندتر از توان سیگنال صحبت تمیز صورت می پذیرد.

برای تخمین نویز از روش Continuous Spectral Minimum Tracking (ردیابی کمینه ها، روش دوم) که در فصل قبل توضیح داده شد، استفاده گردیده است. همانطور که در فصل قبل گفته شد، در این روش مینیمم توان سیگنال نویزی به صورت پیوسته و بدون نیاز به پنجره دنبال می شود. در این روش تخمین نویز، با مقایسه طیف هموار شده سیگنال نویزی در هر فرکانس با تخمین نویز در همان فرکانس و در قاب قبلی، با یک رابطه هموارسازی غیرخطی، به روز می شود. در این روش برای دنبال کردن مینیمم توان سیگنال نویزی به جای پرلودوگرام سیگنال نویزی، از پرلودوگرام هموار شده آن استفاده می شود.

$$P(m, k) = \alpha p(m-1, k) + (1 - \alpha) |Y(m, k)|^2 \quad (46-5)$$

که در آن α فاکتور هموارسازی و در بازه $0.7 \leq \alpha \leq 0.9$ و k فرکانس و m شماره قاب می باشد. رابطه غیر خطی که برای تخمین نویز بر پایه دنبال کردن مینیمم طیف توان سیگنال نویزی، در هر فرکانس استفاده می شود، به صورت زیر است:

If $p_{min}(m-1, k) < p(m, k)$

$$p_{min}(m, k) = \lambda p_{min}(m-1, k) + \frac{1-\lambda}{1-\beta} (p(m, k) - \beta p(m-1, k))$$

Else

$$p_{min}(m, k) = p(m, k)$$

End

(47-5)

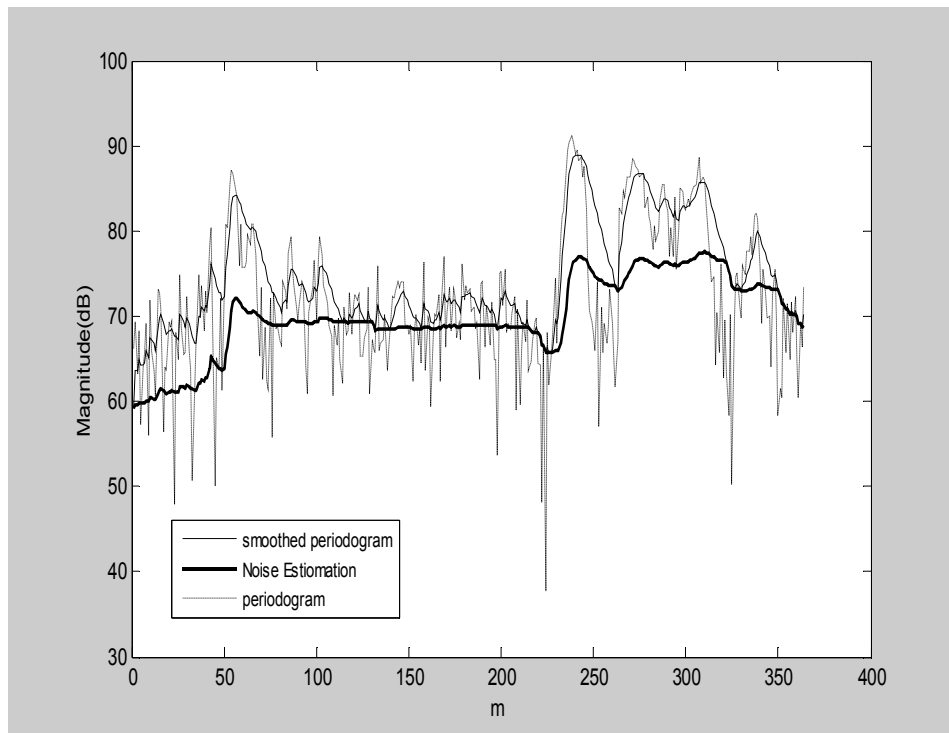
که در آن $p_{min}(m, k)$ تخمین نویز و $p(m, k)$ پرلودوگرام هموار شده سیگنال نویزی می باشد. برای پیاده سازی این روش با توجه به نتایج روش پیشنهادی، مقادیر زیر را پیشنهاد شده است: (این مقادیر بهترین جواب ها را در بخش ارزیابی روش پیشنهادی داشته است.)

$$\alpha = 0.7, \beta = 0.96, \lambda = 0.999$$

شکل (۵-۹) نمایش دهنده پریودوگرام سیگنال نویزی که با نویز $F16$ و در سیگنال به نویز ۵ دسی بل، آغشته شده است، پریودوگرام هموار شده و نویز تخمین زده شده به روش فوق، در فرکانس $k=10$ ، برای کلیه قاب ها می باشد.

۵-۶- ارزیابی روش پیشنهادی

پس از پیاده سازی روش پیشنهادی، برای بررسی و سنجش عملکرد روش، باید آن را با معیارهای مشخص ارزیابی کرد. در این پایان نامه، برای ارزیابی روش پیشنهادی ده جمله از پایگاه داده TIMIT که به وسیله پنج زن و پنج مرد گفته شده بودند، انتخاب شده و در آزمایش های مختلف، سه نویز سفید گوسی، جنگنده $F16$ و همهمه در سه سطح سیگنال به نویز ۰، ۵ و ۱۰ به سیگنال اضافه گردید. (نویز سفید گوسی به وسیله نرم افزار MATLAB تولید شده است) سپس سیگنال نویزی به قاب هایی به طول ۳۲ میلی ثانیه با همپوشانی ۵۰٪ تقسیم و در پنجره همینگ ضرب شد و از آن تبدیل فوریه زمان کوتاه گرفته شد. لازم به ذکر است که علاوه بر پنجره همینگ به طول ۳۲



شکل (۵-۹): خط چین: پرئودوگرام سیگنال نویزی که با نویز F1۶ و در سیگنال به نویز ۵ دسی بل، آغشته شده است. خط یکپارچه نازک: پرئودوگرام هموار شده و خط یکپارچه کلفت: تخمین نویز.

میلی ثانیه، از پنجره هنینگ به طول ۳۲ میلی ثانیه و پنجره همینگ به طول ۱۶ میلی ثانیه نیز برای پیاده سازی روش پیشنهادی استفاده شد که سیگنال بهسازی شده با آنها بر حسب معیارهای ارزیابی که بعداً توضیح داده می شوند، نسبت به همینگ به طول ۳۲ میلی ثانیه قابل قبول نبودند. چون نتایج بهسازی گفتار با آنها خوب نبودند و فقط برای برخی از نویزها و سیگنال به نویزها سنجیده شده بودند، در این بخش از آوردن آنها خودداری شده است. پس از اعمال روش پیشنهادی بر روی طیف زمان کوتاه سیگنال نویزی، از آن تبدیل فوریه زمان کوتاه معکوس گرفته شده و سیگنال بهسازی شده به حوزه زمان انتقال داده می شود. ولی سیگنال به دست آمده به صورت قاب هایی است که باهم همپوشانی دارند. با اعمال روش همپوشانی-جمع بر روی آن، سیگنال بهسازی شده نهایی به دست می آید. حال، برای ارزیابی عینی عملکرد روش پیشنهادی، از سه معیار سیگنال به نویز قطعه ای، PESQ و LLR(Log Likelihood ratio) که در فصل دو توضیح داده شد، استفاده می شود. علت انتخاب سیگنال به نویز قطعه ای به عنوان معیارهای ارزیابی این بود که سیگنال به نویز پارامتر مهمی در نشان دادن موفقیت روش در حذف نویز است به همین جهت در بیشتر روش هایی که در زمینه بهسازی گفتار ارائه می شود، از آن به عنوان معیار ارزیابی استفاده می شود. ولی این معیار همبستگی پایینی (حدود ۰.۳۱) با معیار های ارزیابی ذهنی دارد [۳۱].

همچنین در این پایان نامه استفاده از معیارهای ذهنی یا تست های شنیداری مانند MOS مقدور نبود. زیرا، در این معیارها فرد به سیگنال بهسازی شده گوش داده و به آن امتیازی در بازه مشخص می دهد. استاندارد تست های شنوایی این است که سیگنال صوتی که فرد به آن گوش می دهد و به آن امتیاز می دهد، باید به زبان مادری آن فرد گفته شده باشد. از طرفی، جملات گفته شده در پایگاه داده استفاده شده در این پایان نامه، به زبان انگلیسی است و چون پیدا کردن افرادی که زبان مادریشان انگلیسی باشد، مشکل بود، از دو معیار عینی استفاده شد که همبستگی بالایی با معیارهای

ذهنی داشته باشند. این دو معیار عبارتند از: PESQ و LLR که بیانگر کیفیت سیگنال بهسازی شده هستند و همبستگی بالایی با معیار های ارزیابی ذهنی دارند [۳۱]. همچنین در مقالات معتبر بسیاری از این دو معیار جهت ارزیابی سیگنال بهسازی شده، استفاده شده است. در ابتدا نتیجه سه معیار ارزیابی، با نویزهای ذکر شده و در سه سطح سیگنال به نویز بیان شده، به ازای تغییر تعداد اجزای LMM که با N نشان داده شده است و سپس با توجه به مقادیر به دست آمده از این آزمایشات، مقدار مناسب N انتخاب گردیده است.

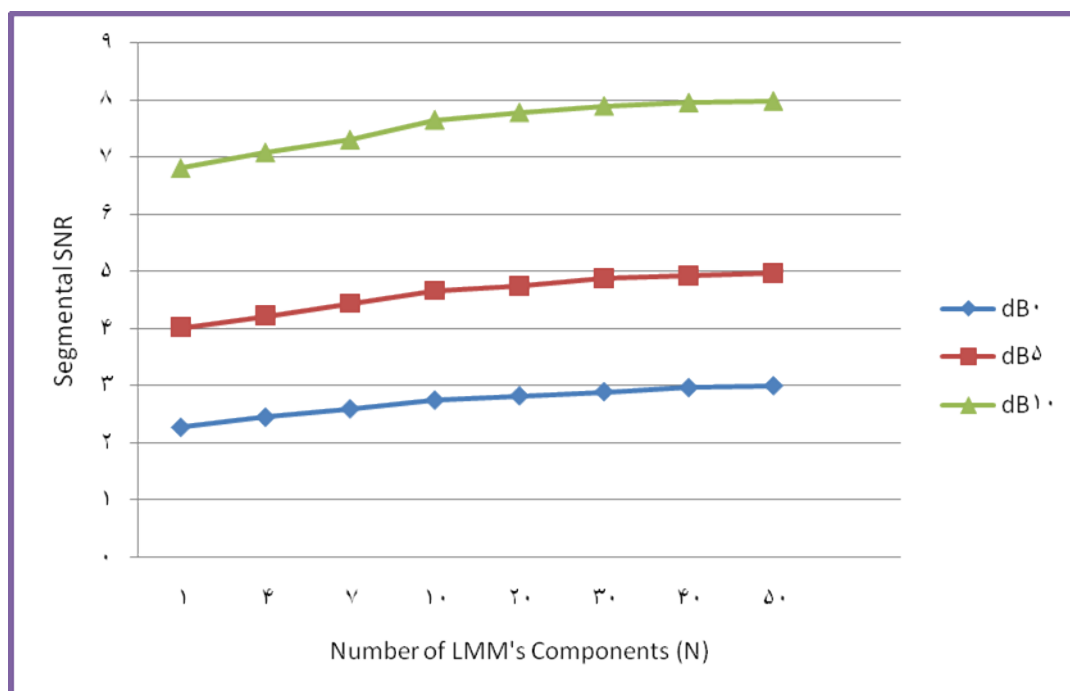
۵-۶-۱- معیار ارزیابی سیگنال به نویز قطعه ای

اولین معیاری که مورد بررسی قرار داده شده است، سیگنال به نویز قطعه ای می باشد. مقادیر این معیار هم در نمودار و هم در جدول ذکر شده است. در سه شکل (۵-۱۰)، (۵-۱۱) و (۵-۱۲) نمودارهای سیگنال به نویز قطعه ای برای سه نویز سفید، همهمه و F_{16} در سه سطح سیگنال به نویز ۰، ۵ و ۱۰ دسی بل به ازای مقادیر مختلف تعداد اجزای مخلوط لاپلاس یعنی N رسم شده است. مقادیر آستانه برای محاسبه سیگنال به نویز قطعه ای -10dB و 35dB در نظر گرفته شده است. [۳۱]. لازم به ذکر است که در این حالت محاسبه سیگنال به نویز قطعه ای ورودی نیز الزامی است. زیرا به ازای یک سیگنال به نویز ورودی، اگر سیگنال به نویزهای قطعه ای ورودی متفاوت باشد، سیگنال به نویزهای قطعه ای خروجی نیز متفاوت خواهند بود. این مطلب درباره دو معیار دیگر نیز صادق است. یعنی در یک سیگنال به نویز معین، هر چه سیگنال به نویز قطعه ای ورودی بیشتر باشد، نتایج به ازای این دو معیار ارزیابی نیز بهتر خواهند بود.

