

دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده علوم مهندسی
گروه هوش مصنوعی



گزارش پروژه نهایی یادگیری ماشین

پروژه نهایی برای دریافت درجهٔ در رشتهٔ مهندسی کامپیوتر
گرایش هوش مصنوعی

سید محمد عرفان موسوی منزله، غزل زمانی نژاد، نیما کمبرانی، سید
مهدی موسوی

اساتید راهنما

دکتر اعرابی و دکتر ابوالقاسمی و دکتر توسلی پور

پاییز و زمستان ۱۴۰۲

فهرست مطالب

پ	فهرست تصاویر
۱	فصل ۱: گزارش کار
۱	۱.۱ پیش پردازش
۱	۱.۱.۱ کوتاه کردن داده
۲	۲.۱.۱ تبدیل فوریه کوتاه زمانی
۲	۳.۱.۱ فیلتر مل
۴	۴.۱.۱ استخراج ویژگی
۴	۲.۱ طبقه بندی
۵	۱.۲.۱ طبقه بند های معرفی شده در کلاس
۵	۱.۱.۲.۱ طبقه بند گاوسی
۶	۲.۱.۲.۱ طبقه بند بردار پشتیبان
۷	۳.۱.۲.۱ درخت تصمیم
۷	۴.۱.۲.۱ طبقه بند آدابوست
۸	۲.۲.۱ سایر طبقه بند ها
۹	۱.۲.۲.۱ طبقه بند رندوم فارست
۱۰	۲.۲.۲.۱ طبقه بند بوست گرادیانی
۱۱	۳.۲.۲.۱ طبقه بند پرسپترون چندلایه
۱۲	۳.۲.۱ نتیجه گیری در طبقه بندی
۱۹	۳.۱ خوشه بندی

۱۳.۱	پردازش ویژگی	۱۹
۲.۳.۱	نرمال سازی ویژگی	۱۹
۳.۳.۱	فشرده سازی ویژگی ها و نمایش توزیع داده ها	۲۰
۱.۳.۳.۱	PCA	۲۰
۴.۳.۱	T-SNE	۲۰
۵.۳.۱	خوشه بندی	۲۸
۱.۵.۳.۱	K-Means	۲۸
۲.۵.۳.۱	Hierarchical Agglomerative	۳۱

۳۷	Automatic Speech Recognition(ASR)	فصل ۲:
۳۷	مقدمه	۱.۲
۳۷	پیش پردازش	۲.۲
۳۸	آموزش	۳.۲
۳۹	ارزیابی	۴.۲

فهرست تصاویر

۱.۱	اسپکتوگرام سیگنال ورودی یکی از نمونه ها	۲
۲.۱	اسپکتوگرام بعد از اعمال فیلتر مل	۳
۳.۱	اسپکتوگرام بعد از اعمال MFCC	۳
۴.۱	ماتریس درهم ریختگی طبقه بند بیزی	۶
۵.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند بیزی	۷
۶.۱	ماتریس درهم ریختگی طبقه بند بردار پشتیبان	۹
۷.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند بردار پشتیبان	۱۰
۸.۱	ماتریس درهم ریختگی طبقه بند درخت تصمیم	۱۱
۹.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند درخت تصمیم	۱۲
۱۰.۱	ماتریس درهم ریختگی طبقه بند آدابوست	۱۴
۱۱.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند آدابوست	۱۴
۱۲.۱	ماتریس درهم ریختگی طبقه بند رندوم فارست	۱۵
۱۳.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند رندوم فارست	۱۵
۱۴.۱	ماتریس درهم ریختگی طبقه بند گرادیان بوست	۱۶
۱۵.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند گرادیان بوست	۱۶
۱۶.۱	ماتریس درهم ریختگی طبقه بند پرسپترون چندلایه	۱۷
۱۷.۱	منحنی مشخصه عملکرد سیستم برای طبقه بند پرسپترون چند لایه	۱۸
۱۸.۱	توزیع نمونه‌ها بر اساس جنسیت گوینده، پس از اعمال PCA . دسته‌های یک جنسیت در کنار هم قرار گرفته اند که ممکن است بخاطر مشابه بودن گوینده باشد.	۲۱

- ۱۹.۱ توزیع نمونه‌ها بر اساس لهجه‌های مختلف، پس از اعمال PCA. در اکثر نمونه‌ها لهجه فارسی است. ۲۲
- ۲۰.۱ توزیع انواع جملات خوانده شده با توجه به تن هر یک، پس از اعمال PCA. ۲۳
- ۲۱.۱ توزیع نمونه‌ها بر اساس طول ویژگی‌ها، پس از اعمال PCA. نمونه‌های با طول کمتر از ۱۵۰ کوتاه، بین ۱۵۰ تا ۲۵۰ متوسط و بیشتر از ۲۵۰ بلند برچسب زده شده‌اند. اکثر نمونه‌های با طول کوتاه در کنار یکدیگر در گوشه راست تصویر قرار گرفته‌اند. ۲۴
- ۲۲.۱ توزیع نمونه‌ها بر اساس جنسیت گوینده، پس از اعمال t-SNE. ۲۵
- ۲۳.۱ توزیع نمونه‌ها بر اساس لهجه‌های مختلف، پس از اعمال t-SNE. ۲۵
- ۲۴.۱ توزیع انواع جملات خوانده شده با توجه به تن هر یک، پس از اعمال t-SNE. ۲۶
- ۲۵.۱ توزیع نمونه‌ها بر اساس طول ویژگی‌ها، پس از اعمال t-SNE. نمونه‌های با طول کمتر از ۱۵۰ کوتاه، بین ۱۵۰ تا ۲۵۰ متوسط و بیشتر از ۲۵۰ بلند برچسب زده شده‌اند. اکثر نمونه‌های با طول کوتاه در کنار یکدیگر در گوشه راست تصویر قرار گرفته‌اند. ۲۷
- ۲۶.۱ روند تغییر مقدار silhouette score با افزایش تعداد دسته‌ها در الگوریتم k-means. ۲۹
- ۲۷.۱ نمایش دسته‌های مختلف در الگوریتم k-means در حالت $k=4$. این حالت بیشتری خروجی معیار را دارد. ۲۹
- ۲۸.۱ نمایش دسته‌های مختلف در الگوریتم k-means در حالت $k=3$ ۳۰
- ۲۹.۱ نمایش دسته‌های مختلف در الگوریتم k-means در حالت $k=17$ ۳۰
- ۳۰.۱ نمودار تغییرات معیار ارزیابی عملکرد خوشه‌بندی با افزایش تعداد دسته‌ها در الگوریتم سلسله‌مراتبی. ۳۱
- ۳۱.۱ نمایش دندروگرام برای داده‌ها در کاهش بعد به ۳ با استفاده از روش PCA با توجه به نمودار دندروگرام و همچنین معیار Silhouette مقادیر ۶ و ۴ برای تعداد نمونه‌ها بهترین عملکرد را خواهند داشت. در این حالات فاصله بین دسته‌ها بیشترین حالت و هر دسته تا حد ممکن یکپارچه خواهد بود. ۳۲

- ۳۲.۱ نمایش دسته‌های مختلف در الگوریتم سلسله مراتبی با تعداد دسته برابر با ۶. با توجه به وجود ۳ بعد برای خوشه بندی خوشه ی ۴ در نمایش تا دو تکه نشان داده شده است. دسته ۵ نیز با توجه به قرارگیری به صورت جدا از بقیه داده به عنوان یک خوشه در نظر گرفته شده است که بخش زیادی از گویندگان در این دسته خانوم هستند. ۳۳
- ۳۳.۱ نمایش دسته‌های مختلف در الگوریتم سلسله مراتبی با تعداد دسته برابر با ۳. در این حالت با توجه به تعداد پایین دسته ها دو دسته آبی و بنفش به دارای همپوشانی بوده اند. ۳۴
- ۳۴.۱ نمایش دسته‌های مختلف در الگوریتم سلسله مراتبی با تعداد دسته برابر با ۱۴. در این حالت تعداد دسته ها زیاد تر از حد نیاز شده است و تعداد اعضای هر دسته تا حدی نامتوازن شده است. هر چند همچنان الگوریتم توانسته تا حدی با توجه به توزیع جنسیت و دیگر کلاس ها مانند طول صوت خوشه بندی را انجام دهد. ۳۵

فصل ۱

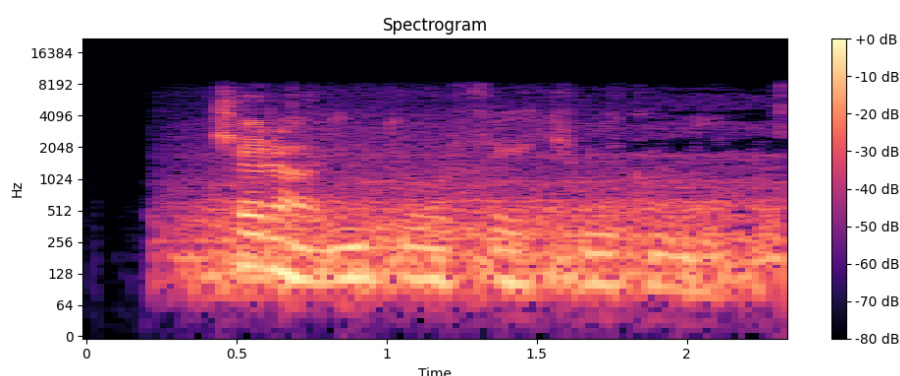
گزارش کار

۱.۱ پیش پردازش

در اولین قسمت از پروژه به تمیز کردن داده های صوتی در دست و همچنین استخراج ویژگی از آن می پردازیم. این کار را در چند قسمت انجام می دهیم.

۱.۱.۱ کوتاه کردن داده

در این مرحله، تلاش می کنیم با استفاده از توابع موجود، بخش هایی از سیگنال ورودی که دارای شدت و انرژی پایین تر از یک حد مشخص هستند را حذف کنیم. این کار باعث می شود که سکوت های میان کلمات یا در ابتدای صداها حذف شوند. بنابراین، سیگنال باقی مانده فقط شامل لحظاتی خواهد بود که در آن ها یک کلمه با معنی گفته شده است.



شکل ۱.۱: اسپکتروگرام سیگنال ورودی یکی از نمونه ها

۲.۱.۱ تبدیل فوریه کوتاه زمانی

در ادامه، به استخراج تبدیل فوریه کوتاه زمانی^۱ می‌پردازیم. این تبدیل، همانطور که از نامش پیداست، سیگنال را در بازه‌های کوتاه زمانی به تبدیل فوریه می‌رساند. دو پارامتری که به این تبدیل داده می‌شود، طول پنجره‌ای است که تبدیل فوریه روی آن انجام می‌شود و همچنین مقدار حرکت این پنجره در هر گام. در نهایت، قدر مطلق خروجی این تبدیل را می‌گیریم تا به اندازه این تبدیل برسیم.

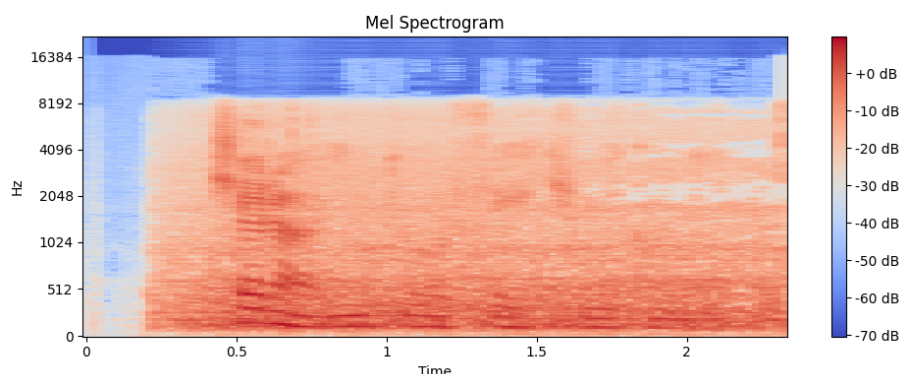
دلیل استفاده از تبدیل فوریه این است که این تبدیل به ما نمایشی از سیگنال در فضای فرکانسی می‌دهد. به عبارت دیگر، این تبدیل و خود سیگنال دوگان هم هستند و با استفاده از فضای فرکانسی، احتمالاً به جدپذیری بهتری خواهیم رسید. همچنین، مفاهیمی که در نمایش سیگنال صوت در زمان نمود ندارند، مثل بم بودن صدا، را می‌توان در این فضا نمایش داد. در این فضا، به راحتی می‌توانیم فیلتر طراحی کنیم یا از فیلترهای تعریف شده استفاده کنیم.

۳.۱.۱ فیلتر مل

در ادامه کار، از فیلتر مل^۲ استفاده می‌کنیم. این فیلتر بر اساس خصوصیات ادراکی انسان طراحی شده است و به خروجی فرکانسی مرحله قبلی اعمال می‌شود. در خروجی این مرحله، به اسپکتروگرام سیگنال می‌رسیم که توصیف سیگنال در طول زمان را نشان می‌دهد. حالا، تنها نیاز است که این خروجی را به مقیاس لگاریتمی ببریم

^۱STFT

^۲MEL

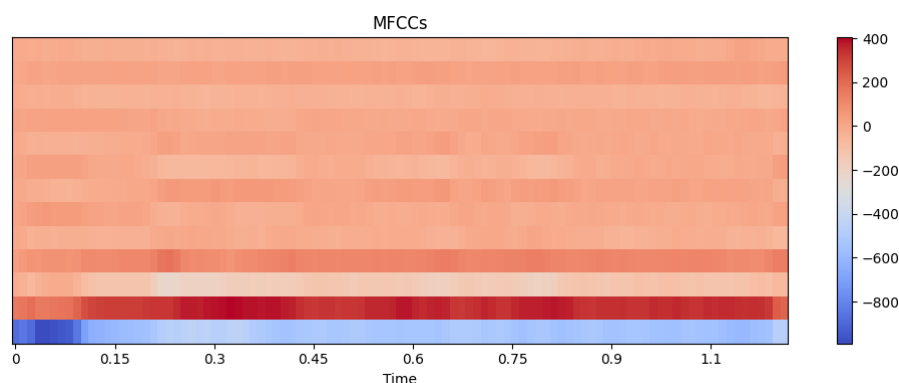


شکل ۲.۱: اسپکتوگرام بعد از اعمال فیلتر مل

تا بتوانیم شنوایی انسان را بهتر مدل کنیم.

فیلتر مل، که نامش از مخفف Mel Frequency Cepstral Coefficients گرفته شده، یکی از متداول‌ترین تکنیک‌ها در پردازش سیگنال‌های صوتی است. این فیلتر بر اساس این فرض که انسان صداهای نزدیک به هم را به خوبی تشخیص نمی‌دهد، طراحی شده است. بنابراین، فیلتر مل با تقسیم بندی فرکانس‌ها به صورت لگاریتمی، سعی در تقلید از شنوایی انسان دارد.

