



دانشگاه آزاد اسلامی

واحد بندرعباس

دانشکده فنی و مهندسی  
رشته مهندسی کامپیوتر – ارشد نرم افزار

عنوان:

تشخیص نویسنده‌ها از روی متن

استاد راهنما:

دکتر عباس عکاسی

نگارندگان:

کمیل آقابابایی، سحر ارزمانی، مریم یوسفی‌زاده، سمیه رجبی، امیر جعفری‌زاده، پری دریاکش، هادی ناظری،  
مرتضی سهرابی، میثم سالاری، امید نیاستی

بهار ۱۳۹۸

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

تقدیم به :

او که آموخت مارا تا بیاموزیم

{استاد گرامی جناب آقای دکتر عکاسی}

تقدیم به آنان که وجودمان جز هدیه وجودشان نیست

{ پدران و مادران عزیزمان }

سپاسگزاری:

از استاد گرامی دکتر عباس عکاسی که با حسن خلق و فروتنی، ما را در این عرصه یاری نمودند کمال تشکر را داریم.



## فهرست

چکیده .....	۱
فصل اول: .....	۲
کلیات تحقیق .....	۲
۱-۱- مقدمه .....	۳
۱-۲- بیان مسأله .....	۳
۱-۳- اهداف تحقیق .....	۴
۱-۴- نوآوری‌های تحقیق .....	۵
۱-۵- تعریف واژگان .....	۵
۱-۶- ساختار پژوهش .....	۶
فصل دوم : .....	۷
ادبیات و پیشینه تحقیق .....	۷
۲-۱- مقدمه .....	۸
۲-۲- داده‌کاوی .....	۸
۲-۳- دسته‌بندی .....	۱۰
۲-۴- الگوریتم‌های رایج دسته‌بندی .....	۱۱
۲-۴-۱- دسته بندی کننده بردار پشتیبان خطی .....	۱۱
۲-۴-۲- شبکه بیزین .....	۱۶
۲-۴-۳- رگرسیون منطقی .....	۱۸
۲-۴-۴- درخت های اضافی .....	۲۰
۲-۵- یادگیری تجمعی .....	۲۱
۲-۵-۱- متدولوژی آموزش دسته جمعی .....	۲۲

۲۳	۲-۶-انواع روش ترکیب دسته بندها
۲۳	۲-۶-۱- ترکیب همه دسته بندها
۲۴	۲-۶-۲- انتخاب رو به جلو
۲۴	۲-۶-۳- حذف رو به عقب
۲۵	۲-۷- معیارهای ارزیابی دسته بندها
۲۵	۲-۷-۱- نگاهی دقیق تر به مسائل دسته بندی
۲۷	۲-۷-۲- دقت - صحت - بازخوانی
۲۹	۲-۷-۳- الگوریتم (مدل) پیشنهادی
۳۳	فصل سوم:
۳۳	روش تحقیق
۳۴	۳-۱- مقدمه
۳۴	۳-۲- روش کار
۳۴	۳-۲-۱- خواندن مجموعه داده
۳۶	۳-۲-۲- تجزیه و تحلیل داده
۳۸	۳-۲-۳- تقسیم بندی دیتاست
۴۱	۳-۲-۴- پاکسازی داده
۴۲	۳-۲-۵- تهیه ورودی برای الگوریتم های پایه با استفاده از TF-IDF و BOW(TF) و BOW(binary)
۴۵	۳-۲-۶- ساخت ۱۰۰ نمونه داده آموزش با رویکرد جایگشتی
۴۷	۳-۲-۷- انتخاب الگوریتم های پایه و ساخت ۱۰۰ دسته بند با استفاده از هر الگوریتم
۵۶	۳-۲-۸- گرفتن میانگین F-score برای هر مدل و اعمال فیلترینگ
۵۸	۳-۲-۹- انتخاب بهترین دسته بند
۵۸	۳-۲-۱۰- رای گیری اکثریت

