





دانشگاه صنعتی امیرکبیر

دانشکده مهندسی برق

پایان نامه کارشناسی ارشد گرایش الکترونیک

مقاوم سازی بازشناسی گفتار با اعمال پردازش زیرباندی

نگارش:

حجت اله یگانه

استاد راهنما:

دکتر سید محمد احدی

بهمن ماه ۱۳۸۷

تقدیم بہ

پدری مہربان،

مادری فداکار

و ہمسری دل آرام

برخود لازم می‌دانم از زحمات بی‌شائبه استاد ارجمندم جناب آقای دکتر احدی تشکر و
قدردانی نمایم. بی‌شک، بدون حمایت‌های دلسوزانه خانواده عزیزم پیمودن راه تاکنون امکان‌پذیر
نمی‌بود. امیدوارم بتوانم قدردان زحمات ایشان باشم.

این پروژه تحت قرارداد پژوهشی شماره ۵۰۰/۸۱۹۹/ت مورخ ۸۷/۵/۲۸ از حمایت های مالی و معنوی مرکز تحقیقات مخابرات ایران بهره مند شده است.

چکیده

در این پروژه مقاوم سازی بازشناسی گفتار در محیط های نویزی بر مبنای پردازش زیرباندی بررسی شده است. مقاوم سازی بازشناسی گفتار یکی از مسائل مهم در این حوزه می باشد که کار بر روی آن همچنان ادامه دارد. از روش های گوناگونی به منظور تحقق یافتن این مهم استفاده می شود و ایده های متنوعی نیز در مقالات و تحقیقات ارائه می گردند. عیب عمده اکثر روش های پیشنهاد شده پیچیدگی زیاد و سرعت کم الگوریتم های آن است. ما در این رساله تلاش نموده ایم به ارائه روش هایی پردازیم که عیوب مذکور را در حد بسیار کمتری داشته باشند و در عین حال هدف ما را که همانا افزایش دقت بازشناسی گفتار در محیط های نویزی است محقق سازند.

از این رو در ابتدا با بررسی روند استخراج ویژگی های مطرح MFCC سعی نمودیم عیب این بردار ویژگی را برطرف نماییم. دلیل این امر آن است که بردار ویژگی MFCC دارای الگوریتم بسیار ساده و سریعی است و در محیط های عاری از نویز یا به اصطلاح تمیز از دقت بازشناسی خیلی خوبی برخوردار است. بنابراین اگر بتوان الگوریتمی پیشنهاد نمود که ساختاری شبیه MFCC داشته باشد و در عین حال بردار ویژگی حاصل نسبت به نویز محیط مقاوم تر نشان دهد، این روش از امتیاز بالایی برخوردار می گردد. با دنبال کردن روند استخراج ویژگی MFCC به این مسئله پی می بریم که برای به دست آمدن هر کدام از مولفه های این بردار، تمام طیف یک فریم در محاسبه تاثیر دارد. این بدان معنی است که آلوده بودن نواحی خاصی از طیف به تمام بردار MFCC سرایت می کند و کارایی این بردار ویژگی را به شدت پایین می آورد.

روند کلی ایده های پیشنهادی ما حول دو گام اصلی قابل بیان است. در ابتدا سعی بر آن داریم تا با فیلتر کردن سیگنال گفتار در حوزه زیرباندی میزان تاثیر نویز را کم کنیم. در ادامه و در گام دوم با اعمال وزن بر خروجی زیرباندهای حوزه مل میزان مشارکت زیرباندهای با کیفیت تر را در به دست آوردن ویژگی های پیشنهادی خود بیشتر نموده و از تاثیر زیرباندهای با کیفیت کمتر می کاهیم.

نتایج آزمایش های ما نشان دادند که تا حد خیلی خوبی به اهداف اصلی خود رسیده ایم. الگوریتم های پیشنهادی ما در عین حالی که ساده هستند در محیط های نویزی بسیار مقاوم می باشند. مقایسه روش های پیشنهادی با روش های مطرح دیگر بر این مطلب صحه می گذارد. روش های ارائه شده منجر به بهبود ۳۲ درصدی روش پایه شده است.