با استفاده از فیلتر مل، می‌توانیم ویژگی‌های مهم و کلیدی سیگنال صوتی را استخراج کنیم. این ویژگی‌ها می‌توانند در کاربردهای مختلفی مانند تشخیص گفتار، تشخیص لهجه، تشخیص هویت از طریق صدا و ... مورد استفاده قرار گیرند. در نهایت، با استفاده از این فیلتر، می‌توانیم سیگنال صوتی را به صورتی که بیشتر با شنوایی انسان همخوانی دارد، مدل کنیم.



شکل ۳.۱: اسپکتوگرام بعد از اعمال MFCC

۴.۱.۱ استخراج ویژگی

در مرحله نهایی، از سیگنال فیلتر شده استفاده می‌کنیم و با استفاده از توابع آماده، ویژگی MFCC را استخراج می‌کنیم. از این ویژگی به دلیل چندین دلیل استفاده می‌کنیم. اولاً، استفاده از این ویژگی باعث کاهش قابل توجه حجم داده‌ها می‌شود. دوماً، این ویژگی به تغییرات خارجی مثل نویز پس زمینه یا تغییرات جنسیت گوینده وابسته نیست. سوماً، این ویژگی به عنوان یکی از ویژگی‌های مهم در زمینه پردازش سیگنال‌های صوتی شناخته شده است.

MFCC یا همان Mel Frequency Cepstral Coefficients، یک روش استخراج ویژگی برای تجزیه و تحلیل سخنرانی و صدا است. این روش سیگنال‌های صوتی خام را به یک نمایش فشرده تبدیل می‌کند که اطلاعات مهم فرکانسی و زمانی را ضبط می‌کند. این ویژگی‌ها می‌توانند در کاربردهای مختلفی مانند تشخیص گفتار، تشخیص لهجه، تشخیص هویت از طریق صدا و ... مورد استفاده قرار گیرند. در نهایت، با استفاده از این فیلتر، می‌توانیم سیگنال صوتی را به صورتی که بیشتر با شنوایی انسان همخوانی دارد، مدل کنیم.

۲.۱ طبقه بندی

برای طبقه بندی داده‌ها بر اساس جنسیت از دو سری طبقه بند مختلف استفاده کردیم. سری اول شامل طبقه بند های معرفی شده در کلاس است. شامل طبقه بند بیزی با فرض توزیع گاوسی و فرض ساده لوحی، طبقه بند بردار پشتیبان، درخت تصمیم و طبقه بند آداپوست است. در دسته دوم طبقه بند هایی را داریم که کمتر در کلاس معرفی شده اند. مانند، رندوم فارست، گردینت بوستینگ و پرسپترون چند لایه. قبل از انجام دسته بندی، از بردار ویژگی بر روی بعد دوم میانگین گرفتیم تا تمام داده ها به اندازه ثابت ۱۳ عدد ویژگی داشته باشند. به عنوان پیش پردازش بیشتر داده ها از PCA و Normalize کردن نیز استفاده کردیم منتها بدلیل کاهش عملکرد مدل بعد از این پیش پردازش ها، از آن ها در مدل نهایی استفاده نکردیم. از آنجا که توضیح داده ها به نفع داده های مرتبط با جنسیت مرد است، آموزش مدل با همین داده ها مدل را به سمت پیش بینی بیشتر داده ها به عنوان مرد می برد. برای جلوگیری از این مشکل از تکنیک آور سمپلینگ برای افزایش تعداد داده های مرتبط با کلاس زن استفاده کردیم. تکنیک SMOTE یا Synthetic Minority Oversampling Technique یک روش برای مقابله با مشکل نامتوازن بودن داده‌ها در مسائل طبقه‌بندی است. در بسیاری از مسائل واقعی، تعداد نمونه‌های یک کلاس

(کلاس اقلیت) بسیار کمتر از کلاس دیگر (کلاس اکثریت) است. این موضوع می‌تواند باعث شود الگوریتم‌های یادگیری ماشین به نامناسب عمل کنند و عملکرد ضعیفی روی کلاس اقلیت داشته باشند. روش SMOTE برای حل این مشکل، نمونه‌های سنتتیک یا مصنوعی برای کلاس اقلیت ایجاد می‌کند. این نمونه‌های سنتتیک با استفاده از بین‌المللی بر روی نمونه‌های موجود کلاس اقلیت تولید می‌شوند.

۱.۲.۱ طبقه بند های معرفی شده در کلاس

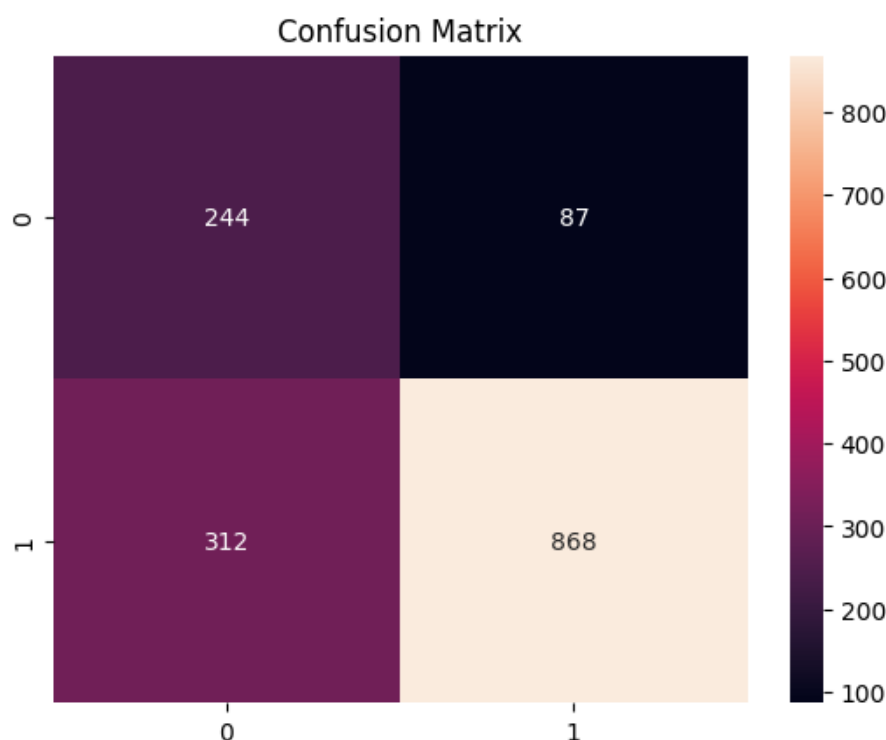
اولین طبقه بند معرفی شده در کلاس، مدل گاوسی مبتنی بر فرض ساده لوحانه بود. این مدل با وجود سادگی توانایی خوبی در دسته بندی داده ها دارد.

۱.۱.۲.۱ طبقه بند گاوسی

طبقه‌بند گاوسی یک الگوریتم یادگیری ماشین برای طبقه‌بندی است. این الگوریتم بر پایه توزیع احتمال گاوسی است و می‌تواند به عنوان پایه برای الگوریتم‌های یادگیری ماشین غیر پارامتریک پیچیده برای طبقه‌بندی و رگرسیون استفاده شود. در زیر نتایج استفاده از طبقه بند بیزی آورده شده است.

جدول ۱.۱: ارزیابی مدل بیزی

Support	F1-score	Recall	Precision	
۳۳۱	۵۵.۰	۷۴.۰	۴۴.۰	۰
۱۱۸۰	۸۱.۰	۷۴.۰	۹۱.۰	۱
۱۵۱۱	۷۴.۰			Accuracy
۱۵۱۱	۶۸.۰	۷۴.۰	۶۷.۰	Macro avg
۱۵۱۱	۷۶.۰	۷۴.۰	۸۱.۰	Weighted avg

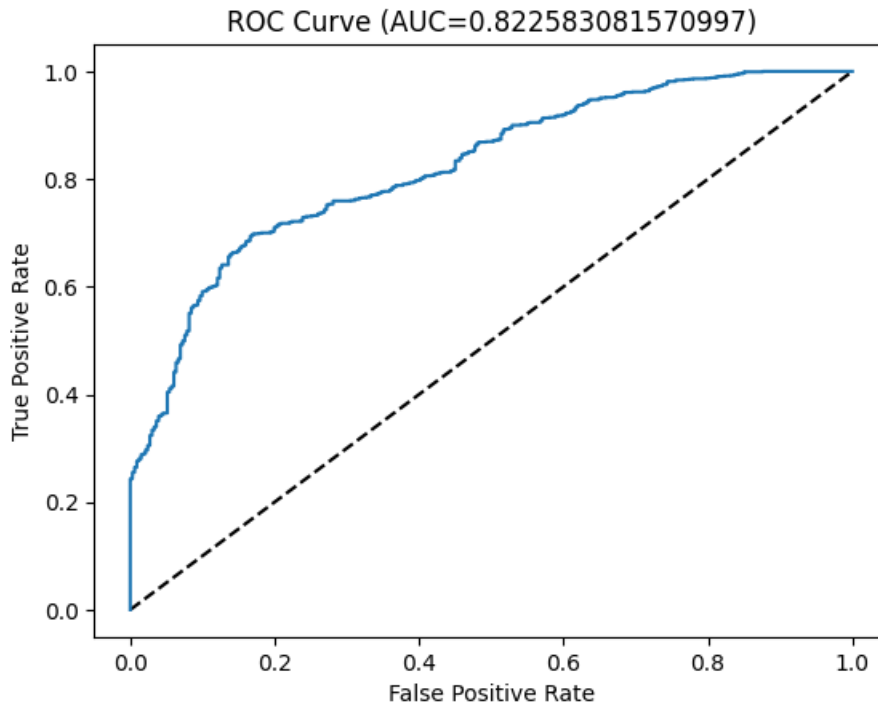


شکل ۴.۱: ماتریس درهم ریختگی طبقه بندی

۲.۱.۲.۱ طبقه بندی بردار پشتیبان

طبقه بندی بردار پشتیبان SVM یک مدل یادگیری ماشینی است که برای حل مسائل طبقه بندی دو گروهی استفاده می شود. پس از ارائه مجموعه هایی از داده های آموزشی برچسب دار به مدل SVM برای هر دسته، آن ها قادر به دسته بندی متن جدید هستند.

SVM دو مزیت اصلی نسبت به الگوریتم های جدیدتر مانند شبکه های عصبی دارد: سرعت بیشتر و عملکرد بهتر با تعداد محدودی از نمونه ها (در هزارها)، که این موضوع باعث می شود الگوریتم برای مسائل طبقه بندی متن بسیار مناسب باشد. SVM با تلاش برای یافتن هایپرپلین که بیشترین فاصله را از نزدیک ترین نقاط کلاس های مختلف دارد، عمل می کند. اگر چنین هایپرپلینی وجود داشته باشد، به آن هایپرپلین حداکثر حاشیه یا حاشیه سخت گفته می شود. SVM همچنین قادر است تا با داده های غیرخطی کار کند. در این موارد، از ترفندهایی مانند "ترفند هسته" استفاده می شود تا داده ها را به فضایی با بعد بالاتر تبدیل کند که در آن می توان یک حاشیه خطی را یافت. برای این پروژه ما از هسته RBF برای تبدیل داده ها به فضای جدید استفاده کردیم.



شکل ۵.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند بیزی

۳.۱.۲.۱ درخت تصمیم

درخت تصمیم یک الگوریتم یادگیری ماشین برای حل مسائل طبقه بندی و رگرسیون است. این الگوریتم یک ساختار درختی مانند جریان نما می سازد که هر گره داخلی آن یک ویژگی را نشان می دهد، هر شاخه قوانین را نشان می دهد و هر گره برگ نتیجه الگوریتم را نشان می دهد.

در طول آموزش، الگوریتم درخت تصمیم بهترین ویژگی را برای تقسیم داده ها بر اساس یک معیار مانند انتروپی یا ناخالصی گینی انتخاب می کند، که سطح ناخالصی یا تصادفی در زیر مجموعه ها را اندازه گیری می کند. هدف یافتن ویژگی ای است که بیشترین افزایش اطلاعات یا کاهش ناخالصی را پس از تقسیم ایجاد کند.

۴.۱.۲.۱ طبقه بند آدا بوست

طبقه بند AdaBoost یک الگوریتم یادگیری ماشین است که برای حل مسائل طبقه بندی و رگرسیون استفاده می شود. AdaBoost یکی از اولین الگوریتم های تقویتی است که معرفی شده است. این الگوریتم با ترکیب

جدول ۲.۱: ارزیابی طبقه بند ماشین بردار پشتیبان

Support	F1-score	Recall	Precision	
۳۳۱	۶۲.۰	۸۴.۰	۴۹.۰	۰
۱۱۸۰	۸۴.۰	۷۵.۰	۹۴.۰	۱
۱۵۱۱	۷۷.۰			Accuracy
۱۵۱۱	۷۳.۰	۸۰.۰	۷۲.۰	Macro avg
۱۵۱۱	۷۹.۰	۷۷.۰	۸۴.۰	Weighted avg

جدول ۳.۱: ارزیابی طبقه بند درخت تصمیم

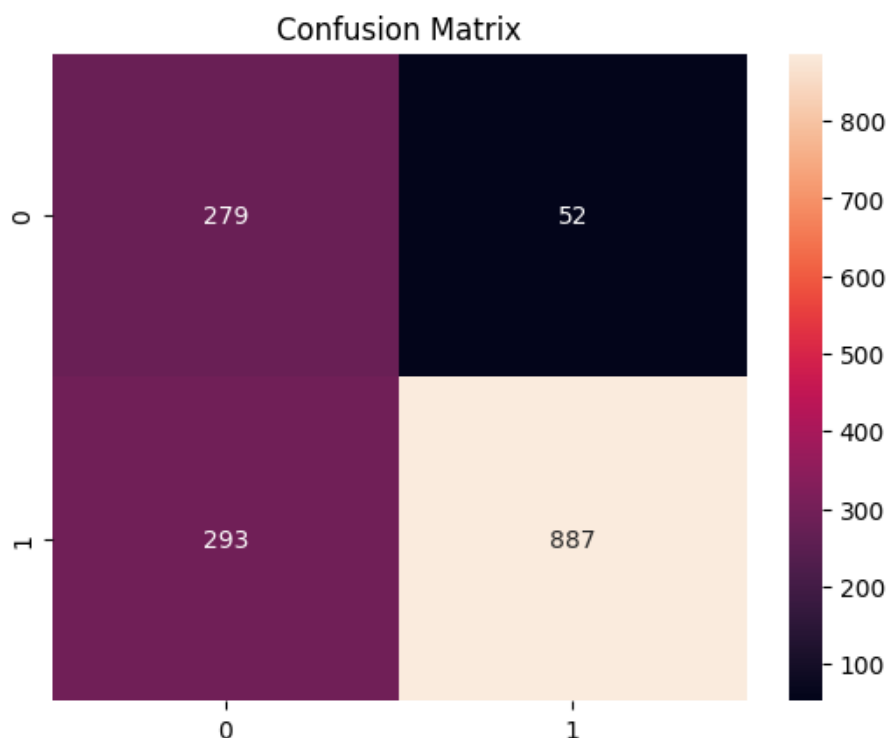
Support	F1-score	Recall	Precision	
۳۳۱	۸۸.۰	۹۰.۰	۸۵.۰	۰
۱۱۸۰	۹۶.۰	۹۶.۰	۹۷.۰	۱
۱۵۱۱	۹۴.۰			Accuracy
۱۵۱۱	۹۲.۰	۹۳.۰	۹۱.۰	Macro avg
۱۵۱۱	۹۴.۰	۹۴.۰	۹۵.۰	Weighted avg

چندین “طبقه‌بند ضعیف” یک “طبقه‌بند قوی” ایجاد می‌کند.

