





دانشگاه صنعتی امیرکبیر

دانشکده مهندسی برق

پایان نامه کارشناسی ارشد مهندسی برق گرایش الکترونیک

استفاده از تکنیک خوشه بندی گوینده در تطبیق گوینده در سیستم بازشناسی گفتار

نگارش:

اولدوز حضرتی یادکوری

استاد راهنما:

دکتر سید محمد احدی

بهمن ماه ۸۶

تاریخ :
شماره :



فرم اطلاعات پایان نامه
کارشناسی ارشد و دکترا

دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

معاونت پژوهشی
فرم پروژه تحصیلات تکمیلی ۷

مشخصات دانشجو

نام و نام خانوادگی : اولدوز حضرتی یادکوری دانشجوی آزاد بورسیه معادل

شماره دانشجویی: ۸۴۱۳۰۷۳ دانشکده : مهندسی برق رشته تحصیلی: الکترونیک

نام و نام خانوادگی استاد راهنما : سید محمد احدی

عنوان به فارسی: استفاده از تکنیک خوشه بندی گوینده در تطبیق گوینده در سیستم بازشناسی گفتار

عنوان به انگلیسی: Speaker Adaptation Using Speaker Clustering in Speech Recognition Systems

نوع پروژه: کارشناسی ارشد کاربردی بنیادی توسعه ای نظری

تاریخ شروع: ۸۵/۸ تاریخ خاتمه : ۸۶/۱۱ تعداد واحد: ۶

سازمان تامین کننده اعتبار :

واژه های کلیدی به فارسی : تطبیق گوینده، خوشه بندی گوینده، تابع کرنل، ماشین های بردار پشتیبان تک کلاسی، K-means ملایم

واژه های کلیدی به انگلیسی : Speaker adaptation, speaker clustering, kernel method, one class support vector machine (OCSVM), soft K-means

نظرها و پیشنهادهای به منظور بهبود فعالیت های پژوهشی دانشگاه:

استاد راهنما: دکتر سید محمد احدی

دانشجو: اولدوز حضرتی یادکوری

امضاء استاد راهنما : تاریخ:

نسخه ۱: معاونت پژوهشی

نسخه ۲: کتابخانه و به انضمام دو جلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز اسناد و مدارک علمی

تقدیم به مادر فداکار و پدر اندیشمندم.

با تشکر فراوان از زحمات استاد گرامی، جناب آقای دکتر احدی که مرا در پیشبرد و تکمیل این پایان‌نامه یاری رساندند. همچنین از آقای مهندس سجادی که مرا در انجام این پروژه مساعدت نمودند، کمال تشکر را دارم.

چکیده

بطور کلی در بازشناسی گفتار، مدل وابسته به گوینده (SD) عملکرد بهتری نسبت به مدل ناوابسته به گوینده (SI) در تشخیص گفتار یک گوینده خاص دارد. یکی از روش های عملی تر برای دستیابی به عملکردی نزدیک به سیستم SD استفاده از تکنیک های تطبیق گوینده است.

خوشه بندی گوینده یکی از تکنیک های اصلی در تطبیق گوینده است. روش خوشه بندی می تواند به دلیل راحتی ترکیب با تکنیک های رایج تطبیق نظیر MAP و MLLR مورد استفاده قرار بگیرد.

در این پروژه، روشی مبتنی بر کرنل تنها در فاز خوشه بندی مورد استفاده قرار می گیرد. پیاده سازی این روش ساده و حجم محاسباتی مطلوبی دارد. ما این روش خوشه بندی بر مبنای توابع کرنل را که الهام گرفته از روش متداول K-means و بر مبنای ماشین های بردار پشتیبان تک کلاسی (OCSVM) است به عنوان مرحله پیش تطبیق قبل از تکنیک های رایج تطبیق نظیر MAP و MLLR در تطبیق گوینده با نظارت سریع مورد استفاده قرار می دهیم.

در بخشی از کار الگوریتم های خوشه بندی مختلفی نظیر خوشه بندی های جنسیت، خوشه بندی K-means و Soft K-means و روشی بر مبنای بردار های پشتیبان تک کلاسی جهت تطبیق مورد استفاده قرار گرفته اند و با روش پیشنهادی به دقتی در حدود ۳٪ بهتر از مدل مبنا با دقت ۹۰/۳۷٪ (دقت مدل SI) دست یافتیم.

