



دانشگاه صنعتی نوشیروانی بابل

جهت اخذ درجه کارشناسی ارشد
رشته مهندسی برق-گرایش الکترونیک

موضوع:

جداسازی صوت تک کانال با روش Bayesian

استاد راهنما:

دکتر محمدرضا کرمی

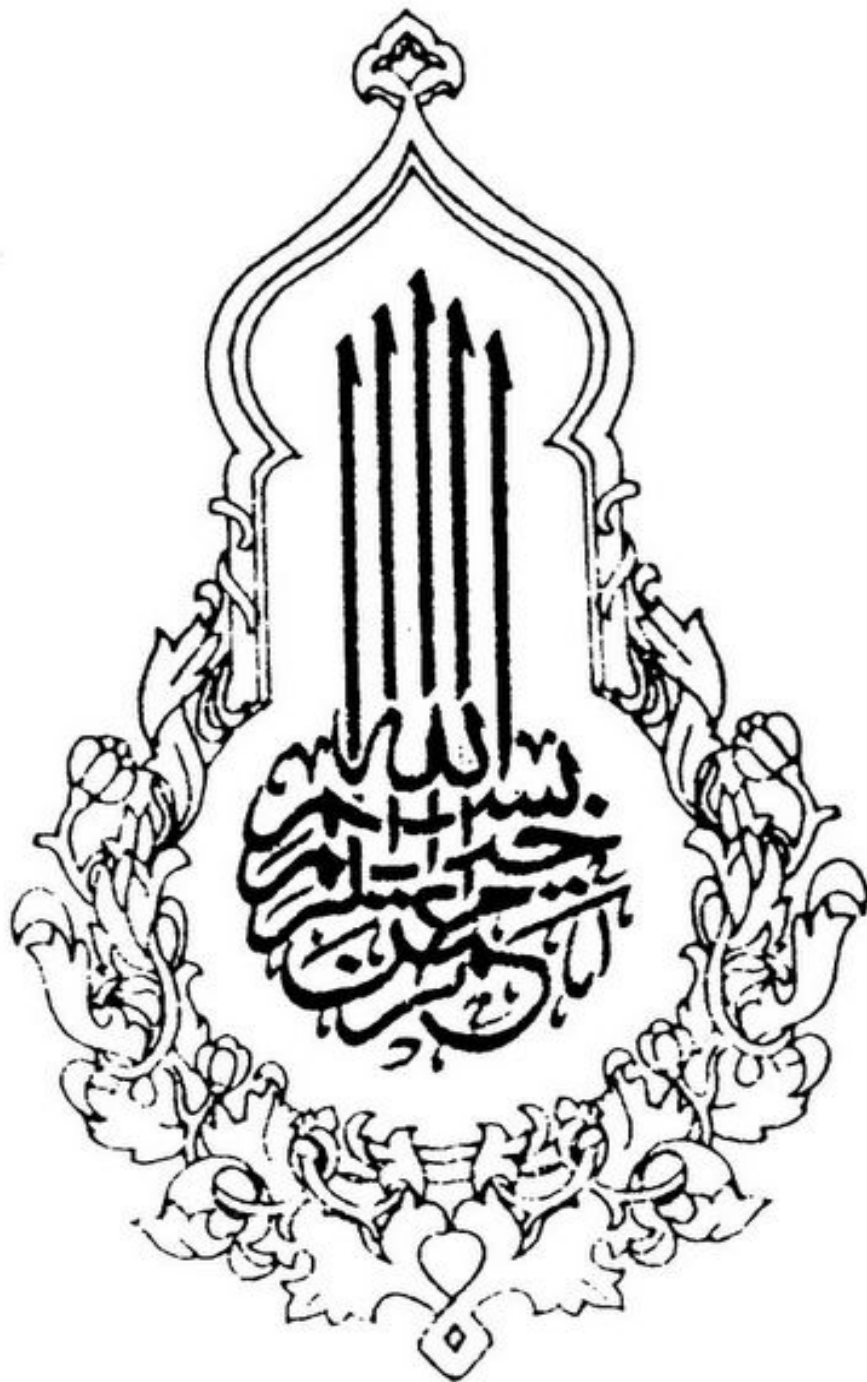
استاد مشاور:

دکتر محمدرضا ذهابی

نگارش:

سونای کمی

تابستان ۹۰





دانشگاه صنعتی نوشیروانی بابل

جهت اخذ درجه کارشناسی ارشد
رشته مهندسی برق-گرایش الکترونیک

موضوع:

جداسازی صوت تک کانال با روش Bayesian

استاد راهنما:

دکتر محمدرضا کرمی

استاد مشاور:

دکتر محمدرضا ذهابی

نگارش:

سونای کمی

تابستان ۹۰

سپاسگزاری

بر خود لازم می دانم از زحمات و راهنمایی های بی دریغ استاد راهنمای گرانقدر، جناب آقای دکتر محمدرضا کرمی که رهنمود های عالمانه ایشان من را در این راه یاری نموده نهایت تشکر و قدردانی را داشته باشم.