AdaBoost با شروع آموزش یک طبقه‌بند روی مجموعه داده اصلی و سپس آموزش نسخه‌های اضافی از طبقه‌بند روی همان مجموعه داده کار می‌کند، اما وزن نمونه‌هایی که به طور نادرست طبقه‌بندی شده‌اند تنظیم می‌شود تا طبقه‌بندهای بعدی بیشتر روی موارد سخت تمرکز کنند.

۲.۲.۱ سایر طبقه‌بندها

علاوه بر طبقه‌بندهای قبلی که در کلاس معرفی شد، سه طبقه‌بند دیگر نیز بررسی شده‌اند که در کلاس زیاد به آنها پرداخته نشده است.



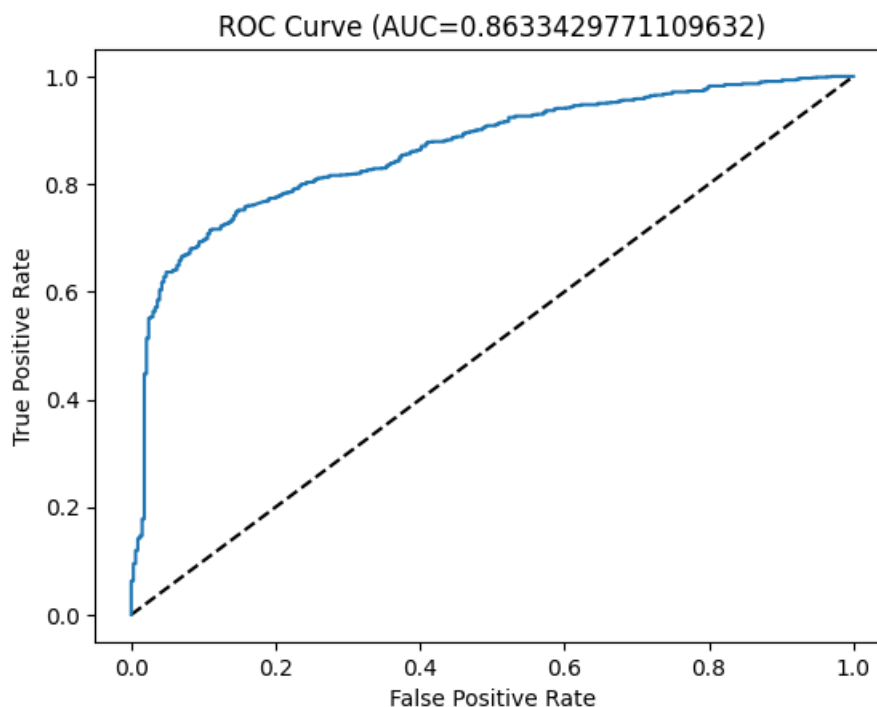
شکل ۶.۱: ماتریس درهم ریختگی طبقه بند بردار پشتیبان

۱.۲.۲.۱ طبقه بند رندوم فارست

Random Forest یک الگوریتم یادگیری ماشین است که برای حل مسائل طبقه‌بندی و رگرسیون استفاده می‌شود. این الگوریتم ترکیبی از چندین درخت تصمیم است که با هم ترکیب شده‌اند تا یک "جنگل" ایجاد کنند. در طول فرآیند آموزش، Random Forest چندین درخت تصمیم را بر اساس مجموعه داده‌های آموزشی می‌سازد. سپس، پیش‌بینی‌های هر درخت ترکیب می‌شوند تا یک پیش‌بینی نهایی ایجاد کنند. این فرآیند باعث می‌شود که Random Forest عملکرد بهتری نسبت به یک درخت تصمیم منفرد داشته باشد.

Random Forest همچنین از روش Bagging یا Bootstrap Aggregation استفاده می‌کند. در این روش، چندین نمونه آموزشی به طور تصادفی از مجموعه داده‌های آموزشی انتخاب می‌شوند (با جایگذاری مجدد) و سپس برای آموزش هر درخت استفاده می‌شوند.

Random Forest برای حل مسائل مختلفی از جمله طبقه‌بندی، رگرسیون، و کاهش ابعاد داده‌ها استفاده می‌شود. این الگوریتم به خاطر قابلیت‌هایی مانند تعدیل بیش‌برازش، کار با داده‌های بزرگ، و قابلیت کار با



شکل ۷.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند بردار پشتیبان

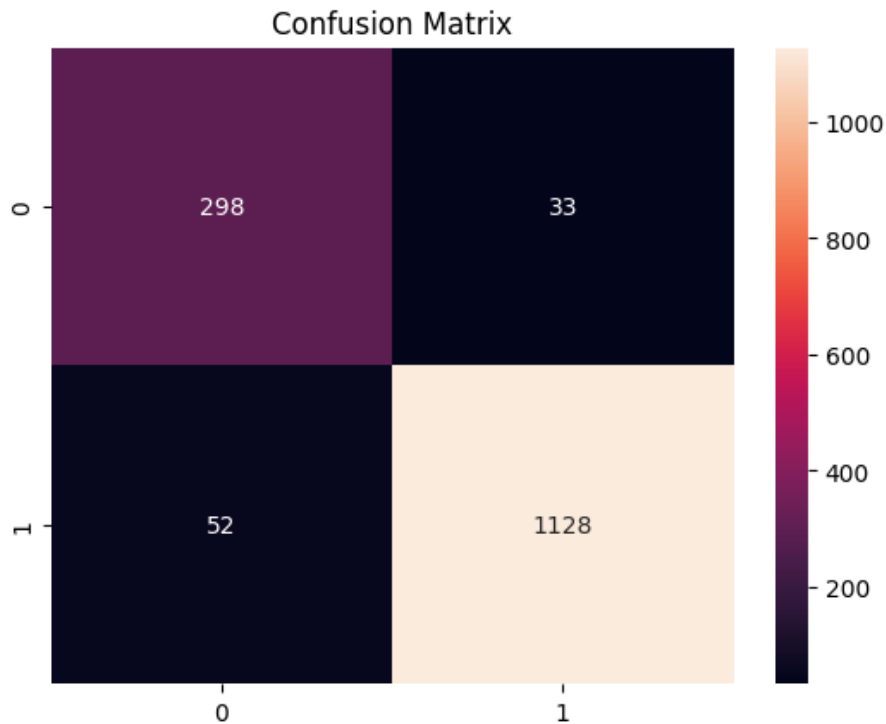
داده‌هایی که دارای ویژگی‌های نامرتبط یا گم‌شده هستند، محبوب است.

۲.۲.۲.۱ طبقه بند بوست گرادیانی

گرادیان بوست یک الگوریتم یادگیری ماشین است که برای حل مسائل طبقه‌بندی و رگرسیون استفاده می‌شود. این الگوریتم با ترکیب چندین "طبقه‌بند ضعیف" (معمولاً درختان تصمیم) یک "طبقه‌بند قوی" ایجاد می‌کند. گرادیان بوست با ساختن مدل‌های پیش‌بینی ساده به صورت توالی کار می‌کند که هر مدل سعی می‌کند خطای باقی مانده از مدل قبلی را پیش‌بینی کند. در اینجا، یک مدل یادگیری ضعیف، مدلی است که کمی بهتر از پیش‌بینی تصادفی عمل می‌کند.

گرادیان بوست برای مسائل مختلف یادگیری ماشین، از جمله رگرسیون و طبقه‌بندی، بسیار مناسب است. این الگوریتم با ترکیب طبقه‌بندهای ضعیف در یک مجموعه، مانند درختان تصمیم، یک پیش‌بین قوی را ایجاد می‌کند.

گرادیان بوست با بهینه‌سازی وزن‌های مدل بر اساس خطاهای تکرارهای قبلی، به تدریج خطای پیش‌بینی را



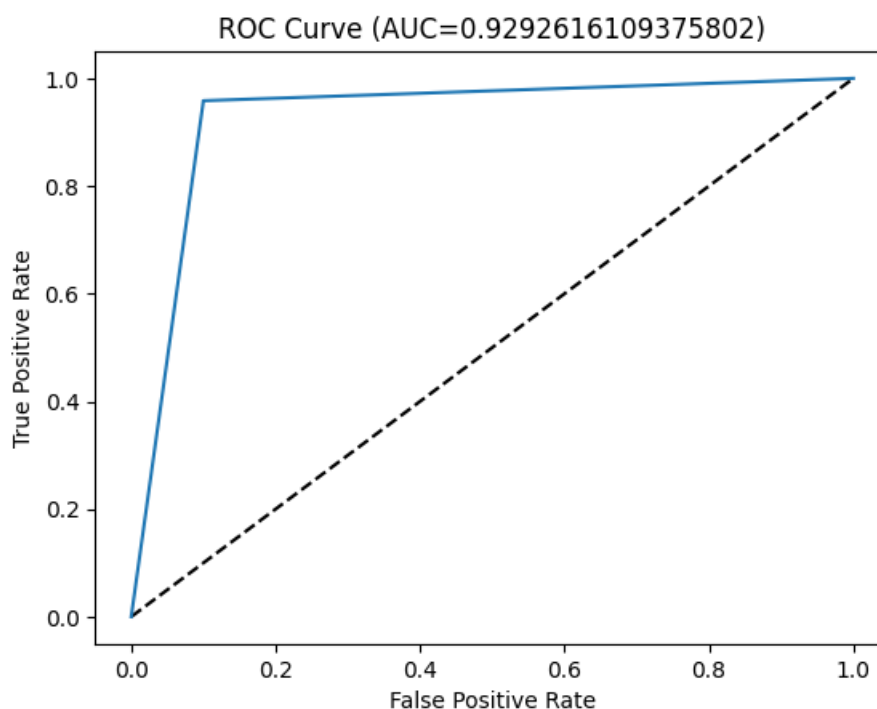
شکل ۸.۱: ماتریس درهم ریختگی طبقه بند درخت تصمیم

کاهش می دهد و دقت مدل را افزایش می دهد.

۳.۲.۲.۱ طبقه بند پرسپترون چندلایه

پرسپترون چند لایه یک نوع شبکه عصبی مصنوعی است که از یک سری نورون های کاملاً متصل با تابع فعال سازی غیرخطی تشکیل شده است. این شبکه ها حداقل سه لایه دارند و قادر به تشخیص داده هایی هستند که به صورت خطی جدا نمی شوند.

MLP از یک لایه ورودی و یک لایه خروجی تشکیل شده است، اما ممکن است چندین لایه پنهان در بین این دو لایه وجود داشته باشد.



شکل ۹.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند درخت تصمیم

۳.۲.۱ نتیجه گیری در طبقه بندی

در میان روش هایی غیر از شبکه عصبی، روش های مبتنی بر Ensemble Learning بهترین عملکرد را در مجموعه داده فعلی ما نشان داده اند. درخت تصمیم و شبکه عصبی پرسپترون عملکرد مشابهی داشته اند که این نشان دهنده قدرت طبقه بندی الگوریتم درخت تصمیم است، زیرا با وجود معماری بسیار ساده تر توانسته عملکردی مشابه با پرسپترون ارائه دهد. ماشین بردار پشتیبان بر خلاف انتظار عملکرد قوی نشان نداد. ما از هسته های خطی، RBF و چند جمله ای استفاده کردیم و بهترین عملکرد که مرتبط با RBF بود هم نتوانست جوابی نزدیک به مدل های دیگر بگیرد. و در نهایت مدل بیزی ساده لوح پایین ترین دقت را ارائه داد که نشان می دهد ویژگی های استخراج شده با فرض ساده لوحانه سازگاری نداشته اند و باید مدل های بیزی دقیق تری تست شود.