واژگان کلیدی: مقاوم سازی بازشناسی گفتار، فیلترکردن در زیرباند، وزن دهی زیرباندهای مل

فهرست مطالب

فصل اول	۱
مقدمه	۱
فصل دوم	۸
مروری بر روشهای مقاوم سازی بازشتاسی گفتار	۸
۱-۲-۱- روش های بهبود گفتار	۱۰
۱-۱-۲- تفاضل طیف	۱۲
۲-۱-۲- فیلتر وینر	۱۴
۲-۲- روش جبران ویژگی	۱۵
۱-۲-۲- روش های اعمال تبدیل بر بردار ویژگی	۱۶
۲-۲-۲- روش های اعمال تبدیل بر طیف	۱۸
۳-۲- روش های تطبیق مدل	۲۷
۱-۳-۲- معیار تصویروزن دار	۲۸
۲-۳-۲- بازگشت خطی با بیشترین درستنمایی (MLLR)	۲۹
۳-۳-۲- ترکیب موازی مدل ها	۳۲
۴-۲- روش های مبتنی بر خواص شنوایی انسان	۳۳
۱-۴-۲- بازشناسی چندباندی گفتار	۳۳
۲-۴-۲- روش ویژگیهای گم شده	۳۶
فصل سوم	۳۹
معرفی دادگان و سیستم بازشناسی مورد استفاده در آزمایشها	۳۹
۱-۳- دادگان AURORA2	۳۹
۲-۳- نرم افزار HTK	۴۳
۳-۳- نرم افزار MATLAB	۴۴
فصل چهارم:	۴۶
روشهای پیشنهادی	۴۶
۱-۴- مقدمه	۴۶
۲-۴- الگوریتم استخراج بردار ویژگی MFCC	۴۷
۳-۴- دورنگاه روشهای پیشنهادی	۵۱

۴-۴-۴-۴	فیلتر کردن سیگنال گفتار در زیرباندها	۵۲
۴-۴-۱-۴	فیلتر تفاضل طیفی زیر بانندی	۵۳
۴-۴-۲-۴	اعمال فیلتر وینر در حوزه زیرباند ها	۵۷
۴-۵-۵-۴	وزن دهی زیرباندها در رابطه تبدیل کسینوسی گسسته	۵۹
۴-۵-۱-۴	وزن دهی بر اساس SNR در هر زیرباند	۶۰
۴-۵-۲-۴	وزن دهی بر اساس نسبت SNR به آنتروپی در هر زیرباند	۶۲
۶۸	فصل پنجم :	
۶۸	آزمایشها و نتایج	
۶۸-۱-۵	نتایج روش فیلتر کردن به روش تفاضل طیفی و وزن دهی بر اساس SNR	۶۸
۶۸-۲-۵	نتایج روش فیلترکردن به روش وینر در زیرباندها و وزن دهی بر اساس نسبت SNR به آنتروپی	۷۱
۶۸-۳-۵	روش فیلتر کردن به روش وینر در زیر باند با استفاده از تخمین نویز چندی گزینی و وزن دهی بر	
۶۸	اساس نسبت SNR به آنتروپی در زیرباندها	۷۸
۸۱	فصل ششم :	
۸۱	جمع‌بندی، نتیجه گیری و پیشنهادات	۸۱
۸۷	مراجع	۸۷
۹۳	واژه نامه انگلیسی به فارسی	۹۳