۵۹	۳-۲-۱۱- انتخاب رو به جلو
۶۱	۳-۲-۱۲- حذف رو به عقب
۶۳	۳-۲-۱۳- ترکیب همه رویکردها:
۶۴	فصل چهارم :
۶۴	محاسبات و یافته های تحقیق
۶۵	۴-۱- داده های مورد استفاده
۶۵	۴-۲- تنظیم پارامترها
۶۵	۴-۲-۱- پارامترهای الگوریتم SVC:
۶۷	۴-۲-۲- پارامترهای الگوریتم Logistic Regression :
۶۷	۴-۲-۳- Extra trees :
۶۷	۴-۲-۴- Naïve Bayes
۶۸	۴-۳- روش های استفاده شده به منظور مقایسه
۷۰	۴-۴- نتایج
۷۱	۴-۵- نتیجه گیری
۷۳	فصل پنجم :
۷۳	نتیجه گیری و مشاهدات
۷۴	۵-۱- خلاصه و نتیجه گیری
۷۴	۵-۲- پیشنهادات
۷۶	مراجع



## فهرست اشکال

- ۲- ۱- مرحله آموزش روش دسته جمعی ..... ۲۲
- ۲- ۲- مرحله آزمایشی روش دسته جمعی ..... ۲۳
- ۲- ۳- ارزیابی کارایی یک مدل ..... ۲۶

## فهرست جداول

- جدول ۴-۱: پارامترهای تنظیم شده توسط کاربر ..... ۶۶
- جدول ۴-۲: پارامترهای تنظیم شده توسط کاربر ..... ۶۷
- جدول ۳-۴: بهترین نتیجه بدست آمده از الگوریتم های ترکیب بر روی مجموعه داده تست ..... ۷۰
- جدول ۴-۴: مقایسه نتایج بدست آمده برای مجموعه داده Spooky\_Author\_Identification  
با سایر روش ها ..... ۷۰



## چکیده

چالش شناسایی نویسنده نیازمند کاربران است که بخش‌هایی از متن نوشته شده توسط نویسندگان را بازایی کنند، یک مدل پیشگویانه داده‌شده با داده را آموزش دهند و قادر باشد تا متن‌های جدید را به نویسندگان درست دسته بندی کنند. این یک مساله یادگیری نظارت شده در حوزه پردازش زبان طبیعی است.

در این پژوهش از الگوریتم‌های متعددی برای دسته بندی ارائه شده‌اند که از آن دسته می‌توان؛ به شبکه‌های بیزین، الگوریتم دسته بندی کننده بردار پشتیبان خطی، الگوریتم رگرسیون منطقی والگوریتم درخت‌های اضافی اشاره کرد که با استفاده از هر الگوریتم ۱۰۰ دسته بند را ساخته و با توجه به معیارهای ارزیابی هر مدل را بر روی داده‌های اعتبارسنجی ارزیابی کرده و بهترین دسته بند را انتخاب می‌کنیم. و در نهایت از روش‌های ترکیب دسته بندها استفاده کرده و بهترین روش را بر روی داده‌های تست، ارزیابی می‌کنیم.

**کلمات کلیدی :** شناسایی نویسنده، الگوریتم شبکه بیزین، الگوریتم دسته بندی کننده بردار پشتیبان خطی، الگوریتم رگرسیون منطقی، الگوریتم درخت‌های اضافی، ترکیب دسته بندها

فصل اول:

کلیات تحقیق

افزایش استفاده از کامپیوترها در فعالیت‌های کسب و کار، منجر به رشد سریع پایگاه‌های اطلاعاتی و اجتماع داده‌ها توسط بیشتر سازمان‌ها شده است. روزانه حجم عظیمی از داده‌ها تولید شده و در پایگاه‌های مختلف داده ذخیره می‌شود. کشف دانش در پایگاه‌های اطلاعاتی (KDD<sup>۱</sup>) یک تحلیل خودکار، اکتشافی و مدل‌سازی از مخازن داده بزرگ است. KDD یک فرآیند سازمان‌یافته شناسایی الگوهای معتبر، جدید، مفید، و قابل فهم از مجموعه‌های داده بزرگ و پیچیده است. داده‌کاوی (DM<sup>۲</sup>) هسته اصلی این فرآیند است که به کشف داده‌ها، توسعه مدل و کشف الگوهای ناشناخته پیشین می‌پردازد [1]. داده‌کاوی، فرآیند کشف الگوهایی است که برای ساخت مدل‌های پیش‌بینی‌کننده استفاده می‌شود، که در بسیاری از رشته‌های علوم مالی گرفته تا مراقبت‌های بهداشتی گسترده شده است [2].

در سال‌های اخیر، صحت اطلاعات (آنالین) توجه زیادی را به خود جلب کرده است، به ویژه در زمینه "اخبار جعلی" در زمینه انتخابات ریاست جمهوری آمریکا. در حال حاضر تاکید زیادی بر روی منشا و صحت اطلاعات صورت گرفته است. در مورد اسناد نوشتاری، یک جنبه مهم از این نوع انتقاد، به نویسندگی مربوط می‌شود. ارزیابی صحت اطلاعات اساساً مربوط به شناسایی نویسنده اصلی این اسناد می‌باشد. شناسایی نویسنده، تلاش می‌کند نویسندگان پشت متن را آشکار کند. این یک حوزه نوظهور از تحقیقات مرتبط با کاربرد در تحقیقات ادبی، امنیت سایبری، قانونی، و تحلیل رسانه‌های اجتماعی است [3].

## ۱-۲- بیان مسأله

رویکردهای کمی به وظایفی مانند تخصیص دادن اسناد، تایید، ثبت مشخصات و یا خوشه‌بندی نویسنده بر این فرض پایه تکیه دارند که سبک نوشتاری اسناد به گونه‌ای مشخص، آموخته و برای ساختن مدل‌های پیش‌بینی استفاده می‌شود. هر سیستم شناسایی موفق نویسنده، می‌تواند در یک تنظیمات اختصاصی یا در یک تنظیمات تایید (بازبینی)، تعیین هویت قدرتمندی در متن در ژانرهای مختلف، رفتار با موضوعات مختلف و یا داشتن مخاطبان هدف مختلف در ذهن ایجاد کند. مدل‌های سنتی برای تخصیص دادن اسناد در مواردی که چندین

<sup>۱</sup> KDD : Knoledge Discovery in Database

<sup>۲</sup> DataMining

نویسنده در یک سند مشارکت دارند، قابل اجرا نیستند. در نتیجه می توان با ابزار های داده کاوی این موضوع را حل کرد [3].

در این پژوهش ابتدا با توجه به ماهیت مسئله با استفاده از کتابخانه پاندا داده های ورودی با سه ویژگی را دریافت می کنیم. روش دسته بندی یک روش یادگیری با نظارت است که داده های ورودی به سه بخش داده های آموزش و داده های آزمون و داده های اعتبارسنجی تقسیم می شوند. هر الگوریتم کاندید، ابتدا با استفاده از مجموعه داده آموزش یک مدل را که نشان دهنده الگوی حاکم بر داده ها می باشد را استخراج می کند و سپس با استفاده از مجموعه آزمون و اعتبارسنجی، دقت مدل ارائه شده برای دسته بندی را بررسی می کند.

الگوریتم های متعددی برای دسته بندی ارائه شده اند که از آن دسته می توان؛ به شبکه های بیزین، الگوریتم دسته بندی کننده بردار پشتیبان خطی، الگوریتم رگرسیون منطقی و الگوریتم درخت های اضافی اشاره کرد.