در بخش دیگری از این پروژه استفاده از روش های خوشه بندی ذکر شده به عنوان یک مرحله پیشین در تطبیق به روش های MAP و MLLR مورد استفاده قرار گرفته است. در این حالت نیز به ۶/۵٪ بهبودی نسبت به مدل مبنا رسیدیم.

در ادامه به جهت بررسی کارایی این روش در امر تطبیق، از چندین روش مختلف برای انتخاب HMM ها استفاده شده است.

در بخش دیگری از کار به منظور مقایسه کارایی HMM و GMM در روش استفاده شده جهت امر خوشه بندی در این پروژه، به جای HMM های بکار رفته در مرحله پیش از خوشه بندی، از GMM استفاده شده است و نتایج بدست آمده مورد بررسی قرار گرفته است.

تمامی آزمایش های انجام شده در این پروژه بر روی دادگان TIDIGITS صورت گرفته و هیچگونه همپوشانی ای بین گوینده های آموزش و تست وجود ندارد.

در انتها نتایج حاصل از تطبیق گوینده به روش eigenvoice و eigenvoice مقاوم آورده شده است که متأسفانه نسبت به مدل مبنا عملکرد ضعیفتری دارند.

کلمات کلیدی: تطبیق گوینده، خوشه بندی گوینده، تابع کرنل، ماشین های بردار پشتیبان تک کلاسی، K-means ملایم.

۲	فصل اول: مقدمه
۲	۱-۱- کلیات
۴	۲-۱- طرح کلی پایان نامه
۷	فصل دوم : کلیاتی بر بازشناسی گفتار
۷	۱-۲- مقدمه
۷	۲-۲- ویژگی های گفتار
۸	۳-۲- بازشناسی گفتار
۱۰	۴-۲- مدل های مارکوف پنهان
۱۰	۲-۴-۱- پارامترهای مدل مارکوف پنهان و تخمین آنها
۱۳	۲-۵-۲- مدل های مخلوط گوسی
۱۵	۲-۵-۱- تخمین پارامترها و آموزش مدل
۱۷	۲-۶- استخراج ویژگی
۲۱	فصل سوم: تطبیق گوینده و روش های موجود
۲۱	۱-۳- مقدمه
۲۲	۳-۲- تطبیق گوینده سریع به روش EV
۲۴	۳-۲-۱- آنالیز مولفه های اساسی (PCA)
۲۶	۳-۲-۲- فضاهای پارامترهای EV
۲۶	۳-۲-۳- تطبیق گوینده به روش EV
۲۹	۳-۲-۴- EV مقاوم
۳۱	۳-۳- روش MAP
۳۵	۳-۴- روش MLLR
۴۱	فصل چهارم: خوشه بندی گوینده و الگوریتم های مورد استفاده
۴۱	۴-۱- مقدمه
۴۱	۴-۲- خوشه بندی گوینده
۴۲	۴-۳- روش های خوشه بندی
۴۲	۴-۳-۱- خوشه بندی جنسیت
۴۲	۴-۳-۲- خوشه بندی K-means
۴۳	۴-۳-۳- الگوریتم تقسیم باینری
۴۴	۴-۳-۳- خوشه بندی Soft K-means
۴۷	۴-۳-۵- خوشه بندی C-means ملایم
۴۸	۴-۳-۶- خوشه بندی گوینده بر مبنای توابع کرنل

فصل پنجم: پیاده سازی و بررسی نتایج ۵۵

- ۵۵-۱-۵ مقدمه
- ۵۵-۲-۵ پایگاه داده
- ۵۶-۳-۵ آموزش مدل ها و تقسیم دادگان
- ۵۶-۱-۳-۵ استخراج ویژگی
- ۵۶-۲-۳-۵ مدل های آکوستیکی
- ۵۷-۳-۳-۵ تقسیم بندی دادگان
- ۵۸-۴-۵ آزمایشات
- ۵۸-۱-۴-۵ روش های خوشه بندی
- ۵۸-۲-۴-۵ آزمایش اول
- ۶۱-۳-۴-۵ آزمایش دوم
- ۶۲-۴-۴-۵ آزمایش سوم
- ۶۳-۵-۴-۵ آزمایش چهارم
- ۶۴-۶-۴-۵ آزمایش پنجم
- ۶۵-۷-۴-۵ آزمایش ششم
- ۶۶-۸-۴-۵ روش EV
- ۶۸-۹-۴-۵ روش EV مقاوم

