



University of Tehran

School of Electrical and Computer Engineering



MACHINE LEARNING

Final Project Report

Sasan Keshavarz

810199253

Mohammad Roudbari

810100364

Reza Souri

810100380

Summer 2022

۴ مقدمه
۴ جمع آوری داده
۴ تمیزسازی داده ها
۶ استخراج ویژگی
۷ Chroma STFT
۷ Chroma CQT
۸ Chroma CENS (Chroma Energy Normalized Statistics)
۸ Mel spectrogram
۸ MFCC
۹ Tonnetz (Tonal Centroid Features)
۹ Spectral Centroid
۱۰ Spectral Bandwidth
۱۰ Spectral Flatness
۱۰ Spectral Contrast
۱۱ Spectral Rolloff
۱۱ RMS
۱۲ Zero Crossing Rate
۱۲ Poly_Features
۱۲ پیش پردازش داده ها :
۱۶ طبقه بندی :
۱۶ طبقه بندی جنسیت :

۱۶	مدل MLP برای طبقه بندی جنسیت:
۲۰	مدل SVM برای طبقه بندی جنسیت:
۲۲	مدل Random Forest برای طبقه بندی جنسیت:
۲۵	طبقه بندی احساسات :
۲۵	مدل MLP برای طبقه بندی احساسات:
۲۸	مدل SVM:
۲۹	مدل Random Forest برای طبقه بندی احساسات :
۳۳	خوشه بندی
۳۳	روش k-means
۳۴	K=2 دو خوشه
۳۶	K=4 دو خوشه
۳۸	K=10 دو خوشه
۴۱	K=400 دو خوشه
۴۲	الگوریتم GMM
۴۳	K=2
۴۵	K=4
۴۷	K=10
۵۱	K=400

مقدمه

هدف نهایی این پروژه طبقه بندی و خوشه بندی صدا بر اساس احساسات و جنسیت آن است. مراحل مختلف پروژه جمع آوری داده‌ها، تحلیل و پیش پردازش آن‌ها و طبقه بندی و خوشه بندی است. ابتدا داده‌های مربوط به افراد متخلف با احساسات خواسته شده را جمع آوری کردیم. پس از تمیز سازی داده آن‌ها را طبقه بندی و خوشه بندی کردیم.

جمع آوری داده

در زمینه جمع آوری داده چالش‌هایی وجود داشت که ممتزین آن این بود که به علت کوتاه بودن عبارات انتقال احساس برای سوژه‌ها دشوار بود. یکی دیگر از چالش‌ها وجود نویز در محیط است که سعی شده کاملاً محیط آرام باشد که مشکل تمیز سازی داده‌ها در مراحل بعدی کمتر شود. جنسیت بر دامنه و فرکانس صدا تأثیر می‌گذارد و چون این‌ها ویژگی‌های اصلی تشخیص احساس هستند. به طور کلی صدای زنان زیر تراست و فرکانس بالاتری دارد اما مردها صدای بم تر با فرکانس پایین تری دارند. به طور میانگین دامنه ی صدای مردها بالاتر است. همچنین ویژگی‌های فردی هم موثر هستند. لحن بیان کردن جملات فارغ از جنسیت و سن، بسته به شرایط اجتماعی فرد و محل زندگی نیز تغییر می‌کند.

تمیزسازی داده‌ها

در ابتدای کار داده‌ها نویزی بودند و صدای ضبط شده فقط شامل صدای سوژه نبود در نتیجه باید داده‌ها را تمیز کردیم. برای این کار ابتدا صدا‌های موجود در فایل را تفکیک کردیم و صدا‌های اضافی را حذف کردیم. همچنین ابتدا و انتهای هر داده که سوژه در آن ساکت است باید بریده شد. بعضی از برچسب داده‌ها تکراری بود و این موارد اصلاح شدند.

در قدم اول تعداد داده‌های null در هر کلاس را شمارش می‌کنیم.

در قدم بعد correlation ستون هارا بدست آورده و میزان وابستگی هر feature به برچسب های جنسیت و احساس را بدست می آوریم .



همانطور که در شکل فوق مشخص است ستون جنسیت وابستگی بیشتری به ویژگی های chroma stft، mfcc، spectral centroid، نسبت به سایر ویژگی ها دارد و همچنین ستون احساسات وابستگی بیشتری به ویژگی های mfcc، poly features، melspectrogram ، rmse دارد .

استخراج ویژگی

هر سیگنال صوتی شامل ویژگی های بسیاری است. در نتیجه، باید مشخصه هایی در تشخیص جنسیت و عاطفه از صدا مرتبط هستند را استخراج کرد. فرآیند استخراج ویژگی برای استفاده از آن ها در تحلیل ها، استخراج ویژگی نامیده می شود. اکنون برخی از ویژگی های سیگنال های صوتی همراه با جزئیات مورد بررسی قرار می گیرد. در این مرحله سیگنال های گفتاری پردازش شده به یک نمایش مختصر اما منطقی که متمایزتر و قابل اعتمادتر از سیگنال واقعی است تبدیل می شوند.

پس از انتخاب ویژگی های مورد نظر با قطعه کد زیر از داده های صوتی ویژگی هارا استخراج کرده و از هر کدام میانگین می گیریم و دیتاست را تشکیل می دهیم .

```
header = 'voice_id chroma_stft rmse spectral_centroid spectral_bandwidth rolloff zero_crossing_rate chroma_cqt chroma_cens
|melspectrogram spectral_contrast spectral_flatness poly_features tonnetz tempogram'
for i in range(1, 21):
    header += f' mfcc{i}'
# header += ' label'
header = header.split()
```

```

file = open('voice_data.csv', 'w', newline='')
with file:
    writer = csv.writer(file)
    writer.writerow(header)

for filename in os.listdir(f'./drive/MyDrive/Voice'):
    songname = f'./drive/MyDrive/Voice/{filename}'
    y0, sr0 = librosa.load(songname, mono=True, duration=30)
    rmse = librosa.feature.rms(y=y0)
    chroma_stft = librosa.feature.chroma_stft(y=y0, sr=sr0)
    spec_cent = librosa.feature.spectral_centroid(y=y0, sr=sr0)
    spec_bw = librosa.feature.spectral_bandwidth(y=y0, sr=sr0)
    rolloff = librosa.feature.spectral_rolloff(y=y0, sr=sr0)
    zcr = librosa.feature.zero_crossing_rate(y0)
    mfcc = librosa.feature.mfcc(y=y0, sr=sr0)

    chroma_cqt = librosa.feature.chroma_cqt(y0, sr0)
    chroma_cens = librosa.feature.chroma_cens(y0, sr0)
    melspectrogram = librosa.feature.melspectrogram(y0, sr0)
    spectral_contrast = librosa.feature.spectral_contrast(y0, sr0)
    spectral_flatness = librosa.feature.spectral_flatness(y0)
    poly_features = librosa.feature.poly_features(y0, sr0)
    tonnetz = librosa.feature.tonnetz(y0, sr0)
    tempogram = librosa.feature.tempogram(y0, sr0)

    to_append = f'{filename} {np.mean(chroma_stft)} {np.mean(rmse)} {np.mean(spec_cent)} {np.mean(spec_bw)} {np.mean(rolloff)} {np.mean(zcr)} {np.mean(chroma_cqt)} {np.mean(chroma_cens)} {np.mean(melspectrogram)} {np.mean(spectral_contrast)} {np.mean(spectral_flatness)} {np.mean(poly_features)} {np.mean(tonnetz)} {np.mean(tempogram)}'

    for e in mfcc:
        to_append += f' {np.mean(e)}'

    file = open('voice_data.csv', 'a', newline='')
    with file:
        writer = csv.writer(file)
        writer.writerow(to_append.split())

data = pd.read_csv('voice_data.csv')
data.head()

```

در قطعه کد فوق ، در ابتدا نام های ویژگی هایی که استفاده کردیم را درون header قرار دادیم و در ادامه با استفاده از پکیج librosa ، ویژگی های مورد نظر را از دیتای صوتی استخراج کردیم .

هر کدام از ویژگی ها در خروجی یک بردار به ما می دهد که نیاز است از آن میانگین بگیریم . لازم به ذکر است در ابتدا تمام ویژگی هایی که در بخش spectral features پکیج librosa وجود دارد را از داده ها استخراج کردیم و طبق دستور العمل پروژه داده هارا به دو دسته ی تقسیم کرده و مدل های مختلف را روی آن ها آموزش می دهیم و نتایج بدست آمده را تحلیل و مقایسه می کنیم. در ادامه توضیح مختصری برای هر ویژگی ارائه شده است.

Chroma STFT

مقدار Chroma یک صدا اساساً نشان دهنده شدت دوازده کلاس صدای متمایز است که برای مطالعه موسیقی استفاده می شود. آنها را می توان در تمایز پروفیل های کلاس گام بین سیگنال های صوتی به کار برد. در واقع یک توصیفگر است که محتوای آهنگی سیگنال صوتی را به صورت فشرده نشان می دهد. بنابراین ویژگی های کروما را می توان به عنوان پیش نیاز مهم برای تحلیل معنایی سطح بالا، مانند تشخیص آکورد یا تخمین تشابه هارمونیک در نظر گرفت.

Chroma STFT مقدار Chroma صدا اساساً نشان دهنده شدت دوازده کلاس صدای متمایز است که برای مطالعه موسیقی استفاده می شود. آنها را می توان در تمایز پروفیل های کلاس گام بین سیگنال های صوتی استفاده کرد. در نتیجه هم برای تشخیص جنسیت و هم تشخیص عاطفه هم ویژگی مهمی است.

Chroma CQT

یکی دیگر از رویکردهای کروماگرام، تبدیل Chroma Constant-Q است که در آن STFT را به گونه ای تنظیم می کند که به صورت لگاریتمی فاصله بین bin های فرکانس ایجاد کند. با این حال، زمانی که فرکانس را در مقیاس لگاریتمی ترسیم می کنید، اندازه bin ثابت برای همه فرکانس ها منجر به مشکلاتی می شود. به دنبال حل این مشکل با افزایش اندازه بافر برای فرکانس های پایین تر و کاهش فشار محاسباتی ناشی از آن با کاهش اندازه بافر مورد استفاده برای فرکانس های بالا است.

CQT از دو جهت با STFT متفاوت است: (۱) مقیاس فرکانس لگاریتمی است نه خطی. (۲) طول پنجره برای هر فرکانس متفاوت است: فرکانس های پایین از پنجره های بلند استفاده می کنند، فرکانس های بالا از پنجره های کوتاه استفاده می کنند.

Chroma CENS (Chroma Energy Normalized Statistics)

توابع CENS را می توان به طور موثر به دلیل وضوح فضایی کم تفسیر کرد. ویژگی های Chroma فقط بر روی ۱۲ ویژگی صدای زیر و بمی متمرکز هستند که در نت موسیقی غربی آشنا هستند که در آن هر متغیر Chroma نشان می دهد که چگونه شدت صدا در میان ده ها باند Chroma پراکنده می شود. این توابع با در نظر گرفتن شاخصه های آمارهای کوتاه مدت بر روی توزیع انرژی در باندهای کروم، ویژگی های CENS (آمار نرمال شده انرژی کروم) را که خانواده ای از ویژگی های صوتی مقیاس پذیر و قوی را تشکیل می دهند، به دست می آورد.

ایده اصلی ویژگی های CENS این است که آمارگیری از پنجره های بزرگ انحرافات محلی در تمپو، بیان و تزیینات موسیقی مانند تریل ها و آکوردهای آرپجیت شده را هموار می کند. CENS به بهترین وجه برای کارهایی مانند تطبیق صدا و شباهت استفاده می شود.

Mel spectrogram

طیف نگار Mel ترکیبی از مقیاس Mel و طیف نگار است که در آن مقیاس Mel تبدیل غیرخطی مقیاس فرکانس را نشان می دهد. در این حالت ابتدا سیگنال صوتی به فریم های کوچکتر تقسیم می شود و یک پنجره

Hamming روی هر فریم اعمال می شود. سپس DFT برای تغییر از حوزه زمان به حوزه فرکانس اعمال می شود. یک لگاریتم در مرحله نهایی برای تولید طیف استفاده می شود که به تولید طیف نگار Mel کمک میکند.

طیف نگار mel مقادیر بر حسب هرتز را به مقیاس mel تبدیل می کند. طیف نگار صوتی خطی برای برنامه هایی که همه فرکانس ها اهمیت یکسانی دارند، ایده آل است، در حالی که طیف نگارهای mel برای برنامه هایی که نیاز به مدل سازی درک شنوایی انسان دارند، مناسب تر هستند.

طیف نگاری تجسم طیف فرکانس یک سیگنال است، که در آن طیف فرکانس سیگنال، محدوده فرکانسی است که توسط سیگنال موجود است. طیف نگار Mel طیفی است که به مقیاس Mel تبدیل می شود. مطالعات نشان داده است که انسان فرکانس ها را در مقیاس خطی درک نمی کند.