همچنین از جناب آقای دکتر محمدرضا ذهابی که مشاوره این پایان نامه را به عهده داشته اند کمال تشکر را دارم.

تقدیم به:

روح پاک پدرم که عالمانه به من آموخت تا چگونه در عرصه زندگی،

ایستادگی را تجربه نمایم

و به مادرم، دریای بی کران فداکاری و عشق که وجودم برایش همه

رنج بود و وجودش برایم همه مهر

چکیده

سیگنال های گفتار به ندرت به صورت خالص برای کاربرد های پردازش گفتار موجود می باشند و اغلب به وسیله ی تداخل آوایی نظیر نویز پس زمینه، اعوجاج، سیگنال گفتار گوینده ی دیگر و ... مخدوش می شوند. در چنین حالتی لازم است که ابتدا سیگنال گفتار از پس زمینه جدا شود. به ویژه عمل جداسازی گفتار چند گوینده که به عنوان جداسازی گفتار شناخته می شود امری چالش برانگیز است زیرا شامل جداسازی سیگنال هایی است که دارای مشخصات آوایی و آماری بسیار مشابهی هستند. چالش برانگیز ترین حالت جداسازی گفتار، جداسازی سیگنال های گفتار ناشی از رکورد تک کاناله است؛ مسئله ای که معمولاً به عنوان جداسازی گفتار تک کانال شناخته می شود. سیستم های جداسازی گفتار تک کانال با کیفیت بالا هنوز از اهمیت زیادی برخوردارند و به عنوان یک مسئله ی حل نشده باقی مانده اند. هدف این پایان نامه، ارائه ی راهکار هایی برای بهبود سیستم های جداسازی گفتار تک کانال می باشد. یک سیستم مؤثر جداسازی گفتار تک کانال، پیش پردازشی مهم در کاربرد های بسیاری نظیر شناسایی اتوماتیک گفتار و تعیین هویت گوینده به شمار می رود. این پایان نامه روی جداسازی مبتنی بر مدل گفتار تک کانال متمرکز میشود که به تکنیک هایی اطلاق می شود که از مدل های آموزش داده شده ی منابع استفاده می کنند تا منابع را از رکورد تک کاناله ی ترکیب خطی آنها جدا کنند. اولین راهکار پیشنهادی در این پایان نامه جداسازی منابع در سطح زیربخش و اعمال این رویکرد به روش مبتنی بر تخمین حداکثر احتمال پسین و روش مبتنی بر کوانتیزاسیون بردار می باشد. در راهکار پیشنهادی دوم، روش مبتنی بر تخمین حداکثر احتمال پسین برای حالتی که در آن سیگنال های تست و آموزشی سطح انرژی مختلفی نسبت به هم دارند (منابع بهره مختلفی دارند) بسط داده می شود. این حالت در بسیاری از روش های مبتنی بر مدل جداسازی گفتار تک کانال در نظر گرفته نشده است. روش هایی که این حالت را در نظر می گیرند به روش های وفق بهره موسوم اند. اگرچه روش های وفق بهره ی موجود در تخمین بهره ی منابع موفق عمل می کنند اما از پیچیدگی محاسباتی و زمان پردازش بالایی برخوردارند. روش پیشنهادی سوم بر مبنای جداسازی زیربخشی و کوانتیزاسیون بردار تلاشی برای غلبه بر این محدودیت است. نتایج تجربی نشان می دهند که راهکار های پیشنهادی، عملکرد سیستم جداسازی گفتار تک کانال را بهبود می دهند.

واژه های کلیدی:

جداسازی گفتار تک کانال، مدل ترکیبی گوسی، کوانتیزاسیون بردار، تخمین حداکثر احتمال پسین، تخمین حداقل میانگین مربعات خطا.