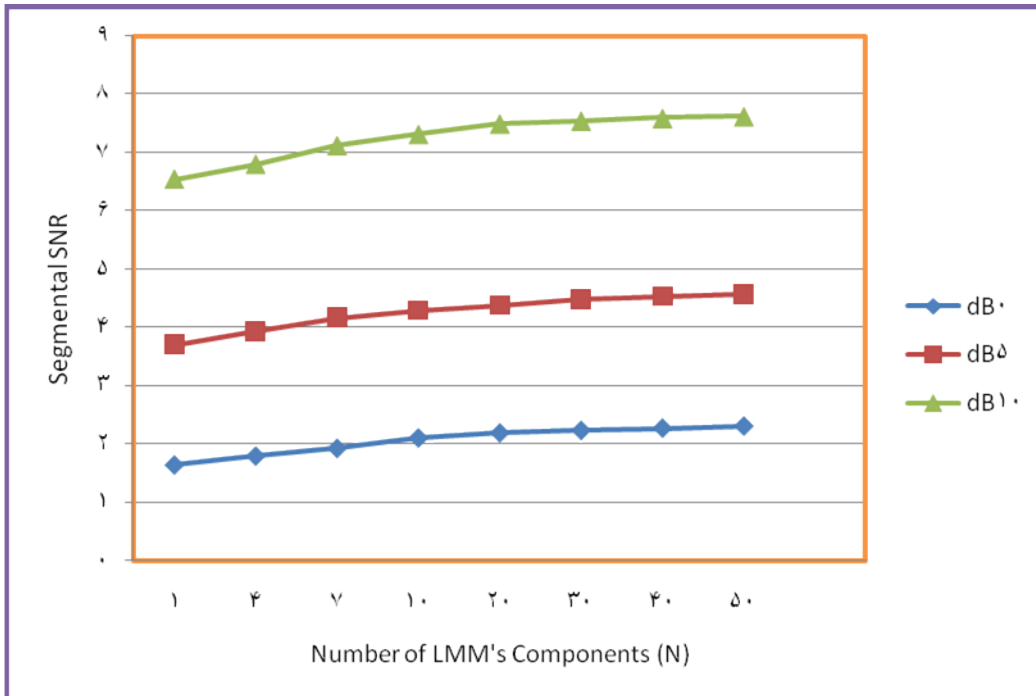
با بررسی این سه نمودار مشخص می شود که میزان سیگنال به نویز قطعه ای سیگنال بهسازی شده با افزایش N، افزایش می یابد، یعنی با افزایش N میزان حذف نویز روش پیشنهادی بیشتر می شود، که منطقی به نظر می رسد. زیرا هر چه N بزرگتر شود مدلی که برای سیگنال در نظر گرفته شده است، با دقت بیشتری به توزیع واقعی سیگنال برآزنده می شود. مثلاً برای حالتی که طیف سیگنال

تمیز با مخلوطی از $N=1$ یا $N=4$ لاپلاس تقریب زده می شود، نتایج خوب نیستند. زیرا این دو حالت تقریب خوبی از توزیع سیگنال نمی باشند.

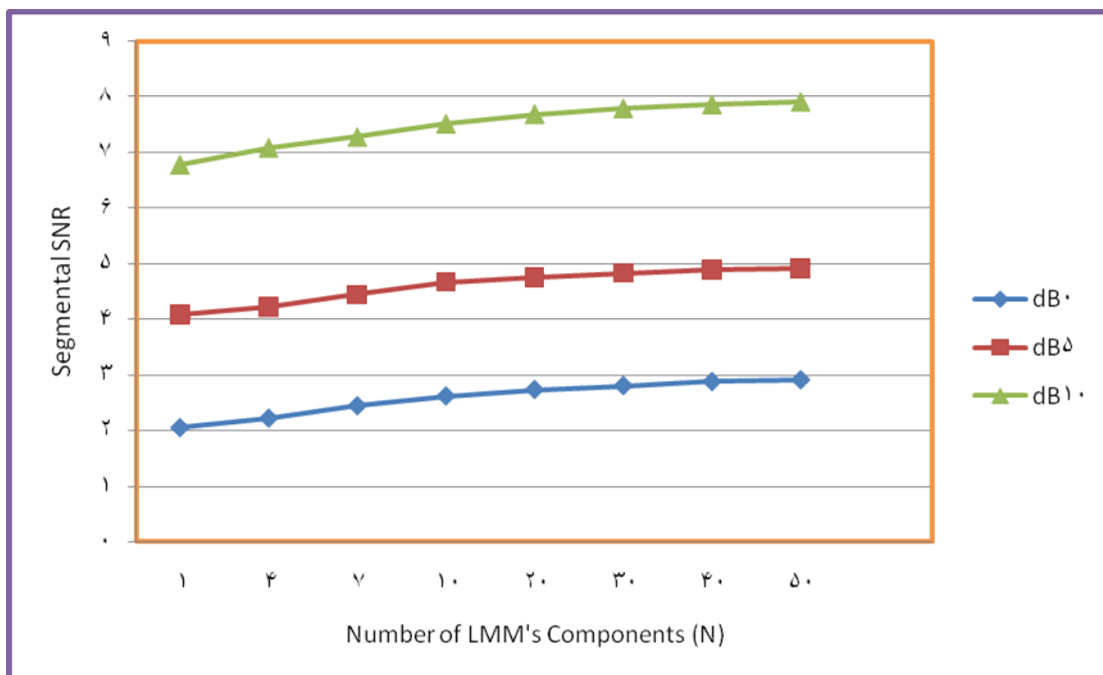
نکته دیگری که از نمودارها مشخص می شود، این است که به ازای تغییر N ، وقتی N ها کوچک هستند، اختلاف مقادیر سیگنال به نویز قطعه ای متناظر با آنها، برای هر سه نویز و سیگنال به نویز بیشتر است. علت این مطلب این است که به ازای مقادیر کوچک N ، با تغییرات N ، تقریب های متناظر با این N ها از سیگنال، با هم اختلاف زیادی دارند (این مطلب در بخش پیاده سازی EM الگوریتم بیان شد)، ولی با زیاد شدن N ، همان طور که قبلا بیان شد، اختلاف بین تقریب ها کمتر می شود، و در نتیجه اختلاف نتایج نیز کم و کمتر می شود. با توجه به نمودارها، $N=30$ پاسخ خوبی داشته و پس از آن به ازای تغییرات N ، نتایج تغییرات زیادی ندارند. این نشان می دهد که تقریب توزیع سیگنال تمیز با توزیع LMM از حدود $N=30$ به بعد، به واقعیت نزدیک و نزدیک تر می شود. منطقی است که با افزایش N نتایج بهتر شود. ولی پیچیدگی محاسباتی نیز به میزان قابل توجهی افزایش یافته و اجرای الگوریتم بسیار زمان بر می شود. به همین دلیل افزایش N در 50 متوقف شده است. زیرا با افزایش N نتایج تغییر چشم گیری نداشتند در حالی که افزایش N به خاطر زمان بر شدن و حجم حافظه اشغالی مقرون به صرفه نبود.



شکل (۵-۱۰): نمودار سیگنال به نویز قطعه ای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز سفید .



شکل (۵-۱۱): نمودار سیگنال به نویز قطعه ای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز F16 .



شکل (۵-۱۲): نمودار سیگنال به نویز قطعه ای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز همهمه .

جدول (۵-۲) : مقایسه سیگنال به نویزهای قطعه ای سیگنال بهسازی شده، برای سه نویز همهمه، سفید و F16، در سه سیگنال به نویز ۰، ۵، ۱۰ و N های متفاوت .

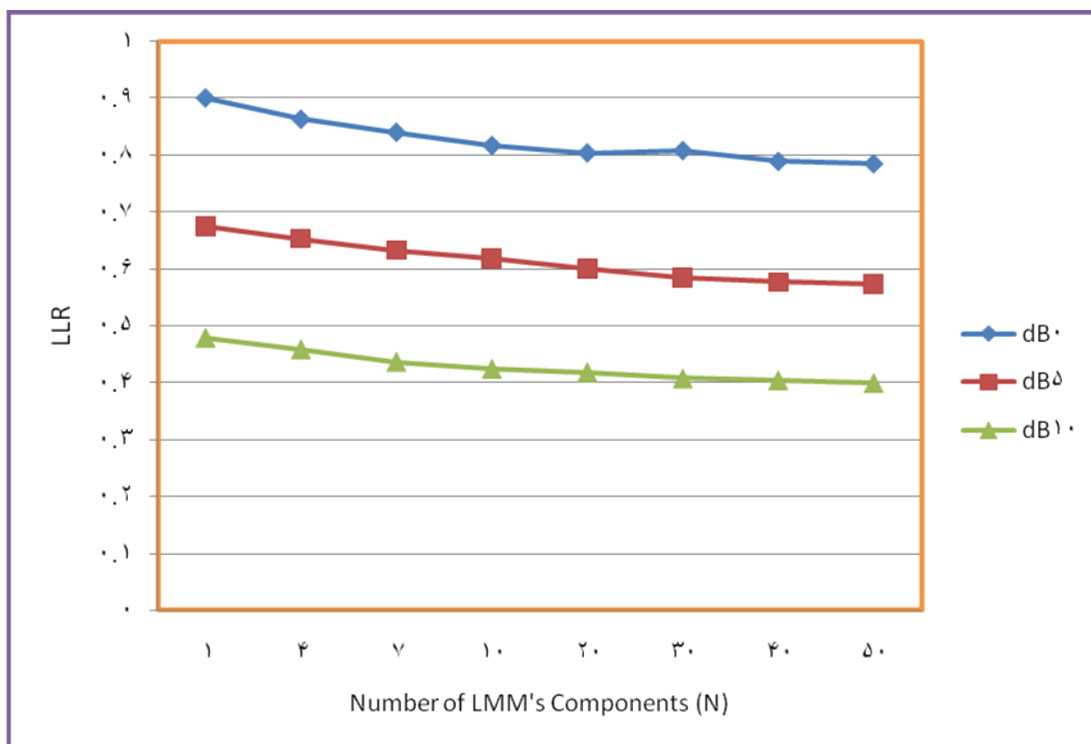
Noise type	White noise			F16 noise			Babble noise		
	-۱۲.۱/۰dB	-۷.۴۵/۵dB	-۲.۲۲/۱۰dB	-۱۲.۴/۰dB	-۷.۸۳/۵dB	-۲.۳/۱۰dB	-۱۱.۹/۰dB	-۷.۱/۵dB	-۲.۰/۱۰dB
N=۱	۲.۲۷۹	۴.۰۲۵	۶.۸۲	۱.۶۴	۳.۷۰۲	۶.۵۴۱	۲.۰۵۵	۴.۰۰۸	۶.۷۸۳
N=۴	۲.۴۵۶	۴.۲۳۶	۷.۰۹۵	۱.۷۹۳	۳.۹۳۲	۶.۸	۲.۲۲۵	۴.۲۱۶	۷.۰۸۹
N=۷	۲.۶۰۴	۴.۴۵۱	۷.۳۱۸	۱.۹۲۶	۴.۱۶۹	۷.۱۲۱	۲.۴۴۷	۴.۴۴۲	۷.۲۸۲
N=۱۰	۲.۷۵۸	۴.۶۷۴	۷.۶۶۵	۲.۱۰۵	۴.۲۸۶	۷.۳۱۵	۲.۶۲۲	۴.۶۶۸	۷.۵۱۷
N=۲۰	۲.۸۳۲	۴.۷۵۵	۷.۷۹۴	۲.۱۹۲	۴.۳۷۹	۷.۴۹۸	۲.۷۳۸	۴.۷۴۴	۷.۶۹۲
N=۳۰	۲.۹۰۴	۴.۸۸۱	۷.۹۰۶	۲.۲۳۵	۴.۴۷۴	۷.۵۴۶	۲.۸۰۷	۴.۸۲۷	۷.۷۹۳
N=۴۰	۲.۹۷۶	۴.۹۲۷	۷.۹۶۴	۲.۲۷۵	۴.۵۳	۷.۵۹۱	۲.۸۷۹	۴.۸۹۱	۷.۸۶۲
N=۵۰	۳.۰۱	۴.۹۸۱	۷.۹۹۳	۲.۳۱	۴.۵۷۱	۷.۶۲۳	۲.۹۱۲	۴.۹۱۵	۷.۹۱۲