MFCC

Mel frequency cepstral coefficients از یک سیگنال، مجموعه ای کوچک از ویژگی ها (معمولاً

حدود ۱۰ الی ۲۰) است که به طور خلاصه شکل کلی یک طیف را توصیف می کند. MFCC ها با تبدیل کسینوس گسسته (DCT) به یک طیف نگار مل فرکانس محاسبه می شوند. این ویژگی مشخصه های صدای انسانی را مدل می کند. ضرایب مغزی فرکانس مل، نمایش های فشرده ای از طیف هستند که معمولاً برای شناسایی خودکار گفتار استفاده می شوند و همچنین به عنوان یک ویژگی اصلی در بسیاری از زمینه های تحقیقاتی که شامل سیگنال های صوتی می شود، استفاده می شود. همچنین می توان مقیاس دهی ویژگی ها را به گونه ای انجام داد که هر «بعد ضریب (coefficient dimension)» دارای میانگین صفر و واریانس واحد باشد.

می توان آن را با نگاشت سیگنال تبدیل فوریه بر روی مقیاس mel با استفاده از پنجره های مثلثی یا کسینوسی به دست آورد. جایی که پس از گرفتن لاگ توان ها در هر یک از فرکانس های Mel و پس از تبدیل کسینوس گسسته توان های مل log دامنه یک طیف را می دهد. لیست دامنه MFCC است.

Tonnetz (Tonal Centroid Features)

Tonnetz با تشخیص تغییرات هارمونیک در کلیپ های صوتی موسیقی، ویژگی های مرکز تونال را محاسبه می کند. این یک نمایش مسطح شناخته شده از روابط زیر و بمی در یک کلیپ صوتی است که به عنوان یک صفحه بی نهایت در نظر گرفته می شود. همانطور که در قسمت ویژگی chroma گفته شد، یک فایل صوتی می تواند شامل ۱۲ کلاس گام صدا باشد، این پدیده با ادغام برخی از کلاس ها با هم در نظر می گیرد که فایل

صوتی از ۶ کلاس گام است. بنابراین نمودار ویژگی کروما ترسیم شده از این روش زمین ها را فقط در ۶ کلاس طبقه بندی می کند که ما آن را Tonnetz می نامیم.

پروفیل های کلاس گام زمانی - که معمولاً به آنها کروماگرام گفته می شود - نمایش سیگنال استاندارد بالفعل برای روش های مبتنی بر محتوا تحلیل هارمونیک موسیقی هستند. داده های صوتی در فضای Tonnetz، با نمایشی هندسی فواصل زیر و بمی همسان مبتنی بر تئوری موسیقی، پخش می شوند. نشان داده شده است که این رویکرد از دقت طبقه بندی نمایش های کرومای قبلی بهتر عمل می کند، در حالی که فضای ویژگی سطح متوسطی را فراهم می کند که چالش های ذاتی کروما را دور می زند.

Spectral Centroid

این ویژگی نشان می دهد که مرکز جرم برای یک صدا در کجا قرار گرفته و به عنوان میانگین وزنی برای فرکانس های موجود در صدا محاسبه می شود. مثلاً بین دو اهنگ یکی از سبک بلوز و یکی از سبک متال، موسیقی بلوز در سراسر طولش یکسان است ولی موسیقی متال فرکانس های بیشتری تا پایان دارد. بنابراین، (Spectral Centroid) برای موسیقی بلوز جایی نزدیک اواسط طیف آن است، در حالیکه برای موسیقی متال در انتهای آن قرار دارد. دستور librosa.feature.spectral_centroid این ویژگی را برای هر فریم در سیگنال را محاسبه می کند.

در بعضی جاها می توان آن را میانه طیف در نظر گرفت اما بین اندازه گیری مرکز طیفی و میانه طیف تفاوت وجود دارد. مرکز طیفی مانند یک میانه وزنی و میانه طیف شبیه به میانگین است. هر دوی آنها تمایل مرکزی سیگنال را اندازه گیری می کنند. در برخی موارد، هر دوی آنها نتایج مشابهی را نشان می دهند.

Spectral Bandwidth

پهنای باند تفاوت بین فرکانس های بالا و پایین در یک باند فرکانس پیوسته است. همانطور که می دانیم سیگنال ها حول یک نقطه نوسان می کنند، بنابراین اگر نقطه مرکز سیگنال باشد، مجموع حداکثر انحراف سیگنال در دو طرف آن نقطه را می توان به عنوان پهنای باند سیگنال در آن فریم زمانی در نظر گرفت.

این مفهوم از مرکز طیفی مشتق شده است. این محدوده طیفی مورد توجه در اطراف مرکز است، یعنی واریانس از مرکز طیفی. همبستگی مستقیمی با صدای درک شده دارد. پهنای باند مستقیماً با انرژی پخش شده در باندهای فرکانسی متناسب است. از نظر ریاضی، میانگین وزنی فواصل باندهای فرکانسی از مرکز طیفی است.

Spectral Flatness

مسطح بودن طیفی که همچنین به عنوان آنتروپی وینر شناخته می شود، معیاری است که در پردازش سیگنال دیجیتال برای مشخص کردن طیف صوتی استفاده می شود. این ویژگی، یکنواختی توزیع انرژی سیگنال در حوزه فرکانس را تخمین می زند. صافی طیفی معمولاً بر حسب دسی بل اندازه گیری می شود و روشی را برای تعیین کمیت یک صدا در مقایسه با نویز مانند بودن، ارائه می دهد.

مسطح بودن طیفی بالا (نزدیک به ۱,۰ برای نویز سفید) نشان می دهد که طیف در همه باندهای طیفی دارای مقدار مشابهی از قدرت است، این صدا شبیه نویز سفید است و نمودار طیف نسبتاً صاف و صاف به نظر می رسد و به معنی یکنواختی بیشتر در توزیع انرژی سیگنال در حوزه فرکانس است. در حالی که آنتروپی وینر پایین دلالت بر یکنواختی کمتر دارد. صافی طیفی کم (نزدیک به ۰,۰ برای یک تن خالص) نشان می دهد که توان طیفی در تعداد نسبتاً کمی از باندها متمرکز است.

این معیار را می توان برای تمایز یک سیگنال باند باریک از سیگنال باند پهن، به عنوان مثال، تمایز بین یک صدا و نویز سفید استفاده کرد. اما نمی توان از آنها برای تشخیص دو سیگنال پهن باند، به عنوان مثال، سیگنال LFM و نویز استفاده کرد.

Spectral Contrast

در یک سیگنال صوتی، کنتراست طیفی اندازه گیری انرژی فرکانس در هر پنجره زمانی است. از آنجایی که اکثر فایل های صوتی حاوی فرکانس هایی هستند که انرژی آن با گذشت زمان در حال تغییر است. اندازه گیری سطح انرژی دشوار می شود. کنتراست طیفی راهی برای اندازه گیری این نوع انرژی است. کنتراست طیفی حداکثر طیفی، حداقل طیفی و نابرابری بین هر فرکانس زیر باند را در نظر می گیرد. مقادیر کنتراست بالا به طور کلی با سیگنال های باند باریک و شفاف مطابقت دارد، در حالی که مقادیر کنتراست پایین مربوط به نویز باند وسیع است. کنتراست انرژی با مقایسه انرژی متوسط در قاب انرژی اوج با قاب انرژی پایین یا دره اندازه گیری می شود.

Spectral Rolloff

این ویژگی پهنای باند سیگنال صوتی را با تعیین یک bin فرکانسی که در آن درصد معینی از کل انرژی وجود دارد، اندازه گیری می کند. در واقع فرکانس roll-off به عنوان فرکانسی تعریف می شود که در آن درصدی (قطع) از کل انرژی طیف قرار می گیرد. برای مثال ۸۰٪ یا هر درصد دلخواه دیگر. فرکانس رول آف را می توان برای تمایز بین صداهای هارمونیک (زیر رول آف) و نویزدار (بالای رول آف) استفاده کرد.

می توان آن را به عنوان عملکرد نوع خاصی از فیلتر تعریف کرد که برای کاهش فرکانس های خارج از محدوده خاص طراحی شده است. دلیل اینکه ما آن را رول آف می نامیم این است که یک رویه تدریجی است. دو نوع فیلتر وجود دارد: Hi-pass و Low Pass و هر دو می توانند فرکانس را از سیگنالی که خارج از محدوده آنها قرار می گیرد حذف کنند.

نقطه rolloff طیفی برای تمایز بین گفتار صدادر و بدون صدا، تمایز گفتار/موسیقی، طبقه بندی ژانر موسیقی، تشخیص صحنه آکوستیک و طبقه بندی حالت موسیقی استفاده شده است.

RMS

مقدار (Root-Mean-Square) همان مقدار مؤثر شکل موج کل است. این برابر با سطح سیگنال DC است که همان توان متوسط سیگنال دوره ای را ارائه می دهد. چون مربوط به ریشه میانگین مربعات است پس با RMS بالاتر به طور کلی به صدای قوی تر است. که برای تشخیص صدای مردان و همچنین تشخیص عاطفه عصبانیت میتواند مفید باشد.

این اندازه گیری عمداً یک اندازه گیری آهسته است که میانگین قله ها و پایین ترین زمان ها را برای انعکاس بلندی درک شده را نشان می دهد. RMS معیار تقریبی روشی است که گوش شما سطوح صدا را درک می کند. مثلاً گوش شما معمولاً قله های تیز را به اندازه واقعی بلند نمی بیند.

Zero Crossing Rate

نرخ تغییر علامت ها در طول یک سیگنال و در واقع نرخ است که در آن سیگنال از مثبت به منفی یا بالعکس تغییر می کند. این ویژگی به طور سنگین هم در «بازشناسی گفتار (Speech Recognition)» و هم در «بازیابی اطلاعات موسیقی (Music Information Retrieval)» مورد استفاده قرار می گیرد. ویژگی مذکور معمولاً دارای مقداری بیشتر برای صداهای بسیار کوبه ای است، مثلاً در هنگام عصبانیت و خوشحالی احتمالاً این ویژگی بیشتر از عواطف غم و خنثی باشد. همچنین می توان از آن برای الگوریتم های تشخیص زمین و همچنین برای تشخیص فعالیت های صوتی استفاده کرد.

Poly_Features

ضرایب برازش یک چند جمله ای مرتبه n را به ستون های یک طیف نگار برمی گرداند. این ویژگی را هم می توان به راحتی با استفاده از Librosa استخراج کرد. ضرایب چند جمله ای برای هر فریم را محاسبه میکند. به این ترتیب که اول ضرایب $[0]$ مربوط به بالاترین درجه، ضرایب $[1]$ مربوط به بالاترین درجه بعدی تا ضرایب

ثابت را به عنوان خروجی برمیگرداند. سپس می توان میانگینی از ضرایب را اعمال کرد یا آماری را برای توصیف سری خروجی به دست آورد، در غیر این صورت مقادیر بیش از حد داده می شود.

پیش پردازش داده ها :

پس از مرحله ی استخراج ویژگی ها و بدست آوردن دیتا ست اولیه از داده ها باید داده ها را پیش پردازش کنیم و به ازای هر دیتا یک برچسب به آن نسبت دهیم.

دیتا ست اولیه داده ها بعد از استخراج ویژگی ها را Voice_data نام گذاری می کنیم که به صورت زیر می باشد:

	voice_id	chroma_stft	rms	spectral_centroid	spectral_bandwidth	rolloff	zero_crossing_rate	mfcc1	mfcc2	mfcc3	...	mfcc11	mfcc12	mfcc13	mfcc14	mfcc15	w
0	1000.wav	0.399225	0.048800	2329.831190	2413.742288	4077.076802	0.118935	-304.132566	65.614784	-0.230034	...	-4.279748	-2.955371	-14.975864	-1.610190	-14.000487	-11.6
1	1000.wav	0.299051	0.041566	2258.563870	2362.174317	4477.544056	0.108119	-306.777816	77.251091	-1.067098	...	-8.351185	-2.557911	-19.813825	2.116489	-12.450156	-10.0
2	1101.wav	0.337110	0.027311	1628.588396	1480.443570	1579.451886	0.049447	-478.123016	82.828812	28.429399	...	-12.248988	-0.775021	-3.362810	-1.868957	-7.215577	-6.7
3	1101.wav	0.428910	0.011069	2284.009634	2696.274514	5271.647135	0.091834	-533.723877	89.416367	27.698174	...	-2.108326	-4.852647	-7.825149	-1.205485	-1.209085	-8.8
4	1101.wav	0.390582	0.016276	2322.139983	2653.368293	5101.143484	0.105826	-480.891223	67.826586	16.232029	...	-2.383848	-1.810100	-0.273352	-1.016275	-6.315440	-6.7
...
17152	10133.wav	0.478882	0.010091	3489.003144	2628.625417	5508.562122	0.092584	-511.809542	79.757819	-4.267803	...	-7.007113	-7.512409	-2.180201	-2.588432	-5.206617	-1.3
17153	10134.wav	0.522728	0.001647	2777.438593	2724.228810	5081.070867	0.151424	-584.925242	79.813018	2.097298	...	-6.757550	-5.675238	-6.795041	1.310687	-3.112588	0.3
17154	10135.wav	0.544817	0.001786	2876.029820	2734.382038	6098.184877	0.185784	-582.591850	75.084543	1.962888	...	-4.828536	-3.738486	-6.230413	-2.142874	-3.748188	-0.8
17155	10136.wav	0.333702	0.027438	2071.005225	1699.756761	4959.060265	0.112778	-478.583740	61.517368	8.261489	...	-11.922656	-2.483714	-6.218111	4.417655	-7.250745	-3.7
17156	10137.wav	0.309188	0.037861	1854.365253	1830.258652	3396.802314	0.078250	-433.368897	69.326766	3.264902	...	-11.350521	-1.512034	-10.716776	-5.300368	-5.786279	-8.8

17157 rows x 27 columns

ستون voice id به صورت String می باشد . به همین منظور باید ابتدا پسوند wav. را از دیتا های این ستون حذف کنیم و آن ها را به مقادیر integer تبدیل کنیم .