۱	۱- مقدمه	
۲	۱-۱- انگیزه پژوهش	
۳	۲-۱- هدف پایان نامه	
۳	۳-۱- ساختار پایان نامه	
۶	۲- پیش زمینه و مروری بر کارهای گذشته	
۶	۱-۲- مقدمه	
۶	۲-۲- مسئله‌ی جداسازی صوت تک کانال	
۱۴	۳-۲- کارهای گذشته	
۱۵	۱-۳-۲- رویکرد های تحلیل محاسباتی صحنه‌ی شنوایی (CASA)	
۲۰	۲-۳-۲- رویکرد های جداسازی کور منابع (BSS)	
۲۲	۳-۳-۲- رویکرد های مبتنی بر مدل	
۲۸	۴-۲- خلاصه	
۳۰	۳- بررسی چند الگوریتم مبتنی بر مدل	
۳۰	۱-۳- مقدمه	
۳۰	۲-۳- جداسازی صوت تک کانال با روش تخمین حداکثر احتمال پسین (MAP)	
۳۰	۱-۲-۳- مدل کردن منابع و مشاهده	
۳۰	۱-۱-۲-۳- مدل کردن منبع	
۳۲	۲-۱-۲-۳- مدل کردن مشاهده	
۳۴	۲-۲-۳- بدست آوردن تخمین گر حداکثر احتمال پسین	
۳۴	۱-۲-۲-۳- فرموله کردن مسئله‌ی تخمین	
۳۴	۲-۲-۲-۳- مرحله‌ی آشکارسازی	
۳۵	۳-۲-۲-۳- مرحله‌ی تخمین	
۳۶	۳-۲-۳- بازسازی سیگنال های گفتار	
۳۷	۳-۳- جداسازی صوت تک کانال بر اساس کوانتیزاسیون بردار (VQ)	
۳۸	۱-۳-۳- مرحله‌ی استخراج ویژگی	
۳۸	۲-۳-۳- مرحله‌ی آموزش	
۳۹	۳-۳-۳- مرحله‌ی جداسازی	
۳۹	۱-۳-۳-۳- فاز آشکارسازی	
۴۰	۲-۳-۳-۳- فاز تخمین	
۴۱	۴-۳-۳- بازسازی سیگنال های گفتار	
۴۲	۴-۳- روش های وفق بهره	
۴۲	۵-۳- روش وفق بهره با تخمین گر حداکثر شباهت (ML)	
۴۳	۱-۵-۳- رابطه‌ی بین بهره و نسبت هدف به تداخل (TIR)	
۴۴	۲-۵-۳- رابطه‌ی سیگنال مشاهده و منابع	
۴۴	۳-۵-۳- تابع چگالی احتمال سیگنال مشاهده	

۴۵	۳-۵-۴- تخمین بهره و منابع
۴۵	۳-۵-۴-۱- فاز آشکارسازی
۴۶	۳-۵-۴-۲- فاز تخمین
۴۶	۳-۶-۶- روش وفق بهره با تخمین گر حداقل میانگین مربعات خطا (MMSE)
۴۷	۳-۶-۱- تعاریف اولیه
۴۷	۳-۶-۲- مدل کردن منابع و سیگنال مشاهده
۴۷	۳-۶-۳- بدست آوردن تخمین گر حداقل میانگین مربعات خطا
۴۷	۳-۶-۳-۱- رویکرد آشکارسازی-تخمین
۴۹	۳-۶-۳-۲- مرحله‌ی آشکارسازی
۵۱	۳-۶-۳-۳- مرحله‌ی تخمین
۵۳	۳-۷- خلاصه
۵۵	۴- معرفی روش های جدید و نتایج شبیه سازی
۵۵	۴-۱- مقدمه
۵۵	۴-۲- روش جداسازی زیربخشی
۵۶	۴-۲-۱- اعمال روش جداسازی زیربخشی به الگوریتم تخمین حداکثر احتمال پسین
۵۸	۴-۲-۲- اعمال روش جداسازی زیربخشی به الگوریتم مبتنی بر کوانتیزاسیون بردار
۶۰	۴-۲-۳- نتایج تجربی
۶۰	۴-۲-۳-۱- داده های استفاده شده برای آزمایش
۶۱	۴-۲-۳-۲- نحوه ایجاد مدل ترکیبی گوسی یا کتاب کد برای هر گوینده
۶۱	۴-۲-۳-۳- معیار ارزیابی عملکرد سیستم
۶۲	۴-۲-۳-۴- نتایج شبیه سازی
۶۸	۴-۳- روش وفق بهره با تخمین گر حداکثر احتمال پسین
۷۰	۴-۳-۱- نتایج تجربی
۷۱	۴-۴- روش وفق بهره مبتنی بر کوانتیزاسیون بردار و جداسازی زیربخشی
۷۳	۴-۴-۱- مرحله‌ی تخمین بهره
۷۴	۴-۴-۲- مرحله‌ی تخمین منابع
۷۵	۴-۴-۳- نتایج تجربی
۷۹	۴-۵- خلاصه
۸۱	۵- جمع بندی و پیشنهادات ادامه کار
۸۲	۵-۱- پیشنهادات ادامه کار
۸۳	ضمیمه‌ی I
۸۶	ضمیمه‌ی II
۸۸	منابع و ماخذ