فصل ششم: جمع بندی، نتیجه گیری و پیشنهادات ۷۱

- ۷۲-۱-۶ جمع بندی و نتیجه گیری
- ۷۴-۲-۶ پیشنهادات

مراجع ۷۵

فهرست اشکال

شکل ۱-۲- الف) فاز آموزش ب) فاز بازشناسی ج) فاز پیش پردازش در یک سیستم بازشناسی گفتار [۱].	۹
شکل ۲-۲- یک مدل مارکوف پنهان پنج حالت	۱۱
شکل ۳-۲- مدل مخلوط گوسی	۱۴
شکل ۴-۲- بلوک استخراج ویژگی	۱۹
شکل ۱-۳- تطبیق به روش EV [13]	۲۳
شکل ۲-۳- بلوک دیاگرام تطبیق گوینده به روش EV	۲۴
شکل ۴-۳- درخت رگرسیون MLLR [13]	۳۷
شکل ۵-۳- شکل کلی EMLLR [13]	۳۹
شکل ۱-۴- الگوریتم K-means در حالت داشتن دو خوشه غیر مشابه، الف) داده توزیع شده، ب) مجموعه ای از تخصیص ها و میانگین ها. چهار داده از خوشه وسیعتر اشتباهاً به خوشه کوچکتر نسبت داده شده اند [35].	۴۵
شکل ۲-۴- الف) دو خوشه باریک، ب) حل پایدار بوسیله الگوریتم K-means	۴۶
شکل ۳-۴- نگاشت ویژگی طبقه بندی را ساده تر می سازد.	۵۰
شکل ۱-۵- استفاده از مدل های خوشه در بازشناسی گوینده جدید [24].	۵۹
شکل ۲-۵- الگوریتم آزمایش اول	۵۹

فهرست جداول

- جدول ۵-۱- دقت بازشناسی برای الگوریتم های خوشه بندی مختلف و تست با مدل های خوشه. ۶۰
- جدول ۵-۲- دقت بازشناسی برای الگوریتم های خوشه بندی مختلف و استفاده از مدل های وابسته به خوشه به عنوان مدل های اولیه در تطبیق MAP و MLLR [49]، [50]. ۶۱
- جدول ۵-۳- دقت بازشناسی برای الگوریتم های خوشه بندی مختلف و برداشتن مدل هر کلمه از خوشه ای که بیشترین درستنمایی را دارد. ۶۲
- جدول ۵-۴- دقت بازشناسی در حالت اعمال الگوریتم های خوشه بندی مختلف بر گوینده های آموزشی به تعداد مدل های موجود و برداشتن مدل هر کلمه از خوشه ای که بیشترین درستنمایی را دارد. ۶۳
- جدول ۵-۵- دقت بازشناسی بدست آمده با افزایش داده تطبیق برای روش های MAP و MLLR و استفاده از مدل های SI به عنوان مدل های اولیه (بدون خوشه بندی). ۶۴
- جدول ۵-۶- دقت بازشناسی برای الگوریتم های خوشه بندی مختلف اعمال شده بر GMM های گوینده های آموزشی و تست با مدل های خوشه در حالتیکه فاصله ابربردارهای گوینده ها محاسبه می شود. ۶۵
- جدول ۵-۷- دقت بازشناسی روش eigenvoice با دو عدد eigenvoice و ۲ تکرار و تغییر همزمان همه مدل ها. ۶۷
- جدول ۵-۸- دقت بازشناسی روش eigenvoice با دو عدد eigenvoice و ۲ تکرار و تغییر مدل داده تطبیق. ۶۸
- جدول ۵-۹- دقت بازشناسی روش eigenvoice مقاوم با دو عدد eigenvoice و ۲ تکرار و تغییر همزمان همه مدل ها. ۶۹
- جدول ۵-۱۰- دقت بازشناسی روش eigenvoice مقاوم با دو عدد eigenvoice و ۲ تکرار و تغییر مدل داده تطبیق. ۷۰