فهرست شکل ها

- شکل (۱-۲) : دسته بندی روش های مقاوم سازی ۱۰
- شکل (۲-۲) : بلوک دیاگرام عملیات RASTA ۲۰
- شکل (۳-۲) : پاسخ ضربه و پاسخ فرکانسی فیلتر RASTA ۲۱
- شکل (۱-۳) : مشخصه فرکانسی فیلترهای G712 و MIRS ۴۰
- شکل (۲-۳) : مشخصات فرکانسی نویزهای اضافه شده به دادگان AURORA2 ۴۲
- شکل (۱-۴) : فیلترپیش تاکید ۴۸
- شکل (۲-۴) : مشخصه زمانی و فرکانسی پنجره مستطیلی ۴۹
- شکل (۳-۴) : مشخصه زمانی و فرکانسی پنجره همینگ ۴۹
- شکل (۴-۴) : بانک فیلتر استفاده شونده در فرایند استخراج MFCC ۵۰
- شکل (۵-۴) : اندازه انرژی طیف فرکانس های ۲k، ۴k و ۳۰۰ هرتز در SNR=10dB ۵۶
- شکل (۶-۴) : اندازه انرژی طیف فرکانس های ۲k، ۴k و ۳۰۰ هرتز در سیگنال تمیز ۵۶
- شکل (۷-۴) : اندازه انرژی طیف فرکانس های ۲k، ۴k و ۳۰۰ هرتز در SNR=-5dB ۵۷
- شکل (۸-۴) : رابطه بین SNR_i ، ξ_i و W_i برای حالتی که $\gamma = 1$ است ۶۱
- شکل (۹-۴) : رابطه بین W_i و γ ۶۲
- شکل (۱۰-۴) : رابطه بین SNR_i ، $Entropy_i$ ، R_i و W_i ۶۶
- شکل (۱-۵) : مقایسه روش پیشنهادی ۱-۵ با روش های مختلف در مجموعه A دادگان AURORA2 ۶۹
- شکل (۲-۵) : مقایسه روش پیشنهادی ۱-۵ با روش های مختلف در مجموعه B دادگان AURORA2 ۶۹
- شکل (۳-۵) : مقایسه روش پیشنهادی ۱-۵ با روش های مختلف در مجموعه C دادگان AURORA2 ۷۰
- شکل (۴-۵) : مقایسه روش پیشنهادی ۲-۵ با روش های مختلف در مجموعه A دادگان AURORA2 ۷۲
- شکل (۵-۵) : مقایسه روش پیشنهادی ۲-۵ با روش های مختلف در مجموعه B دادگان AURORA2 ۷۲
- شکل (۶-۵) : مقایسه روش پیشنهادی ۲-۵ با روش های مختلف در مجموعه C دادگان AURORA2 ۷۳

- شکل (۷-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز همهمه ۷۴
- شکل (۸-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز مترو ۷۴
- شکل (۹-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز اتومبیل ۷۵
- شکل (۱۰-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز نمایشگاه ۷۵
- شکل (۱۱-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز رستوران ۷۶
- شکل (۱۲-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز خیابان ۷۶
- شکل (۱۳-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز هواپیما ۷۷
- شکل (۱۴-۵) : مقایسه روش پیشنهادی ۲-۵ با روش‌های مختلف در نویز ایستگاه قطار ۷۷
- شکل (۱۵-۵) : مقایسه روش پیشنهادی ۳-۵ با روش‌های مختلف در مجموعه A دادگان AURORA2 ۷۹
- شکل (۱۶-۵) : مقایسه روش پیشنهادی ۳-۵ با روش‌های مختلف در مجموعه B دادگان AURORA2 ۸۰
- شکل (۱۷-۵) : مقایسه روش پیشنهادی ۳-۵ با روش‌های مختلف در مجموعه C دادگان AURORA2 ۸۰

فهرست جدول‌ها

جدول (۱-۵) : مقایسه متوسط درصد بازشناسی روش ۱-۵ با چند روش دیگر در مجموعه های A، B و C.... ۷۰

جدول (۲-۵) : مقایسه متوسط درصد بازشناسی روش ۲-۵ با چند روش دیگر در مجموعه های A، B و C.... ۷۸

جدول (۳-۵) : مقایسه متوسط درصد بازشناسی روش ۳-۵ با چند روش دیگر در مجموعه های A، B و C.... ۷۹

فصل اول

مقدمه

فصل اول

مقدمه

امروزه سیستم های بازشناسی گفتار کاربردهای زیادی پیدا کرده اند. در این میان سهم کشورهای پیشرفته از جنبه تکنولوژی به مراتب بیشتر از سایر کشورها بوده است. گستردگی دنیای ماشینی و رایانه هایی که اطراف ما را پوشانده اند محققین را به سمت استفاده از این ابزار سوق داده است. بشر امروز سعی دارد تا با حداقل زمان ممکن بیشترین کارایی روزانه را داشته باشد. یکی از مشکلات پیش روی این مهم رابطه میان انسان و ابزار مورد نیاز اوست. سیستم های بازشناس گفتار یکی از راه حل هایی است که توجه محققین را به خود معطوف داشته است. از ساده ترین کاربردهای این سیستم می توان به تایپ کردن در محیط های نوشتاری رایانه ای با بیان متن مورد نظر و بدون تایپ دستی اشاره نمود. از طرف دیگر انسان امروز تلاش می کند به دلایل اقتصادی تا حد ممکن حضور خود را در محیط های کاری کم کند. از اینرو سیستم های بازشناسی گفتار در مراکز ارائه دهنده اطلاعات تلفنی نیز گسترش یافته اند. این گونه مراکز به سیستم های رایانه ای مجهز شده اند که در این رایانه ها سیستم های بازشناس گفتار تعبیه نموده اند. با تماس با این مراکز می توان اطلاعات مورد نیاز خود را از این رایانه ها دریافت کرد به گونه ای که بدون حضور فرد خاصی این سیستم بازشناسی گفتار است که پیام پرسنده را تشخیص می دهد و پس از تشخیص آن، داده های مورد نیاز را در اختیار طرف مقابل می گذارد. کاربردهای این گونه

سیستم‌ها به این موارد ختم نمی‌شوند و امروزه همچنان محققین به دنبال زمینه‌های کاری این رشته می‌باشند.