در قطعه کد موارد فوق انجام شده است و دیتا برحسب ستون voice id به صورت صعودی مرتب شده است .

```
for i in range(0,17157):
    voice_data['voice_id'][i] = voice_data['voice_id'][i].replace('.wav', '')

voice_data['voice_id']=voice_data['voice_id'].astype(str).astype(int)

voice_data=voice_data.sort_values('voice_id')
```

در قدم بعد دیتاست متشکل از برچسب های داده ها را import می کنیم (datav2) که به صورت زیر می باشد:

	emotionID	textID	sex	age	voice id
0	1	1	m	21	15997
1	1	2	m	21	16001
2	1	3	m	21	16005
3	1	4	m	21	16009
4	1	5	m	21	16013
...
16860	4	6	f	54	10563
16861	4	7	f	54	10567
16862	4	8	f	54	10571
16863	4	9	f	54	10575
16864	4	10	f	54	10539

16865 rows x 5 columns

همانطور که مشاهده می کنید تعداد دیتا های موجود در دیتا ست برچسب ها ۱۶۸۶۵ و دیتاست داده ها ۱۷۱۵۷ می باشد ، در نتیجه دیتا باید پیش پردازش شوند و دیتاهای نادرست حذف شوند.

در قدم اول ، سطرهایی که در دیتاست متشکل از برچسب های داده ها، دارای voice id یکسان هستند را ، تشخیص داده و به صورت زیر حذف می کنیم .

```
same_id = []
for i in range(0,16865):
    for j in range(1,16865):
        if (dataV2['voice id'][i] == dataV2['voice id'][j]) & (i != j) :
            same_id.append(i)

# same_id
dataV2.drop(same_id , axis=0, inplace=True)
```

در قدم بعد ، متوجه شدیم که تعدادی از سطر های دیتاست داده ها و برچسب ها دارای voice id هستند که در دیتاست دیگر وجود ندارد ، در نتیجه با مقایسه ستون voice id دیتاست داده ها را با دیتاست برچسب ها ان ها را حذف می کنیم.

```
for j in voice_data['voice_id'] :
    if not ((dataV2.voiceId == j).any()) :
        voice_data.drop(voice_data[voice_data['voice_id']== j].index, axis=0, inplace=True)

for j in dataV2['voiceId'] :
    if not ((voice_data.voice_id == j).any()) :
        dataV2.drop(dataV2[dataV2['voiceId']== j].index, axis=0, inplace=True)
```

تعدادی از داده های موجود در دیتاست برچسب ها ، چندین بار تکرار شده اند و نیاز است از سطر هایی که بیش از یک بار تکرار شده اند یکی از ان ها نگه داشته شود و باقی ان ها حذف شوند. برای این منظور به صورت زیر عمل می کنیم .

```
for j in datav2['voiceId'] :|
    if datav2['voiceId'].value_counts()[j] > 1 :
        for i in range(datav2['voiceId'].value_counts()[j]-1) :
            datav2.drop(datav2[datav2['voiceId']== j].index[i], axis=0, inplace=True)
```

پس از اتمام تمیز سازی داده ها نیاز است دیتاست داده ها و برچسب ها را یکی کنیم . برای این منظور دیتاست هارا بر حسب ستون voice id مرتب کرده و سپس دو دیتاست را به هم پیوند می زنیم .

```
voice_data.reset_index(inplace=True)
voice_data=voice_data.drop(['index'],axis=1)
```

```
datav2.reset_index(inplace=True)
datav2=datav2.drop(['index'],axis=1)
```

```
Data_final = pd.concat([voice_data, datav2], axis=1, join='inner')
Data_final=Data_final.drop(['voiceId'],axis=1)
```

می دانیم ستون sex باید متشکل از دو نوع String باشد ، یکی 'f' برای داده های دارای جنسیت زن و دیگری 'm' برای داده های دارای جنسیت مرد . اما این ستون دارای عناصر دیگری به جز f , m می باشد که نیاز است اصلاح شود . همچنین با استفاده از کتابخانه ی LabelEncoder به ازای f ، ۰ و به ازای m ، ۱ قرار می دهیم. برای این منظور به صورت زیر عمل می کنیم.

```
Data_final['sex'].unique()
array(['f', 'm', 'H', 'F', 'w', 'f '], dtype=object)
```

```
for i in range (0 , 16784 ) :
    if (Data_final['sex'][i] == 'f') :
        Data_final['sex'][i] = 'f'
    if (Data_final['sex'][i] == 'f ') :
        Data_final['sex'][i] = 'f'
    if (Data_final['sex'][i] == 'F') :
        Data_final['sex'][i] = 'f'
    if (Data_final['sex'][i] == 'M') :
        Data_final['sex'][i] = 'm'
    if (Data_final['sex'][i] == 'm') :
        Data_final['sex'][i] = 'm'
    if (Data_final['sex'][i] == 'w') :
        Data_final['sex'][i] = 'm'

LE=LabelEncoder()
Data_final['sex']=LE.fit_transform(Data_final['sex'])
```

ستون emotionID نیز دارای عناصری غیر از ۱و۲و۳و۴ می باشد. برای این منظور دیتا هایی که دارای emotionID غیر از موارد ذکر شده هستند را به صورت زیر حذف می کنیم.

```
for j in [5,6,7,8,9,10] :
    if ((Data_final.emotionID == j).any()) :
        Data_final.drop(Data_final[Data_final['emotionID']== j].index, axis=0, inplace=True)
```

در نتیجه دیتاست نهایی ما پس از مراحل پیش پردازش و تمیز سازی داده ها به صورت زیر خواهد.

[illegible]

طبقه بندی :

طبق دستور عمل ویژگی های استخراج شده را به دو دسته تقسیم میکنیم. سعی شده است ویژگی های شبیه به هم در یک دسته قرار بگیرند.

دسته ویژگی های اول :

ویژگی‌هایی که در دسته‌ی اول استفاده شده‌اند عبارت‌اند از :

0	chroma_stft	16664	non-null	float64
1	rmse	16664	non-null	float64
2	zero_crossing_rate	16664	non-null	float64
3	poly_features	16664	non-null	float64
4	chroma_cqt	16664	non-null	float64
5	chroma_cens	16664	non-null	float64

دستہ ویژگی، های دوم :

ویژگی هایی که در دسته ی دوم استفاده شده اند عبارت اند از :

0	spectral_centroid	10664	non-null	float64
1	spectral_bandwidth	10664	non-null	float64
2	rolloff	10664	non-null	float64
3	mel_spectrogram	10664	non-null	float64
4	spectral_contrast	10664	non-null	float64
5	spectral_flatness	10664	non-null	float64
6	mfcc1	10664	non-null	float64
7	mfcc2	10664	non-null	float64
8	mfcc3	10664	non-null	float64
9	mfcc4	10664	non-null	float64
10	mfcc5	10664	non-null	float64
11	mfcc6	10664	non-null	float64
12	mfcc7	10664	non-null	float64
13	mfcc8	10664	non-null	float64
14	mfcc9	10664	non-null	float64
15	mfcc10	10664	non-null	float64
16	mfcc11	10664	non-null	float64
17	mfcc12	10664	non-null	float64
18	mfcc13	10664	non-null	float64
19	mfcc14	10664	non-null	float64
20	mfcc15	10664	non-null	float64
21	mfcc16	10664	non-null	float64
22	mfcc17	10664	non-null	float64
23	mfcc18	10664	non-null	float64
24	mfcc19	10664	non-null	float64
25	mfcc20	10664	non-null	float64

طبقه بندی جنسیت :

برای این منظور در ابتدا دیتاست y , X تعیین کرده و با استفاده از Standard Scaler ، داده ها را نرمالایز می کنیم. در قدم بعد داده را به دو دسته ی train و test تقسیم کرده و ۳۰ درصد داده ها را به تست اختصاص می دهیم.

مدل MLP برای طبقه بندی جنسیت:

دسته ای از شبکه های عصبی مصنوعی پیشخور است. یک MLP شامل حداقل سه لایه گره است: یک لایه ورودی، یک لایه پنهان و یک لایه خروجی. به جز گره های ورودی، هر گره یک نورون است که از یک تابع فعال سازی غیر خطی استفاده می کند. MLP از تکنیک یادگیری نظارت شده به نام بازپرداخت برای آموزش استفاده می کند. لایه های متعدد آن و فعال سازی غیر خطی آن MLP را از یک پرسپترون خطی متمایز می کند. در واقع می تواند داده هایی را متمایز کند که به صورت خطی قابل تفکیک نیستند. این شبکه ها شامل سه یا تعداد بیشتری از لایه ها است که از گره های غیرخطی فعال کننده هستند. از آنجا که MLP ها به طور کامل متصل شده اند، هر گره در یک لایه با وزن مشخص W_{ij} در هر گره در لایه بعدی متصل می شود.

نودهای شبکه عصبی که به آن ها نورون نیز گفته می شود، واحدهای محاسباتی در یک شبکه عصبی محسوب می شوند. در این شبکه عصبی، از خروجی های لایه اول (ورودی)، به عنوان ورودی های لایه بعدی (پنهان) استفاده می شود؛ این کار به همین شکل ادامه پیدا می کند، تا زمانی که، پس از تعداد خاصی از لایه ها، خروجی های آخرین لایه پنهان به عنوان ورودی های لایه خروجی مورد استفاده قرار می گیرد. به لایه هایی که بین لایه ورودی و لایه خروجی قرار می گیرند، لایه های پنهان گفته می شود. شبکه های پرسپترون چند لایه، مانند شبکه های عصبی پرسپترون تک لایه، حاوی مجموعه ای از وزن ها نیز هستند که باید برای آموزش و یادگیری شبکه عصبی تنظیم شوند.

مدل طراحی شده برای شبکه های عصبی mlp به صورت زیر می باشد . این مدل از ۴ لایه تشکیل شده است که تابع فعال ساز سه لایه اول relu و لایه اخر softmax انتخاب شده است . تعداد نورون های در نظر گرفته شده برای هر لایه در قطعه کد زیر نمایش داده شده است .

```
model = Sequential()
model.add(layers.Dense(256, activation='relu', input_shape=(X_train.shape[1],)))
model.add(layers.Dense(128, activation='relu'))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(2, activation='softmax'))
model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

classifier = model.fit(X_train,
                      y_train,
                      epochs=20,
                      batch_size=128, validation_split = 0.2)
```

اعمال مدل MLP بر روی دسته ویژگی های اول :

دقت و خطای بدست آمده برای داده های train و validation در 5 epochs اخر به شرح زیر است.

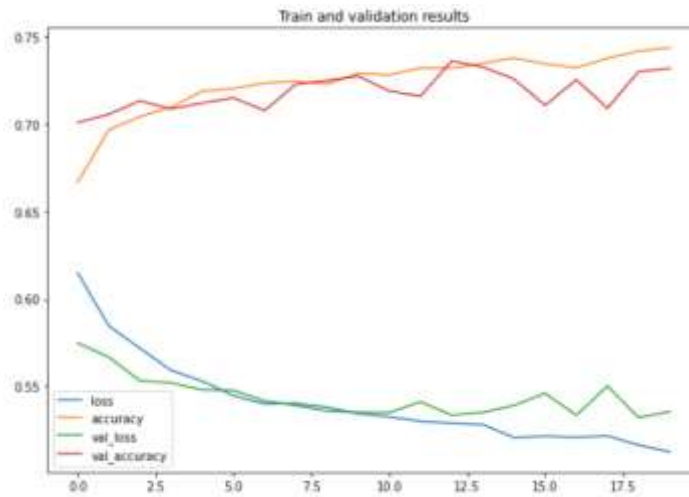
```
Epoch 16/20
73/73 [=====] - 0s 5ms/step - loss: 0.5215 - accuracy: 0.7343 - val_loss: 0.5461 - val_accuracy: 0.7107
Epoch 17/20
73/73 [=====] - 0s 5ms/step - loss: 0.5210 - accuracy: 0.7322 - val_loss: 0.5335 - val_accuracy: 0.7252
Epoch 18/20
73/73 [=====] - 0s 5ms/step - loss: 0.5216 - accuracy: 0.7375 - val_loss: 0.5502 - val_accuracy: 0.7090
Epoch 19/20
73/73 [=====] - 0s 5ms/step - loss: 0.5165 - accuracy: 0.7416 - val_loss: 0.5321 - val_accuracy: 0.7300
Epoch 20/20
73/73 [=====] - 0s 5ms/step - loss: 0.5126 - accuracy: 0.7435 - val_loss: 0.5356 - val_accuracy: 0.7317
```

دقت و خطای مدل روی داده های تست برابر است با :