فهرست علائم اختصاری

<i>DFT</i>	<i>Discrete Fourier Transform</i>
<i>DTW</i>	<i>Dynamic Time Warping</i>
<i>eKEV</i>	<i>Embedded Kernel Eigenvoice</i>
<i>EMAP</i>	<i>Extended MAP</i>
<i>EV</i>	<i>EigenVoice</i>
<i>FFT</i>	<i>Fast Fourier Transform</i>
<i>GMM</i>	<i>Gaussian Mixture Model</i>
<i>HMM</i>	<i>Hidden Markov Models</i>
<i>KMLLR</i>	<i>Kernel Maximum Likelihood Linear Regression</i>
<i>MAP</i>	<i>Maximum a Posteriori</i>
<i>MFCC</i>	<i>Mel-Frequency Cepstral Coefficients</i>
<i>ML</i>	<i>Maximum Likelihood</i>
<i>MLE</i>	<i>Maximum Likelihood Estimation</i>
<i>MLLR</i>	<i>Maximum Likelihood Linear Regression</i>
<i>OCSVM</i>	<i>One Class Support Vector Machine</i>
<i>SD</i>	<i>Speaker Dependent</i>
<i>SI</i>	<i>Speaker Independent</i>
<i>SVM</i>	<i>Support Vector Machine</i>
<i>VQ</i>	<i>Vector Quantization</i>

فهرست نمادها

T	تعداد بردارهای مشاهده
n	بعد بردارهای مشاهده
o	بردارهای مشاهده
s_1, \dots, s_N	حالت های مدل مارکوف پنهان
M	تعداد عناصر مخلوط
c_{jk}	وزن عنصر مخلوط
μ	بردار میانگین
Σ	ماتریس کواریانس
σ	واریانس
λ	مدل گوینده
e	Eigenvoice
z_j	مرکز خوشه
A	مرکز کره

فصل اول

مقدمه

فصل اول

مقدمه

۱-۱- کلیات

همانگونه که استفاده از گفتار و کلام یکی از راه های ابتدایی برقراری ارتباط بین انسانهاست، بکارگیری آن در اداره امور مختلف به جای دست و پا می تواند تا حد زیادی باعث سهولت گردد. بکارگیری این امر در تسهیل انجام کارهای کنترلی و استفاده از ابزارهایی که برای اداره آنها به حرکت دست و پا نیاز است، برای همه به ویژه افراد ناتوان بسیار حائز اهمیت است.

این امر مستلزم بازشناسی گفتار^۱ افراد مختلف است که یکی از دشوارترین مسائل موجود در مبحث شناسایی الگوست. تمهیدات فراوانی به منظور کاهش نرخ خطای^۲ بازشناسی گفتار صورت گرفته است که منجر به سیستم های بازشناسی مقاومتری شده اند.

مشکلی که معمولاً در بکارگیری یک سیستم بازشناسی مطرح است، استفاده از یک سیستم آموزش دیده برای گوینده های جدیدی است که به علت نداشتن سیستم وابسته به گوینده مناسب معمولاً دقت بازشناسی مطلوبی ندارند. پس از آموزش یک سیستم با داده آموزشی کافی جمع آوری شده از گوینده های مختلف سیستم حاصله یک سیستم مستقل از گوینده^۳ (SI) خواهد بود و چنین سیستمی مسلماً نسبت به سیستمی که تنها با گفتار یک گوینده خاص آموزش دیده و برای همان گوینده نیز تست می

^۱ Speech Recognition

^۲ Error Rate

^۳ Speaker Independent

شود (سیستم وابسته به گوینده^۴ (SD))، دقت بازشناسی پایین تری خواهد داشت. از آنجائیکه همیشه جمع آوری داده زیاد از یک گوینده به منظور آموزش سیستم SD او ممکن نیست، بنابراین سعی بر این است که پس از یکبار آموزش سیستم SI، برای گوینده هایی غیر از گوینده های آموزشی نیز بتوان از سیستم SI با دقت خوبی استفاده کرد. به این منظور از روش های تطبیق گوینده استفاده می گردد. در این روش ها با داشتن تنها مقدار محدودی داده از گوینده جدید، سیستم SI به گونه ای تغییر داده می شود که برای گوینده جدید عملکردی نظیر سیستم SD آن گوینده داشته باشد.