اگر بخواهیم نمایی کلی از یک سیستم بازشناسی گفتار ارائه دهیم باید از معرفی نوع سیگنال گفتار آغاز کنیم. همان گونه که می‌دانیم در سیگنال‌های تصادفی مشخصات آماری مختلفی وجود دارد که با تحلیل آنها می‌توان به اطلاعات مهمی در مورد یک سیگنال دست یافت. اگر این مشخصات آماری مانند میانگین و واریانس داده‌ها در طول زمان ثابت باشند این سیگنال را ایستاد^۱ می‌گویند. اما اگر چنین نباشد و این اطلاعات در طول زمان متغیر باشد سیگنال را نایستاد^۲ می‌نامند. معمولاً تحلیل سیگنال‌های نایستاد مشکلات مخصوص به خود را دارد. کافی است به عنوان مثال به این مطلب اشاره نماییم که تمام سواد حوزه فرکانسی ما و تبدیل فوریه و ابزارهای دیگر برای سیگنال‌های ایستاد تعریف شده‌اند. یکی از راه حل‌های این مشکل بریدن سیگنال به قطعاتی است که مشخصات آماری سیگنال در آن‌ها تقریباً ثابت باشد. اصطلاحاً به این قطعات فریم گفته می‌شود. در برخی از کاربردها مانند تجزیه و تحلیل سیگنال‌های راداری طول این فریم‌ها متفاوت است. به بیان دیگر بررسی می‌کنند که سیگنال را در چه بازه‌ای می‌توان ایستاد در نظر گرفت و بنابر ماهیت این گونه سیگنال‌ها بازه ایستاد بودن متغیر است. سیگنال گفتار نیز یک سیگنال غیر ایستاد است و از اینرو برای تحلیل آن نیاز به فریم بندی کردن آن داریم. معمولاً طول فریم‌ها را در مورد سیگنال گفتار ثابت در نظر می‌گیرند و گفتار را در آن بازه‌ها پردازش می‌کنند. با این فرض که سیگنال گفتار در فریم‌ها ایستاد است می‌توان این فریم‌ها را به حوزه فرکانس منتقل کرد و در آن حوزه نیز تغییراتی را به سیگنال اعمال کرد. در ادامه، مسئله‌ای که وجود دارد حجم زیاد داده‌ها و اطلاعات در هر فریم است. برای مثال اگر یک سیگنال گفتار با فرکانس ۸ کیلوهرتز نمونه برداری شده باشد و طول فریم‌ها را ۲۰ میلی ثانیه در نظر بگیریم - که بازه معمولی هم هست - در آن صورت برای هر فریم مقدار ۱۶۰ نمونه را در اختیار داریم. کافی است حجم اطلاعات را برای یک

^۱ Stationary

^۲ Non Stationary

گفتار ۱۰ دقیقه ای محاسبه کنیم تا به ابعاد مشکل پی ببریم. در این حالت راه چاره ای که دانشمندان جسته اند، مشکل را تا حد زیادی برطرف نموده است و این راهکار، استخراج بردار ویژگی از فریم ها است.

محققین راه ها و الگوریتم هایی را ارائه نموده اند که یک سری نماینده تحت عنوان ویژگی از هر فریم گفتار استخراج شود به گونه ای که این ویژگی ها نماینده ای برای مشخصات و خصوصیات یک فریم گفتار باشد. بردارهای ویژگی متعددی در این راستا معرفی شده اند که از مهمترین آنها در بحث سیگنال گفتار می توان به بردار ضرایب کپسترال حوزه فرکانس مل یا MFCC^۱ اشاره نمود.