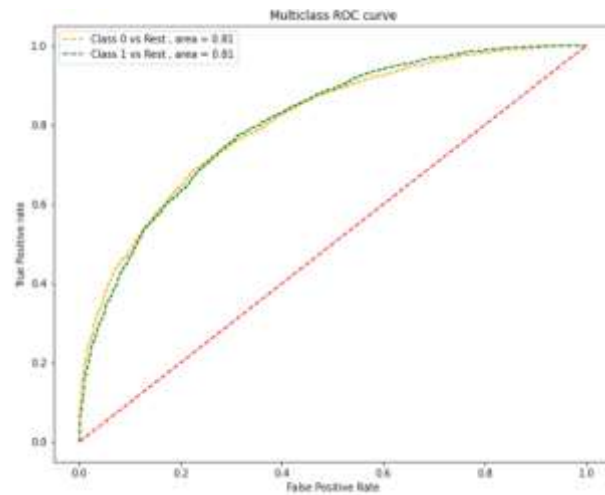
```
results = model.evaluate(X_test, y_test, verbose = 0)
print('test loss, test accuracy:', results)

test loss, test accuracy: [0.5321431756019592, 0.7310000061988831]
```

نمودار دقت و خطا مدل :



نمودار ROC مدل :



اعمال مدل MLP بر روی دسته ویژگی های دوم :

دقت و خطای بدست آمده برای داده های train و validation در 5 epochs آخر به شرح زیر است.

```

Epoch 16/20
73/73 [=====] - 0s 5ms/step - loss: 0.0107 - accuracy: 0.9886 - val_loss: 0.1840 - val_accuracy: 0.9477
Epoch 17/20
73/73 [=====] - 0s 5ms/step - loss: 0.0097 - accuracy: 0.9984 - val_loss: 0.1548 - val_accuracy: 0.9559
Epoch 18/20
73/73 [=====] - 0s 5ms/step - loss: 0.0046 - accuracy: 0.9998 - val_loss: 0.1680 - val_accuracy: 0.9546
Epoch 19/20
73/73 [=====] - 0s 5ms/step - loss: 0.0035 - accuracy: 0.9999 - val_loss: 0.1678 - val_accuracy: 0.9567
Epoch 20/20
73/73 [=====] - 0s 5ms/step - loss: 0.0024 - accuracy: 1.0000 - val_loss: 0.1726 - val_accuracy: 0.9559

```

دقت و خطای مدل روی داده های تست برابر است با :

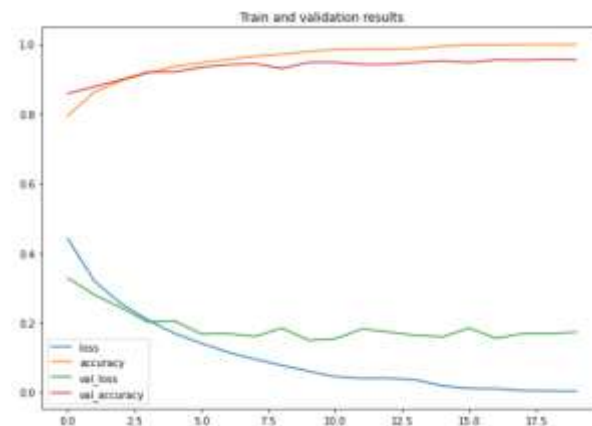
```

results = model.evaluate(X_test, y_test, verbose = 0)
print('test loss, test accuracy:', results)

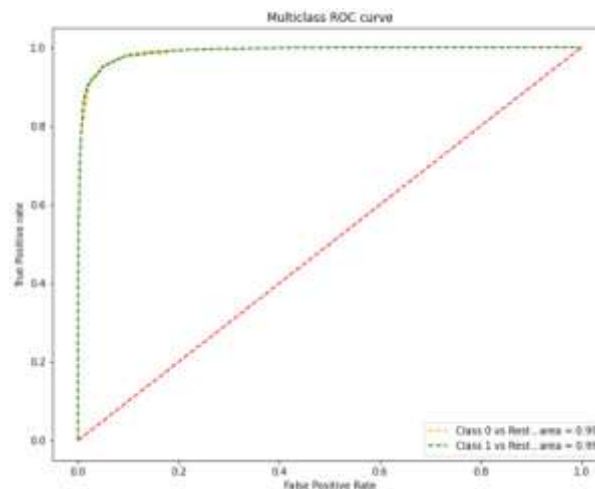
test loss, test accuracy: [0.17778711020946585, 0.9491999745368958]

```

نمودار دقت و خطا مدل :



نمودار ROC مدل :



تحلیل نتایج بدست آمده :

دقت مدل روی دسته اول برای داده آموزش ۷۴٪ و برای داده تست ۷۳٪ بوده است. پس بیش برآزش اتفاق نیفتاده است. دقت مدل روی دسته دوم به نحو چشمگیری افزایش افتاده است. روی دسته دوم دقت مدل برای داده آموزش ۱۰۰٪ و برای داده تست ۹۵٪ به دست آمده است. پس دسته ویژگی های دوم که شامل mfcc و چند ویژگی spectral بوده است (منظور از ویژگی spectral ویژگی های spectral rolloff, spectral centroid, spectral bandwidth و در کل ویژگی هایی است که با کلمه spectral شروع میشوند) دقت بسیار بهتری به دست داده است. دسته اول ویژگی ها شامل ZCR, RMSE, poly_features و chroma ها است دقت کمتری داده است. علت این اس که داده های دسته اول مشخصه های طیفی و زیر و بم را بسیار بهتر نمایان میکنند. همچنین ویژگی های mfcc مدل بسیار نزدیکی به صدای انسان تولید میکنند که برای آموزش داده ها بهتر عمل کرده اند. در دسته دوم بیشتر ویژگی ها از نوع زمانی و معمولی بوده اند و برای تفکیک داده صوتی کافی نبوده اند. همان عملکرد نسبتاً خوب دسته ویژگی دوم هم مربوط به chroma است که گام صدا را استخراج میکند.

مدل SVM برای طبقه بندی جنسیت:

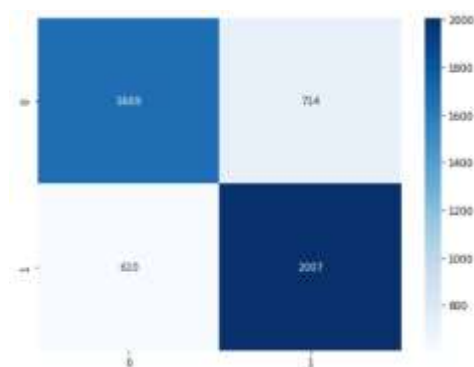
این روش از جمله روش های نسبتاً جدیدی است که در سال های اخیر کارایی خوبی نسبت به روش های قدیمی تر برای طبقه بندی نشان داده است. مبنای کاری دسته بندی کننده ی SVM دسته بندی خطی داده ها است و در تقسیم خطی داده ها سعی می کنیم خطی را انتخاب کنیم که حاشیه اطمینان بیشتری داشته باشد. حل معادله پیدا کردن خط بهینه برای داده ها به وسیله روش های QP که روش های شناخته شده ای در حل مسائل محدودیت دار هستند صورت می گیرد. قبل از تقسیم خطی برای اینکه ماشین بتواند داده های با پیچیدگی بالا را دسته بندی کند داده ها را به وسیله ی تابع phi به فضای با ابعاد خیلی بالاتر می بریم. برای اینکه بتوانیم مسئله ابعاد خیلی بالا را با استفاده از این روش ها حل کنیم از قضیه دوگانی لاگرانژ برای تبدیل مسئله ی مینیمم سازی مورد نظر به فرم دوگانی آن که در آن به جای تابع پیچیده ی phi که ما را به فضایی با ابعاد بالا می برد، تابع ساده تری به نام تابع هسته که ضرب برداری تابع phi است ظاهر می شود استفاده می کنیم. از توابع هسته مختلفی از جمله هسته های نمایی، چند جمله ای و سیگموئید می توان استفاده نمود.

```
classifier = SVC(kernel='rbf',C=150)
classifier.fit(X_train,y_train)
y_pred_test = classifier.predict(X_test)
```

اعمال مدل SVM بر روی دسته ویژگی های اول :

مقادیر دقت روی داده های تست و آموزش و f1 score و jaccard score و همچنین ماتریس confusion به شرح زیر می باشد.

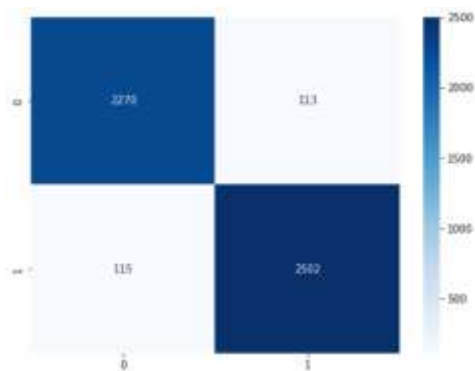
Train_accuracy : 76.8004%
Test_accuracy : 73.5200%
f1_score : 73.398%
jaccard_score : 58.008%



اعمال مدل SVM بر روی دسته ویژگی های دوم :

مقادیر دقت روی داده های تست و آموزش و f1 score و jaccard score و همچنین ماتریس confusion به شرح زیر می باشد.

Train_accuracy : 99.9914%
Test_accuracy : 95.4400%
f1_score : 95.430%
jaccard_score : 91.261%



تحلیل نتایج بدست آمده SVM برای جنسیت روی دو دسته ویژگی :

دقت مدل روی دسته اول برای داده آموزش ۷۶/۸٪ و برای داده تست ۷۳/۵٪ بوده است. پس بیش برآزش اتفاق نیفتاده است. دقت مدل روی دسته دوم به نحو چشمگیری افزایش افتاده است. روی دسته دوم دقت مدل برای داده آموزش ۱۰۰٪ و برای داده تست ۹۵/۴٪ به دست آمده است. برای دسته دوم ویژگی هم بیش برآزش اتفاق نیفتاده است. برای مدل SVM هم مانند مدل قبل دسته دوم بسیار بهتر عمل میکند. دلایل مشابه قبل است.

مدل Random Forest برای طبقه بندی جنسیت:

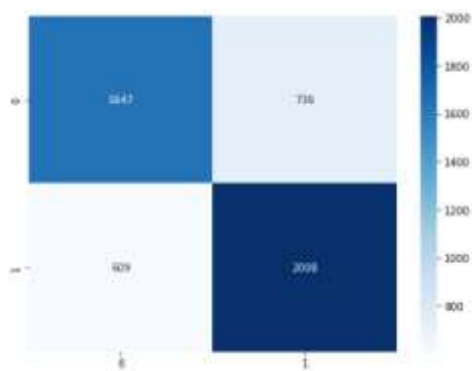
جنگل تصادفی یک روش یادگیری آنسمل است که برای دسته بندی و رگرسیون به کار میرود. این الگوریتم بر اساس ساختاری متشکل از تعداد زیادی درخت تصمیم، مدل مورد نظر به صورت زیر می باشد .

```
clf = RandomForestClassifier(n_estimators=1000 , random_state=42)
clf.fit(X_train,y_train)
y_pred_test=clf.predict(X_test)
```

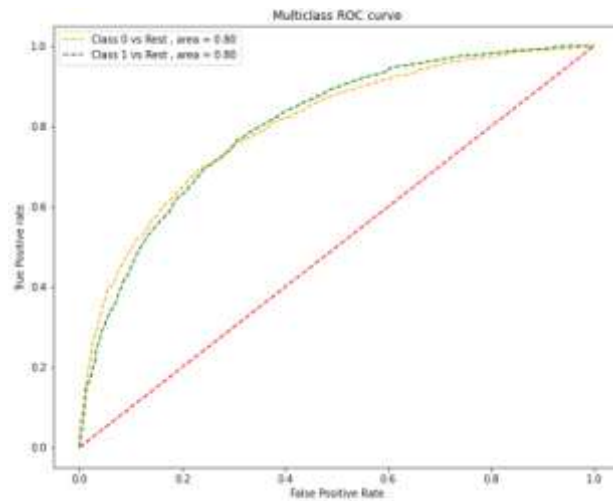
مقادیر دقت روی داده های تست و آموزش و f1 score و jaccard score و همچنین ماتریس confusion به شرح زیر می باشد.