۸	شکل ۱-۲- مسئله‌ی جداسازی صوت تک کانال
۹	شکل ۲-۲- مؤلفه های اصلی جداسازی گفتار تک کانال
۱۱	شکل ۳-۲- کاربرد های جداسازی گفتار تک کانال
۱۲	شکل ۴-۲- دقت تشخیص کلمه برای شنوندگان انسانی به صورت تابعی از تعداد صدا های رقابت کننده و شدت تداخل. سطح گفتار هدف در ۹۵ dB ثابت شده است
۱۳	شکل ۵-۲- مقایسه‌ی دقت شناسایی سیستم ASR با استفاده از گفتار تک کانال پردازش نشده (خط مشخص شده با ستاره) و گفتار هدف جدا شده بوسیله‌ی روش ارائه شده در [۹] (خط مشخص شده با دایره)
۱۴	شکل ۶-۲- نرخ تعیین درست سیستم SID با و بدون استخراج گفتار قابل استفاده
۱۶	شکل ۷-۲- تعدادی از معیار های گروه بندی مورد استفاده در رویکرد های CASA
۱۷	شکل ۸-۲- دیاگرام شماتیک یک سیستم CASA معمولی
۲۱	شکل ۹-۲- دیاگرام شماتیک سیستم جداسازی کور منابع برای حل مسئله‌ی جداسازی گفتار تک کانال
۲۳	شکل ۱۰-۲- شماتیک کلی روش های مبتنی بر مدل (۱)
۲۴	شکل ۱۱-۲- شماتیک کلی روش های مبتنی بر مدل (۲)
۲۵	شکل ۱۲-۲- نحوه‌ی استخراج لگاریتم اندازه‌ی تبدیل فوریه
۳۳	شکل ۱-۳- تخمین MIXMAX برای دو فریم مصوت
۳۳	شکل ۲-۳- تخمین MIXMAX برای دو فریم غیر مصوت
۳۷	شکل ۳-۳- بلوک دیاگرام سطح بالای سیستم جداسازی به روش تخمین MAP
۳۸	شکل ۴-۳- بازنمایی سیگنال در حوزه‌ی زمان، فرکانس و لگاریتم فرکانس
۴۱	شکل ۵-۳- بلوک دیاگرام تکنیک جداسازی گفتار براساس کوانتیزاسیون بردار
۴۴	شکل ۶-۳- میانگین $G_0^2(g_x^2 + g_y^2)$ و G_z^2 ده جمله‌ی ترکیبی بر حسب θ (dB)
۴۹	شکل ۷-۳- $p(s_x = i, s_y = j Z^r, \theta)$ بر حسب تمام زوج های (i, j) برای نه فریم که به طور اتفاقی انتخاب شده اند
۵۱	شکل ۸-۳- $Q(\theta)$ بر حسب θ برای چهار سیگنال ترکیبی
۵۵	شکل ۱-۴- شماتیک روش جداسازی زیربخشی
۵۶	شکل ۲-۴- روش مبتنی بر تخمین MAP با جداسازی زیربخشی
۵۸	شکل ۳-۴- روش مبتنی بر کوانتیزاسیون بردار با جداسازی زیربخشی
۶۲	شکل ۴-۴- (a) سیگنال اولیه‌ی گوینده‌ی ۱، (b) سیگنال اولیه‌ی گوینده‌ی ۲، (c) سیگنال ترکیبی، d و (e) به ترتیب تخمین سیگنال گوینده‌ی ۱ با روش تخمین MAP با جداسازی زیربخشی و بدون جداسازی زیربخشی، f و g) به ترتیب تخمین سیگنال گوینده‌ی ۲ با روش تخمین MAP با جداسازی زیربخشی و بدون جداسازی زیربخشی
۶۴	شکل ۵-۴- میانگین SNR حاصل از تخمین پنج سیگنال هدف از روی سیگنال های ترکیبی با استفاده از روش تخمین MAP با جداسازی زیربخشی (خط لوزی) و بدون جداسازی زیربخشی (خط مربع)