معمولاً هدف از بازشناسی گفتار تشخیص این مطلب است که گفتار مورد نظر به اصطلاح چه می گوید. به عبارت دیگر می توان گفت هدف از بازشناسی گفتار تشخیص مطلب ادا شده در این گفتار است. اگر بتوان این مسئله را تشخیص داد آنوقت این گفتار می تواند به متن تبدیل شود یا به یک دستور برای راه اندازی یک سیستم الکتریکی یا کاربردهای دیگری که قبلاً اشاره شد. هنگامی که صحبت از تشخیص مطلب یا اطلاعاتی از یک سیگنال آن هم توسط رایانه باشد، چگونگی تبدیل این تشخیص به الگوریتم های ریاضی و به تبع آن رایانه ای سوالی است که باید جواب داده شود. ساختار پیچیده مغز انسان برای تشخیص و پردازش اطلاعات به مرور زمان آموزش دیده است. به این خاطر است که انسان حاوی چنین قدرت بالای تشخیص و تصمیم گیری می باشد. اگر این انتظار را از رایانه نیز داریم باید شبیه اتفاقی که در ساختار مغز انسان می افتد را شبیه سازی نماییم.

انسان برای تشخیص یک مطلب سعی می کند ماهیت داده جدیدی که می بیند را با اطلاعات گذشته مقایسه نماید و بعد از این مقایسه تصمیم گیری کند. پس اگر قرار است الگوریتم های طبقه بندی کننده^۲ نیز این عمل را انجام دهند باید در ابتدا با یک سری از داده ها و اطلاعات آموزش ببینند تا ساختار آن ها را با داده های جدید مقایسه کنند و آنگاه تشخیص دهند که ماهیت و اطلاعات یک ورودی جدید

^۱ Mel Frequency Cepstral Coefficients

^۲ Classifier

چيست. با اين مقدمه مشخص شد که در بحث بازشناسی گفتار نیز ما با اين مسئله روبرو هستيم که گفتار ورودی باید با داده های آموزش و از طريق یک طبقه بندی کننده مقایسه شود تا مشخص شود گوینده چه مطلبی را ادا کرده است.

در بازشناسی گفتار سعی بر آن است که سیستم را با داده های تمیز یا به عبارت دیگر سیگنال های گفتاری که در محیط با نویز بسیار کم ادا شده اند آموزش دهند. طبیعی است همانگونه که بیان شد آموزش از طریق بردارهای ویژگی صورت می پذیرد. در ادامه و به هنگام ورود یک داده جدید این سیستم ها نیاز به یک طبقه بندی کننده دارند که تشخیص دهند آیا این گفتار ورودی شبیه آن چه با آن آموزش دیده اند می باشند یا خیر. طبقه بندی کننده های زیادی در بحث شناسایی آماری الگو مطرح می شوند که از جمله آنها می توان به شبکه های عصبی مصنوعی (ANN^1) یا مدل مخلوط گوسی (GMM^2) اشاره نمود. اما در بحث گفتار طبقه بندی کننده بسیار مطرح و کارآیی که مورد استفاده قرار می گیرد مدل های مخفی مارکف یا HMM^3 است. دلیل استفاده از این طبقه بندی کننده ساختار خاص آن است که به یک ویژگی بسیار مهم گفتار منطبق است. این ویژگی پویا بودن سیگنال گفتار در طول زمان می باشد. به بیان ساده تر مشخصه مهم گفتار این است ادای یک کلمه حتی از یک گوینده در دو بار پی در پی مشخصه زمانی و فرکانسی مختلفی دارد. از این رو برای تشخیص یک کلمه آنچه اهمیت دارد تشخیص این مطلب است که کلمه از کجا تا کجا ادا شده و احتمال اینکه این سیگنال زمانی کلمه خاصی باشد چقدر است. در واقع این پویا بودن در حوزه زمان است که ما را به سمت استفاده از مدل مخفی مارکف برای بازشناسی می برد. مدل مخفی مارکف می تواند احتمال یک مسیر خاص را در یک داده گفتاری و فارغ از طول داده محاسبه نماید و با مقایسه این احتمال با آنچه آموزش دیده است تشخیص دهد که چه کلمه یا حتی جمله ای ادا شده است. لازم به ذکر است از HMM می توان به منظور استخراج ویژگی نیز استفاده نمود.