جدول ۴.۱: ارزیابی طبقه بند آدابوست

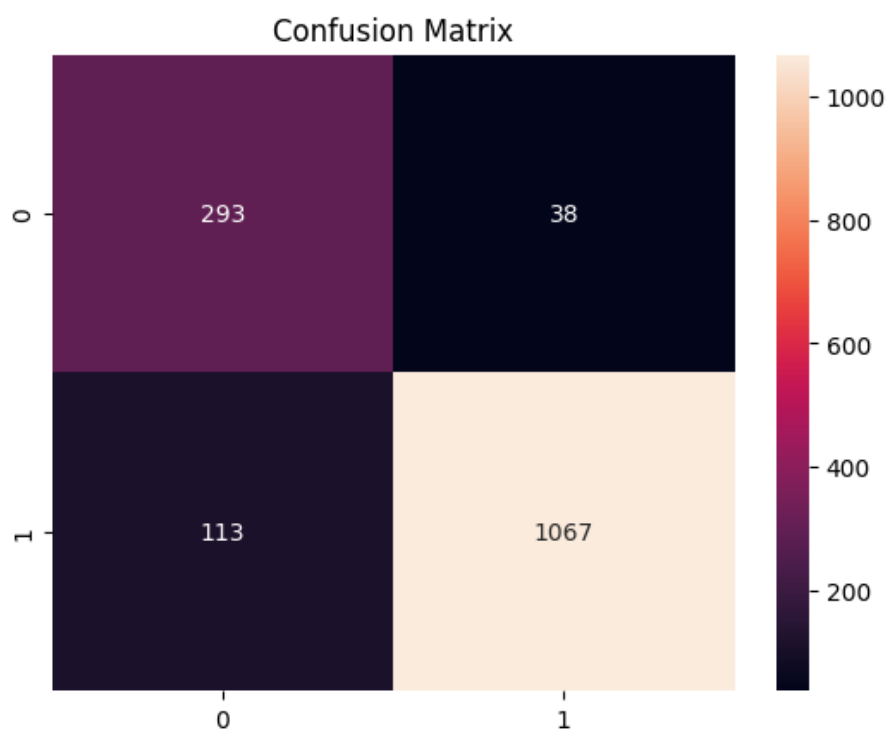
Support	F1-score	Recall	Precision	
۳۳۱	۸۰.۰	۸۹.۰	۷۲.۰	۰
۱۱۸۰	۹۳.۰	۹۰.۰	۹۷.۰	۱
۱۵۱۱	۹۰.۰			Accuracy
۱۵۱۱	۸۶.۰	۸۹.۰	۸۴.۰	Macro avg
۱۵۱۱	۹۰.۰	۹۰.۰	۹۱.۰	Weighted avg

جدول ۵.۱: ارزیابی طبقه بند رندوم فارست

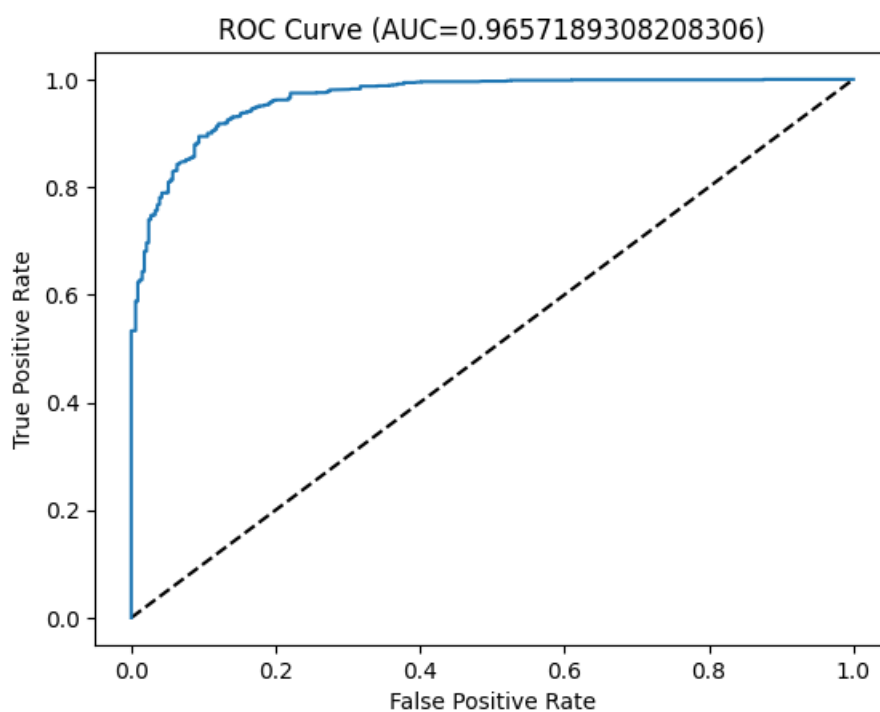
Support	F1-score	Recall	Precision	
۳۳۱	۹۶.۰	۹۵.۰	۹۷.۰	۰
۱۱۸۰	۹۹.۰	۹۹.۰	۹۹.۰	۱
۱۵۱۱	۹۸.۰			Accuracy
۱۵۱۱	۹۷.۰	۹۷.۰	۹۸.۰	Macro avg
۱۵۱۱	۹۸.۰	۹۸.۰	۹۸.۰	Weighted avg

جدول ۶.۱: ارزیابی طبقه بند گرادیان بوستینگ

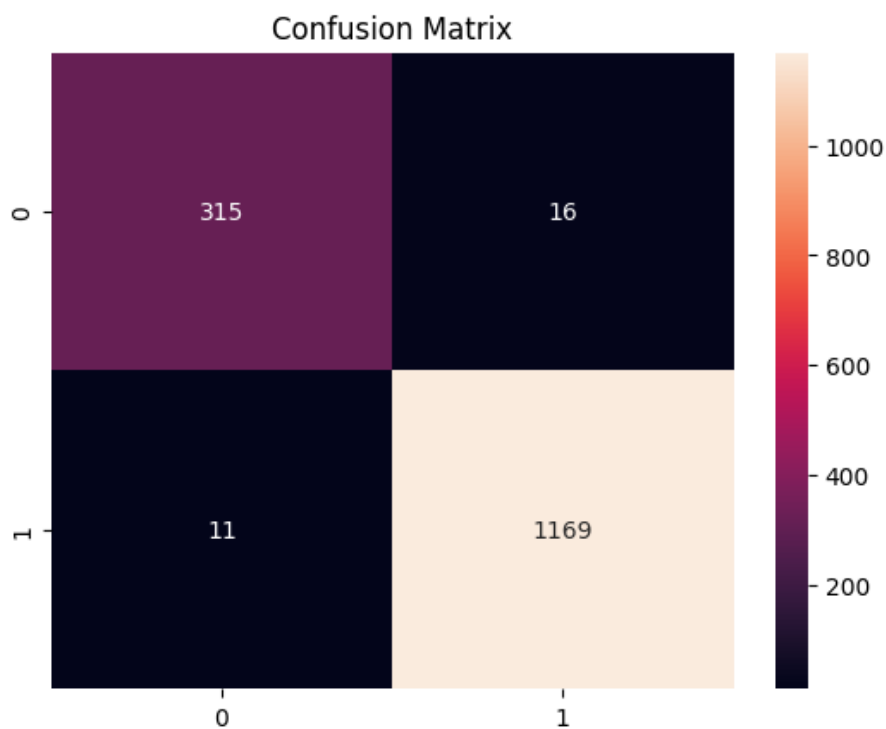
Support	F1-score	Recall	Precision	
۳۳۱	۸۸.۰	۹۳.۰	۸۴.۰	۰
۱۱۸۰	۹۷.۰	۹۵.۰	۹۸.۰	۱
۱۵۱۱	۹۵.۰			Accuracy
۱۵۱۱	۹۲.۰	۹۴.۰	۹۱.۰	Macro avg
۱۵۱۱	۹۵.۰	۹۵.۰	۹۵.۰	Weighted avg



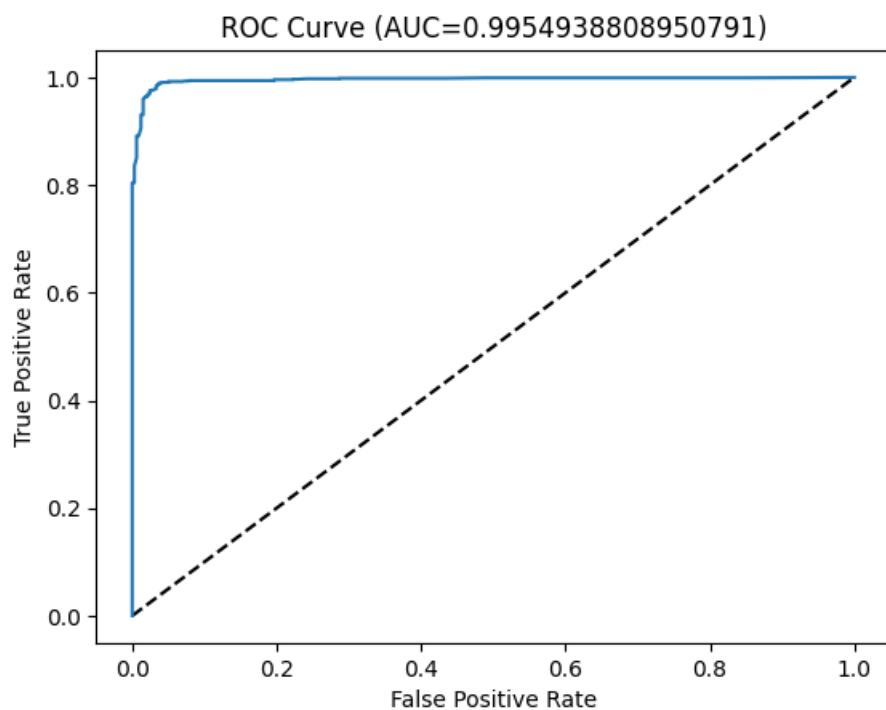
شکل ۱۰.۱: ماتریس درهم ریختگی طبقه بند آدابوست



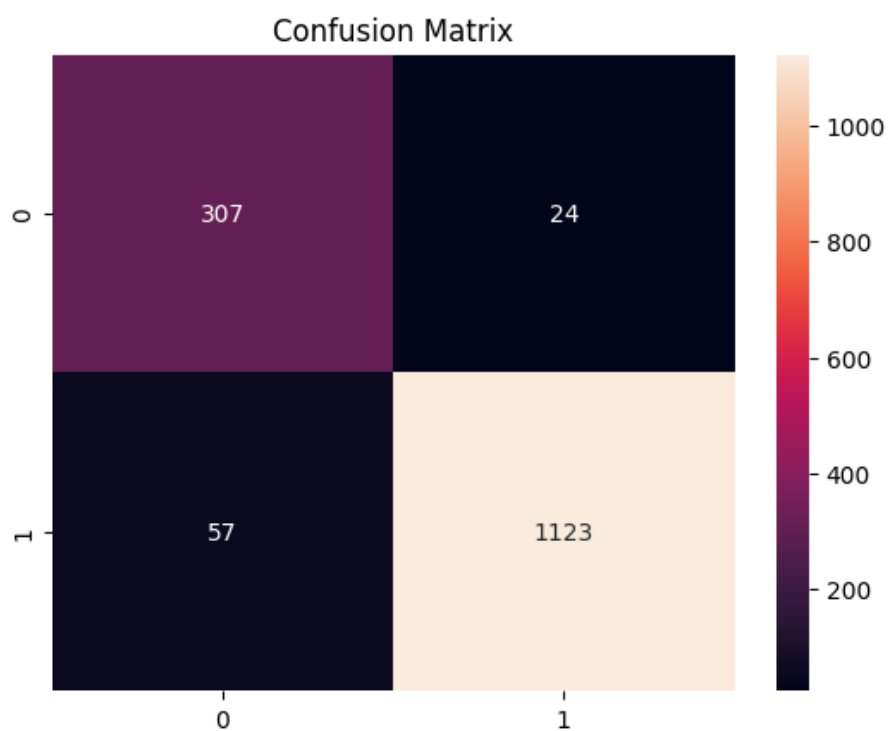
شکل ۱۱.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند آدابوست



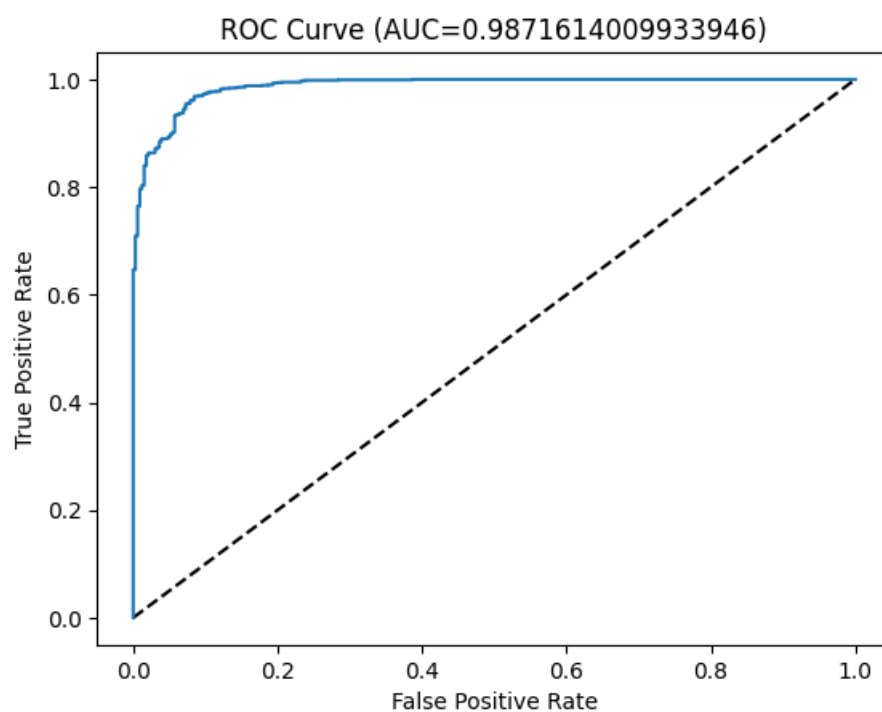
شکل ۱۲.۱: ماتریس درهم ریختگی طبقه بند رندوم فارست



شکل ۱۳.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند رندوم فارست



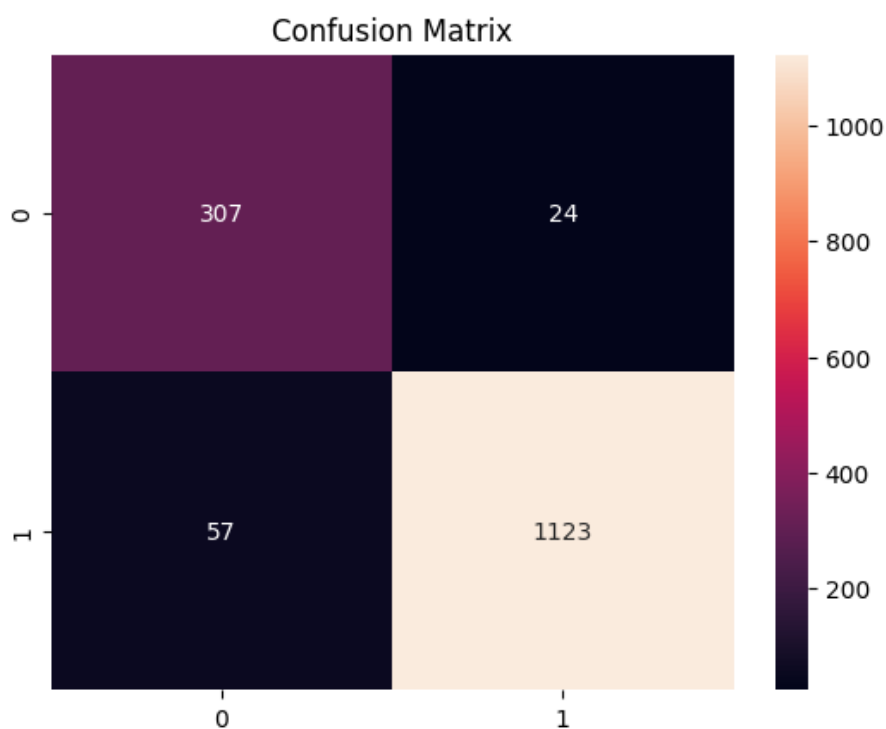
شکل ۱۴.۱: ماتریس درهم ریختگی طبقه بند گرادیان بوست