- شکل ۴-۶- میانگین SNR حاصل از تخمین پنج سیگنال تداخل از روی سیگنال های ترکیبی با استفاده از روش تخمین MAP با جداسازی زیربخشی (خط لوزی) و بدون جداسازی زیربخشی (خط مربع) ۶۴
- شکل ۴-۷- (a) سیگنال اولیه ی گوینده ی ۱، (b) سیگنال اولیه ی گوینده ی ۲، (c) سیگنال ترکیبی، d و e) به ترتیب تخمین سیگنال گوینده ی ۱ با روش مبتنی بر کوانتیزاسیون بردار با جداسازی زیربخشی و بدون جداسازی زیربخشی، f و g) به ترتیب تخمین سیگنال گوینده ی ۲ با روش مبتنی بر کوانتیزاسیون بردار با جداسازی زیربخشی و بدون جداسازی زیربخشی ۶۵
- شکل ۴-۸- میانگین SNR حاصل از تخمین پنج سیگنال هدف از روی سیگنال های ترکیبی با استفاده از روش مبتنی بر کوانتیزاسیون بردار با جداسازی زیربخشی (خط لوزی) و بدون جداسازی زیربخشی (خط مربع) ۶۷
- شکل ۴-۹- میانگین SNR حاصل از تخمین پنج سیگنال تداخل از روی سیگنال های ترکیبی با استفاده از روش مبتنی بر کوانتیزاسیون بردار با جداسازی زیربخشی (خط لوزی) و بدون جداسازی زیربخشی (خط مربع) ۶۷
- شکل ۴-۱۰- میانگین SNR حاصل از تخمین ده سیگنال هدف از روی سیگنال های ترکیبی با استفاده از روش تخمین وقتی MAP (خط لوزی) و روش تخمین وقتی ML (خط مربع) ۷۰
- شکل ۴-۱۱- میانگین SNR حاصل از تخمین ده سیگنال تداخل از روی سیگنال های ترکیبی با استفاده از روش تخمین وقتی MAP (خط لوزی) و روش تخمین وقتی ML (خط مربع) ۷۱
- شکل ۴-۱۲- شماتیک روش وفق بهره مبتنی بر کوانتیزاسیون بردار و جداسازی زیربخشی ۷۲
- شکل ۴-۱۳- $Q(\theta)$ بر حسب θ برای شش سیگنال ترکیبی ۷۶
- شکل ۴-۱۴- میانگین SNR حاصل از تخمین ده سیگنال هدف از روی سیگنال های ترکیبی با استفاده از روش وفق بهره مبتنی بر کوانتیزاسیون بردار و جداسازی زیربخشی (خط لوزی)، روش تخمین وقتی MMSE (معادل تخمین وقتی MAP) (خط مربع)، روش مبتنی بر کوانتیزاسیون بردار (خط مثلث) و روش تخمین MAP (خط دایره) ۷۸
- شکل ۴-۱۵- میانگین SNR حاصل از تخمین ده سیگنال تداخل از روی سیگنال های ترکیبی با استفاده از روش وفق بهره مبتنی بر کوانتیزاسیون بردار و جداسازی زیربخشی (خط لوزی)، روش تخمین وقتی MMSE (معادل تخمین وقتی MAP) (خط مربع)، روش مبتنی بر کوانتیزاسیون بردار (خط مثلث) و روش تخمین MAP (خط دایره) ۷۸
- شکل II-۱- فرم معمولی $Q(\theta)$ با سه نقطه متناظر با $A = Q(\theta_i)$ ، $C = Q(\theta_c)$ و $B = Q(\theta_r)$ ۸۶

۵۲	جدول ۳-۱- مقایسه‌ی θ واقعی (ستون دوم از سمت چپ) و θ^* بدست آمده از معادله‌ی (۳-۵۳) با استفاده از الگوریتم بهینه سازی درجه دومی برای بیست سیگنال ترکیبی
۶۳	جدول ۴-۱- SNR بدست آمده از سیگنال های تخمین زده شده با روش تخمین MAP با جداسازی زیربخشی و بدون جداسازی زیربخشی برای ده سیگنال ترکیبی
۶۶	جدول ۴-۲- SNR بدست آمده از سیگنال های تخمین زده شده با روش مبتنی بر کوانتیزاسیون بردار با جداسازی زیربخشی و بدون جداسازی زیربخشی برای ده سیگنال ترکیبی
۷۷	جدول ۴-۳- مقایسه‌ی θ واقعی (ستون دوم از سمت چپ) و θ^* بدست آمده از معادله‌ی (۴-۲۶) با استفاده از الگوریتم بهینه سازی درجه دومی برای بیست سیگنال ترکیبی
۸۷	جدول II-۱- الگوریتم بهینه سازی درجه دومی

لیست علائم و اختصارات

ACF	تابع خود همبستگی (Autocorrelation Function)
AM	مدولاسیون دامنه (Amplitude Modulation)
ASA	تحلیل صحنه شنوایی (Auditory Scene Analysis)
ASR	شناسایی اتوماتیک گفتار (Automatic Speech Recognition)
BSS	جداسازی کور منابع (Blind Source Separation)
CASA	تحلیل محاسباتی صحنه شنوایی (Computational Auditory Scene Analysis)
CSM	مدل منبع مرکب (Composite Source Modeling)
EM	بیشینه سازی امید (Expectation Maximization)
GMM	مدل ترکیبی گوسی (Gaussian Mixture Model)
HMM	مدل مخفی مارکوف (Hidden Markov Model)
ICA	تحلیل مؤلفه های مستقل (Independent Component Analysis)
IDFT	تبدیل فوریه معکوس گسسته (Inverse Discrete Fourier Transform)
MAP	حداکثر احتمال پسین (Maximum A Posteriori)
ML	حداکثر شباهت (Maximum Likelihood)
MMSE	حداقل میانگین مربعات خطا (Minimum Mean Squared Error)
MSE	میانگین مربعات خطا (Mean Squared Error)
OLA	افزودن همپوشانی (Overlap Add)
PDF	تابع چگالی احتمال (Probability Density Function)
SACF	تابع خود همبستگی اختصاری (Summary Autocorrelation Function)
SDR	نسبت سیگنال به اعوجاج (Signal to Distortion Ratio)
SE	بهبود گفتار (Speech Enhancement)
SID	تعیین هویت گوینده (Speaker Identification)
SNR	نسبت سیگنال به نویز (Signal to Noise Ratio)
SPL	سطح شدت صدا (Sound Pressure Level)
SRR	نسبت سیگنال به مانده (Signal to Residual Ratio)
TIR	نسبت هدف به تداخل (Target to Interference Ratio)
VQ	کوانتیزاسیون بردار (Vector Quantization)