¹ Artificial Neural Networks

² Gaussian Mixture Models

³ Hidden Markov Models

آموزش در بحث بازشناسی گفتار در واقع تنظیم پارامترهای طبقه بندی کننده HMM است. HMM احتمال وقوع بردارها را در حالت ها محاسبه می نماید اما آنچه برای ما مهم است این مطلب میباشد که در بین این احتمالات مسیر بهینه کدام است و در انتها احتمال وقوع این مسیر بهینه چقدر است و به چه گفتاری شبیه میباشد. پیدا کردن این مسیر بهینه بر عهده الگوریتم های جستجو است. در بحث گفتار الگوریتم جستجویی که بسیار مورد استفاده قرار میگیرد الگوریتم ویتربی¹ می باشد.

مسئله مهمی که گریبان سیستم های بازشناسی گفتار را گرفته است این است که مدل های HMM معمولاً با دادگانی آموزش می بینند که در شرایط خاصی ضبط شده اند. مثلاً در شرایط با نویز بسیار کم. بنابراین این هنگامی که داده هایی که در شرایط دیگری ضبط شده اند را به سیستم میدهیم، HMM با مشکل مواجه می گردد. زیرا جنس سیگنال گفتاری که به آن نشان می دهیم متفاوت از سیگنالی است که در حافظه دارد. این مسئله در حضور نویز محیط به شدت حاد می شود و دانشمندان مدتهاست به فکر چاره ای برای این مشکل هستند. بهبود دقت بازشناسی گفتار در محیطهای نویزی را مقاوم سازی بازشناسی گفتار می گویند. زمینه کاری این رساله نیز رفع این مشکل است.

روش های گوناگونی برای مقاوم سازی بازشناسی گفتار ارائه شده اند که می توان آنها را به چهار دسته کلی تقسیم بندی کرد. دسته اول روش های بهسازی سیگنال گفتار است. هدف از این روش ها کم کردن نویز سیگنال به منظور کاهش عدم تطبیق بین داده های آموزش و آزمایش می باشد. روش های جبران ویژگی یا استخراج ویژگی های مقاوم، روش های تطبیق مدل و روشهای مبتنی بر ساختار گوش انسان از دیگر راهکارهای ارائه شده توسط محققین است. ما نیز توجه خود را معطوف به ارائه یک روش پیشنهادی نموده ایم که با بررسی زیرباندهای فرکانسی بردار ویژگی مقاوم تری را نسبت به بردارهای موجود از سیگنال گفتار استخراج نماید.

در ادامه و در فصل دوم مروری بر روش های مطرح مقاوم سازی بازشناسی گفتار خواهیم داشت. فصل سوم به توضیحات در باب دادگان و نرم افزارهای مورد استفاده در این پروژه می پردازد. روش ها و

¹ Viterbi

ایده های پیشنهادی خود را در فصل چهارم تشریح نموده ایم و در ادامه نتایج آزمایش ها را در فصل پنجم مشاهده خواهید کرد. در انتها جمع بندی، نتیجه گیری و پیشنهادات برای ادامه کار در فصل ششم آورده شده است.

فصل دوم

مروری بر روش‌های مقاوم سازی بازشناسی گفتار

فصل دوم

مروری بر روش‌های مقاوم‌سازی بازشناسی گفتار

یکی از مشکلات اصلی در سیستم‌های بازشناسی گفتار افت کارایی آن در محیط‌های نویزی است. از آنجاییکه معمولاً آموزش سیستم‌های بازشناس گفتار در محیط‌های به اصطلاح تمیز یا با نویز بسیار کم انجام می‌شود لذا این سیستم‌ها هنگام کار در محیط‌های نویزی با مشکل روبرو می‌شوند. نویزهای محیطی را معمولاً به دو دسته نویز کانال و نویز جمع پذیر تقسیم می‌کنند. از طرفی نویز می‌تواند ایستاد¹ یا نا ایستاد باشد. تلاش محققین بر آن است روش‌هایی ارائه دهند که جامعیت مطلوبی را دارا باشد به این معنا که تاثیر نویز را جدا از ایستاد یا نا ایستاد بودن و جمع پذیر یا کانال بودن کاهش دهند.

در طی سالیان اخیر ایده‌ها و روش‌های متعددی برای مقاوم‌سازی بازشناسی گفتار ارائه شده‌اند. این رویکردها را می‌توان به چهار دسته کلی تقسیم کرد. گروه اول مربوط به روش‌های کاهش نویز یا به عبارتی بهسازی² گفتار است. هدف در این گونه روش‌ها کم کردن نویز سیگنال گفتار به منظور کاهش عدم تطبیق شرایط آموزش و آزمایش قبل از استخراج ویژگی است. در دسته دوم سعی بر آن است که

¹ Stationary

² Speech Enhancement