نکته دیگری که می‌توان در این سه نمودار ملاحظه کرد، این است که با تغییر N بیشترین تغییرات سیگنال به نویز قطعه ای برای هر سه نویز، مربوط به حالت سیگنال به نویز ۱۰dB و سپس ۵dB و پس از آن ۰dB است. این امر بیانگر این است که با افزایش سیگنال به نویز عملکرد روش پیشنهادی بهتر می‌شود. بهترین جواب ها در سه سیگنال به نویز ۰ و ۵ و ۱۰ مربوط به نویز سفید است. پس از آن نویز همهمه و در نهایت نویز F16 قرار می‌گیرد. با مشاهده و بررسی طیف توان این سه نویز مشخص شد که تغییرات توان نویز سفید از دو نویز دیگر کمتر است. این تغییرات برای نویز همهمه کمتر از نویز F16 می‌باشد. همچنین در بخش تخمین نویز بیان شد که فرض اساسی روش‌های تخمین نویز، کند بودن تغییرات توان نویز است. احتمالاً به این دلیل نتایج برای نویز سفید از دو نویز دیگر بهتر است و نویز همهمه نیز نتایج بهتری از نویز F16 دارد.

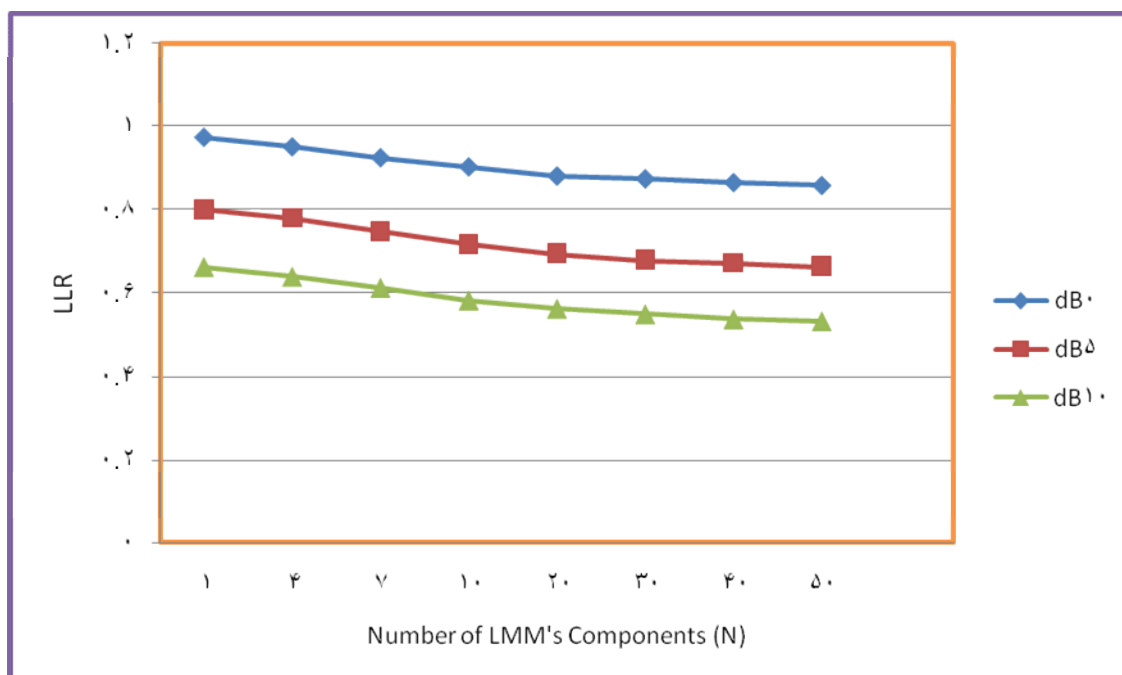
برای اینکه مقادیر رسم شده در نمودارهای بالا به طور دقیق تر مشاهده شود، در جدول (۲-۵) نیز نمایش داده شده اند.

۵-۶-۲- معیار ارزیابی LLR

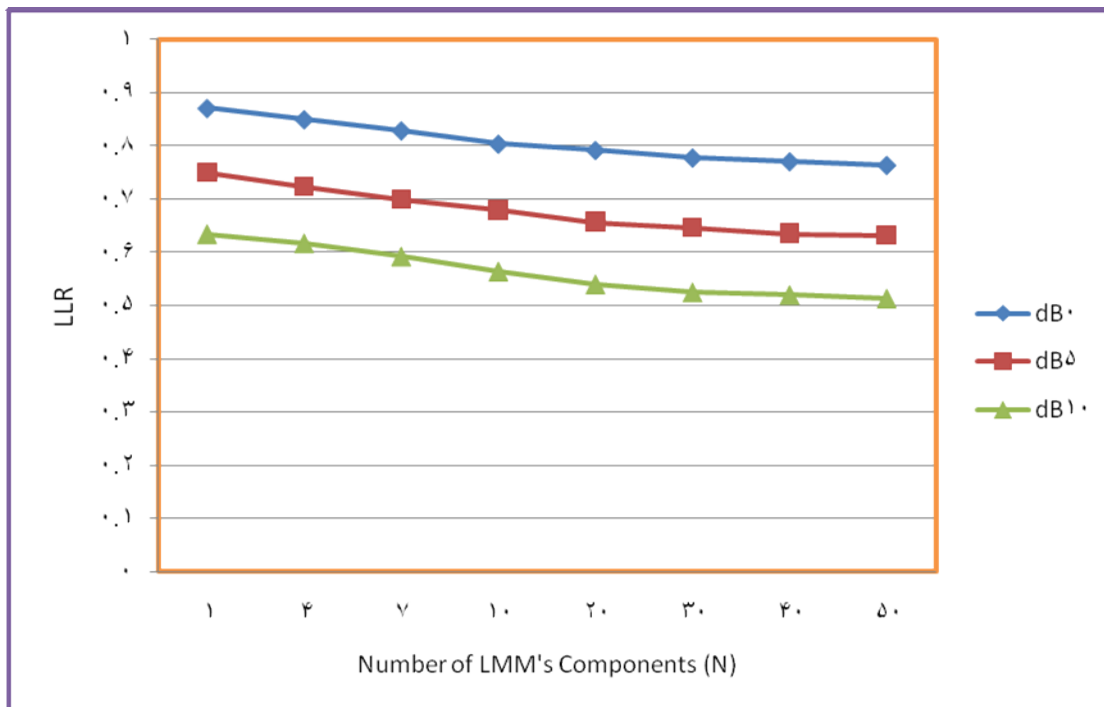
پس از بررسی معیار سیگنال به نویز قطعه ای، معیار LLR برای سه نویز مذکور و سه سطح سیگنال به نویز قبل بررسی شده است. مقادیر معیار LLR برای این سه نویز و سه سطح سیگنال به نویز، در نمودارهای (۵-۱۳)، (۵-۱۴) و (۵-۱۵) رسم شده است. قبلاً بیان شد که کمتر بودن مقدار LLR به معنی بهتر بودن کیفیت سیگنال خروجی است. از نمودارها مشخص است که مقدار LLR با افزایش N کم می شود، یعنی کیفیت سیگنال خروجی با افزایش N زیاد می شود که همان نتیجه ای است که در بخش قبل نیز به دست آمد. در این بخش نیز بهترین جواب ها مربوط به حذف نویز سفید و سپس همه می باشد. به ازای یک N مشخص، در هر سه نویز، کمترین LLR مربوط به سیگنال به نویز ۱۰ است، که نشان می دهد در سیگنال به نویزهای بالاتر، که دامنه سیگنال خیلی از دامنه نویز بزرگتر است، کیفیت سیگنال خروجی بالاتر است. در این بخش نیز مشاهده می شود که مانند بخش قبل، تغییرات نتایج در ابتدای نمودارها که مربوط به مقادیر پایین N است، زیاد می باشد. اما این



شکل (۵-۱۳): نمودار LLR برای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز سفید.



شکل (۵-۱۴): نمودار LLR برای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز F16.



شکل (۵-۱۵): نمودار LLR برای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز هممه .

جدول (۵-۳): مقایسه LLR سیگنال بهسازی شده، برای سه نویز هممه، سفید و F16 در سه سیگنال به نویز 0، 5، 10 و N های متفاوت .

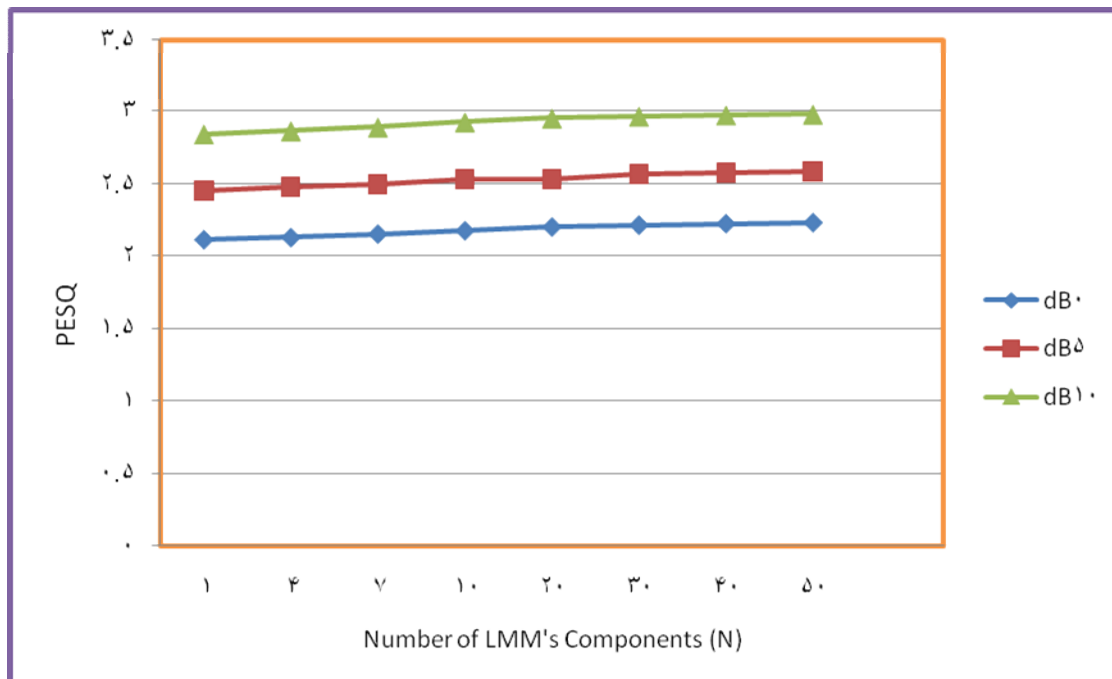
Noise type	White noise			F16 noise			Babble noise		
	0dB	5dB	10dB	0dB	5dB	10dB	0dB	5dB	10dB
N=1	0.9	0.675	0.479	0.993	0.782	0.691	0.891	0.73	0.634
N=4	0.863	0.653	0.459	0.97	0.758	0.668	0.87	0.704	0.617
N=7	0.84	0.633	0.437	0.944	0.727	0.642	0.849	0.683	0.592
N=10	0.817	0.618	0.425	0.922	0.696	0.61	0.824	0.659	0.564
N=20	0.804	0.601	0.419	0.9	0.674	0.592	0.818	0.637	0.54
N=30	0.808	0.585	0.408	0.883	0.659	0.579	0.809	0.626	0.525
N=40	0.789	0.577	0.405	0.872	0.651	0.567	0.8	0.615	0.52
N=50	0.785	0.573	0.4	0.861	0.643	0.562	0.794	0.612	0.513

تغییرات در انتهای نمودار و از حدود $N = 30$ به بالا بسیار کم می شود که علت این امر در قسمت قبل توضیح داده شد. به همین خاطر افزایش N در $N = 50$ متوقف شده است. جدول (۳-۵) نیز مقادیر عددی LLR را برای سه نویز مذکور، در سه سیگنال به نویز آورده شده، نشان می دهد.