فصل اول

مقدمه

۱- مقدمه

۱-۱- انگیزه پژوهش

در بسیاری از کاربرد های پردازش گفتار نظیر شناسایی اتوماتیک گفتار^۱ (ASR)، تعیین هویت گوینده^۲ (SID) و بهبود کیفیت سیگنال گفتار^۳ (SE)، سیگنال ورودی اغلب به نویز آوایی محیطی آغشته می شود که این عمل سبب کاهش کیفیت سیگنال گفتار و در نتیجه باعث کاهش کارایی کلی الگوریتم پردازش گفتار می شود. هنگامی که تداخل آوایی، شامل سیگنال های صحبت گوینده های دیگر باشد (معمولاً به عنوان اثر مهمانی شلوغ^۴ شناخته می شود)، آنگاه به علت مشابهت در ماهیت سیگنال های خواسته و ناخواسته، این کاهش کارایی بیشتر می شود. بنابراین در چنین شرایطی یک الگوریتم جداسازی صوت، پیش پردازشی ضروری به شمار می رود تا کیفیت سیگنال گفتار را برای پردازش بعدی بهبود ببخشد. اگر ما بتوانیم سیگنال گفتار مورد نظر را قبل از پردازش آن از سیگنال اولیه جدا کنیم، این عمل در بهبود کارایی کلی الگوریتم پردازش گفتار کمک شایانی خواهد کرد.

در بعضی از شرایط عملی که تنها یک کانال موجود است تکنیک های جداسازی تک کانال باید مورد استفاده قرار گیرد. این حالت ممکن است به وسیله سیستم مورد استفاده (مانند کاربردهای مبتنی بر تلفن) یا به علت فراهم بودن سیگنال مورد نظر (مانند کاربردهای از پیش ضبط شده) اعمال شود. این شرایط به علت سهولت در نصب میکروفن جالب هستند اما محدودیت اصلی روش های تک کانال این است که هیچ سیگنال مرجعی برای تداخل^۵ موجود نیست. بنابراین چگالی طیف قدرت^۶ سیگنال تداخلی باید تنها بر اساس سیگنال گفتار تک کانال موجود تخمین زده شود و این چیزی است که آن را به امری چالش برانگیز تبدیل می کند. این مسئله معمولاً به مسئله جداسازی صوت تک کانال^۷ اطلاق می شود.

ایدهی جداسازی صوت تک کانال، پردازش اتوماتیک سیگنال ترکیبی به منظور بازیابی سیگنال گفتار اولیهی هر گوینده است. حداقل کردن آرتیفکت در سیگنال گفتار پردازش شده امری کلیدی است به ویژه اگر هدف نهایی استفاده از گفتار بازیابی شده در کاربردهای مبتنی بر ماشین نظیر سیستم های شناسایی اتوماتیک گفتار و تشخیص گوینده باشد. بنابراین هدف الگوریتم جداسازی صوت تک کانال عبارت است از:

¹ Automatic speech recognition

² Speaker identification

³ Speech enhancement

⁴ Cocktail party effect

⁵ Interference

⁶ Power spectral density

⁷ Monaural speech separation problem

* بهبود جنبه های ادراکی سیگنال گفتار مخدوش شده

* بهبود عملکرد سیستم نهایی پردازش گفتار

* افزایش توانمندی سیستم های پردازش گفتار مبتنی بر ماشین

اگرچه کاملا واضح است که سیستم شنوایی انسان در تمرکز روی یک یا چند گوینده‌ی خاص در میان جمعی از گویندگان که همزمان در حال صحبت هستند بسیار تواناست [۱]، الگوریتم های کامپیوتری که برای انجام همان کار طراحی شده اند درجه‌ی محدودی از موفقیت را نشان داده اند.

۱-۲- هدف پایان نامه

هدف این پایان نامه ارائه‌ی راهکار هایی برای بهبود جداسازی گفتار تک کانال می باشد. در میان رویکرد های موجود برای حل مسئله‌ی جداسازی گفتار تک کانال، رویکرد های مبتنی بر مدل^۱ توانایی بیشتری از خود نشان داده اند به طوریکه در چند سال گذشته مورد توجه محققین قرار گرفته اند. در این پایان نامه ما بر روی رویکرد های مبتنی بر مدل متمرکز می شویم. راهکار های ارائه شده در این پایان نامه به صورت زیر می باشد

۱- بهبود جداسازی گفتار تک کانال با استفاده از جداسازی زیربخشی

۲- وفقی کردن بهره^۲ در روش تخمین حداکثر احتمال پسین

۳- ارائه‌ی روشی مؤثر و سریع بر مبنای جداسازی زیربخشی و کوانتیزاسیون بردار^۳ برای جبران تفاوت بهره^۴ سیگنال های ورودی

۱-۳- ساختار پایان نامه

ادامه‌ی این پایان نامه به صورت زیر سازمان دهی شده است

فصل دوم مسئله‌ی جداسازی گفتار تک کانال را معرفی می کند و به مرور سریع رویکرد های مختلفی می پردازد که برای حل این مسئله ارائه شده اند.