اعمال مدل Random Forest بر روی دسته ویژگی های اول :

```
Train_accuracy : 100.0000%
Test_accuracy : 73.1000%
f1_score : 72.959%
jaccard_score : 57.467%
```

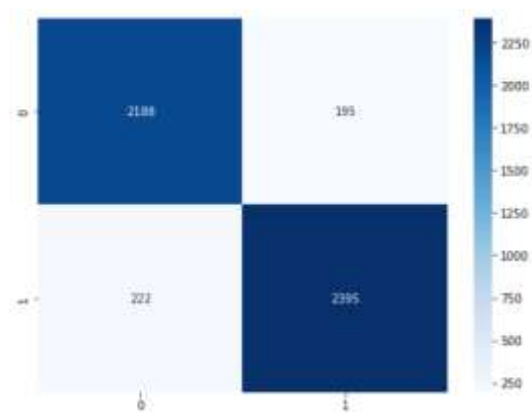


نمودار ROC مدل :

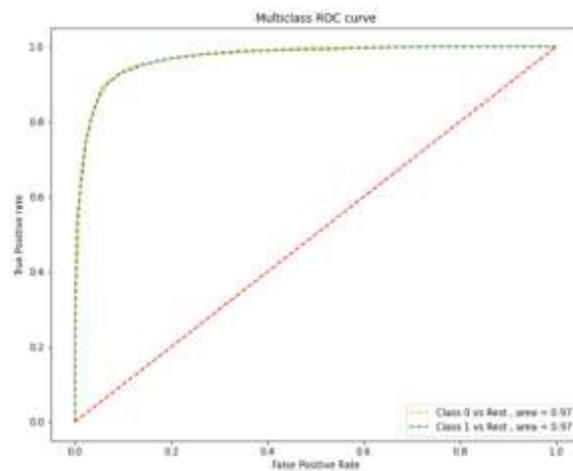


اعمال مدل Random Forest بر روی دسته ویژگی های دوم :

Train_accuracy : 100.0000%
 Test_accuracy : 91.6600%
 f1_score : 91.646%
 jaccard_score : 84.582%



نمودار ROC مدل :



تحلیل نتایج بدست آمده جنگل تصادفی روی دسته اول و دوم:

دقت مدل روی دسته اول برای داده آموزش ۱۰۰٪ و برای داده تست ۷۳٪ بوده است. پس بیش برآزش اتفاق افتاده است. با وجود اینکه روش جنگل تصادفی که امکان رخ دادن بیش برآزش بسیار کم است باز هم این اتفاق افتاده است. دقت مدل روی دسته داده دوم به نحو چشمگیری افزایش یافته است. روی دسته دوم دقت مدل برای داده آموزش ۱۰۰٪ و برای داده تست ۹۱/۶٪ به دست آمده است. دسته دوم ویژگی به قدری برای تفکیک صدای زن و مرد خوب عمل میکند که حتی روی داده آموزش هم دقت عالی گرفته ایم. ویژگی های بر حسب طیف و mfcc برای امر حداسازی صدای زن و مرد بسیار خوب عمل میکند.

مقایسه مدلها برای جنسیت:

همانطور که در جدول زیر نمایش داده شده است ، همانطور که انتظار داشتیم، مدل SVM بر روی دسته ویژگی های دوم دارای عملکرد بهتری نسبت به سایر مدل ها دارا است. البته اختلاف میان دقت مدلها روی داده تست بسیار نزدیک بوده است به هر حال random forest ضعیف تر از بقیه مدلها عمل کرده است.

مدل	دقت روی دسته ویژگی اول		دقت روی دسته ویژگی دوم	
	Train	test	Train	Test
MLP	74.35%	73.1 %	100%	94.9%
SVM	76.89%	73.52%	99.99%	95.44%
Random Forest	100%	73.1%	100%	91.66%

طبقه بندی احساسات :

برای این منظور در ابتدا دیتاست X , y تعیین کرده و با استفاده از Standard Scaler ، داده ها را نرمالایز می کنیم. در قدم بعد داده را به دو دسته ی train و test تقسیم کرده و ۳۰ درصد داده ها را به تست اختصاص می دهیم.

مدل MLP برای طبقه بندی احساسات:

مدل طراحی شده برای شبکه های عصبی mlp به صورت زیر می باشد . این مدل از ۴ لایه تشکیل شده است که تابع فعال ساز سه لایه اول relu و لایه اخر softmax انتخاب شده است . تعداد نورون های در نظر گرفته شده برای هر لایه در قطعه کد زیر نمایش داده شده است .

```
model = Sequential()
model.add(layers.Dense(256, activation='relu', input_shape=(X_train.shape[1],)))
model.add(layers.Dense(128, activation='relu'))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(4, activation='softmax'))
model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])
```

اعمال مدل MLP بر روی دسته ویژگی های اول :

دقت و خطای بدست آمده برای داده های train و validation در 5 epochs اخر به شرح زیر است.

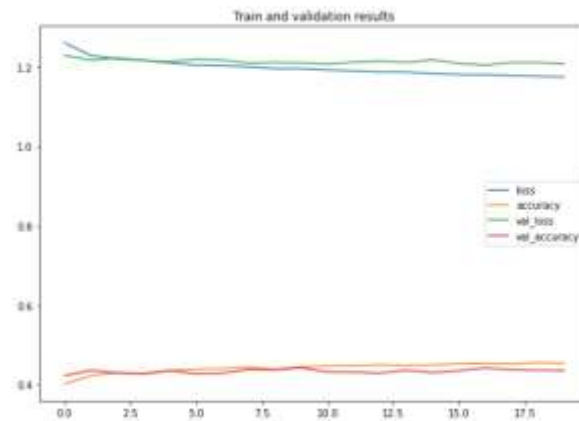
```
Epoch 16/20
73/73 [-----] - 0s 4ms/step - loss: 1.1810 - accuracy: 0.4529 - val_loss: 1.2085 - val_accuracy: 0.4351
Epoch 17/20
73/73 [-----] - 0s 5ms/step - loss: 1.1806 - accuracy: 0.4539 - val_loss: 1.2045 - val_accuracy: 0.4423
Epoch 18/20
73/73 [-----] - 0s 5ms/step - loss: 1.1781 - accuracy: 0.4524 - val_loss: 1.2108 - val_accuracy: 0.4378
Epoch 19/20
73/73 [-----] - 0s 5ms/step - loss: 1.1761 - accuracy: 0.4559 - val_loss: 1.2109 - val_accuracy: 0.4383
Epoch 20/20
73/73 [-----] - 0s 5ms/step - loss: 1.1752 - accuracy: 0.4542 - val_loss: 1.2074 - val_accuracy: 0.4359
```

دقت و خطای مدل روی داده های تست برابر است با :

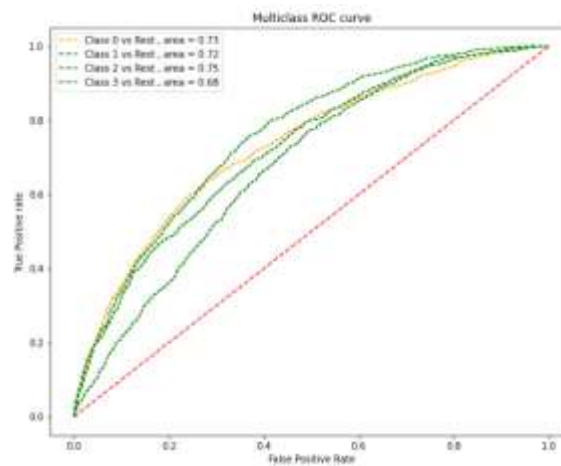
```
results = model.evaluate(X_test, y_test, verbose = 0)
print('test loss, test accuracy:', results)

test loss, test accuracy: [1.2038077116012573, 0.4390000104904175]
```

نمودار دقت و خطا مدل :



نمودار ROC مدل :



اعمال مدل MLP بر روی دسته ویژگی های دوم :

دقت و خطای بدست آمده برای داده های train و validation در 5 epochs آخر به شرح زیر است.

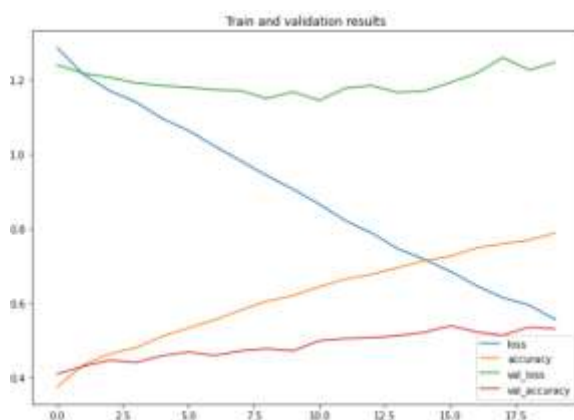
```
Epoch 16/20
73/73 [=====] - 0s 6ms/step - loss: 0.6848 - accuracy: 0.7252 - val_loss: 1.1919 - val_accuracy: 0.5384
Epoch 17/20
73/73 [=====] - 0s 5ms/step - loss: 0.6461 - accuracy: 0.7477 - val_loss: 1.2183 - val_accuracy: 0.5221
Epoch 18/20
73/73 [=====] - 0s 5ms/step - loss: 0.6142 - accuracy: 0.7592 - val_loss: 1.2591 - val_accuracy: 0.5131
Epoch 19/20
73/73 [=====] - 0s 6ms/step - loss: 0.5943 - accuracy: 0.7682 - val_loss: 1.2255 - val_accuracy: 0.5349
Epoch 20/20
73/73 [=====] - 0s 5ms/step - loss: 0.5554 - accuracy: 0.7879 - val_loss: 1.2464 - val_accuracy: 0.5382
```

دقت و خطای مدل روی داده های تست برابر است با :

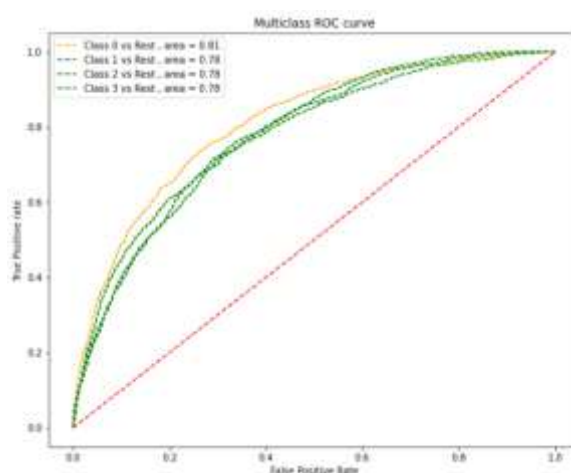
```
results = model.evaluate(X_test, y_test, verbose = 0)
print('test loss, test accuracy:', results)
```

test loss, test accuracy: [1.2326005697250366, 0.5266000032424927]

نمودار دقت و خطا مدل :



نمودار ROC مدل :



تحلیل نتایج بدست آمده mlp روی دسته ویژگی اول برای تفکیک احساسات:

دقت مدل روی دسته اول برای داده آموزش ۴۵٪ و برای داده تست ۴۳٪ بوده است. پس بیش برآزش اتفاق نیفتاده است. اما دقت مدل روی این دسته ویژگی برای تفکیک احساسات حتی از دقت رندوم هم پایینتر است و اصلا مناسب نیست. دقت مدل برای دسته ویژگی دوم بهبود یافته است چون ویژگی های این دسته حتی برای تفکیک احساس هم بهتر هستند. روی دسته دوم دقت مدل برای داده آموزش ۷۸/۸٪ و برای داده تست ۵۳٪ به دست آمده است. روی قسمت آموزش دقت قابل قبول گرفته ایم ولی روی داده تست باز هم تفاوت چشمگیر

نیست و از دسته بدی رندوم تصادفی تنها ۳٪ بهتر شده است. در مورد این دسته بیش برآزش اتفاق افتاده است که در انتهای این بخش در این مورد بحث شده است.

مدل SVM :

```
classifier = SVC(kernel='rbf',C=3.8 , gamma=0.1)
classifier.fit(X_train,y_train)
y_pred_test = classifier.predict(X_test)
```

اعمال مدل SVM بر روی دسته ویژگی های اول :

مقادیر دقت روی داده های تست و آموزش و f1 score و jaccard score و همچنین ماتریس confusion به شرح زیر می باشد.

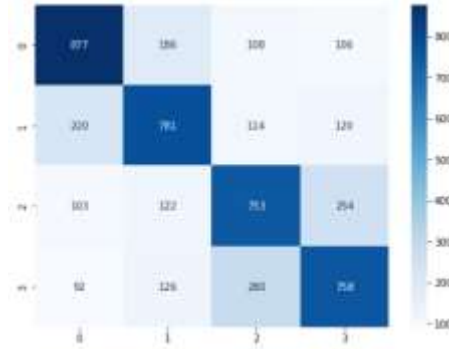
Train_accuracy : 45.7562%
Test_accuracy : 44.4600%
f1_score : 44.381%
jaccard_score : 28.622%



اعمال مدل SVM بر روی دسته ویژگی های دوم :

مقادیر دقت روی داده های تست و آموزش و f1 score و jaccard score و همچنین ماتریس confusion به شرح زیر می باشد.

Train_accuracy : 96.3820%
Test_accuracy : 63.3800%
f1_score : 63.343%
jaccard_score : 46.429%



تحلیل نتایج بدست آمده :

دقت مدل روی دسته اول برای داده آموزش ۴۵/۷٪ و برای داده تست ۴۴/۴٪ بوده است. پس بیش برآزش اتفاق نیفتاده است. اما دقت این مدل هم مثل مدل قبل روی این دسته ویژگی برای تفکیک احساسات اصلا مناسب نیست. دقت مدل برای دسته دوم بهبود یافته است چون ویژگی های این دسته حتی برای تفکیک احساس هم بهتر هستند. روی دسته دوم دقت مدل برای داده آموزش ۹۶/۳٪ و برای داده تست ۶۳/۴٪ به دست آمده است. دقت روی داده آموزش و تست بسیار بهتر از مدل mlp است. اما باز هم بیش برآزش اتفاق افتاده است که در این مورد هم در انتها بحث خواهد شد. دقت روی داده تست قابل قبول است.