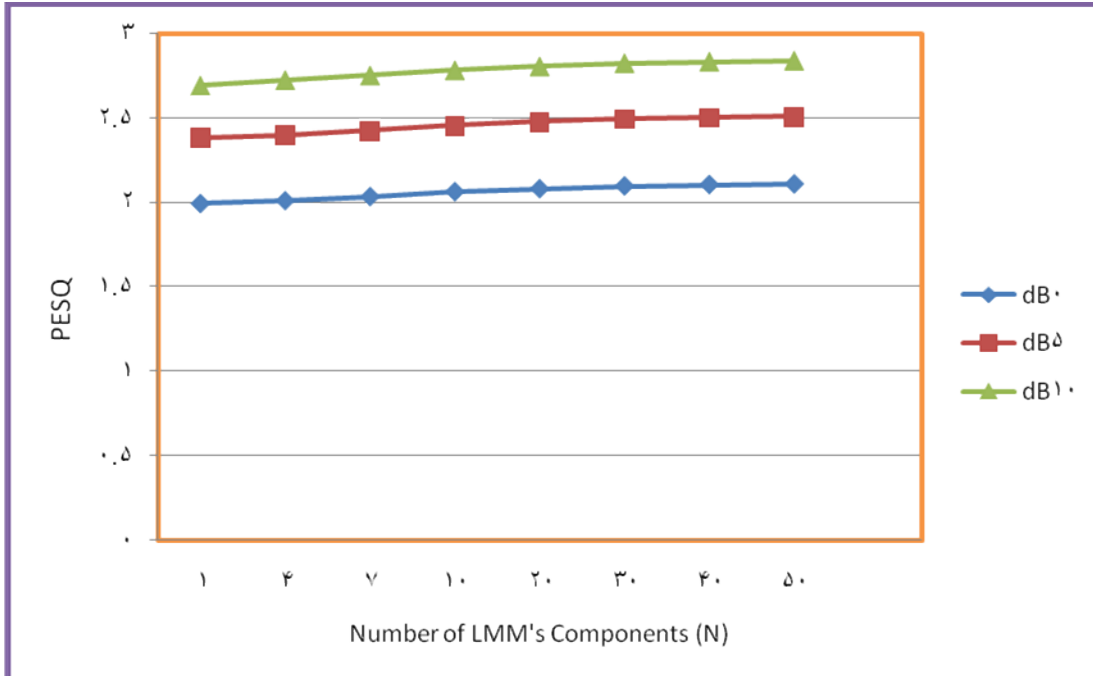
۳-۶-۵- معیار ارزیابی PESQ

همانطور که در فصل های قبل گفته شده PESQ برای پیش بینی میزان امتیاز معیار MOS طراحی شده و عددی بین ۰.۵ - تا ۴.۵ را می دهد [۱۰]. هر چه این عدد بزرگتر باشد بیانگر کیفیت بالاتر سیگنال صوتی است. در نمودارهای (۵-۱۶)، (۵-۱۷) و (۵-۱۸) مقادیر این معیار برای سه نویز قبل و سه سیگنال به نویز مذکور آورده شده است. با توجه به نمودارها مشخص می شود که با افزایش N مقادیر PESQ افزایش یافته و در واقع نتایج بهبود می یابند. بالاترین مقادیر PESQ برای هر سه نویز، در سیگنال به نویز ۱۰ اتفاق افتاده است. بر اساس این معیار نیز بهترین نتایج برای نویز سفید و پس از آن همهمه و $F16$ است. روند تغییرات نمودارها نیز مانند نمودارهای دو بخش قبل است. یعنی

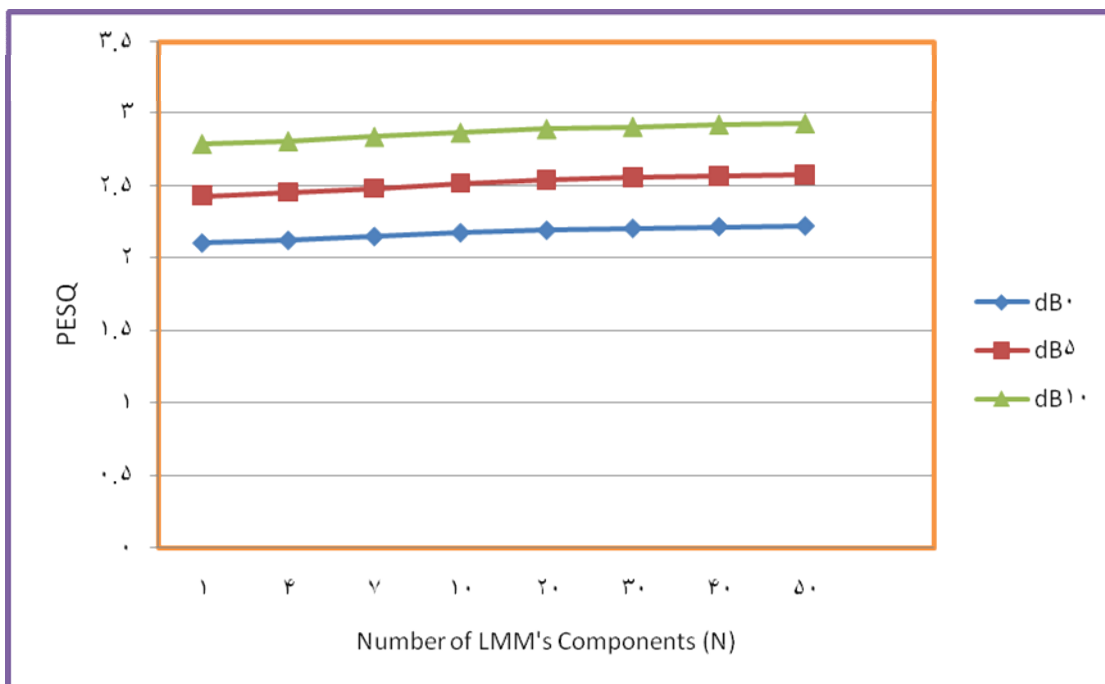
در



شکل (۵-۱۶): نمودار PESQ برای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز سفید.



شکل (۵-۱۷): نمودار PESQ برای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز F16.



شکل (۵-۱۸): نمودار PESQ برای سیگنال بهسازی شده، پس از اعمال روش پیشنهادی روی سیگنال آغشته به نویز همهمه .

جدول (۵-۴): مقایسه PESQ سیگنال بهسازی شده، برای سه نویز همهمه، سفید و F16 در سه سیگنال به نویز ۰، ۵، ۱۰ و N های متفاوت .

Noise Type	white noise			F16 noise			Babble noise		
	۰dB	۵dB	۱۰dB	۰dB	۵dB	۱۰dB	۰dB	۵dB	۱۰dB
N=۱	۲.۱۱۷	۲.۴۵۳	۲.۸۴۵	۱.۹۹۴	۲.۳۸۴	۲.۶۹۴	۲.۱۰۸	۲.۴۳۱	۲.۷۹۲
N=۴	۲.۱۳۱	۲.۴۷۸	۲.۸۶۷	۲.۰۱۱	۲.۴۰۲	۲.۷۲۶	۲.۱۲۶	۲.۴۵۵	۲.۸۱۱
N=۷	۲.۱۵۵	۲.۵	۲.۸۹۳	۲.۰۳۵	۲.۴۲۵	۲.۷۵۳	۲.۱۵۲	۲.۴۸۳	۲.۸۴
N=۱۰	۲.۱۷۸	۲.۵۲۹	۲.۹۲۸	۲.۰۶۴	۲.۴۵۶	۲.۷۸۴	۲.۱۷۸	۲.۵۱۶	۲.۸۶۹
N=۲۰	۲.۲۰۴	۲.۵۳	۲.۹۵۴	۲.۰۸۲	۲.۴۷۸	۲.۸۰۷	۲.۱۹۵	۲.۵۳۹	۲.۸۹۳
N=۳۰	۲.۲۱۷	۲.۵۶۸	۲.۹۶۸	۲.۰۹۸	۲.۴۹۷	۲.۸۲۶	۲.۲۰۸	۲.۵۵۵	۲.۹۰۸
N=۴۰	۲.۲۲۷	۲.۵۷۷	۲.۹۷۶	۲.۱۰۶	۲.۵۰۵	۲.۸۳۵	۲.۲۲	۲.۵۶۸	۲.۹۲۳
N=۵۰	۲.۲۳۴	۲.۵۸۴	۲.۹۸۱	۲.۱۱	۲.۵۱	۲.۸۴۱	۲.۲۲۵	۲.۵۷۶	۲.۹۳

بخش انتهایی نمودارها و پس از $N=30$ تغییر چشمگیری در نتایج دیده نمی شود. افزایش N نیز به دلایلی که قبلاً گفته شد، در $N=50$ متوقف شده است. جدول (۵-۴) مقادیر PESQ را برای هر سه نویز و سیگنال به نویز مذکور، نمایش می دهد.

۵-۷- مقایسه روش پیشنهادی با روش های دیگر

برای بررسی بهبود یا عدم بهبود روش پیشنهادی، سه روش معروف دیگر بهسازی گفتار پیاده سازی شدند و روش پیشنهادی با این سه روش مقایسه گردیده است. نتایج این مقایسه ها در این بخش آورده شده است.

روش هایی که برای مقایسه پیاده سازی شده اند، عبارت اند از:

- تخمین MMSE دامنه مبتنی بر توزیع گوسی [۲۲] که با MMSE نمایش داده شده است.
- تخمین LogMMSE دامنه مبتنی بر توزیع گوسی [۲۳] که با Log-MMSE نمایش داده شده است. (در LogMMSE یک تخمین بهینه با کمینه کردن میزان خطای مربعات لگاریتم دامنه به دست می آید).
- تخمین MMSE طیف مبتنی بر توزیع لاپلاس [۴۳] که با Lap-MMSE نمایش داده شده است.
- تخمین گر پیشنهادی نیز با LMM-MMSE نمایش داده شده است.

می توان گفت که این سه روش، موارد مناسبی برای مقایسه می باشند. زیرا هر سه روش آماری محض هستند و در هر سه روش از تخمین گر MMSE استفاده شده است. دو روش از این سه روش بر پایه فرض توزیع گوسی برای سیگنال تمیز و روش دیگر بر پایه فرض توزیع لاپلاس برای سیگنال تمیز می باشند. در واقع هدف از مقایسه روش پیشنهادی با این سه روش، بررسی مناسب یا نامناسب بودن فرض LMM برای توزیع سیگنال تمیز نسبت به توزیع های لاپلاس و گوسی است.

برای مقایسه این چهار روش از دو نویز F_{16} و همهمه و سه سطح سیگنال به نویز قبل استفاده شده است. لازم به ذکر است که برای پیاده سازی سه روش دیگر نیز از همان ده جمله از پایگاه داده TIMIT استفاده شد که برای پیاده سازی روش پیشنهادی مورد استفاده قرار داده شد. برای مقایسه نتایج، از مقادیر روش پیشنهادی در $N=50$ استفاده شده است.