روش هایی نظیر روش های مبتنی بر خوشه بندی گوینده، MAP و MLLR، سالهاست که برای این منظور مورد استفاده قرار می گیرند.

خوشه بندی گوینده یکی از تکنیک های اصلی در تطبیق گوینده است. در این روش، گوینده هایی که خصوصیات آکوستیکی مشابه یا حداقل ویژگی های مشابهی در فضای ویژگی دارند خوشه بندی شده و سپس مدل های وابسته به خوشه برای هر خوشه تشکیل می گردد. در مرحله بازشناسی، خوشه ای که بیشترین مطابقت را با گفتار ورودی داراست، خوشه ای که بیشترین درستنمایی را دارد، انتخاب شده و از مدل های مارکوف پنهان (HMM) این خوشه برای بازشناسی استفاده می شود. این HMM ها همچنین می توانند در تطبیق گوینده به عنوان مدل های ابتدایی به جای HMM های مدل SI مورد استفاده قرار بگیرند. روش خوشه بندی همچنین می تواند به دلیل راحتی ترکیب با تکنیک های رایج تطبیق نظیر MAP و MLLR مورد استفاده قرار بگیرد.

روش های تطبیق گوینده موجود را می توان به دو دسته تقسیم بندی کرد:

روش های مبتنی بر ویژگی و روش های مبتنی بر مدل. نرمالیزاسیون لوله صوتی (VTLN) یکی از روش های تطبیق بر مبنای ویژگی است که تغییرات گوینده ای ایجاد شده به علت طول های مختلف لوله صوتی را اصلاح می کند.

تکنیک های مبتنی بر مدل را می توان به سه مجموعه تقسیم بندی کرد: تکنیک های بیژن، تکنیک های مبتنی بر تبدیل و تکنیک های مبتنی بر فضای ویژه^۵. این روش ها به مقدار مشخصی داده تطبیق از گوینده جدید احتیاج دارد تا با استفاده از این اطلاعات عملکرد بازشناسی گوینده جدید را بهبود بخشند.

⁴ Speaker Dependent

⁵ Eigen Space

در تطبیق گوینده، در صورتیکه نسخه آوایی از داده تطبیق موجود باشد تطبیق با نظارت و در غیر اینصورت بدون نظارت نامیده می شود. همچنین در این پروژه منظور از تطبیق گوینده سریع، تطبیق گوینده با داده تطبیقی کمتر از ۱۰ ثانیه می باشد.

۱-۲- طرح کلی پایان نامه

در این پروژه به منظور بهبود بازشناسی گفتار گوینده های جدید، از روش های خوشه بندی استفاده شده است. با استفاده از داده تطبیق کمتر از ۱۰ ثانیه، با بکارگیری تکنیک خوشه بندی به عنوان یک مرحله ابتدایی پیش از تطبیق های رایج MAP و MLLR عملکرد سیستم را به میزان قابل توجهی بهبود داده ایم.

در فصل دوم از این پایان نامه ابتدا کلیاتی در مورد ویژگی های گفتار، بازشناسی گفتار، استخراج ویژگی، مدل های مارکوف پنهان و مدل های مخلوط گوسی آورده شده است.

در فصل سوم روش های متداول موجود در تطبیق گوینده شامل روش MAP^۶ و MLLR^۷ و روش نسبتاً جدید Eigenvoice معرفی شده اند.

در فصل چهارم چند روش خوشه بندی شامل خوشه بندی جنسیت، باینری، K-means و C-means و Soft K-means توضیح داده شده اند. همچنین در انتهای این فصل روش خوشه بندی ای بر مبنای بردارهای پشتیبان تک کلاسی^۸ (OCSVM) معرفی شده است.