مدل Random Forest برای طبقه بندی احساسات :

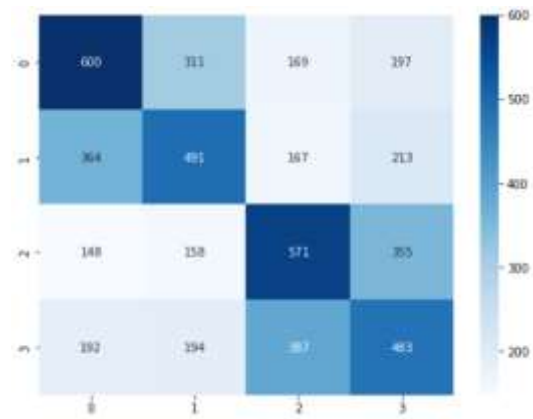
مدل مورد نظر به صورت زیر می باشد .

```
clf = RandomForestClassifier(n_estimators=2000 , random_state=42 , max_depth=50 )
clf.fit(X_train,y_train)
y_pred_test=clf.predict(X_test)
```

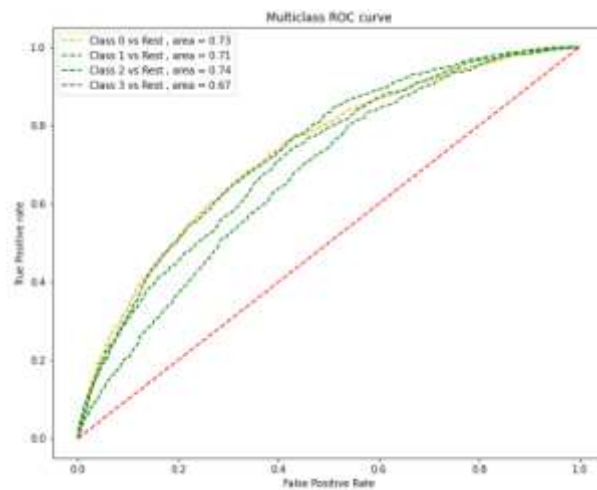
مقادیر دقت روی داده های تست و آموزش و f1 score و jaccard score و همچنین ماتریس confusion به شرح زیر می باشد.

اعمال مدل Random Forest بر روی دسته ویژگی های اول :

```
Train_accuracy : 99.9228%
Test_accuracy : 42.9000%
f1_score : 42.847%
jaccard_score : 27.316%
```

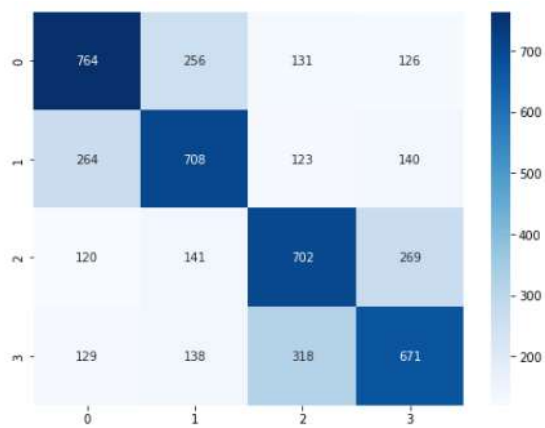


نمودار ROC مدل :

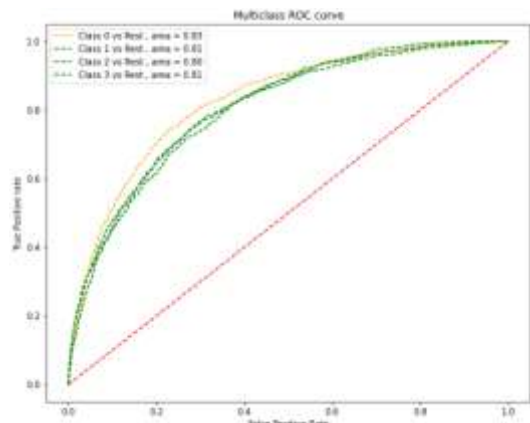


اعمال مدل Random Forest بر روی دسته ویژگی های دوم :

Train_accuracy : 99.9228%
 Test_accuracy : 56.9000%
 f1_score : 56.876%
 jaccard_score : 39.765%



نمودار ROC مدل :



تحلیل نتایج بدست آمده :

دقت مدل روی دسته اول برای داده آموزش ۹۹/۹٪ و برای داده تست ۴۲/۹٪ بوده است. با وجود استفاده از مدل جنگل تصادفی بیش برآزش زیادی اتفاق افتاده است. دقت مدل روی این داده تست دسته ویژگی برای تفکیک احساسات حتی از دقت رندوم هم پایینتر است و اصلا مناسب نیست. روی دسته دوم دقت مدل برای داده آموزش ۹۹/۹٪ و برای داده تست ۵۶/۹٪ به دست آمده است. روی دسته دوم ویژگی هم بیش برآزش زیادی رخ داده اما حداقل دقت روی داده تست بیشتر از طبقه بند رندوم است و میتوان به آن اتکا کرد.

در مورد بیش برآزش در داده احساسات میتوان تحلیل کرد که در روشهای انتخاب ویژگی ۱۰ صدای هر سوژه با یک عاطفه مشخص برچسب گذاری شده است. پس همبستگی این وویسها برای یک شخص زیاد است. ممکن است دو شخص که یک جمله را با یک عاطفه را گفته اند همبستگی آنچنان زیادی نداشته باشند. پس داده های یک عاطفه خاص برای هر شخص شباهت زیادی دارند. این ممکن است منجر به overfit شود. تعداد داده های مشابه در قسمت آموزش است و همبستگی بالایی وجود دارد و انگار چند داده کاملا مثل هم را وارد مدل کردیم و مدل دارد به جای یادگیری حفظ میکند. به ازای هر نفر ۴۰ داده داریم ولی اگر داده های ورودی متفاوت بودند و مثلا از هر فرد فقط ۴ داده میرفتیم و از هر عاطفه برای هر فرد یک وویس در داده موجود بود دیگر این شباهت و اورفیتینگ پیش نمی آمد. اورفیتینگ در جایی اتفاق میافتد که دیتاست بالانس نباشد. مدل های ساده احتمال اورفیتینگ بیشتری دارد. ولی مثلا cnn اگر میزدیم اینطور نمیشد.

مقایسه مدلها برای تفکیک احساسات :

همانطور که در جدول زیر نمایش داده شده است ، مدل SVM بر روی دسته ویژگی های دوم دارای عملکرد بهتر قابل توجهی نسبت به سایر مدل ها دارا است. دقت مدل جنگل تصادفی از mlp هم بهتر است. باز هم این مورد انتظار بود چون قبل از ابداع روش cnn و روشهای یادگیری عمیق بهترین مدل همین svm بوده است.

مدل	دقت روی دسته ویژگی اول		دقت روی دسته ویژگی دوم	
	train	test	train	test
MLP	45.42%	43.9%	78.79%	٪۵۲,۶۶
SVM	45.75%	٪۴۴,۴۶	96.38%	٪۶۳,۳۸
Random Forest	92.92%	٪۴۲,۹	92.92%	٪۵۶,۹

خوشه‌بندی

با اینکه دیتاست پروژه برچسب‌گذاری شده است در این بخش قصد داریم برچسبها را نادیده بگیریم و داده‌ها را بر اساس دو روش k-means و GMM خوشه‌بندی کنیم. در نهایت با توجه به معیارهای ارزیابی عملکرد خوشه‌بند را بررسی خواهیم کرد. برای ارزیابی عملکرد خوشه‌بند از معیار silhouette استفاده کردیم. این معیار هم به پیوستگی (Cohesion) درون خوشه‌ها و هم به میزان تفکیک‌پذیری آن‌ها بستگی دارد. مقدار نیم‌رخ برای هر نقطه، میزان تعلق آن را به خوشه‌اش در مقایسه با خوشه مجاور اندازه می‌گیرد. محدوده سیلوئت از -۱ تا +۱ است، که در آن مقدار زیاد نشان می‌دهد که شی به خوبی با خوشه خود مطابقت دارد و با خوشه‌های همسایه همسان نیست. اگر بیشتر اشیا از مقدار بالایی برخوردار باشند، ساختار خوشه‌بندی مناسب است. اگر بسیاری از نقاط دارای مقدار کم یا منفی باشند، در این صورت ممکن است ساختار خوشه‌بندی دارای خوشه‌های بسیار زیاد یا بسیار کم باشد.

مشابه قسمت طبقه‌بندی در این بخش هم دو دسته ویژگی در نظر گرفتیم و هر بار یک دسته را به خوشه‌بند مورد نظر دادیم.

روش k-means

برای ورودی این قسمت همه فضای ویژگی‌های استخراج شده را در نظر می‌گیریم. همانطور که خواسته شده الگوریتم را برای ۲، ۴، ۱۰ و تعداد کل شرکت‌کننده‌ها که ۴۰۰ تخمین زده شده است پیاده‌سازی می‌کنیم. برای این کار از دستور KMeans استفاده می‌کنیم و تعداد خوشه را مشخص می‌کنیم. سپس پارامترهای max_iter که مربوط به تعداد ماکسیمم تکرار الگوریتم است را ۳۰۰ قرار می‌دهیم و n_init که مربوط به دفعات تکرار کل الگوریتم برای در نظر گرفتن انواع نقاط شروع رندوم برای شروع الگوریتم است را ۱۰ قرار می‌دهیم.

در این پروژه از الگوریتم ++K-Means استفاده شده است. در الگوریتم ++K-Means، کاربر می‌خواهد که مرکزوارها در مقداردهی اولیه تا حد ممکن از هم دور باشند. ایده نهفته در پس این روش آن است که مرکزوارها به مرکز خوشه‌های واقعی نزدیک‌تر باشند و بنابراین سریع‌تر به همگرایی برسند. در واقع یک الگوریتم تخمین برای مسائل KMeans، NP-hard است.

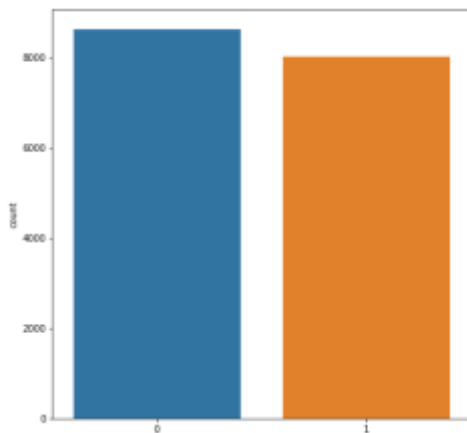
K=2 دو خوشه

برای دسته ویژگی اول:

توقع داریم این ۲ دسته متناظر با توزیع برچسب جنسیت زن و مرد در داده اولیه باشد.

```
model = KMeans(n_clusters = 2, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X1)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X1,y_clusters)
silhouette
```

0.2349702368059524

در جدول زیر فراوانی جنسیت داده ها در هر خوشه نمایش داده شده است.

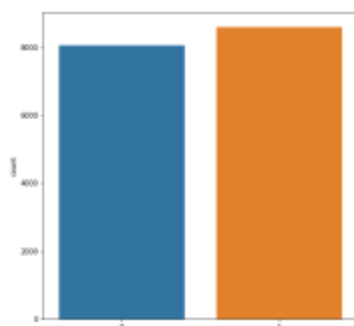
Cluster \ Actual label	Actual label	
	0	1
0	3669	4964
1	4278	3753

برای دسته ویژگی دوم:

مدل آموزش داده به صورت زیر می باشد.

```
model = KMeans(n_clusters = 2, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette
0.10781839934972111
```

در جدول زیر فراوانی جنسیت داده ها در هر خوشه نمایش داده شده است.

Cluster \ Actual label	Actual label	
	0	1
0	3549	4513
1	4398	4204

مقدار silhouette برای دسته ویژگی اول بسیار بهتر است و این برخلاف دقت در طبق بندها میباشد. دسته دوم ویژگی ابعاد بسیار بزرگتری دارد و چون الگوریتم بر حسب فاصله کار میکند و فاصله در ابعاد بالا خیلی زیاد میشود نتیجه خوشه بندی بدتر میشود.

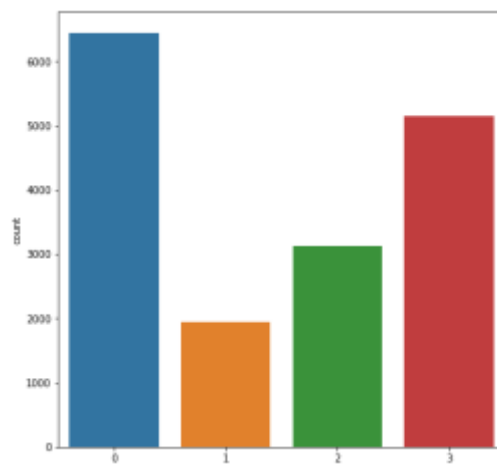
K=4 دو خوشه

برای دسته ویژگی اول:

توقع داریم این ۴ دسته متناظر با توزیع برجسب ۴ احساس خشم، غم، شادی و خنثی در داده اولیه باشد. مدل آموزش داده به صورت زیر می باشد.

```
model = KMeans(n_clusters = 4, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X1)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X1,y_clusters)
silhouette
```

0.24170465861133994

در جدول زیر فراوانی نوع احساس داده ها در هر خوشه نمایش داده شده است.