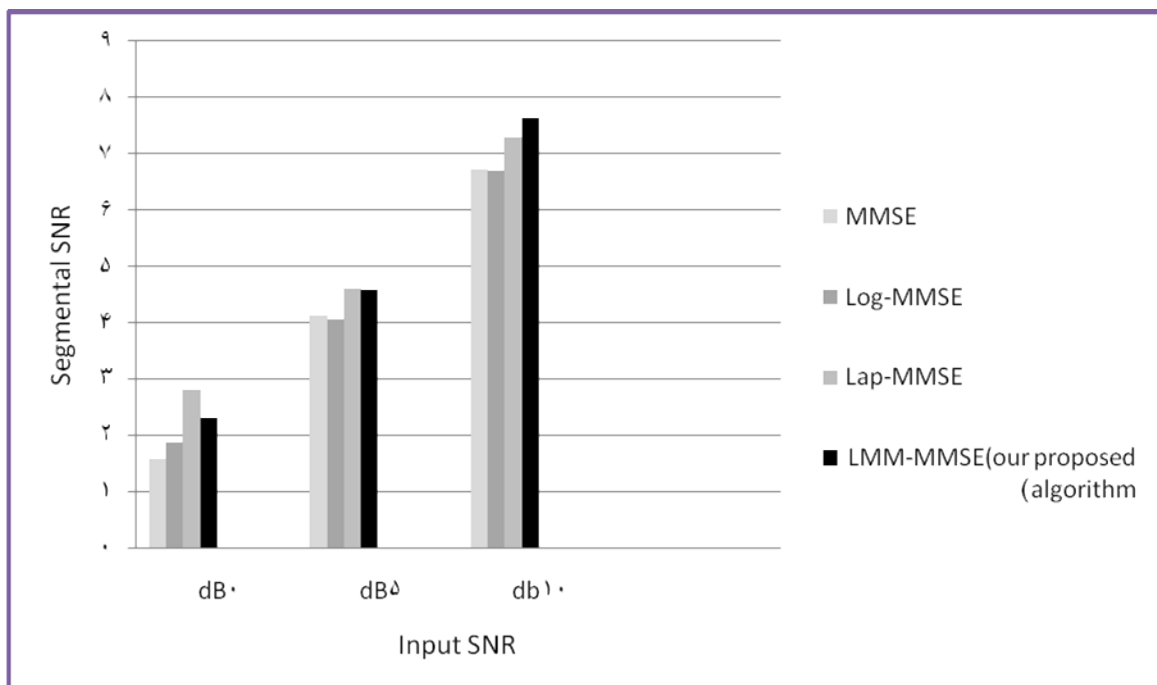
۵-۷-۱- معیار سیگنال به نویز قطعه ای

شکل های (۵-۱۹) و (۵-۲۰) نمودار سیگنال به نویز قطعه ای را برای چهار روش مذکور، در سه سطح سیگنال به نویز ۵ و ۱۰ و برای نویزهای F_{16} و همهمه نشان می دهد.

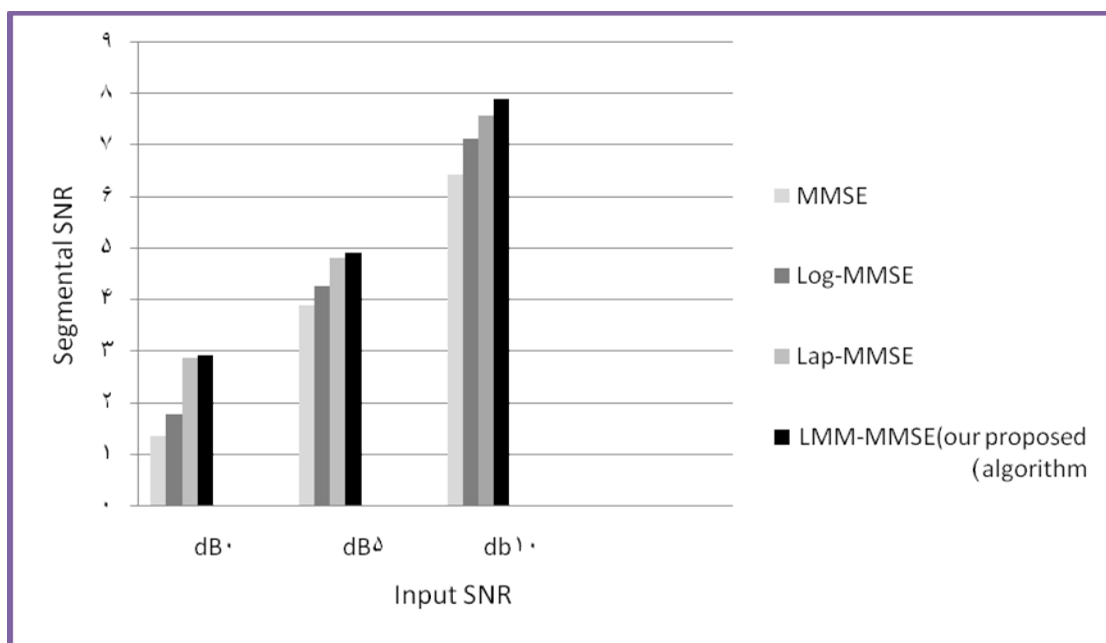
با توجه به نمودار شکل (۵-۱۹)، برای نویز F_{16} ، در سیگنال به نویز صفر و پنج عملکرد روش پیشنهادی از دو روش MMSE و Log-MMS بهتر می باشد ولی عملکرد آن از روش Lap-MMSE بهتر نیست. در حالی که در سیگنال به نویز ۱۰dB روش پیشنهادی از سه روش دیگر موفق تر عمل

می کند. در واقع برای نویز F_{16} و در سیگنال به نویزهای پایین روش Lap-MMSE که در آن سیگنال با یک لاپلاس مدل می شود و پارامترهای لاپلاس، از روی سیگنال نویزی تعیین می شود، از روش پیشنهادی بهتر عمل می کند. ولی با افزایش سیگنال به نویز عملکرد روش پیشنهادی بهتر می شود و از این روش پیشی می گیرد.

با توجه به نمودار (۵-۱۹)، دو روش MMSE و Log-MMSE که در آن سیگنال با توزیع گوسی مدل می شود، مقادیر سیگنال به نویز پایین تری از دو روش دیگر دارند. این مطلب بیانگر این است که مدل کردن سیگنال با توزیع لاپلاس و مخلوطی از لاپلاس از مدل کردن سیگنال با توزیع گوسی مناسب تر می باشد. نکته دیگر که از جدول (۵-۵) مشخص می شود، این است که در همه سیگنال به نویزها، روش پیشنهادی بیش از یک دسی بل بهبود روی سیگنال به نویز قطعه ای، نسبت به این دو روش MMSE و Log-MMSE ایجاد می کند که مقدار قابل توجهی است.



شکل (۵-۱۹): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE بر حسب سیگنال به نویز قطعه ای و برای نویز F_{16} .



شکل (۵-۲۰): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE بر حسب سیگنال به نویز قطعه ای و برای نویز همهمه.

جدول (۵-۵): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE بر حسب سیگنال به نویز قطعه ای و برای نویز همهمه و F۱۶.

Noises	Babble noise			F۱۶ noise		
	-۱۱.۹/۰dB	-۷.۱/۵dB	-	-۱۲.۴/۰dB	-۷.۸۳/۵dB	-۲.۳/۱۰dB
Estimators						
MMSE [۲۲]	۱.۳۴۱	۳.۸۹۲	۶.۴۲	۱.۵۷	۴.۱۳۲	۶.۷۳۱
Log-MMSE [۲۳]	۱.۷۷۳	۴.۲۵۱	۷.۱۲۳	۱.۸۷۳	۴.۰۵۲	۶.۶۹
Lap-MMSE [۴۳]	۲.۸۷۴	۴.۸۰۸	۷.۵۸۲	۲.۸۱	۴.۶۱۳	۷.۲۸۸
LMM-MMSE (Proposed method)	۲.۹۱۲	۴.۹۱۵	۷.۹۱۲	۲.۳۱	۴.۵۷۱	۷.۶۲۳

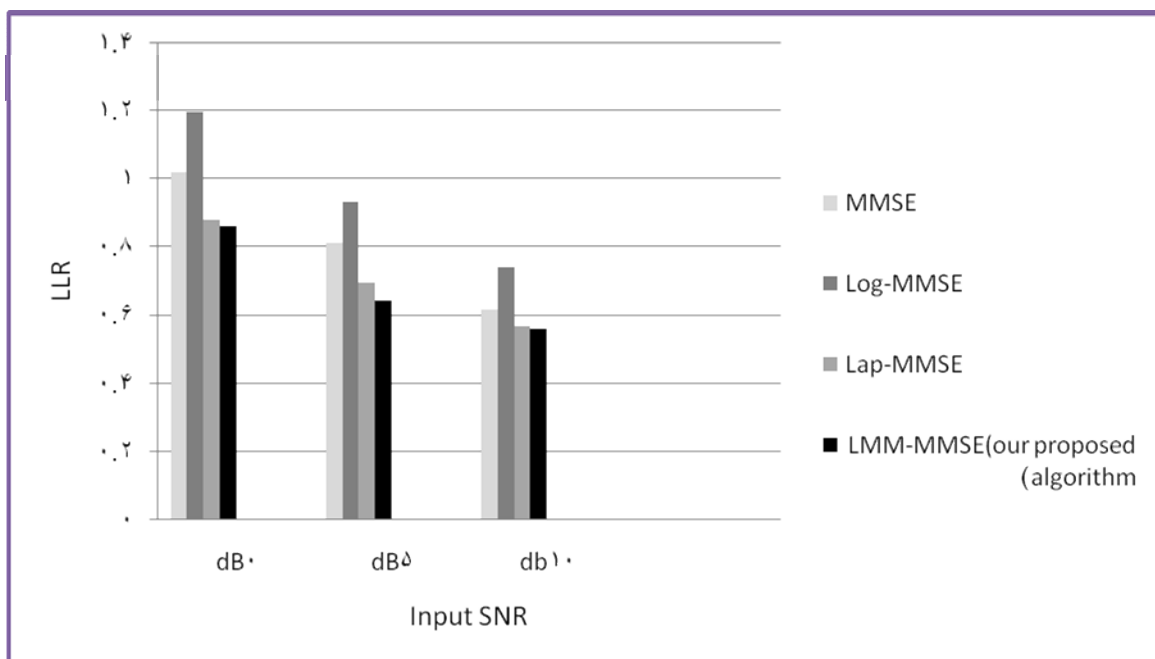
برای نویز همهمه نیز با توجه به نمودار شکل (۵-۲۰)، روش پیشنهادی در همه سیگنال به نویزها از سه روش پیشی گرفته است. البته با توجه به نمودار، با افزایش سیگنال به نویز اختلاف روش پیشنهادی با سه روش دیگر زیاد می شود یعنی عملکرد آن بهتر می شود. جدول (۵-۵) نیز نشان دهنده مقادیر نمودارهای شکل (۵-۱۹) و (۵-۲۰) می باشد.

۵-۷-۲- معیار LLR

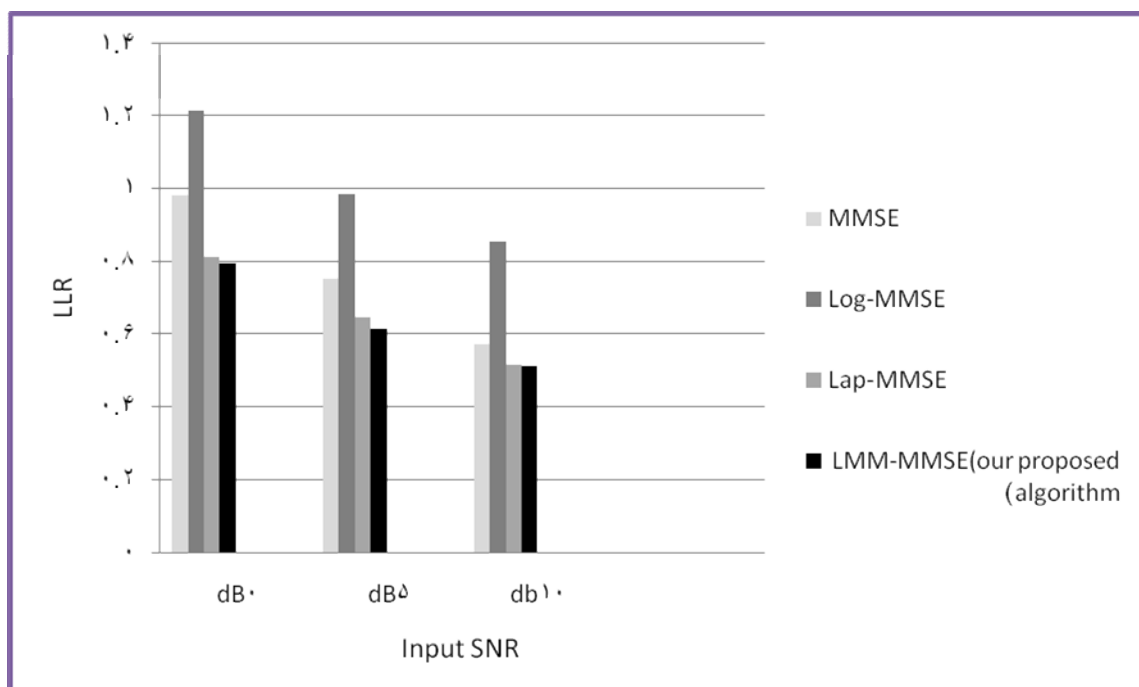
در این بخش معیار LLR جهت ارزیابی روش پیشنهادی با سایر روش های ذکر شده در بخش قبل به کار می رود. نتایج حاصل از این مقایسه، در شکل های (۵-۲۱) و (۵-۲۲) آورده شده است. با توجه به نمودار، در همه سیگنال به نویزها و برای هر دو نویز، روش پیشنهادی دارای LLR کمتر و در نتیجه سیگنال بهسازی شده با کیفیت بالاتر می باشد. در واقع مشخص می شود که LMM گزینه مناسبی برای مدل کردن سیگنال تمیز می باشد. با توجه به نمودارها، اختلاف روش پیشنهادی با دو روش مبتنی بر توزیع گوسی زیاد می باشد؛ اما اختلاف آن با روش مبتنی بر توزیع لاپلاس زیاد نیست. جدول (۵-۶) مقادیر عددی دو نمودار شکل (۵-۲۱) و (۵-۲۲) را نمایش می دهد.