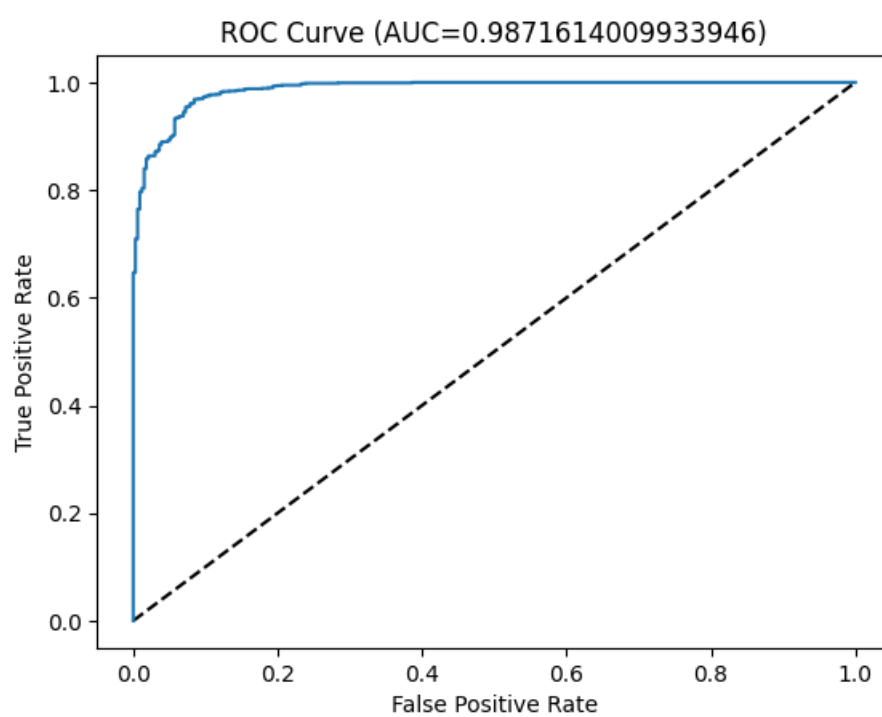
شکل ۱۵.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند گرادیان بوست

جدول ۷.۱: ارزیابی طبقه بند پرسپترون چند لایه

Support	F1-score	Recall	Precision	
۳۳۱	۹۱.۰	۸۷.۰	۹۶.۰	۰
۱۱۸۰	۹۸.۰	۹۹.۰	۹۶.۰	۱
۱۵۱۱	۹۶.۰			Accuracy
۱۵۱۱	۹۴.۰	۹۳.۰	۹۶.۰	Macro avg
۱۵۱۱	۹۶.۰	۹۶.۰	۹۶.۰	Weighted avg



شکل ۱۶.۱: ماتریس درهم ریختگی طبقه بند پرسپترون چند لایه



شکل ۱۷.۱: منحنی مشخصه عملکرد سیستم برای طبقه بند پرسپترون چند لایه

۳.۱ خوشه بندی

خوشه بندی به فرایند دسته بندی نمونه های مشابه در یک دسته بدون استفاده برچسب های متناظر گفته می شود. در فرایند دسته بندی نمونه ها با توجه به تفاوت طول فایل های صوتی مختلف و همچنین وجود نویز در داده های ورودی ابتدا نیاز است تا داده های ورودی تا حد ممکن تمیز شده و سپس از آن ها استخراج ویژگی صورت گیرد. برای این مسئله در بخش پیشین با استفاده از تابع mfcc در کتابخانه librosa ویژگی های مختلفی با توجه به صوت و تبدیل فوریه صوت استخراج شده است که در بخش های دیگر مورد استفاده قرار می گیرد. سپس ویژگی های بدست آمده با استفاده از روش های کاهش بعد فشرده سازی می شوند. در نهایت الگوریتم های خوشه بندی بر روی ویژگی های بدست آمده اعمال شده و نتایج آن نمایش داده می شود.

۱.۳.۱ پردازش ویژگی

با توجه به ساختار داده های صوتی، برای استخراج ویژگی در هر مرحله پنجره کوچکی از صوت انتخاب شده و برای آن بخش ویژگی های متناظر محاسبه می شوند. ویژگی استخراج شده دارای ۱۳ بعد برای هر پنجره است. در نتیجه با توجه به طول اولیه صوت، تعداد پنجره ها و متعاقبا طول ویژگی ها متفاوت خواهد بود. با توجه به اینکه برای اعمال الگوریتم های خوشه بندی نیاز است که طول ویژگی ها برای همه نمونه ها برابر باشد، باید استفاده از روشی برای همه نمونه ها ویژگی با طول ثابت محاسبه شود. برای محاسبه ویژگی واحد در اینجا ما از عملگر میانگین استفاده می کنیم. در نتیجه تمام مقادیر مختلف برای صوت با یکدیگر میانگین گرفته شده و به ازای هر صوت ۱۳ ویژگی خواهیم داشت.

۲.۳.۱ نرمال سازی ویژگی

در مرحله بعد با توجه به حساسیت الگوریتم های خوشه بندی به فاصله و مقیاس، نیاز داریم که این ویژگی ها نرمال سازی شوند. به عنوان مثال، در مورد ویژگی های صوتی، شدت یک ویژگی ممکن است مقیاسی بسیار بزرگتر از دیگری داشته باشد. دیگر اینکه بعضی از الگوریتم های روش های خوشه بندی زمانی که ویژگی ها به مقیاس مشابهی باشند، سریع تر همگرا می شوند. این موضوع به خصوص برای الگوریتم هایی مانند K-Means مهم است که فاصله بین نقاط داده عامل مهمی است. همچنین اعمال نرمال سازی، تفسیر نتایج خوشه بندی را

آسان‌تر می‌کنند. بدون نرمال‌سازی، ویژگی‌های با مقیاس بزرگتر ممکن است نتایج خوشه‌بندی را تسلط کنند و اهمیت نسبی هر ویژگی در فرآیند خوشه‌بندی را مشخص کنند. در نرمال‌سازی تلاش می‌شود تا برای هر ویژگی با استاندارد سازی میانگین ۰ و واریانس ۱ داشته باشیم.

۳.۳.۱ فشرده سازی ویژگی ها و نمایش توزیع داده‌ها

گاهی اوقات یک مجموعه داده ممکن است دچار Curse of Dimensionality شود که در آن تعداد ویژگی‌ها بسیار بیشتر از تعداد مشاهدات است. اساساً در فضاهای با ابعاد بالا با مشاهدات کم، تفکیک مشاهدات در ابعاد فوق‌دشوار می‌شود. همچنین، با توجه به اینکه داده‌ها به صورت دو بعدی نمایش داده می‌شوند نیاز است که ترتیب اهمیت ویژگی‌ها محاسبه و دو ویژگی پر اهمیت تر جهت نمایش توزیع داده‌ها بدست آورده شود. بدین دلایل، جهت کاهش ابعاد دو الگوریتم PCA و t-SNE مورد استفاده قرار گرفته است. در این حالت تعداد ابعاد با استفاده از دو الگوریتم گفته شده از ۱۳ به ۳ کاهش پیدا کرده و سپس دو بعد اول برای نمایش در نظر گرفته شده است.

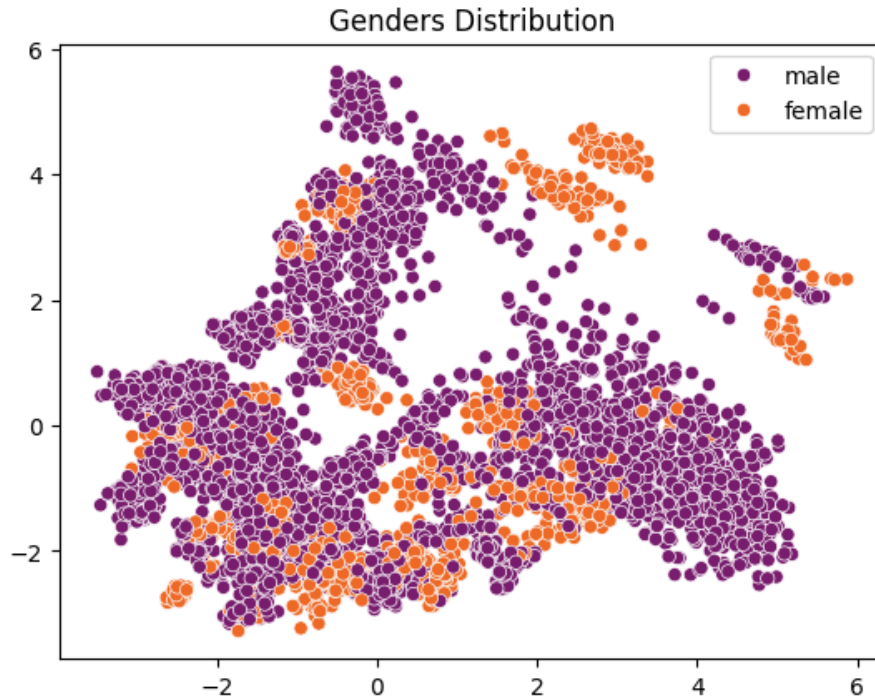
۱.۳.۳.۱ PCA

این روش کاهش بعد با یافتن یک تبدیل خطی بین ویژگی‌های ورودی با استفاده از بردار ویژه و مقدار ویژه سعی می‌کند تا نمونه‌های ورودی را بر روی بردارهای با بیشترین واریانس نگاشت کند. این کار می‌تواند به حذف ویژگی‌های با تاثیر کمتر در جدایی پذیری و همچنین حذف نویز کمک کند. در ادامه، نمایش نمونه‌ها در فضای دوبعدی برای دسته‌های مختلف برچسب‌ها قرار داده شده است. در این حالات تعداد ابعاد با استفاده از روش PCA کاهش بعد داده شده اند.

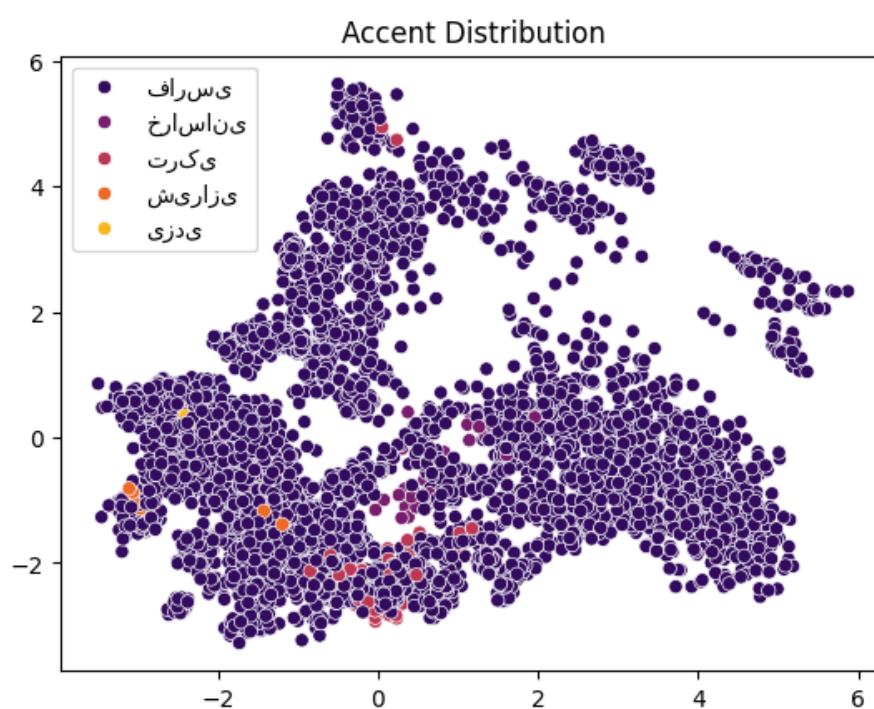
۴.۳.۱ T-SNE

این روش با یادگیری پارامترهایی سعی می‌کند تا یک تبدیل غیر خطی جهت کاهش بعد ارائه کند که ساختار فاصله نمونه‌ها را در ابعاد پایین تر نسبت حالت اولیه حفظ کند. با توجه به غیر خطی بودن این تبدیل برای یادگیری روابط پیچیده‌تر در داده‌ها مناسب است. این تبدیل باعث با توجه پارامتر تعداد همسایگی مورد توجه جهت حفظ

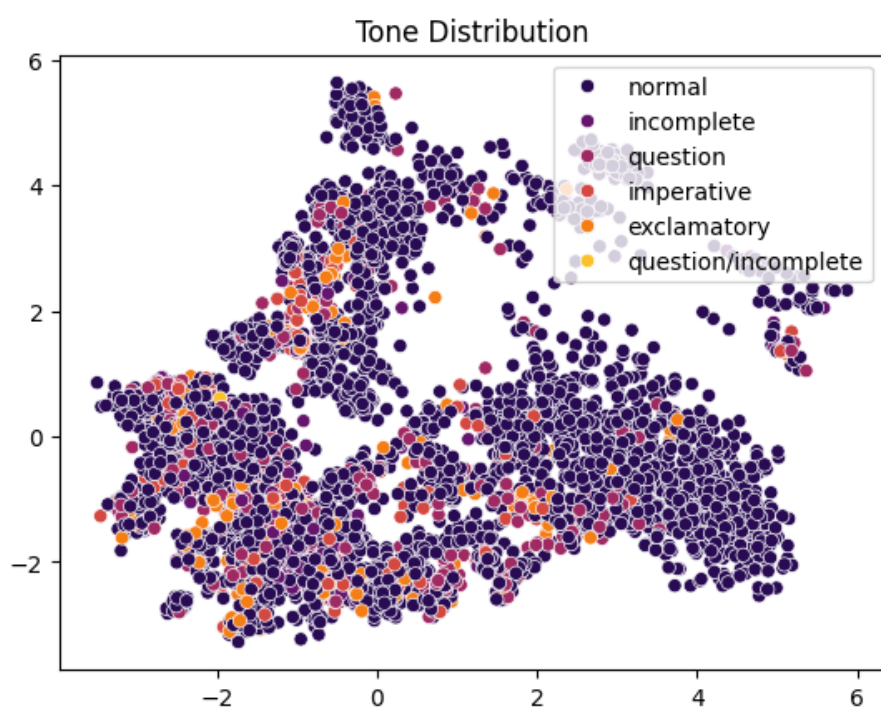
رابطه فاصله می‌تواند باعث ایجاد نقاط جدا از یکدیگر یا دسته‌های متعدد از نمونه‌های نزدیک به یکدیگر شود. در ادامه توزیع داده‌ها پس از کاهش ابعاد با استفاده از این روش قرار داده شده است.