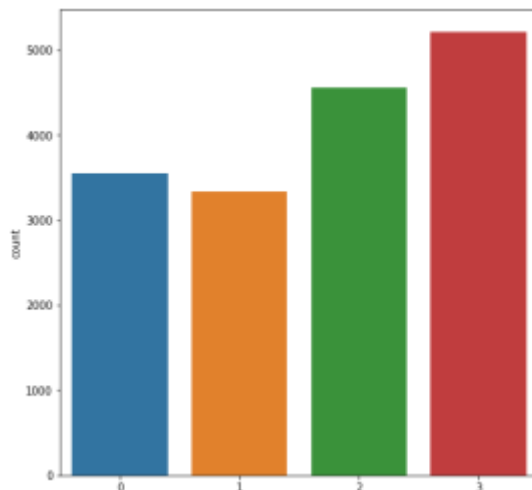
Actual label Cluster label				
	1	2	3	4
0	1439	1265	1858	1885
1	575	350	504	511
2	1366	1165	285	310
3	784	1364	1532	1471

برای دسته ویژگی دوم:

مدل آموزش داده به صورت زیر می باشد.

```
model = KMeans(n_clusters = 4, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette
0.07514463033931845
```

در جدول زیر فراوانی نوع احساس داده ها در هر خوشه نمایش داده شده است.

Actual label Cluster label				
	1	2	3	4
0	827	673	1065	990
1	1138	993	614	593
2	1227	1409	928	989
3	972	1069	1572	1605

برای تعداد ۴ خوشه هم مقدار silhouette برای دسته ویژگی اول بسیار بهتر است و این برخلاف دقت در طبق بندها میباشد. تحلیل این رخداد هم مانند حالت دو خوشه ای به خاطر ابعاد بالای دسته ویژگی دوم است.

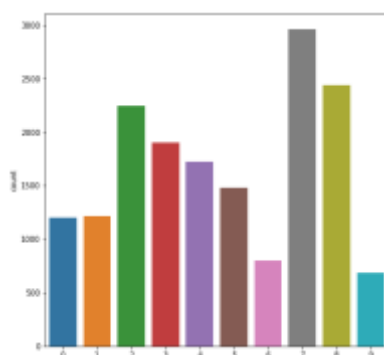
K=10 دو خوشه

برای دسته ویژگی اول:

توقع داریم این ۱۰ دسته متناظر با توزیع برچسب ۱۰ TEXT_ID در داده اولیه باشد. مدل آموزش داده به صورت زیر می باشد.

```
model = KMeans(n_clusters = 10, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X1)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score(X1,y_clusters)
silhouette
0.1858674891168136
```

در جدول زیر فراوانی text id داده ها در هر خوشه نمایش داده شده است.

Actual label \ Cluster label	1	2	3	4	5	6	7	8	9	10
0	129	108	149	92	155	248	80	99	65	75
1	164	220	108	88	96	76	152	106	103	104
2	222	161	185	251	262	234	223	212	255	237

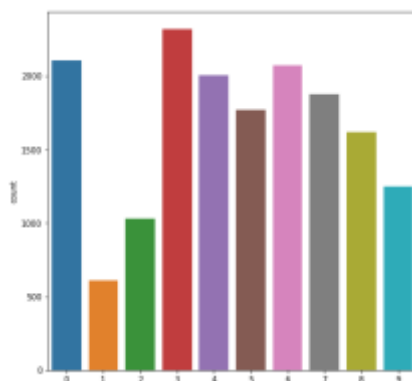
3	199	234	197	157	189	148	182	199	190	202
4	154	106	233	255	132	263	137	197	134	109
5	170	219	121	119	149	74	134	136	145	186
6	67	76	30	82	108	62	107	65	87	116
7	260	246	336	342	267	275	291	309	343	292
8	249	260	258	218	234	149	260	265	287	260
9	73	36	89	84	62	127	57	65	43	52

برای دسته ویژگی دوم:

مدل آموزش داده به صورت زیر می باشد.

```
model = KMeans(n_clusters = 10, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette

0.86152897869811592
```


در جدول زیر فراوانی text id داده ها در هر خوشه نمایش داده شده است.

Actual label \ Cluster label	1	2	3	4	5	6	7	8	9	10
0	208	172	315	210	188	334	170	218	155	134
1	58	67	66	58	73	63	57	53	56	57
2	97	107	43	101	113	68	150	103	102	148
3	204	202	321	240	244	259	219	221	206	202
4	213	254	125	203	241	135	224	167	226	214
5	194	231	102	137	159	89	181	177	256	242
6	214	189	245	249	215	243	197	191	166	164
7	197	184	179	200	133	175	184	226	212	183
8	190	139	148	160	135	143	163	196	162	184
9	112	121	162	130	153	147	108	101	111	105

باز هم مقدار silhouette برای دسته ویژگی اول بسیار بهتر است و این هم برخلاف دقت در طبق بندیها میباشد. علت مشابه قبل است.

K=400 دو خوشه

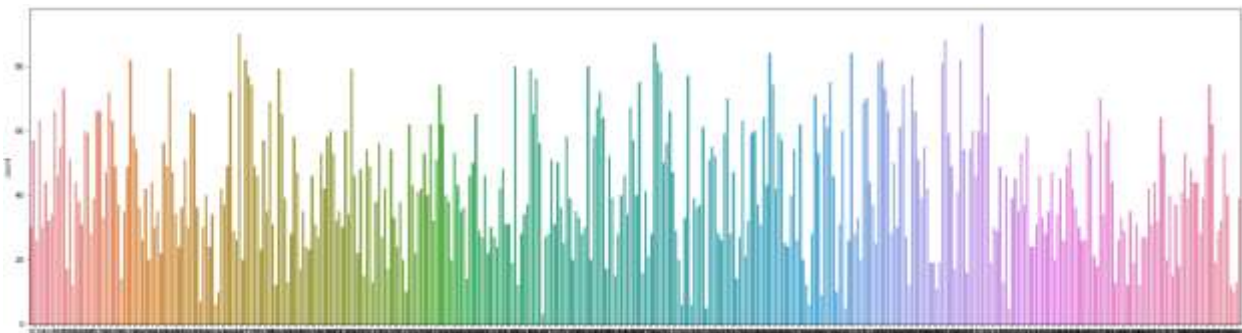
برای دسته ویژگی اول:

توقع داریم این ۱۰ دسته متناظر با توزیع برچسب ۴۰۰ شرکت کننده در داده اولیه باشد. مدل آموزش داده

به صورت زیر می باشد.

```
model = KMeans(n_clusters = 400, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X1)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X1,y_clusters)
silhouette
```

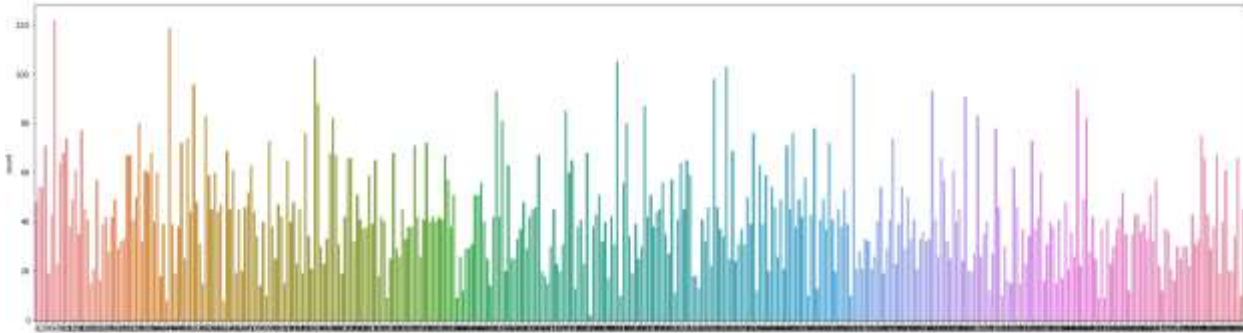
0.14528973950497076

برای دسته ویژگی دوم:

. که این مدل آموزش داده به صورت زیر می باشد.

```
model = KMeans(n_clusters = 400, init = "k-means++", max_iter = 300, n_init = 10, random_state = 0)
y_clusters = model.fit_predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette
```

0.1231195596116675

برای این تعداد خوشه مقدار silhouette برای دسته ویژگی اول کمی بهتر است و این هم برخلاف دقت در طبق بندها میباشد.

به طور تقریبی مقدار silhouette با افزایش تعداد خوشه کاهش یافته است. k-means برای داده های کروی خوش شکل مناسب است و بر اساس فاصله عمل میکند. برای همین برای خوشه بندی این دیتاست انتخاب بهینه ای نیست. میدانیم با افزایش بیش از حد تعداد خوشه ها این امر قابل انتظار است.

الگوریتم GMM

در این شیوه خوشه بندی، فرض بر این است که هر خوشه از داده هایی با توزیع نرمال (گوسی) تشکیل شده و در حالت کلی نیز داده ها نمونه ای از توزیع آمیخته نرمال هستند. هدف از خوشه بندی مدل آمیخته گوسی یا نرمال، برآورد پارامترهای توزیع هر یک از خوشه ها و تعیین برچسب برای مشاهدات است. به این ترتیب مشخص می شود که هر مشاهده به کدام خوشه تعلق دارد. چنین روشی را در یادگیری ماشین، خوشه بندی بر مبنای مدل می نامند. روش های خوشه بندی مانند kmeans قادر به شناسایی خوشه های کروی هستند و اگر ساختار خوشه ها خارج از این شکل باشد، عمل خوشه بندی توسط این الگوریتم ها به خوبی صورت نمی گیرد. به عنوان یک راه حل در چنین مواردی می توان از خوشه بندی بر مبنای مدل یا همان مدل آمیخته گوسی استفاده کرد. البته با توجه به ساختار خوشه های می توان از توزیع های دیگری نیز در این میان کمک گرفت ولی به طور معمول فرض می شود که مشاهدات مربوط به خوشه ها دارای توزیع نرمال هستند.

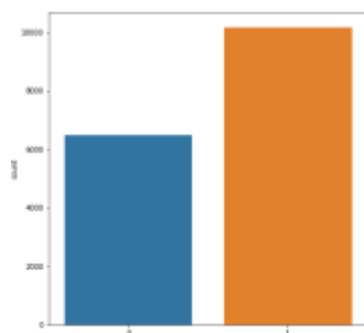
K=2

برای دسته ویژگی اول:

توقع داریم این ۲ دسته متناظر با توزیع برچسب ۲ جنسیت زن و مرد در داده اولیه باشد. مدل آموزش داده به صورت زیر می باشد.

```
gmm = GaussianMixture(n_components=2, covariance_type='full').fit(X1)
y_clusters = gmm.predict(X1)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X1,y_clusters)
silhouette
```

0.20176181348124786

در جدول زیر فراوانی جنسیت داده ها در هر خوشه نمایش داده شده است.

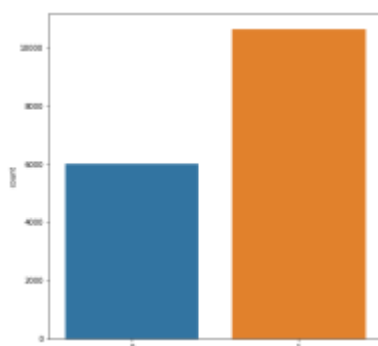
Cluster \ Actual label	Actual label	
	0	1
0	۲۸۶۰	۳۶۳۱
1	۵۰۸۷	۵۰۸۶

برای دسته ویژگی دوم:

تعداد ۸۶۳۳ در خوشه اول قرار میگیرد و تعداد ۸۰۳۱ خوشه در دسته دوم. در جدول زیر اطلاعات مربوط به این دو خوشه گزارش شده است. معیار silhouette برای روش GMM دو خوشه ای برابر با ۰/۲۳ شد. که این مقدار بالاترین مقدار در روش GMM بود. مدل آموزش داده به صورت زیر می باشد.

```
gmm = GaussianMixture(n_components=2, covariance_type='full').fit(X2)
y_clusters = gmm.predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score(X2,y_clusters)
silhouette
0.09620584599600727
```

در جدول زیر فراوانی جنسیت داده ها در هر خوشه نمایش داده شده است.

Cluster \ Actual label	label	
	0	1
0	5278	5375
1	2669	3342

مانند خوشه بندی kmeans برای GMM مقدار silhouette برای دسته ویژگی اول بسیار بهتر است و این برخلاف دقت در طبق بندها میباشد. دسته دوم ویژگی ابعاد بسیار بزرگتری دارد و دچار نحسی ابعاد میشویم.

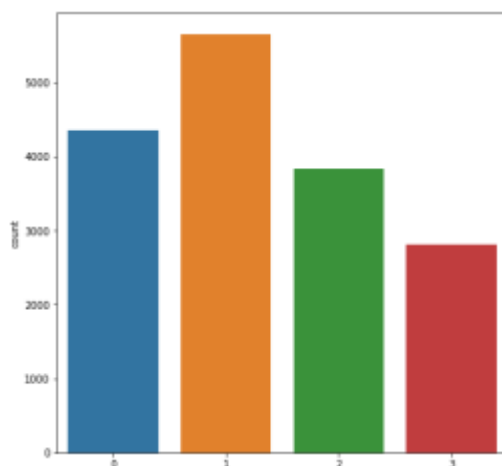
K=4

برای دسته ویژگی اول:

توقع داریم این ۴ دسته متناظر با توزیع برچسب ۴ احساس خشم، غم، شادی و خنثی در داده اولیه باشد. مدل آموزش داده به صورت زیر می باشد.

```
gmm = GaussianMixture(n_components=4, covariance_type='full').fit(X1)
y_clusters = gmm.predict(X1)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score(X1,y_clusters)
silhouette
0.11373828056357054
```

در جدول زیر فراوانی نوع احساس داده ها در هر خوشه نمایش داده شده است.