۵-۷-۳- معیار PESQ

در این بخش معیار PESQ جهت مقایسه کارایی روش پیشنهادی با سه روش دیگر به کار می رود. نتایج این مقایسه، در نمودارهای (۵-۲۳) و (۵-۲۴) آورده شده است. با توجه به نمودار شکل (۵-۲۳) در سیگنال به نویز صفر و برای نویز $F16$ روش پیشنهادی از دو روش مبتنی بر توزیع گوسی بهتر ولی از روش مبتنی بر توزیع لاپلاس بدتر می باشد. (مقادیر بالاتر بیانگر سیگنال بهسازی شده با کیفیت بالاتر می باشند) اما برای سیگنال به نویزهای دیگر، کارایی روش پیشنهادی از هر سه روش دیگر بهتر بوده و با افزایش سیگنال به نویز، اختلاف مقادیر PESQ برای این روش با روش های دیگر، بیشتر هم می گردد، یعنی عملکرد روش بهتر می شود. برای نویز همهمه در همه سیگنال به نویزها، روش پیشنهادی، بر اساس این معیار، نسبت به سه روش دیگر، بهتر عمل می کند. در این حالت نیز با افزایش سیگنال به نویز اختلاف مقادیر PESQ برای روش پیشنهادی و بقیه روش ها بیشتر می شود. جدول (۵-۶) مقادیر نمودارهای شکل (۵-۲۳) و (۵-۲۴) را نشان می دهد.



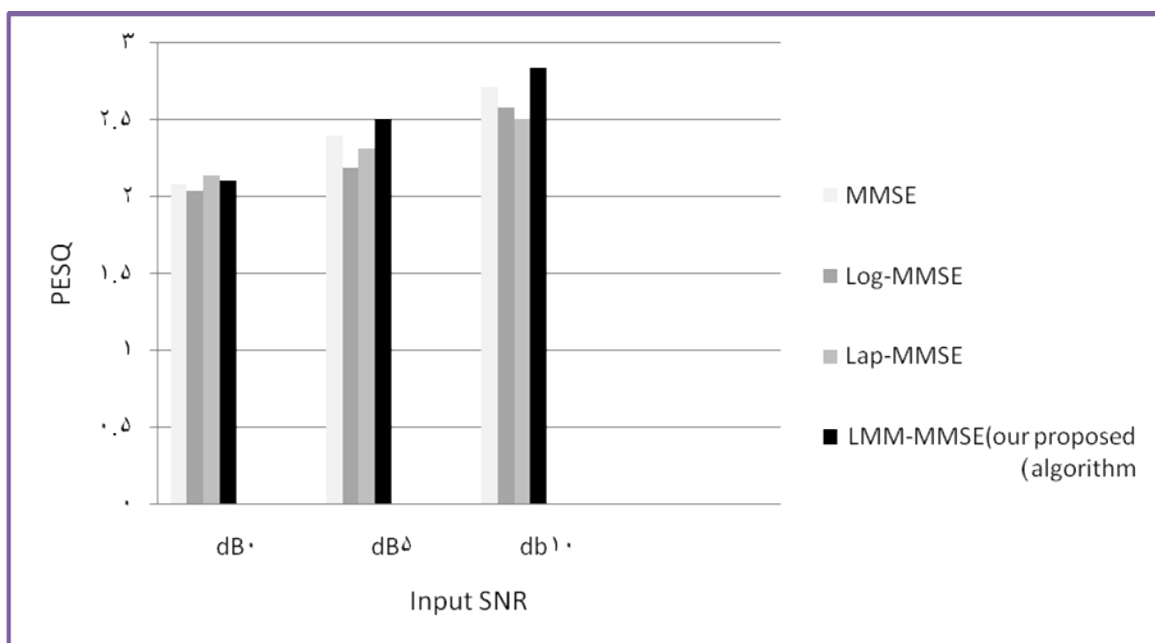
شکل (۵-۲۱): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE, بر حسب LLR و برای نویز ۱۶.F



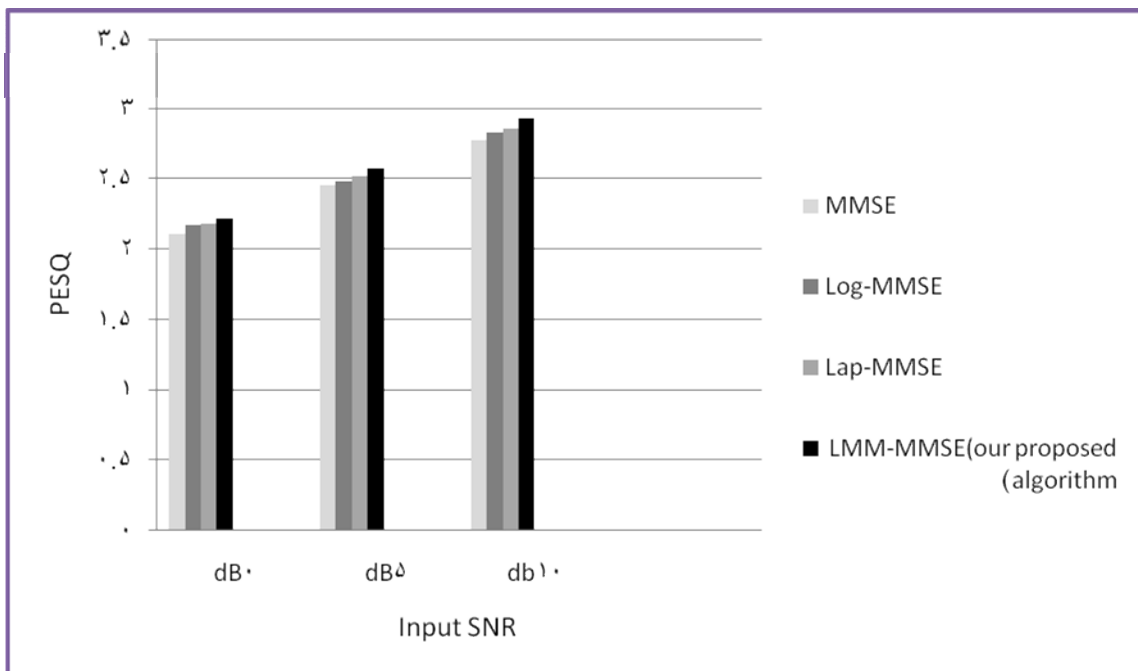
شکل (۵-۲۲): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE, بر حسب LLR و برای نویز مهمه.

جدول (۵-۶): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE, بر حسب LLR و برای نویز همهمه و F۱۶.

Noises	Babble noise			F۱۶ noise		
	۰ dB	۵ dB	۱۰ dB	۰ dB	۵ dB	۱۰ dB
MMSE [۲۲]	۰.۹۸۱	۰.۷۵۱	۰.۵۷۲	۱.۰۲۱	۰.۸۱۲	۰.۶۱۹
Log-MMSE [۲۳]	۱.۲۱۳	۰.۹۸۴	۰.۸۵۲	۱.۱۹۶	۰.۹۳۵	۰.۷۴۳
Lap-MMSE[۴۳]	۰.۸۱	۰.۶۴۴	۰.۵۱۶	۰.۸۷۹	۰.۶۹۶	۰.۵۷
LMM-MMSE (Proposed method)	۰.۷۹۳	۰.۶۱۲	۰.۵۱۳	۰.۸۶۱	۰.۶۴۳	۰.۵۶۲



شکل (۵-۲۳): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE, بر حسب PESQ و برای نویز F۱۶.



شکل (۵-۲۴): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE, بر حسب PESQ و برای نویز همهمه.

جدول (۵-۷): مقایسه روش پیشنهادی با روش های MMSE, Log-MMSE, Lap-MMSE, بر حسب PESQ و برای نویز همهمه و F16. F16

Noises	Babble noise			cockpit noise		
	0 dB	5 dB	10 dB	0 dB	5 dB	10 dB
Estimators						
MMSE [۲۲]	۲.۱۱۴	۲.۴۵۶	۲.۷۸۱	۲.۰۸۱	۲.۳۹۶	۲.۷۱۲
Log-MMSE [۲۳]	۲.۱۷۲	۲.۴۸۷	۲.۸۳۳	۲.۰۳۵	۲.۱۹	۲.۵۸۵
Lap-MMSE [۴۳]	۲.۱۸۳	۲.۵۲	۲.۸۵۷	۲.۱۳۶	۲.۳۱۴	۲.۵۰۸
LMM-MMSE (Proposed method)	۲.۲۲۵	۲.۵۷۶	۲.۹۳	۲.۱۱	۲.۵۱	۲.۸۴۱

با توجه به مجموعه نمودارها و جداول بالا نکاتی مشخص می شوند که آنها را می توان به طور خلاصه و به صورت زیر بیان کرد:

- ۱- با افزایش تعداد اجزای مخلوط لاپلاس N ، عملکرد روش پیشنهادی بهبود می یابد. که به دلیل دقیق تر شدن تقریب LMM از طیف سیگنال می باشد.
- ۲- مقدار بهینه برای N با توجه به مقادیر نمودارها و جداول، ۵۰ انتخاب می گردد. بهینه بودن این مقدار از این جهت مطرح می شود که باید بین تقریب دقیق تر - نتایج بهتر و پیچیدگی - زمان اجرای برنامه - حجم حافظه اشغالی کمتر، مصالحه شود.
- ۲- چهار روشی که با هم مقایسه شدند، برای حذف نویز همهمه عملکرد بهتری نسبت به حذف نویز F1۶ دارند.
- ۳- روش پیشنهادی در بیشتر موارد، علی الخصوص در سیگنال به نویزهای بالا و برای نویز همهمه، از جهت حذف نویز و کیفیت سیگنال بهسازی شده، کارایی بهتری از سه روش دیگر دارد. بنابراین می توان نتیجه گرفت که مدل کردن سیگنال تمیز با LMM مناسب تر از مدل کردن سیگنال با لاپلاس تنها یا گوسی تنها می باشد.
- ۴- با افزایش سیگنال به نویز، عملکرد روش پیشنهادی بهتر می شود. ولی این بدان معنا نیست که روش پیشنهادی در سیگنال به نویزهای پایین عملکرد قابل قبولی نداشته باشد.
- ۵- به ازای N های کوچک مثل ۱ و ۴، روش پیشنهادی برای هر دو نویز و هر سه سیگنال به نویز، کارایی مناسبی نسبت به روش مبتنی بر لاپلاس ندارد. زیرا در روش پیشنهادی به ازای N های کوچک، تقریب خوبی از سیگنال نداشته و همچنین پارامترهای توزیع سیگنال تمیز به صورت برون خطی^۱ تخمین زده می شوند. در حالی که در روش مبتنی بر لاپلاس، شاید یک لاپلاس برای تخمین طیف سیگنال تمیز مناسب نباشد ولی پارامترهای آن به صورت درون خطی^۲، تخمین زده می شوند.