در فصل پنجم آزمایش های انجام شده بر روی دادگان TIDIGITS، نحوه انجام هر یک و نتایج حاصله آورده شده است.

در فصل ششم نیز به جمع بندی، نتیجه گیری و پیشنهاد پرداخته ایم.

^۶ Maximum a Posteriori

^۷ Maximum Likelihood Linear Regression

^۸ One Class Support Vector Machines (OCSVM)

فصل دوم :

کلیاتی بر بازشناسی گفتار

فصل دوم

کلیاتی بر بازشناسی گفتار

۲-۱- مقدمه

در این بخش توضیحاتی در رابطه با بازشناسی گفتار آورده شده که به صورت زیر سازماندهی شده است: ابتدا ویژگی های سیگنال گفتار توضیح داده شده است، سپس بازشناسی گفتار در سیستم های بازشناسی گفتار کنونی توضیح داده شده است. همچنین مدل های مارکوف پنهان^۱ (HMM) و مدل های مخلوط گوسی^۲ (GMM) مورد بررسی قرار گرفته است و در انتها روش استخراج ویژگی سیگنال گفتار توضیح داده شده است.

۲-۲- ویژگی های گفتار

در بازشناسی گفتار سعی بر استخراج اطلاعات موجود در سیگنال و تشخیص آن گفتار است. درانتقال سیگنال از فرستنده به گیرنده ممکن است سیگنال دچار اغتشاشاتی نظیر اغتشاش کانال، صدای دستگاه های اطراف و یا گفتار سایر گوینده ها شود. هر تبدیلی سیگنال را تغییر داده و بازشناسی آن را مشکل تر می سازد.

سیگنال گفتار دارای تغییر پذیری های ویژه ای است که بیان ریاضی آن را مشکل می سازد. برای مفید بودن، یک بازشناس گفتار باید کاری کند که این اثرات حذف شوند و یا ساختاری را متضمن شود که تاثیر این تغییرات را در نظر بگیرد[1].

^۱ Hidden Markov Model

^۲ Gaussian Mixture Model

- در این بخش منبع چنین تغییراتی را بیان می کنیم. برای سادگی این تغییرات را به چندین دسته تقسیم بندی می کنند، البته باید در نظر داشت که این بخش ها ممکن است همپوشانی هم داشته باشند [2]:
- (۱) تغییرات شیوه^۳ صحبت کردن: این تغییرات توسط گوینده قابل کنترل هستند که شامل مواردی از قبیل: دقت، وضوح، شمرده سخن گفتن و ... می باشد.
 - (۲) زمینه: شرایطی که در آن تولید سیگنال صورت می گیرد تاثیراتی بر گفتار دارد. این شرایط ممکن است بر سرعت صحبت، تغییرات شیوه صحبت و تاکید تاثیر داشته باشد.
 - (۳) استرس: شامل فاکتورهای احساسی و تغییرات تحمیلی محیط می باشد. نمونه های متداول آن ترس و واکنش لمبارد^۴ است.
 - (۴) کیفیت صدا: این بخش شامل تاثیراتی نظیر صدای ناراحت و عصبی، سوت زدن و ... است.
 - (۵) آهنگ صحبت: آهنگی که صحبت با آن ادا شده است که در قابلیت فهم نیز موثر است.
- علاوه بر این موارد، تفاوت های فیزیولوژیکی نظیر جنسیت نیز بر گفتار تاثیر می گذارد. به عنوان مثال خانم ها لوله صوتی کوتاهتر و فرکانس گام^۵ بالاتری نسبت به آقایان دارند. همچنین احتمال اینکه خانم ها حجم^۶ صدای پایین تری نسبت به آقایان داشته باشند بیشتر است.

۲-۳- بازشناسی گفتار

در سیستم های بازشناسی متداول دو مرحله وجود دارد:

- (۱) فاز آموزش: در این مرحله سیستم بازشناسی راه اندازی می شود.
- (۲) فاز بازشناسی: در آن از سیستم بازشناسی برای تشخیص صحبت بیان شده استفاده می شود.

³ Style

⁴ Lombard reflex

⁵ Pitch

⁶ volume