¹ Model based approaches

² Gain adaptation

³ Vector quantization

⁴ Gain

فصل سوم به تفصیل به بررسی برخی از روش های مبتنی بر مدل که برای حل مسئله‌ی جداسازی گفتار تک کانال ارائه شده اند می پردازد. به کمک این روش ها با شیوه های مختلف مدل سازی¹ و تخمین آشنا می شویم. در نهایت، روش های پیشنهادی با روش های ارائه شده در این فصل مقایسه خواهند شد.

در فصل چهارم راهکار هایی برای بهبود جداسازی گفتار تک کانال ارائه می شود. این راهکارها شامل جداسازی زیربخشی، وفقی کردن بهره در روش تخمین حداکثر احتمال پسین و ارائه‌ی یک روش وفق بهره بر مبنای جداسازی زیربخشی می باشد.

¹ Modeling

فصل دوم

پیش زمینه و مروری بر کارهای گذشته

۲- پیش زمینه و مروری بر کارهای گذشته

۲-۱- مقدمه

صوت تک کانال هنگامی اتفاق می افتد که سیگنال های گفتار مربوط به چند گوینده از طریق یک کانال با هم ترکیب می شوند. فرایند استخراج یا جداسازی سیگنال گفتار مورد نظر از سیگنال ترکیبی معمولاً به عنوان جداسازی صوت تک کانال شناخته می شود. استفاده از یک الگوریتم جداسازی صوت تک کانال به عنوان پیش پردازش در سیستم های پردازش گفتار می تواند در بهبود سیگنال گفتار هدف^۱ در بسیاری از کاربردها تاثیر زیادی داشته باشد. این کاربردها شامل تکنیک های شناسایی اتوماتیک گفتار (ASR)، تعیین هویت گوینده (SID) و بهبود کیفیت گفتار (SE) هستند. به طور قابل ملاحظه ای، سیستم شنوایی انسان می تواند مسئله جداسازی صوت تک کانال را با دقت بیشتری نسبت به سیستم های جداسازی مبتنی بر ماشین فعلی حل کند. از لحاظ تاریخی، سیستم های جداسازی صوت تک کانال با استفاده از الگوریتم هایی توسعه یافته اند که سیگنال های تداخلی را حذف می کنند، سیگنال هدف را بهبود می دهند، یا هر دو سیگنال را به طور همزمان تخمین می زنند. تاکنون، تکنیکی که به طور رضایت بخش در تمام حالات کار کند وجود نداشته است. در این فصل دورنمای مختصری از مسئله جداسازی صوت تک کانال و مروری بر کارهای انجام شده برای حل این مسئله ارائه می شود به این ترتیب که بخش ۲-۲ به توضیح درباره ای این مسئله و کاربردهای اصلی آن می پردازد و مروری بر رویکردهای گذشته برای حل این مسئله در بخش ۲-۳ ارائه می شود.

۲-۲- مسئله جداسازی صوت تک کانال

صوت تک کانال به صورت سیگنال ترکیبی حاصل از دو یا چند گوینده تعریف می شود. این پدیده معمولاً به علت ترکیب سیگنال های گفتار منابع مستقل و همزمان به صورت یک سیگنال منفرد در گیرنده (میکروفن منفرد) اتفاق می افتد. شرایط متداولی که صوت تک کانال ممکن است اتفاق بیافتد می توان به صورت زیر خلاصه کرد:

۱- هنگام ضبط صدای دو نفر که به طور همزمان نسبت به یک میکروفن در حال صحبت هستند. یک مثال در این مورد، سیگنال های گفتار ضبط شده در جعبه ی ضبط صدا در کابین خلبان می باشد.

۲- هنگامی که سیگنال گفتار هدف همزمان با صداهای پس زمینه شنیده شدنی گوینده های دیگر به سمعک افراد کم شنوا می رسد.

^۱ Target speech signal

۳- هنگامی که به علت انتقال سیگنال گفتار از طریق یک کانال ارتباطی معیوب، تداخل صحبت^۱ اتفاق می افتد.