^۱ - offline
^۲ - online

فصل ششم

نتیجه گیری و

پیشنهادات

۶-۱- مقدمه

در این فصل ابتدا یک جمع بندی کلی در مورد کارهای انجام شده در این پایان نامه آمده است و در ادامه پیشنهاداتی برای کارهای آینده ارائه خواهد شد.

۶-۲- نتیجه گیری و جمع بندی

در این پایان نامه به بررسی بهسازی گفتار، مفاهیم مربوط به آن و روش های مختلف بهسازی گفتار، علی الخصوص روش آماری محض پرداخته شد. سپس با ارائه مدل آماری مخلوط لاپلاس برای سیگنال گفتار تمیز و با استفاده از تخمین گر MMSE روش جدیدی برای تخمین سیگنال تمیز از روی سیگنال نویزی (نویز با سیگنال جمع شده و مستقل از آن است) ارائه شد.

در این پایان نامه، برای تخمین پارامترهای سیگنال تمیز، از EM الگوریتم و برای تخمین پارامترهای نویز از روش مبتنی بر ردیابی کمینه ها استفاده شد.

برای پیاده سازی روش پیشنهادی، سیگنال گفتار تمیز با سه نویز سفید، همهمه و $F16$ و در سه سطح سیگنال به نویزهای ۰، ۵ و ۱۰ جمع گردید و سپس از سه معیار سیگنال به نویز قطعه ای، LLR و PESQ برای ارزیابی روش استفاده شد. در نهایت نیز عملکرد روش پیشنهادی به ازای افزایش تعداد اجزای مخلوط لاپلاس بررسی گردید، و نتایج این روش با سه روش معروف دیگر بهسازی گفتار (MMSE, Log-MMSE, Lap-MMSE) مقایسه شد.

بررسی نمودارها و جداول موجود در فصل پنجم، نشان می دهد که به ازای افزایش تعداد اجزای مخلوط لاپلاس N ، عملکرد روش پیشنهادی، با توجه به مقادیر سه معیار ارزیابی ذکر شده، برای هر سه نویز و سیگنال به نویز بهتر می شود. دلیل این امر را می توان اینگونه بیان کرد که هرچه N بزرگتر شود LMM متناظر با آن با خطای کمتری بر تابع چگالی احتمال داده ها برازنده می شود و به بیان ساده تر مدل به واقعیت نزدیکتر می شود. مثلاً برای حالت های $N=1$ یا $N=4$ که تقریب

خوبی از سیگنال موجود نمی باشد، نتایج هر سه معیار ارزیابی (در مقایسه با مقادیر سه روش بهسازی دیگر) مناسب و قابل قبول نمی باشد.

نکته دیگری که از مشاهده نمودارها و جداول نتیجه می شود، این است که به ازای N های کوچک، با تغییر N ، اختلاف نتایج نسبتاً زیاد است. ولی با زیاد شدن N اختلاف نتایج کم و کمتر می شود. علت این امر این است که وقتی N کوچک است، به ازای تغییرات N تقریب های متناظر با N ها از طیف سیگنال تمیز اختلاف زیادی با هم دارند، ولی با افزایش N اختلاف بین تقریب های متناظر با N های متفاوت، کم می شود. (این مطلب با استفاده از معیار I_{KL} در فصل پنج نشان داده شد) با توجه به نمودارها، $N=30$ پاسخ خوبی داشته و پس از آن به ازای تغییرات N ، نتایج تغییرات زیادی ندارند. این نشان می دهد که تقریب توزیع سیگنال تمیز با توزیع LMM از حدود $N=30$ به بعد، به واقعیت نزدیک و نزدیک تر می شود. به همین جهت افزایش N در ۵۰ توقف می یابد، زیرا در حالیکه با افزایش N پیچیدگی محاسباتی و زمان اجرای برنامه به شدت افزایش می یابد، نتایج تغییرات زیادی ندارد. در واقع، مقدار بهینه برای N با توجه به مقادیر نمودارها و جداول، ۵۰ در نظر گرفته شد. بهینه بودن این مقدار از این جهت مطرح می شود که باید بین تقریب دقیق تر - نتایج بهتر و پیچیدگی - زمان اجرای برنامه - حجم حافظه اشغالی کمتر، مصالحه شود.

مطلب مهم دیگری که می توان به عنوان نتیجه به آن اشاره کرد، این است که؛ بهترین جواب ها براساس هر سه معیار، در سه سیگنال به نویز ۰ و ۵ و ۱۰ دسی بل، مربوط به نویز سفید گوسی است. پس از آن نویز همهمه و در نهایت نویز F_{16} قرار می گیرد. با مشاهده و بررسی طیف توان این سه نویز مشاهده شد که تغییرات توان نویز سفید از دو نویز دیگر کمتر است. همچنین این تغییرات برای نویز همهمه کمتر از نویز F_{16} می باشد. در بخش تخمین نویز بیان شد که فرض اساسی روش های تخمین نویز، کند بودن تغییرات توان نویز است. احتمالاً به این دلایل نتایج برای نویز سفید از دو نویز دیگر بهتر است و نویز همهمه نیز نتایج بهتری از نویز F_{16} دارد.

در مقایسه روش پیشنهادی با سه روش MMSE, Log-MMSE, Lap-MMSE، برای سیگنال های جمع شده با نویزهای F_{16} و همهمه در سه سطح سیگنال به نویزهای ۰ و ۵ و ۱۰ دسی بل، براساس سه معیار سیگنال به نویز قطعه ای، LLR و PESQ، این نتیجه حاصل شد که روش پیشنهادی برای حذف نویز همهمه در هر سه سطح سیگنال به نویز ذکر شده، از سه روش دیگر بهتر عمل می کند. اختلاف نتایج روش پیشنهادی و دو روش مبتنی بر توزیع گوسی خیلی زیاد است در حالیکه اختلاف نتایج روش پیشنهادی و روش Lap-MMSE زیاد نمی باشد. برای حذف نویز F_{16} ، روش پیشنهادی در همه سیگنال به نویزها از دو روش مبتنی بر توزیع گوسی موفق تر است ولی در مقایسه با روش مبتنی بر توزیع لاپلاس، این روش در سیگنال به نویزهای بالاتر (۵ و ۱۰ دسی بل) بهتر عمل می کند و در سیگنال به نویز ۰ دسی بل روش مبتنی بر لاپلاس نتایج بهتری دارد. برای هر دو نویز، با افزایش سیگنال به نویز، اختلاف مقادیر به دست آمده برای روش پیشنهادی با روش های دیگر بیشتر می شود که نشانگر بهبود عملکرد روش پیشنهادی با افزایش سیگنال به نویز است.

نکته آخر این است که با مقایسه جدول های مربوط به بخش ارزیابی روش پیشنهادی و جدول های بخش مقایسه روش پیشنهادی با روش های دیگر، مشخص می شود که به ازای N های کوچک مثل ۱ و ۴ و ۷ مقادیر روش پیشنهادی برای هر دو نویز و هر سه سطح سیگنال به نویز، کارایی مناسبی نسبت به روش مبتنی بر توزیع لاپلاس ندارد. زیرا در روش پیشنهادی به ازای N های کوچک، هم تقریب خوبی از سیگنال نداریم و هم پارامترهای توزیع سیگنال تمیز به صورت برون خطی تخمین زده می شوند. در حالی که در روش مبتنی بر لاپلاس، شاید یک لاپلاس برای تخمین طیف سیگنال تمیز مناسب نباشد ولی پارامترهای توزیع سیگنال تمیز به صورت درون خطی، تخمین زده می شوند.

در کل مقایسه و بررسی نتایج بیانگر عملکرد قابل قبول روش پیشنهادی است و این مطلب نشان می دهد که تقریب سیگنال تمیز با توزیع LMM مناسب تر از تقریب آن با توزیع های گوسی یا

لاپلاس است. البته این مناسب تر بودن از جهت بهبود نتایج می باشد. زیرا از لحاظ پیچیدگی محاسباتی، زمان اجرای برنامه ها و حجم حافظه اشغالی، سه روش دیگر مناسب تر هستند.

۳-۶- پیشنهادات برای کارهای آینده

- در روش پیشنهادی در این پایان نامه، مدل LMM برای سیگنال تمیز و مدل گوسی را برای نویز در نظر گرفته شد. پیشنهاد دیگر این است، که نویز را هم با توزیع لاپلاس، مخلوطی از لاپلاس یا مخلوطی از گوسی مدل کرد.
- در این پایان نامه، برای تخمین سیگنال تمیز از تخمین گر MMSE استفاده شد. استفاده از تخمین گر MAP شاید باعث کاهش محاسبات و بهبود نتایج شود.
- در این پایان نامه، مقدار دهی اولیه به EM الگوریتم به صورت تصادفی انجام شد. می توان مقدار دهی اولیه را با روش های دیگر مثل الگوریتم ژنتیک انجام داد که این کار احتمالاً منجر به ایجاد نتایج بهتری خواهد شد.
- در روش پیشنهادی، از تخمین گر MMSE برای تخمین بخش حقیقی و موهومی طیف سیگنال تمیز استفاده شد. در حالی که می توان دامنه و فاز سیگنال تمیز را تخمین زد. در این حالت باید توزیع دامنه و فاز ضرایب طیف را با فرض اینکه توزیع بخش حقیقی و موهومی LMM باشد، محاسبه کرد.
- در این پایان نامه، از روش ردیابی کمینه ها برای تخمین نویز استفاده شده است. می توان به جای این روش از روش های دیگر تخمین نویز نیز استفاده کرد.
- عملکرد روش را می توان علاوه بر تغییر N با تغییر پارامتر دیگری مثل تغییر طول قاب نیز سنجید.
- روش پیشنهادی را می توان با فرض عدم قطعیت گفتار هم به دست آورد. روش پیشنهادی در حوزه طیف عمل می کند. شاید بتوان این روش را به حوزه های دیگر مثل حوزه زمان، کپسترال و گسترش داد.

فهرست منابع

- [١] Accardi. A. J and Cox. R. V, (١٩٩٩), A Modular Approach to Speech Enhancement with an Application to Speech Coding, In IEEE International Conference on Acoustics, Speech and Signal Processing, pp. ٢٠١-٢٠٤, Arizona, USA
- [٢] Akbari.A and Lebouquin.R,(١٩٩٦),Speech enhancement using a Wiener Filtering under signal presence uncertainty, European signal processing Conference, pp ٩٧١-٩٧٤, Italy
- [٣] Andrianakis.I, (٢٠٠٧), Bayesian Algorithm for Speech Enhancement, Ph. D Thesis, Univecity os Southampton
- [٤] Berouti. M and Schwartz. R, Makhoul. J and Cambridge . M. A, (١٩٧٩) , Enhancement of Speech Corrupted by Acoustic Noise, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. ٢٠٨-٢١١, Washington,DC, USA
- [٥] Bilmes. J. A, (١٩٩٨), A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models, Technical report, Department of Electrical Engineering and Computer Science, California
- [٦] Boll. S ,(١٩٧٩), Suppression of Acoustic Noise in Speech Using Spectral Subtraction, **IEEE Transaction on Aocostic, Speech and Signal Processing**, Vol. ٢٧, No. ٢, pp. ١١٣-١٢٠
- [٧] Breithaupt. C and Martin. R, (٢٠٠٣), MMSE Estimation of Magnitude-Squared DFT Coefficients with Super Gausssian Priors, ICASSP, pp. ٨٩٦-٨٩٩, Hong Kong, China .
- [٨] Cappe. O, (١٩٩٤), Elimination of the Musical Noise Phenomenon with the Ephraim and Malah Noise Suppressor, **IEEE Transaction on Aocostic, Speech and Signal Processing**, Vol. ٢, No. ٢, pp. ٣٤٥-٣٤٩
- [٩] Chen. B, (٢٠٠٥), Speech Enhancement Using A Laplacian-Based MMSE Estimator, Ph. D. Thesis, Department of Electrical Engineering, Texas, USA, University of Texas
- [١٠] Chen. B, Loizou. P. C, (٢٠٠٧), A Laplacian-Based MMSE Estimator for Speech Enhancement, **Speech Communication**, Vol. ٤٩, No. ٢, pp. ١٣٤-١٤٣
- [١١] Cohen. I, (٢٠٠٢), Optimal Speech Enhancement Under Signal Presence Uncertainty Under Spectral Amplitude Estimator, **IEEE Signal Processing Letters**, Vol. ٩, No. ٤, pp.١١٣-١١٦