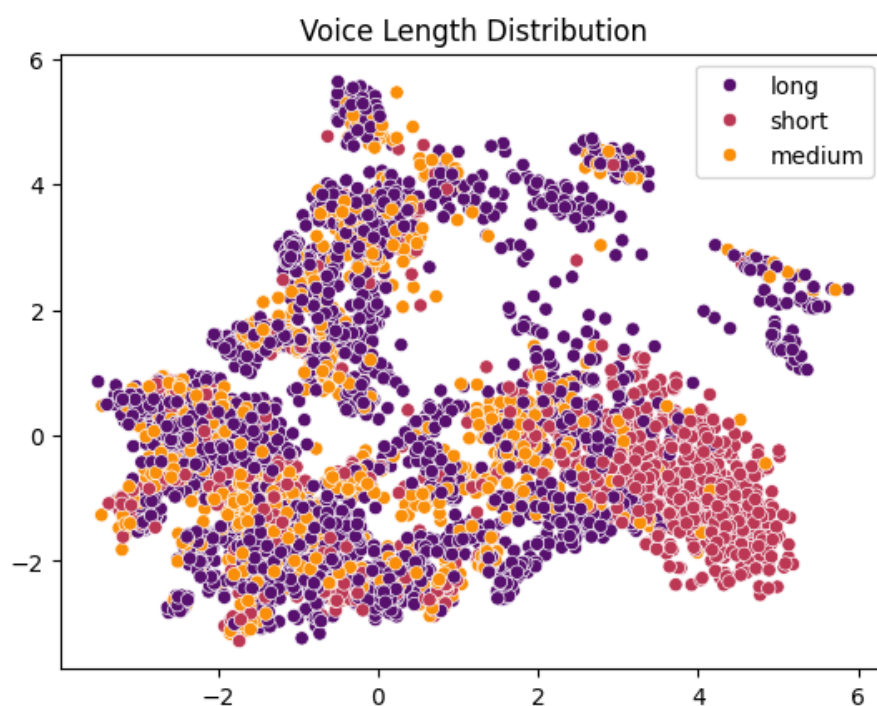
شکل ۱۸.۱: توزیع نمونه‌ها بر اساس جنسیت گوینده، پس از اعمال PCA. دسته‌های یک جنسیت در کنار هم قرار گرفته اند که ممکن است بخاطر مشابه بودن گوینده باشد.



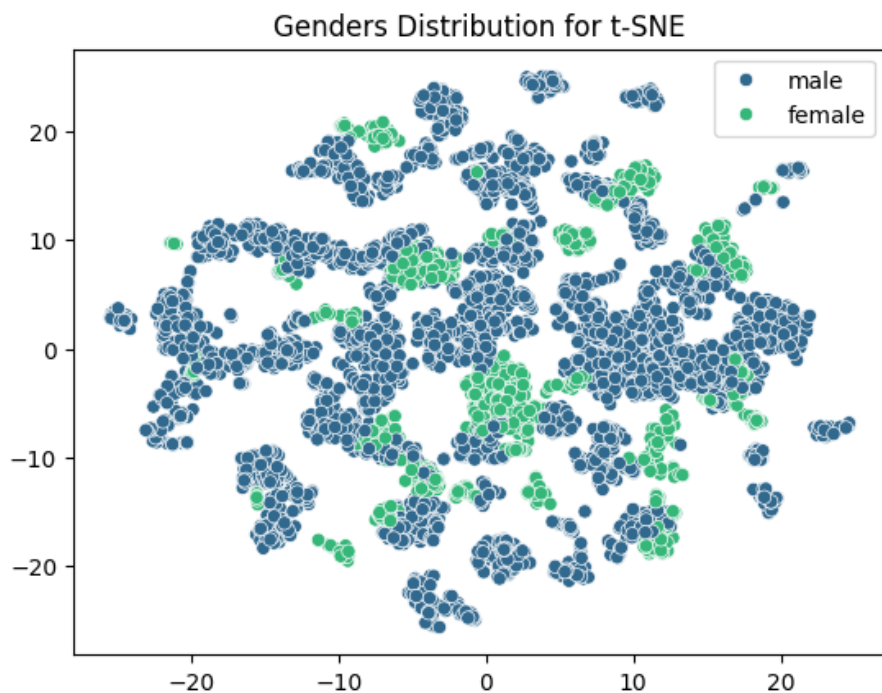
شکل ۱۹.۱: توزیع نمونه‌ها بر اساس لهجه‌های مختلف، پس از اعمال PCA. در اکثر نمونه‌ها لهجه فارسی است.



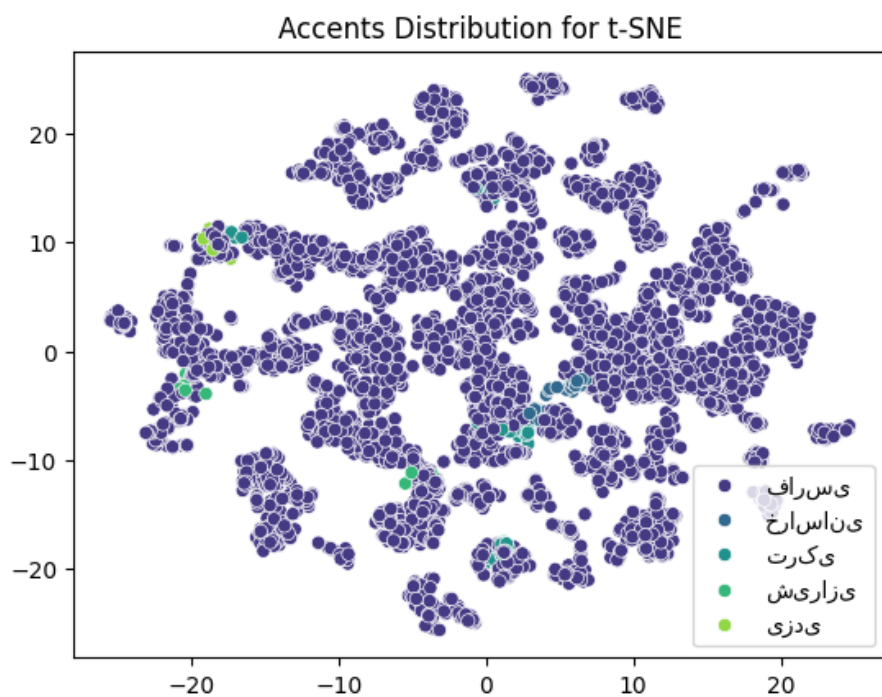
شکل ۲۰.۱: توزیع انواع جملات خوانده شده با توجه به تن هر یک، پس از اعمال PCA



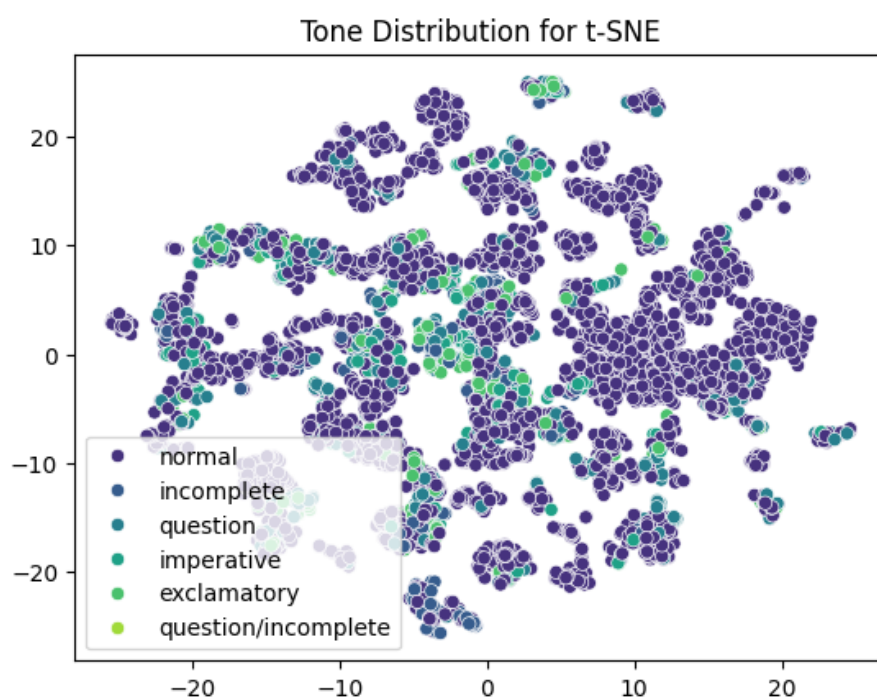
شکل ۲۱.۱: توزیع نمونه‌ها بر اساس طول ویژگی‌ها، پس از اعمال PCA. نمونه‌های با طول کمتر از ۱۵۰ کوتاه، بین ۱۵۰ تا ۲۵۰ متوسط و بیشتر از ۲۵۰ بلند برچسب زده شده‌اند. اکثر نمونه‌های با طول کوتاه در کنار یکدیگر در گوشه راست تصویر قرار گرفته‌اند.



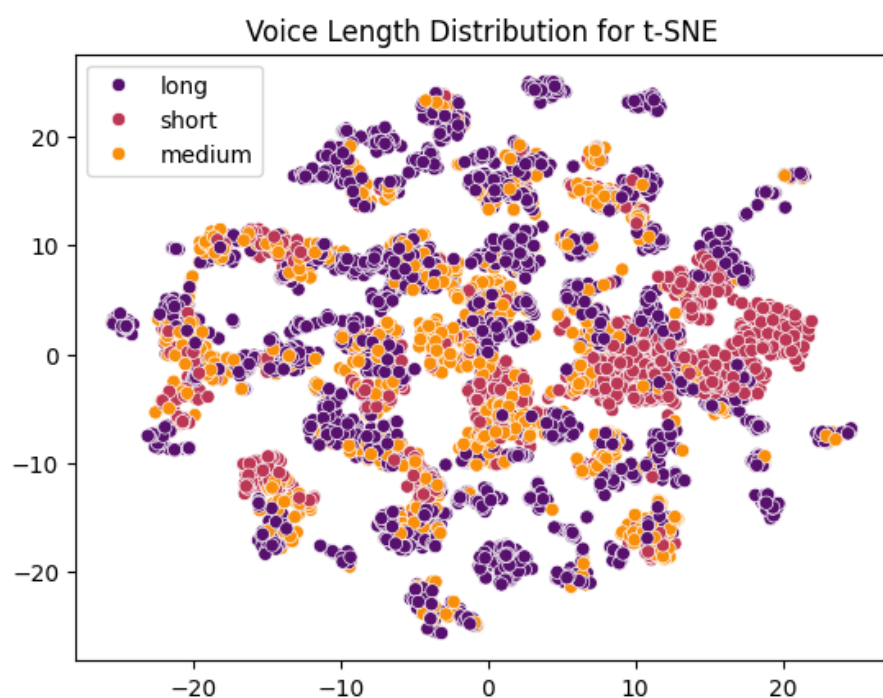
شکل ۲۲.۱: توزیع نمونه‌ها بر اساس جنسیت گوینده، پس از اعمال t-SNE



شکل ۲۳.۱: توزیع نمونه‌ها بر اساس لهجه‌های مختلف، پس از اعمال t-SNE



شکل ۲۴.۱: توزیع انواع جملات خوانده شده با توجه به تن هر یک، پس از اعمال t-SNE



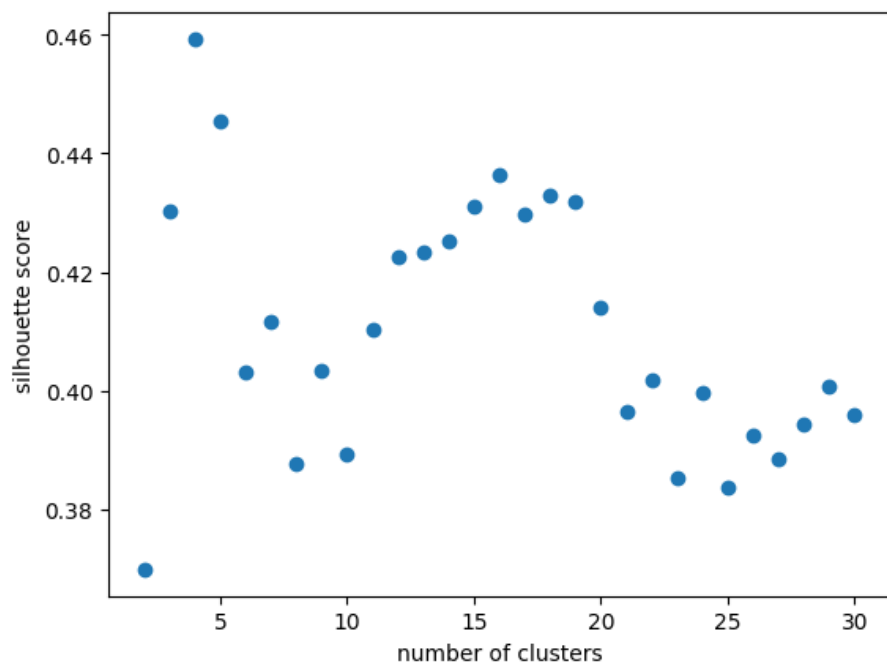
شکل ۲۵.۱: توزیع نمونه‌ها بر اساس طول ویژگی‌ها، پس از اعمال t-SNE. نمونه‌های با طول کمتر از ۱۵۰ کوتاه، بین ۱۵۰ تا ۲۵۰ متوسط و بیشتر از ۲۵۰ بلند برچسب زده شده‌اند. اکثر نمونه‌های با طول کوتاه در کنار یکدیگر در گوشه راست تصویر قرار گرفته‌اند.