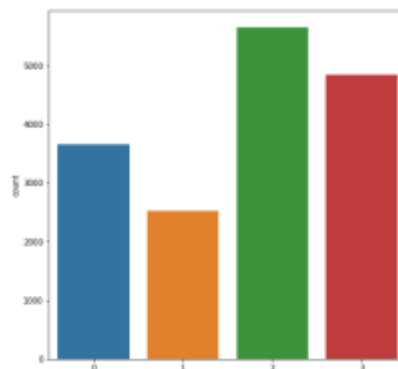
Actual label Cluster label	1	2	3	4
0	۸۵۳	۱۲۹۰	۱۱۱۰	۱۱۰۱
1	۸۲۸	۹۰۰	۲۰۳۳	۱۸۹۹
2	۱۳۵۳	۹۴۲	۶۸۶	۸۵۹
3	۱۱۳۰	۱۰۱۲	۳۵۰	۳۱۸

برای دسته ویژگی دوم:

تعداد ۸۶۳۳ در خوشه اول، تعداد ۸۰۳۱ عدد در خوشه دوم، تعداد * در خوشه سوم و تعداد در خوشه چهارم قرار میگیرد. در جدول زیر اطلاعات مربوط به این دو خوشه گزارش شده است. معیار silhouette برای روش GMM دو خوشه ای برابر با ۰/۲۳ شد. مدل آموزش داده به صورت زیر می باشد.

```
gmm = GaussianMixture(n_components=4, covariance_type='full').fit(X2)
y_clusters = gmm.predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette
```

0.012251905721506032

در جدول زیر فراوانی نوع احساس داده ها در هر خوشه نمایش داده شده است.

Actual label Cluster label	1	2	3	4
0	1300	1269	585	537
1	1224	1045	1247	1346
2	889	988	1906	1796
3	751	842	441	498

برای این تعداد خوشه هم مقدار silhouette برای دسته ویژگی اول بسیار بهتر است و این برخلاف دقت در طبق بندها میباشد. معیار ارزیابی خوشه بندی برای دسته دوم نزدیک به صفر است و خوشه بندی اصلا قابل قبل نیست. دسته دوم ویژگی ابعاد بسیار بزرگتری دارد و دچار نحسی ابعاد میشویم.

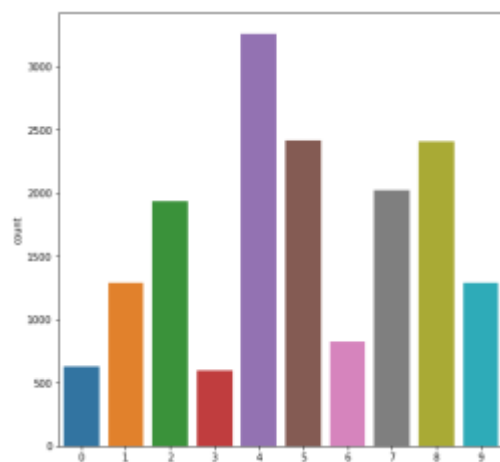
K=10

برای دسته ویژگی اول:

توقع داریم این ۱۰ دسته متناظر با توزیع برچسب ۱۰ TEXT_ID در داده اولیه باشد مدل آموزش داده به صورت زیر می باشد.

```
gmm = GaussianMixture(n_components=10, covariance_type='full').fit(X1)
y_clusters = gmm.predict(X1)
```


تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X1,y_clusters)
silhouette
```

0.05637521668995065

در جدول زیر فراوانی text id داده ها در هر خوشه نمایش داده شده است.

Actual label \ Cluster label	1	2	3	4	5	6	7	8	9	10
0	۸۱	۱۰۵	۶۳	۴۲	۶۰	۴۶	۵۹	۶۷	۴۱	۶۱
1	۱۳۱	۶۹	۱۷ .	۱۷۹	۹۸	۲۲۷	۱۰۵	۱۲۳	۱۰۶	۸۱
2	۲۰۵	۲۵۶	۱۵ .	۱۵۷	۱۷۴	۸۰	۲۳۷	۲۱۰	۲۱۲	۲۵۴
3	۵۸	۳۳	۷۹	۶۸	۵۹	۱۱۶	۴۸	۵۶	۳۷	۴۴
4	۳۲۳	۳۳۱	۳۶ ۸	۳۳۰	۳۱۱	۳۲۸	۳۳۰	۳۵۰	۳۱۲	۲۷۳
5	۲۲۷	۲۱۹	۲۶ ۷	۲۶۸	۲۲۴	۲۱۳	۲۲۲	۲۵۳	۲۶۵	۲۵۵

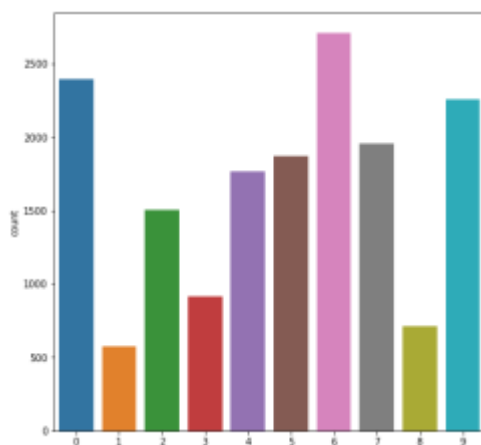
6	۸۱	۸۸	۴۵	۷۶	۹۴	۷۲	۱۰۹	۷۶	۶۸	۱۱۳
7	۲۰۰	۲۵۱	۱۸ ۸	۱۷۹	۱۷۷	۹۶	۲۳۲	۲۰۲	۲۴۶	۲۴۷
8	۲۳۹	۱۸۴	۲۱ ۵	۲۸۰	۲۹۰	۲۵۸	۲۱۰	۲۰۶	۲۹۲	۲۲۹
9	۱۴۲	۱۳۰	۱۶ ۱	۱۰۹	۱۶۷	۲۲۰	۱۰۱	۱۱۰	۷۳	۷۶

برای دسته ویژگی دوم:

مدل آموزش داده به صورت زیر می باشد.

```
gmm = GaussianMixture(n_components=10, covariance_type='full').fit(X2)
y_clusters = gmm.predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette
```

0.0019084484354461647

در جدول زیر فراوانی text id داده ها در هر خوشه نمایش داده شده است.

Actual label Cluster label	1	2	3	4	5	6	7	8	9	10
0	۲۳۱	۲۵۶	۲۳۲	۲۴۳	۲۴۹	۱۸۱	۲۵۱	۲۳۴	۲۵۲	۲۶۷
1	۶۱	۶۰	۵۳	۵۶	۵۶	۵۳	۶۰	۵۵	۶۰	۶۰
2	۱۴۰	۱۵۹	۱۱۵	۱۳۹	۱۶۵	۱۳۲	۱۷۲	۱۴۹	۱۵۳	۱۷۹
3	۹۴	۸۷	۱۰۷	۹۶	۸۴	۹۰	۸۱	۹۹	۸۵	۹۵
4	۱۹۶	۱۵۲	۱۵۳	۱۹۵	۱۶۳	۱۶۸	۱۷۸	۲۰۷	۱۶۳	۱۸۵
5	۱۸۱	۱۵۱	۲۳۶	۱۷۹	۱۷۸	۲۱۱	۱۸۲	۱۹۳	۱۹۰	۱۷۱
6	۲۷۳	۲۷۸	۲۹۶	۲۹۷	۲۶۰	۲۶۵	۲۵۰	۲۶۶	۲۶۷	۲۵۹
7	۲۱۲	۲۲۶	۱۸۳	۱۷۶	۲۰۸	۲۳۲	۱۹۰	۱۵۲	۱۹۰	۱۸۵
8	۷۰	۶۷	۸۳	۷۲	۷۰	۷۹	۶۹	۶۹	۶۹	۶۲
9	۲۲۹	۲۳۰	۲۴۸	۲۳۵	۲۲۱	۲۴۵	۲۲۰	۲۲۹	۲۲۳	۱۷۰

ایین برای این تعداد خوشه silhouette برای هر دو دسته ویژگی . معیار ارزیابی خوشه بندی برای دسته دوم نزدیک به صفر است ولی باز هم برای دسته ویژگی اول بهتر عمل کرده است. برای هر دسته خوشه بندی قابل قبول نیست.

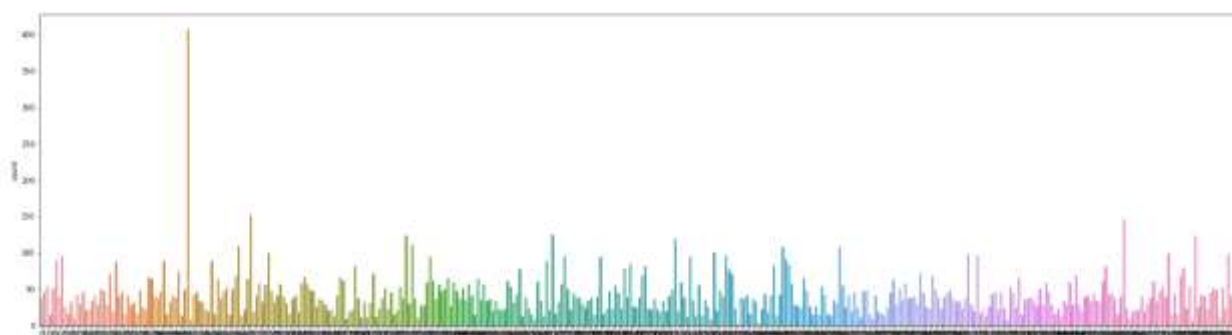
K=400

برای دسته‌بندی اولیه:

توقع داریم این ۱۰ دسته متناظر با توزیع برچسب ۴۰۰ شرکت کننده در داده اولیه باشد. مدل آموزش داده به صورت زیر می‌باشد.

```
gmm = GaussianMixture(n_components=400, covariance_type='full').fit(X1)
y_clusters = gmm.predict(X1)
```

تعداد داده‌های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

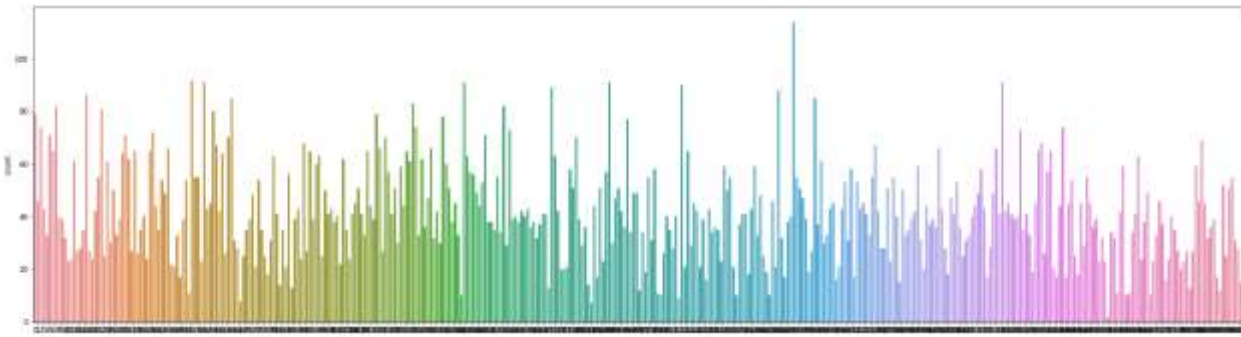
```
silhouette = silhouette_score(X1, y_clusters)
silhouette
-0.04428584678315022
```

برای دسته‌بندی دوم:

. معیار silhouette برای روش GMM دو خوشه‌ای برابر با ۰/۲۳ شد. مدل آموزش داده به صورت زیر می‌باشد.

```
gmm = GaussianMixture(n_components=400, covariance_type='full').fit(X2)
y_clusters = gmm.predict(X2)
```

تعداد داده های موجود در هر خوشه :



مقدار silhouette بدست آمده برای مدل :

```
silhouette = silhouette_score (X2,y_clusters)
silhouette
```

0.11932491431225853

برای این خوشه بندی معیار ارزیابی خوشه بندی برای دسته اول منفی است که اصلا قابل قبول نیست. silhouette دسته دوم مثبت شد ولی بسیار نزدیک به صفر بود. به طور کلی روی این تعداد خوشه الگوریتم اصلا خوب عمل نکرد.

در این الگوریتم خوشه بندی هم به طور تقریبی با افزایش تعداد خوشه ها مقدار silhouette کاهش یافته است. میدانیم با افزایش بیش از حد تعداد خوشه ها این امر قابل انتظار است.

به طور کلی نتایج k-means برای همه تعداد خوشه ها از GMM بهتر بود. علت احتمالا این است که الگوریتم GMM تعدادی گاوسی را برای توزیع کل نمونه ها در نظر میگیرد و داده های ما به خوبی از این توزیع های گاوسی پیروی نمیکنند. پس روش مبتنی بر فاصله k-means جواب بهتر به دست داده است.