- [12] Cohen. I, (2004), On the Decision-Directed Estimation Approach of Ephraim and Malah, ICASSP, pp. 293-6, Montreal, Canada .
- [13] Cohen. I, (2005b), Relaxed Statistical Model for Speech Enhancement and a priori SNR Estimation, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 13, No. 5 part 2, pp. 336-350.
- [14] De Moor. B, (1993), The Singular Value Decomposition and Long and Short Space of Noisy Matrices, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 41, No. 9, pp. 2826-2838
- [15] Dempster. A. P, Laird. N. M and Rubin. D. B, (1977). Maximum Likelihood from Incomplete Data via the EM Algorithm, **Journal of the Royal Statistical Society**, Vol. 39, No. 1, pp. 1-38
- [16] Dendrinos.M, Bakamidis.S, and Carayannis. G, (1991), Speech Enhancement from Noise: a regenerative approach, **speech communication**, Vol. 10, No. 1, pp. 40-46
- [17] Drucker. H, (1968), Speech Processing in a High Ambient Noise Environment, **IEEE Transaction on Audio and Electroacoustics**, Vol. 16, No. 2, pp. 160-8
- [18] Ephraim. Y, (1992a), A Bayesian Estimation Approach for Speech Enhancement Using Hidden Markov Models, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 40, No. 4, pp. 720-730
- [19] Ephraim. Y, (1992b), Gain Adapted Hidden Markov Models for Recognition of Clean and Noisy Speech , **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 40, No. 6, pp. 1303-1316
- [20] Ephraim. Y, (1992c), Statistical Model-Based Speech Enhancement System , **Proceeding of the IEEE** , Vol. 80, No. 10, pp. 1026-1000
- [21] Ephraim. Y, Cohen. I, (2000), Recent Advancements in Speech Enhancement, In the Electrical Engineering Handbook, CRC Press
- [22] Ephraim. Y and Malah. D, (1984), Speech Enhancement Using a Minimum Mean Square Error Short Time Spectral Amplitude Estimator, IEEE International Conference on Acoustics, Speech and Signal Processing , pp. 1109-1121, San Diego, California
- [23] Ephraim. Y and Malah. D, (1980), Speech Enhancement Using a Minimum Mean

Square Error Log Spectral Amplitude Estimator, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 443-6, Tampa, Florida, USA

[24] Ephraim. Y and Malah. D, Juang. B. H, (1989), On the Application of Hidden Markov Models for Enhancing Noisy Speech, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 37, No. 12, pp. 1846-1856

[25] Ephraim. Y and Van Trees. H. L, (1990), A Signal Subspace Approach for Speech Enhancement, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 38, No. 12, pp. 202-216

[26] Eekelens. J. S, Jensen. J and Heusdens. R, (2007b). Speech Enhancement Based on Rayleigh Mixture Modeling of Speech Spectral Amplitude Distributions, In European Signal Processing Conference, pp. 60-9

[27] Haigh. J. A and Mason. J. S, (1993), Robust Voice Activity Detection Using Cepstral Feature, In IEEE International Conference on Computer, Communication, Control and Power Engineering, pp. 321-4

[28] Hendriks. R. C, (2008), Advances in DFT-Based Single-Microphone Speech Enhancement, Ph. D Thesis, Technischm Univerity, Delft

[29] Hirsch. H. G and Ehrlicher. C, (1990), Noise Estimation Technique for Robust Speech Recognition, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 103-106, Detroit, Michigan, USA

[30] Hu. Y and Loizou. P, (2006), Subjective Comparison of Speech Enhancement Algorithms, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 103-6, Toulouse, France

[31] Hu. Y and Loizou. P, (2008), Evaluation of Objective Quality Measures for Speech Enhancement, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 16, No. 1, pp. 202-216

[32] Junqua. J. C, Reaves. B and Mark. B (1991), A Study of Endpoint Detection Algorithm in Adverse Conditions: Incidence on a DTW and HMM Recognizer, In European on Speech Communication and Technology, pp. 1371-1374, Genova, Italy

[33] Kondoz. A. M, (2004), Digital Speech: Coding for low Bit Rate Communication System, Wiley .

[34] Laird.D, Dempster.A.P, Rubin.D.B,(1997), Maximum Likelihood from Incomplete Data via the EM Algorithm, **Journal of the Royal Statistical Society**, Vol. 39, No.1,

pp. 1-38.

[30] Lim. J and Oppenheim. A, (1979), Enhancement and Bandwidth Compression of Noisy Speech, **Proceeding of the IEEE**, Vol. 67, No. 12, pp. 1086-1104

[36] Lin. L, Holmes. W. H and Ambikairajah. E, (2003a), Adaptive Noise Estimation Algorithm for Enhancement, **Electrical Letters**, , Vol. 39, No. 9, pp. 704-706

[37] Loizou. P. C, (2007), **Speech Enhancement: Theory and Practice**, CRC Press

[38] Lotter. T, Vary. P, (2008), Speech Enhancement by MAP Spectral Amplitude Estimation Using a Super-Gaussian Speech Model, **EURASIP Journal on Applied Signal Processing**, Vol. 2008, No. 7, pp. 1110-1126

[39] Malah. D, Cox. R, V and Accardi. A. J, (1999), Tracking Speech Presence Uncertainty to Improve Speech Enhancement in Non-Stationary Noise Environment, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 789-792, Arizona, USA

[40] Martin . R, (1994), Spectral Subtraction Based on Minimum Statistic, European signal processing Conference, pp 1182-86, Edinburgh, Scotland

[41] Martin . R, (2001), Noise Power Spectral Density Estimation Based on Optimal Smoothing and Minimum Statistic, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 9, No. 6, pp. 54-62

[42] Martin . R, (2003), Statistical Method for Enhancement of Noisy Speech, In International Workshop on Acoustic Echo and noise control, pp. 43-60, Kyoto, Japan .

[43] Martin . R and Breithaupt. C, (2003), Speech Enhancement in the DFT domain using Laplacian Speech Priors, In International Workshop on Acoustic Echo and noise control, pp. 87-90, Kyoto, Japan .

[44] Martin.R,(2008), Speech Enhancement Based on Minimum Mean Square Error and Supergaussian Priors, **IEEE Transaction on Speech and Audio Processing**, Vol.13, No.6

[45] McAulay. R and Malpass. M, (1980), Speech Enhancement Using a Soft-Decision Noise Suppression Filter, **IEEE Transaction on Acoustic, Speech and Signal Processing**, Vol. 9, No. 6, pp. 54-62

[46] Mitianoudis. N and Stathaki. T, (2004), Overcomplete source separation using Laplacian Mixture Model, IEEE Signal Processing Letters, Vol. 1, No. 5 .

[47] Paliwal. K. K and Basu. A, (1987), A Speech Enhancement Method Based on Kalman Filtering, IEEE International Conference on Acoustics, Speech and Signal

Processing, pp. 117-118, Dallas, Texas, USA

[18] Paliwal. K and Alsteris. L, (2000), On the Usefulness of STFT Phase Spectrum in Human Listening Tests, **Speech Communication**, Vol. 40, No. 2, pp. 103-110.

[19] Papoulis, A and Pillai. S, (2002), **Probability, Random Variables and Stochastic Processes**, New York, McGraw-Hill

[20] Porter. J and Boll. S, (1988), Optimal Estimator for Spectral Restoration of Noisy Speech, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 11A. 2. 1 - 11A. 2. 4, San Diego, California, USA

[21] Rabiner. L and Sambur. M, (1977), Voiced-Unvoiced Detection Using the Itakura

LPC Distance Measure, ICASSP, pp. 323-6

[22] Rangachari. S, (2004), Noise Estimation Algorithms For Highly NON-Stationary Environment, Ph. D Thesis, University of Texas at Dallas

[23] Ris. C and Dupont. S, (2001), Assessing Local Noise Level Estimation Methods: Application to Noise Robust ASR. **Speech Communication**, Vol. 43, No. 1, pp. 144-108

[24] Sohn. J and Sung. W, (1998), A Voice Activity Detector Employing Soft Decision Based Noise Spectrum Adaptation, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 360-8, Seattle, USA

[25] Soon. I and Yeo. C.K, (1999), Improved Noise Suppression Filter Using Self-Adaptive Estimator of probability of Speech Absence, **Signal Processing**. Vol. 79, No. 2, pp. 101-109

[26] Sudirga. R, (2009), A Speech Enhancement System Based on Statistical and Acoustic-Phonetic Knowledge, M. S Thesis, Queen's University, Ontario, Canada .

[27] Tucker. R, (1992), Voice Activity Detection Using periodicity Measure, **IEEE Proceeding on Communications, Speech and Vision**, Vol. 139, No. 4, pp. 377-380

[28] Vary. P, (1980), Noise Suppression by Spectral Magnitude Estimation Mechanism and Theoretical Limits, **Signal Processing**. Vol. 8, No. 4,

[29] Veisi. H, (2000), Model-Based Method for Noise Robust Speech Recognition Systems, M.S. Thesis, Computer Engineering Department, Sharif University of Technology, Tehran, Iran

[30] Wiener. N, (1949), **Extrapolation, Interpolation and Smoothing of Stationary Time Series: Engineering Application**, MIT Press, Cambridge, Mass

[۶۱] Wolf. P. J and Godsill. S. J, (۲۰۰۱), Simple Alternative to the Ephraim and Malah Suppression Rule for Speech Enhancement, **IEEE Workshop on Statistical Signal Processing**, Vol. ۲, pp. ۴۹۶-۴۹۹

[۶۲] Wolf. P. J and Godsill. S. J, (۲۰۰۳a), Efficient Alternative to the Ephraim and Malah Suppression Rule for Audio Signal Enhancement, **EURASIP Journal on Applied Signal Processing**, Vol. ۲۰۰۳, No. ۱۰, pp. ۱۰۴۳-۱۰۵۵

[۶۳] You. C. H and Rahardja. S, (۲۰۰۳), Adaptive B-Order MMSE Estimator for Speech Enhancement, IEEE International Conference on Acoustics, Speech and Signal Processing, pp. ۸۲۵-۸۵۵, Hong Kong, China

[۶۴] <http://WWW.Wikipedia.org/LaplaceDistribution.Html>

[۶۵] <http://WWW.Wikipedia.org/TIMIT.Html>

[۶۶] اسفندیان. ن، (۱۳۸۶)، استفاده از آنالیز LPC در روش تبدیل موجک، جهت غنی سازی

سیگنال گفتار، پایان نامه کارشناسی ارشد، دانشکده فنی مهندسی، دانشگاه مازندران.

Summary

The name " speech enhancement " refers to larg group of methods that improve the quality and intelligibility of noisy speech by suppressing the background noise from noisy signal.

There are a lot of methods and papers about speech enhancement, but estimation of the clean signal from the noisy signal is still one of the insoluble problem in speech processing field .

One of the most important methods for speech enhancement, are statistical methods, because they have better performance.

In these methods, a distribution for clean speech and noise are assumed. Then an estimator are used for estimating the clean signal from noisy signal.

In this thesis, an estimator for speech enhancement in Discrete Fourier Transform (DFT) domain is proposed. It is shown that the complex DFT coefficients of clean speech can be modeled more accurate, by Laplacian Mixture Model than Gaussian, Laplacian and GMM distributions. Then, an analytical solution for estimating the complex DFT coefficients with the MMSE (Minimum Mean Square Error) estimator is derived, when the clean speech DFT coefficients are mixture of Laplacians distributed and the DFT coefficients of noise are Gaussian distributed. The derived MMSE estimator is non-linear and it is shown that, this estimator has better performance than the estimators which are based on Laplacian or Gaussian models.

For estimating LMM's parameter the Expectation Maximization (EM) Algorithm is used. Indeed, the parameters of LMM are estimated offline with clean speech from TIMIT data base and the parameters of noise are estimated online with minimal tracking based method.

Finally, the proposed algorithm is evaluated in term of Segmental SNR, LLR(Log Likelihood Ratio) and PESQ (Perceptual Evaluation Of Speech Quality) and then the proposed method is compared with Laplacian and Gaussian based methods.

The results of this comparison shows that the proposed method has an acceptable performance.

Important words:

Speech enhancement, Laplacian Mixture Model, Expectation Maximization, MMSE estimator, Voice Activity Detector, Evaluation methods.