۵.۳.۱ خوشه بندی

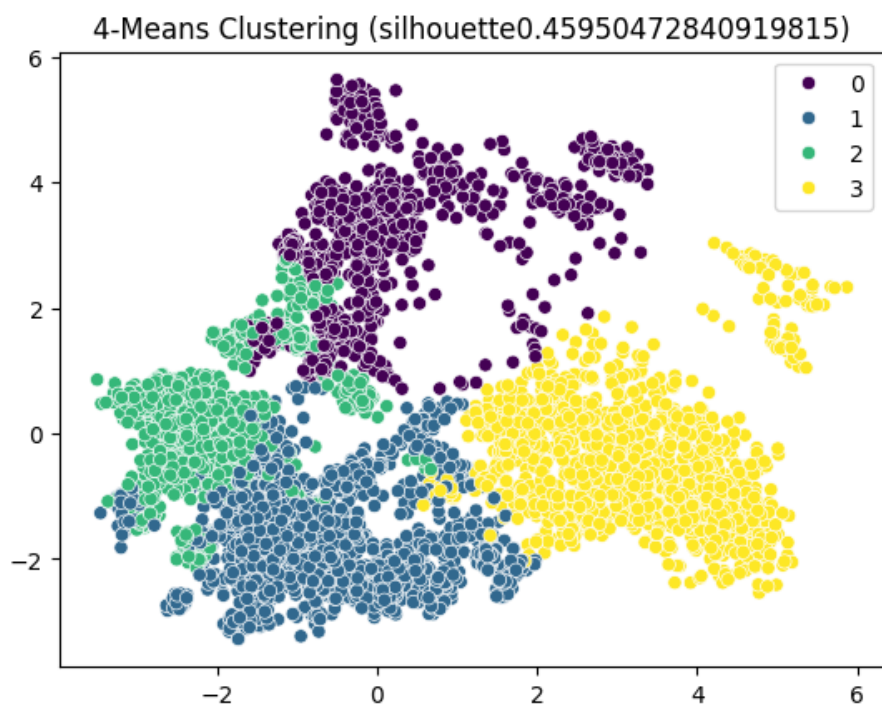
برای خوشه بندی داده ها ۲ روش k-means و Agglomerative Hierarchical مورد بررسی قرار گرفته اند. K-means یک روش ساده و از نظر محاسباتی بهینه است و می تواند برای مجموعه داده بزرگ به کار رود. خوشه های حاصل از این روش با مرکزشان نمایش داده می شوند و تفسیرپذیری راحتی نیز دارند. در این الگوریتم تنها پارامتر تعداد خوشه ها باید مشخص شود که آن هم از طریق امتیاز silhouette و یا روش elbow قابل محاسبه است. عیب k-means این است که خوشه ها را به صورت کروی در نظر می گیرد و ممکن است اگر توزیع داده ها به گونه ای باشد که خوشه ها شکل غیر معمول داشته باشند، چندان عملکرد مناسبی نداشته باشد. Agglomerative Hierarchical که به عنوان روش دوم مورد استفاده قرار گرفته، یک سلسله مراتب از ترکیب خوشه ها را در قالب دندروگرام نمایش می دهد و بدین ترتیب دید خوبی از روند تجمیع به ما می دهد. در این روش نیازی به مشخص کردن تعداد خوشه ها از ابتدای الگوریتم نیست و می توانیم تعداد خوشه ها را در انتهای کار، با جدا کردن نمودار دندروگرام از یک ارتفاع خاص انتخاب کنیم. از مزایای دیگر می توان به منعطف بودن linkage، یعنی روش محاسبه فاصله بین خوشه ها برای ترکیب، اشاره کرد. این الگوریتم نسبت به شکل خوشه ها فرضی ندارد و می تواند برای خوشه ها با شکل غیر معمول مناسب باشد. ایراد این روش آن است که از نظر محاسباتی پیچیده است، مخصوصاً برای مجموعه داده های بزرگ. اکنون به توضیح هر یک از روش ها می پردازیم.

۱.۵.۳.۱ K-Means

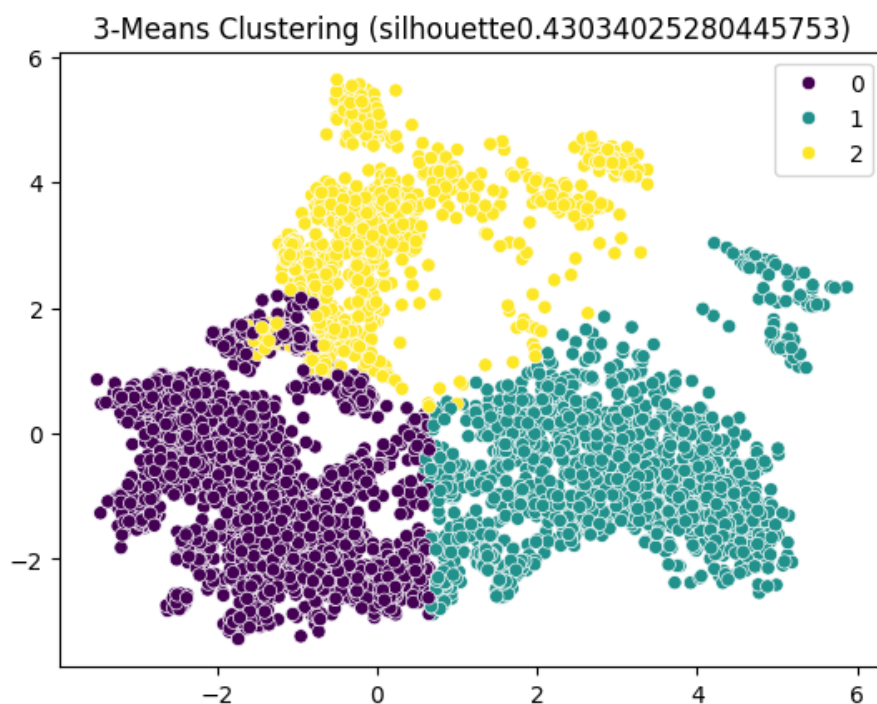
این روش با تعیین تعداد خوشه ها به اندازه k سعی می کند تا برای هر خوشه یک مرکز که از میانگین نمونه های داخل دسته محاسبه می شود قرار دهد. سپس هر نمونه به خوشه ای که به مرکز آن نزدیکتر است اختصاص داده می شود. برای محاسبه عملکرد الگوریتم خوشه بندی از معیار silhouette score استفاده می کنیم. در این حالت ابتدا الگوریتم خوشه بندی بر روی داده ها با تعداد خوشه ها در بازه ۲ تا ۳۰ اجرا شده و برای هر یک از اجراها معیار کیفیت محاسبه می شود. در ادامه روند تغییر مقدار silhouette و نمایش ۳ آزمایش مختلف قرار داده شده است. با توجه به اینکه برای عملیات خوشه بندی از ۳ ویژگی استفاده شده است اما در نمایش تنها دو بعد نمایش داده شده است خوشه ها کمی در نمایش تا حدی در هم تنیدگی دارند.



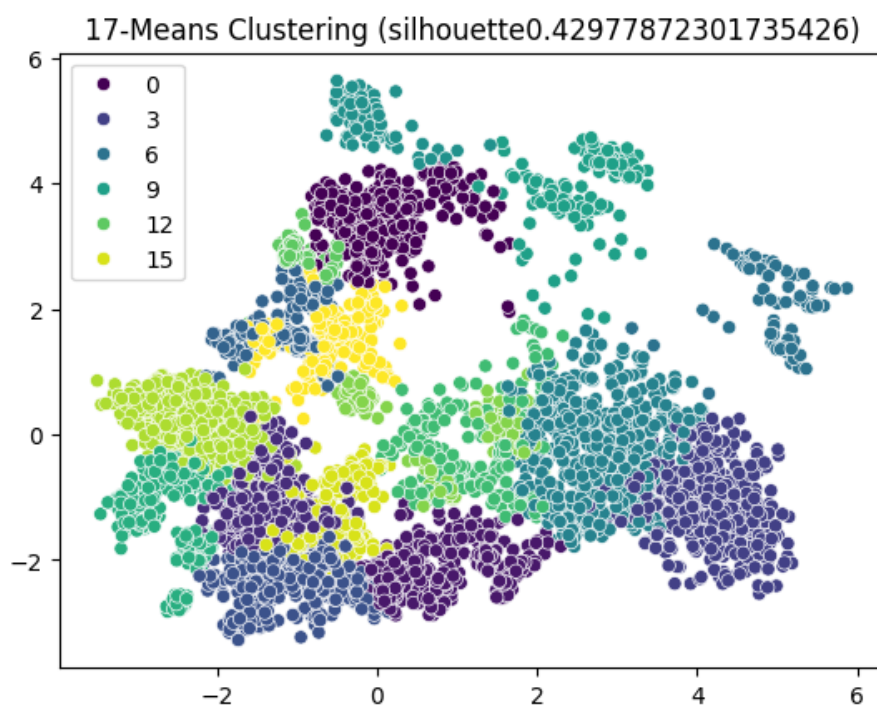
شکل ۲۶.۱: روند تغییر مقدار silhouette score با افزایش تعداد دسته‌ها در الگوریتم k-means



شکل ۲۷.۱: نمایش دسته‌های مختلف در الگوریتم k-means در حالت $k=4$. این حالت بیشتری خروجی معیار را دارد.



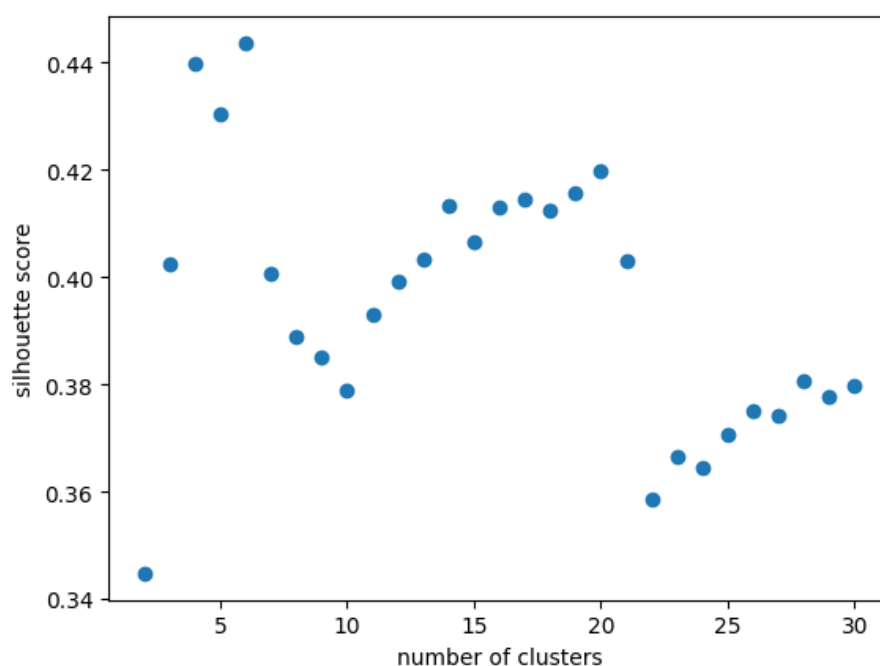
شکل ۲۸.۱: نمایش دسته‌های مختلف در الگوریتم k-means در حالت $k=3$.



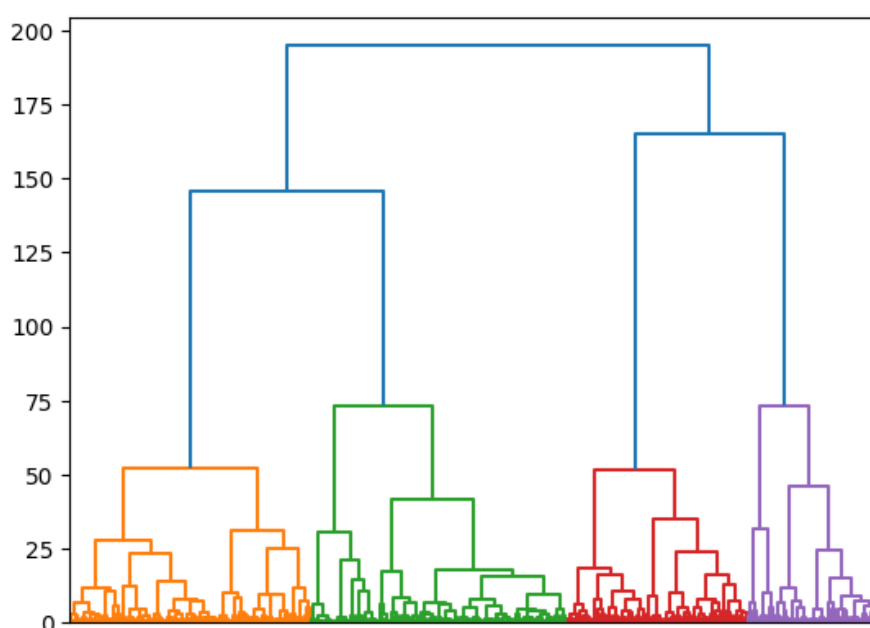
شکل ۲۹.۱: نمایش دسته‌های مختلف در الگوریتم k-means در حالت $k=17$.

۲.۵.۳.۱ Hierarchical Agglomerative

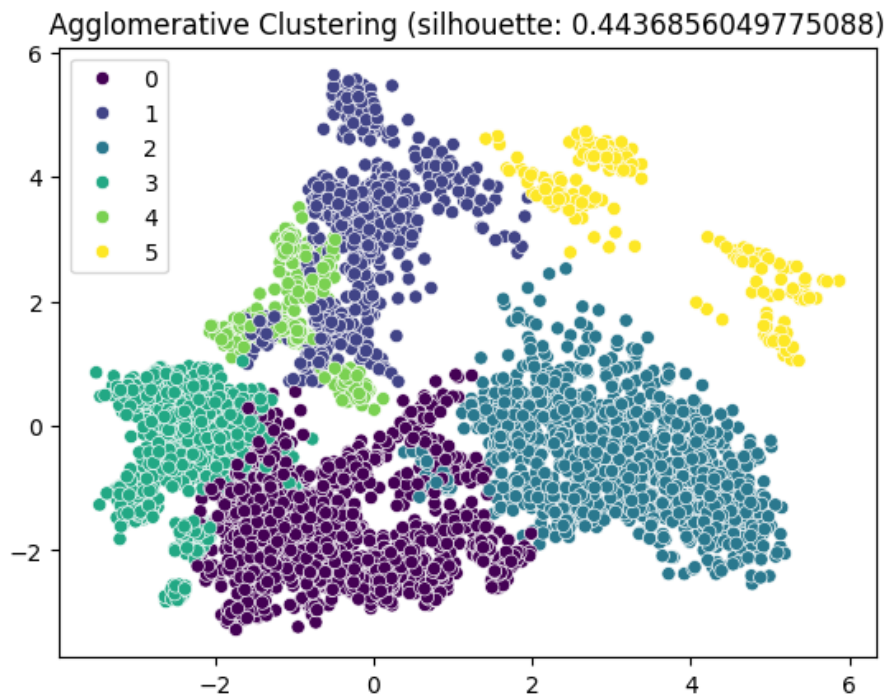
این الگوریتم خوشه‌بندی به صورت سلسله‌مراتبی نمونه‌ها با کمترین فاصله و بیشترین شباهت را با یک دیگر ترکیب می‌کند تا در نهایت پس از مراحل مختلف به یک دسته واحد برسد. با توجه به نیاز می‌توان تعداد دسته‌ها را در مقدار خاصی موقوف کرد. برای نمایش سلسله‌مراتب انتخاب این در روش از رسم دندروگرام استفاده کرده ایم. در این حالت ترتیب ادغام و شباهت دسته‌ها در هنگام ترکیب شدن مشخص می‌شود. در ادامه نمودار تغییر معیار silhouette و دسته‌بندی داده برای تعداد دسته‌های مختلف نمایش داده شده است.