سیستم شنوایی انسان در تمرکز روی یک یا چند گوینده‌ی خاص در جمعی از گویندگان (معمولا به عنوان اثر مهمانی شلوغ شناخته می شود) توانایی های قابل ملاحظه ای از خود نشان می دهد. به طرز شگفت آوری، شنوندگان انسانی حتی هنگامی که فقط از یک گوش استفاده می کنند یا با دو گوش به صدای ضبط شده از طریق یک کانال منفرد^۲ گوش می دهند قادر به انجام این کار هستند. از طرف دیگر، الگوریتم های کامپیوتری طراحی شده برای انجام این کار، درجه‌ی محدودی از موفقیت را از خود نشان داده اند. آزمایش های های میلر [۲] و بروکس و نوتبوم [۳] نشان داده اند که شنوندگان انسانی می توانند به درجه‌ی بالایی از جداسازی صداهای ترکیب شده از طریق یک کانال برسند. در حقیقت، انسان شامل فاکتورهای زیادی برای جداسازی صداهای هم زمان است که ممکن است برای ماشین موجود نباشد. بعضی از این فاکتورها همانطور که به وسیله ی چری [۴] بیان شده اند عبارتند از:

۱- اطلاعات فضایی درباره‌ی جهت منبع صوتی

۲- اطلاعات بصری نظیر اشارات در هنگام سخن گفتن و حرکت لب

۳- اطلاعات درباره‌ی آهنگ صدا (بر اساس سرعت گفتار، جنس گوینده و...)

۴- اطلاعات درباره لهجه ها

۵- اطلاعات درباره‌ی احتمالات گذرا^۳ (بر اساس درک گفتار، علم نحو و ...)

به جز آخرین فاکتور، فاکتورهای دیگر می توانند در حالتی از صوت تک کانال که، به عنوان مثال، دو پیام از یک گوینده به طور همزمان روی یک نوار ضبط می شوند حذف شوند. با این وجود، گوش انسان هنوز قادر به جداسازی گفتار در چنین حالت حدی است و علت آن حافظه‌ی وسیع درباره‌ی احتمالات گذرا است که به انسان اجازه می دهد تا دنباله‌ی کلمات را پیشگویی کند.

هدف جداسازی صوت تک کانال بازیابی یک یا هر دو سیگنال گفتار از سیگنال ترکیبی است. مدل دو گوینده ای مسئله‌ی جداسازی صوت تک کانال که مدل مورد بررسی در این پایان نامه است در شکل (۱-۲)

¹ Cross talk

² Single channel

³ Transition probabilities

نشان داده شده است. همانطور که از شکل پیداست، سیگنال گفتار تک کانال ($z(t)$) مجموع سیگنال گفتار هدف ($x(t)$) و سیگنال گفتار تداخل ($y(t)$) است.

$$z(t) = x(t) + y(t) \quad , \quad t = 1, 2, \dots, T \quad (1-2)$$

که t نشان دهنده‌ی نمونه در حوزه‌ی زمان است. نسبت هدف به تداخل^۱ ورودی (TIR) به صورت نسبت توان $x(t)$ به $y(t)$ تعریف می‌شود.

$$\text{TIR [dB]} = 10 \log_{10} \frac{\sum_t x^2(t)}{\sum_t y^2(t)} \quad (2-2)$$

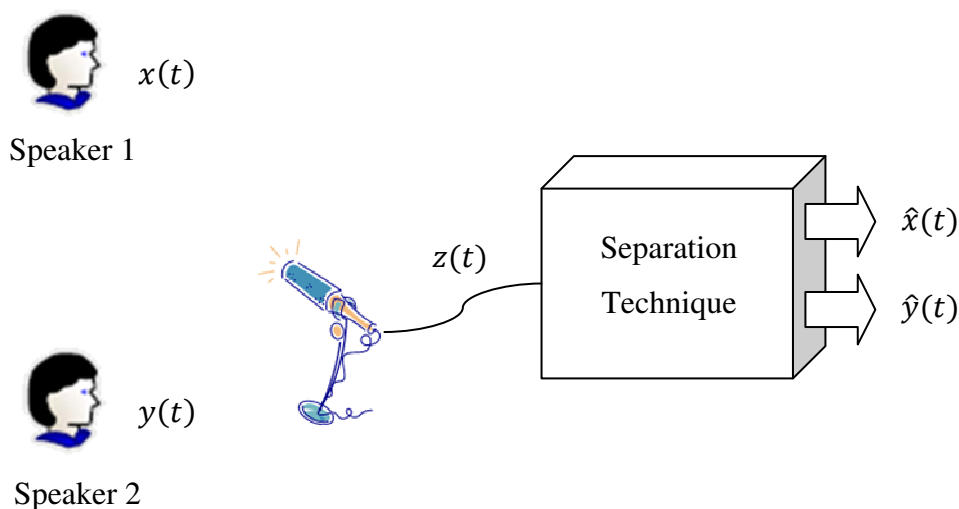
در زمینه‌ی پردازش سیگنال می‌توان به جداسازی صوت تک کانال از چهار منظر مختلف نگریست

۱- بهبود سیگنال گفتار هدف

۲- استخراج سیگنال گفتار هدف

۳- حذف سیگنال گفتار تداخلی

۴- تخمین هر دو سیگنال گفتار



شکل ۱-۲- مسئله‌ی جداسازی صوت تک کانال

¹ Target to interference ratio