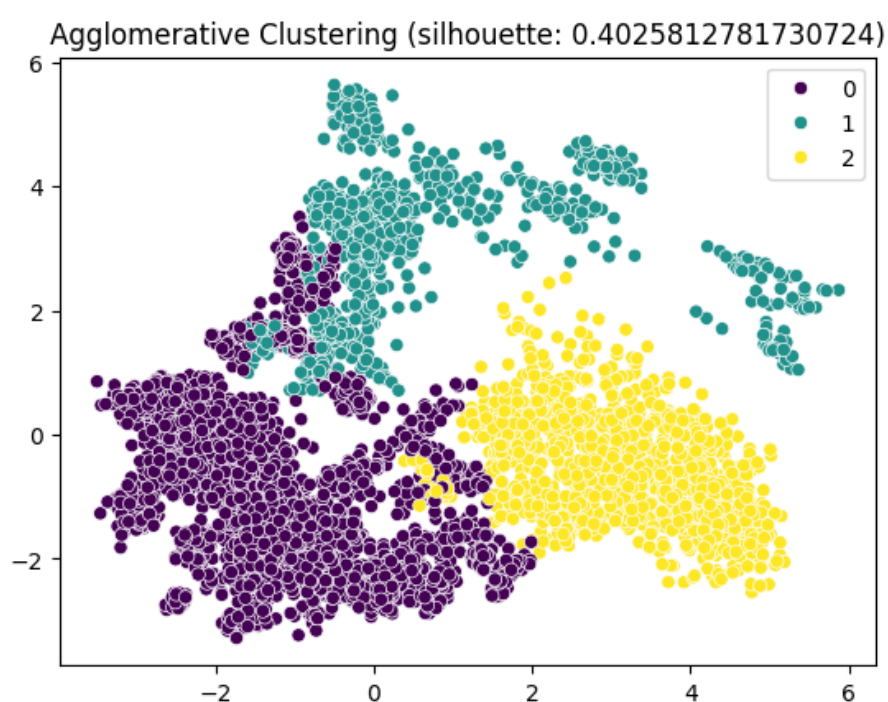
شکل ۳۰.۱: نمودار تغییرات معیار ارزیابی عملکرد خوشه‌بندی با افزایش تعداد دسته‌ها در الگوریتم سلسله‌مراتبی



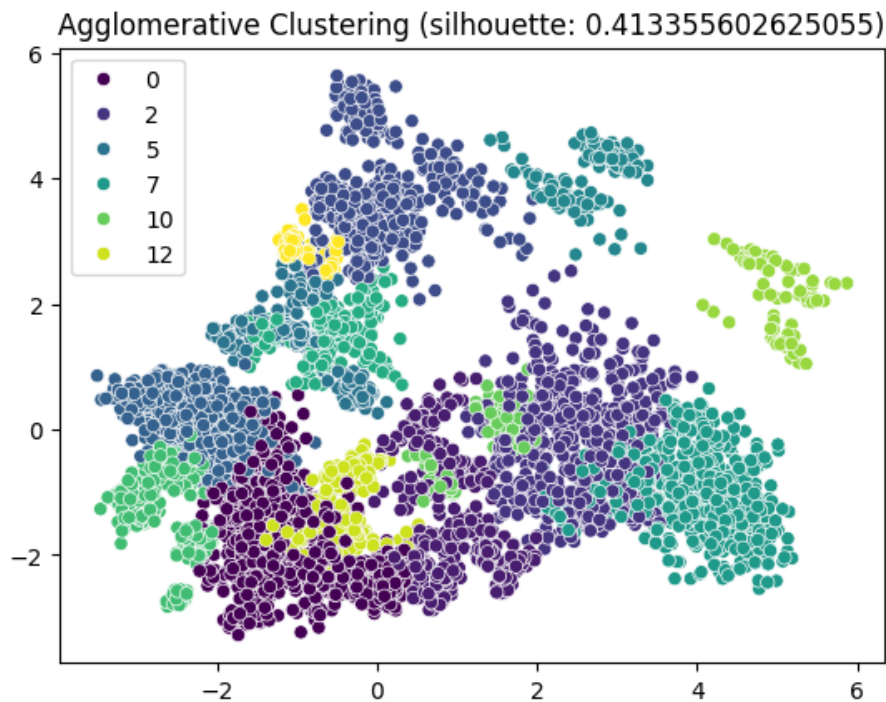
شکل ۳۱.۱: نمایش دندروگرام برای داده‌ها در کاهش بعد به ۳ با استفاده از روش PCA با توجه به نمودار دندروگرام و همچنین معیار Silhouette مقادیر ۴ و ۶ برای تعداد نمونه‌ها بهترین عملکرد را خواهند داشت. در این حالات فاصله بین دسته‌ها بیشترین حالت و هر دسته تا حد ممکن یکپارچه خواهد بود.



شکل ۳۲.۱: نمایش دسته‌های مختلف در الگوریتم سلسله مراتبی با تعداد دسته برابر با ۶. با توجه به وجود ۳ بعد برای خوشه بندی خوشه ی ۴ در نمایش تا دو تکه نشان داده شده است. دسته ۵ نیز با توجه به قرارگیری به صورت جدا از بقیه داده به عنوان یک خوشه در نظر گرفته شده است که بخش زیادی از گویندگان در این دسته خانوم هستند.



شکل ۳۳.۱: نمایش دسته‌های مختلف در الگوریتم سلسله مراتبی با تعداد دسته برابر با ۳. در این حالت با توجه به تعداد پایین دسته‌ها دو دسته آبی و بنفش به دارای همپوشانی بوده اند.



شکل ۳۴.۱: نمایش دسته‌های مختلف در الگوریتم سلسله مراتبی با تعداد دسته برابر با ۱۴. در این حالت تعداد دسته‌ها زیاد تر از حد نیاز شده است و تعداد اعضای هر دسته تا حدی نامتوازن شده است. هر چند همچنان الگوریتم توانسته تا حدی با توجه به توزیع جنسیت و دیگر کلاس‌ها مانند طول صوت خوشه بندی را انجام دهد.

فصل ۲

Automatic Speech Recognition(ASR)

۱.۲ مقدمه

تشخیص خودکار صحبت به فرایند تبدیل فایل صوتی از گفتار یک فرد و تبدیل آن به متن قابل خواندن و معادل با فایل صوتی اولیه گفته می‌شود. در اینجا تبدیل صوت به متن در سه مرحله پیش‌پردازش، آموزش و ارزیابی انجام شده است که در ادامه به توضیح هر کدام می‌پردازیم.

۲.۲ پیش‌پردازش

اولین مرحله از پیش‌پردازش یکسان‌سازی حروف موجود در متون مجموعه داده است. تعداد کاراکترهای موجود در مجموعه داده قبل از انجام هرگونه پردازشی، ۱۶۳ عدد بود. برای عملکرد بهتر مدل، لازم است کاراکترهای عربی را به فارسی تبدیل کنیم؛ به طور مثال داده «نتیجه تطهیرکننده سفرهای زیارتی» به «نتیجه ی تطهیرکننده ی سفرهای زیارتی» تبدیل می‌شود. همچنین بعضی از کاراکترها از جمله علائم نگارشی و کاراکترهای اضافه را حذف می‌کنیم. با انجام این پیش‌پردازش‌ها، تعداد کاراکترهای منحصر بفرد به ۴۷ می‌رسد. الفبا را برای استفاده‌های بعدی در یک فایل ذخیره می‌کنیم. از توکنایزر و استخراج‌کننده ویژگی کتابخانه Wav2Vec استفاده می‌کنیم و داده‌ها را با اعمال processor آماده می‌کنیم. پس از حذف داده‌های بیشتر از ۱۵ ثانیه، داده پردازش شده

را در درایو ذخیره می‌کنیم.

۳.۲ آموزش

ابتدا داده‌های پردازش شده را از درایو لود می‌کنیم. از مدل از پیش آموخته Vec2Wav استفاده می‌کنیم ولی قسمت feature_encoder آن را فریز می‌کنیم تا وزن‌های این بخش به روزرسانی نشود. حین آموزش، مقدار validation loss بی نهایت می‌شد که با اضافه کردن خط زیر این مشکل را برطرف کردیم:

`model.config.ctc_zero_infinity = True`

مدل را در ۲ اپیک، با نرخ آموزش e-4، و با اندازه batch ۱ آموزش می‌دهیم. فرایند آموزش ۳ ساعت به طول انجامید. در ابتدا مقدار خطای آموزش ۱۷/۰۷ و خطای ارزیابی ۱۱/۹۰ بود. WER معیار دیگریست که برای ارزیابی مدل مورد استفاده قرار گرفت. این معیار، اختلاف بین دنباله مرجع (متن اصلی صوت) و دنباله خروجی (پیش‌بینی شده) را به ازای تعداد کلمات اندازه‌گیری می‌کند. فرمول آن به صورت زیر است:

$$WER = \frac{(S + D + I)}{N} \quad (۱.۲)$$

S تعداد خطاهای جایگزینی

D تعداد خطاهای حذف

I تعداد خطاهای اضافه

N تعداد کلمات کل در مرجع

معیار WER از گام رو به کاهش رفت و در نهایت به مقدار ۰/۵۰۶ رسید. در جدول ۱.۲ به صورت خلاصه گزارش مقدار خطا و WER مشاهده می‌شود.

WER	Validation Loss	Training Loss	Step
1.000000	11.907510	17.074800	۲۰
1.000000	3.109851	5.678500	۴۰
...
0.999813	2.356913	2.608600	۴۲۰
0.991504	1.798767	2.243500	۴۴۰
...
0.521427	0.552796	0.696100	۹۲۰
0.506934	0.553852	0.747700	۹۴۰

جدول ۱.۲: معیارهای آموزش و ارزیابی

۴.۲ ارزیابی

در قسمت آخر به صورت تصادفی ۲۰ داده از داده‌های آزمون را انتخاب می‌کنیم و پیش‌بینی مدل را بدست می‌آوریم. در جدول ۲.۲، ۵ نمونه از آنها را مشاهده می‌کنیم. به نظر می‌رسد مدل می‌تواند تا حد خوبی خروجی را پیش‌بینی کند اما در برخی نمونه‌ها خروجی مدل با مرجع متفاوت است. دلایل احتمالی عبارتند از:

- ابتدا یا انتهای صوت به طور کامل ضبط نشده و به همین دلیل مدل کلمه را ناقص پیش‌بینی کرده (نمونه ۱).
- در جمع‌آوری مجموعه داده اشتباهاتی رخ داده و مرجع برخی از سطرهای داده با صوت منطبق نیست (نمونه ۲).
- لهجه فرد منجر به سخت‌تر شدن کار مدل شده و نتوانسته بعضی حروف را به درستی تخمین بزند (نمونه ۳، فرد با لهجه ترکی).
- در زبان فارسی تعداد زیادی از حروف آوای یکسان دارند ولی با کاراکترهای متفاوت نوشته می‌شوند. بزرگترین دلیل پیش‌بینی اشتباه مدل، این حروف هستند (نمونه ۴ و ۵).

شماره	پیش‌بینی	مرجع
۱	ن میتوانم با کداملاتی که به تور تصادفی انتخاب می‌کنم	من میتوانم با کلماتی که به طور تصادفی انتخاب می‌کنم
۲	هر چند بازار کارق برق در ایران رو به رشته است اما همچنان چالش‌هایی نیز در این حوزه وجود دارد	افرادی که بتوانند این چالش‌ها را مدیریت کنند در بازار کار برق در ایران موفقیت‌آمیز خواهند بود
۳	همیشه محربان خواهم ماند حتتا اگر کسی قدر مهربنیم را نداند	همیشه مهربان خواهم ماند حتی اگر کسی قدر مهربانیم را نداند
۴	آدم قروب خورشید را دوست می‌دارد	ادم غروب خورشید را دوست می‌دارد
۵	من صدای قدم خاهش را می‌شنوم	من صدای قدم خاهش را می‌شنوم

جدول ۲.۲: پیش‌بینی و مرجع برای چند نمونه

در آخر، بر روی کل داده آزمون، پیش‌بینی مدل را بدست می‌آوریم. سپس مقادیر ویژگی‌های جنسیت، لهجه و تن صدا را یکدست می‌کنیم. اکنون همبستگی میان هر یک از ویژگی‌ها با مقدار WER را حساب می‌کنیم. به طور کلی برای محاسبه همبستگی با یک ویژگی، WER را برای هر یک از مقادیر آن ویژگی حساب می‌کنیم و سپس از معیار Point-Biserial Correlation استفاده می‌کنیم. این ضریب زمانی استفاده می‌شود که یکی از متغیرها پیوسته و دیگری categorical باشد. مقدار آن از ۱- تا ۱ متغیر است. مقدار ۱- نشان‌دهنده همبستگی منفی کامل، ۱ نشان‌دهنده همبستگی مثبت کامل و ۰ نشان‌دهنده عدم همبستگی است. مقدار همبستگی‌ها در جدول ۳.۲ قابل مشاهده است.

برای جنسیت، مرد را به ۰ و زن را به ۱ تبدیل کردیم. مقدار WER برای مرد 0.44 و برای زن 0.39 است. WER با جنسیت رابطه عکس دارد؛ بدین معنی که با افزایش مقدار جنسیت، میزان خطا کاهش یافته. برای لهجه، فارسی را به ۰، شیرازی به ۱، ترکی به ۲، خراسانی به ۳ و یزدی به ۴ نگاشت کردیم. WER با لهجه رابطه مستقیم دارد؛ بدین معنی که با افزایش مقدار لهجه (دور شدن از لهجه فارسی)، میزان خطا افزایش یافته است. در مجموعه داده آزمون، ۱۱۲۴ عدد از دادگان لهجه فارسی دارند و در مجموع ۲۹ داده مربوط به ۴ لهجه دیگر هستند. چون مدل حین آموزش نیز از این دادگان بسیار کمتر از فارسی دیده، انتظار می‌رود در پیش‌بینی لهجه‌هایی غیر از فارسی خطای بیشتری داشته باشد و مقدار همبستگی نیز شاهی بر همین ادعاست.

برای تن صدا، نرمال را به ۰، سوالی به ۱، تعجبی به ۲، دستوری به ۳ و ناقص به ۴ نگاشت کردیم. WER با تن رابطه عکس دارد؛ بدین معنی که با افزایش مقدار تن، میزان خطا کاهش یافته.

ویژگی	همبستگی با WER
جنسیت	-1
لهجه	0.8
تن	-0.44

جدول ۳.۲: همبستگی ویژگی‌ها با WER