در این نوشتار ابتدا با استفاده از توابع کتابخانه پاندا داده های مربوط به نویسندگان تقسیم بندی می شوند تا به صورت داده های آموزش و آزمون و اعتبارسنجی درآیند. سپس بعد از پاکسازی داده و تهیه ورودی برای الگوریتم های پایه که در بالا ذکر شد، ۱۰۰ دسته بند را با استفاده از هر الگوریتم میسازیم. سپس با استفاده از معیار های ارزیابی بر روی داده های اعتبارسنجی، هر مدل را ارزیابی کرده و بهترین دسته بند را انتخاب می کنیم و در مرحله بعد با روش های ترکیب دسته بند ها از جمله حذف پس رونده و انتخاب پیش رونده و انتخاب تک دسته بند بهترین، بهترین روش را انتخاب کرده و بر روی داده های آزمون تست می کنیم. در نهایت سعی شده با در نظر گرفتن نقاط ضعف و قوت روش های مختلف داده کاوی یک الگوریتم ترکیبی برای تشخیص نویسنده ارائه شود.

### ۳-۱- اهداف تحقیق

هیچ راهنمای مشخصی برای آنکه پژوهشگران یا تحلیلگران بدانند چگونه یک الگوریتم داده کاوی را انتخاب کنند وجود ندارد. در نتیجه انتخاب یک الگوریتم مشخص امری بسیار پیچیده است، لذا در این پژوهش برای ارتقای نتایج داده کاوی از چندین الگوریتم استفاده شده و هر کدام را جداگانه توضیح و تشریح داده و پردازش ها را با الگوریتم های مختلف تکرار و سرانجام از تکنیک رای گیری اکثریت استفاده می شود. هدف از روش های ارائه شده افزایش دقت تشخیص نویسنده با استفاده از کشف الگوها در میان مجموعه داده ها می باشد.

#### ۴-۱- نوآوری های تحقیق

- استفاده از ترکیب دسته بندها
- ساخت تعداد مدل های زیاد

#### ۵-۱- تعریف واژگان

داده کاوی : استخراج داده فرآیند کشف الگوهای جالب و دانش از مقادیر زیاد داده ها است [4]. دسته بندی<sup>۳</sup> : جهت اختصاص دادن یک برچسب به مجموعه ای از داده ها که دسته بندی نشده اند، استفاده می شود. در حقیقت دسته بندی یک وظیفه آنالیز داده است ، به عبارتی فرآیند یافتن مدلی است که کلاس های داده و مفاهیم را توصیف و تمایز می دهد . در دسته بندی داده ها بر اساس ویژگی هایشان به دسته هایی که نام آن ها از قبل مشخص می باشد، تخصیص داده می شوند [5].

یادگیری تجمعی<sup>۴</sup> : یادگیری تجمعی یک تکنیک یادگیری ماشین است که با استفاده از یک الگوریتم یادگیری متفاوت ، چندین یادگیرنده را با همان داده ها آموزش می دهد . یادگیری تجمعی اغلب از یک الگوریتم یادگیری واحد بهتر عمل می کند. این امر به خصوص زمانی صادق است که گروه شامل الگوریتم های متنوعی است که هر کدام یک رویکرد کاملاً متفاوتی را اتخاذ می کنند . هر الگوریتم یادگیری به عنوان یک یادگیرنده پایه است . الگوریتم مورد استفاده برای ترکیب نتایج با نام گروه<sup>۵</sup> شناخته می شود .

زمانی که باید از یادگیری گروهی استفاده کنیم به شرح ذیل می باشد :

- مجموعه داده ها خیلی بزرگ یا خیلی کوچک است .
- داده های پیچیده ( غیرخطی<sup>۶</sup> )
- قابلیت اعتماد بالا [6][7]

---

<sup>۳</sup> Classification

<sup>۴</sup> Ensemble Learning

<sup>۵</sup> Ensemble

<sup>۶</sup> Nonlinear

## ۶-۱- ساختار پژوهش

مطالبی که در فصول بعدی ارائه خواهد شد به شرح زیر خواهد بود:

- در فصل دوم مفاهیم پایه‌ای مانند داده‌کاوی، کلیات مربوط به الگوریتم‌های دسته‌بندی، مفاهیم یادگیری تجمعی، الگوریتم‌های رایج دسته‌بندی، معیارهای ارزیابی این الگوریتم‌ها و انواع روش ترکیب دسته‌بندها مورد بحث قرار می‌گیرد.
- در فصل سوم حاوی کارهای انجام شده و تحقیقات مرتبط با موضوع می‌باشد، همچنین فضای کلی مسأله معرفی می‌شود.
- در فصل چهارم مدل پیشنهادی برای دسته‌بندی بر روی مجموعه داده‌های نویسندگان اعمال و نتایج روش پیشنهادی در این زمینه مورد مقایسه و ارزیابی قرار گرفته است.
- فصل پنجم نیز حاوی خلاصه، نتیجه‌گیری و پیشنهادات می‌باشد.



## فصل دوم :

### ادبیات و پیشینه تحقیق

## ۱-۲- مقدمه

در دنیای امروزی، ما با حجم عظیمی از داده ها مواجه هستیم و تجزیه و تحلیل مجموعه داده های بزرگ، اساس کلیدی رقابت، زیربنای رشد بهره وری و نوآوری خواهد بود. افزایش حجم و جزئیات اطلاعات به دست آمده توسط شرکت ها، رشد چند رسانه ای، رسانه های اجتماعی و اینترنت اشیاء، رشد نمایی در داده را برای آینده قابل پیش بینی خواهد داشت. در نتیجه با توجه به این حجم از داده ها، چالش هایی برای توانایی کنترل داده ها بوجود می آید [1]. دست پیدا کردن به اطلاعات نهفته در پایگاه داده شرکت ها، دانشگاه ها، مؤسسات دولتی و سایر مراکز نیازمند مدیریتی جدید است و با به کارگیری سیستم های سنتی این امر تحقق نمی یابد.

در این فصل ابتدا داده کاوی و مراحل آن شرح داده می شود و سپس دسته بندی داده ها و الگوریتم های رایج آن و معیار های ارزیابی آن ها مورد بحث قرار میگیرد و در پایان توضیحاتی در مورد یادگیری تجمعی و انواع روش ترکیب دسته بندها بصورت کامل شرح داده می شود.

## ۲-۲- داده کاوی

در سال ۱۹۶۰ آماردانان اصطلاح "Data Fishing" یا "Data Dredging" به معنای "صید داده" را جهت کشف هر گونه ارتباط در حجم بسیار بزرگی از داده ها بدون در نظر گرفتن هیچگونه پیش فرضی بکار بردند. بعد از سی سال و با انباشته شدن داده ها در پایگاه داده اصطلاح داده کاوی در حدود سال ۱۹۹۰ رواج بیشتری یافت. این فرآیند یک مرحله فراتر از بازیابی ساده داده ها است و این اجازه را می دهد که دانش را در میان حجم انبوه داده ها کشف کرد.

داده کاوی، به مفهوم استخراج اطلاعات نهان یا الگوها و روابط مشخص در حجم زیادی از داده ها در یک یا چند بانک اطلاعاتی بزرگ گفته می شود. بسیاری از مردم داده کاوی را مترادف واژه های رایج کشف دانش از داده ها (KDD) می دانند. داده کاوی، پایگاه ها و مجموعه حجیم داده ها را در پی کشف و استخراج، مورد تحلیل قرار می دهد. در حقیقت داده کاوی یک مرحله ضروری از فرآیند بزرگ تر کشف دانش می باشد که شامل مراحل زیر می باشد:

۱. پاک‌سازی داده‌ها<sup>۷</sup>: در این مرحله ، داده‌های نویز و ناسازگار حذف می‌شوند .
  ۲. یکپارچگی داده<sup>۸</sup>: در این مرحله , چند منبع داده ترکیب می‌شوند .
  ۳. انتخاب داده<sup>۹</sup>: در این مرحله ، داده‌های مربوط به کار آنالیز از پایگاه‌داده بازیابی می‌شوند
  ۴. تبدیل داده ها<sup>۱۰</sup>: در این مرحله ، داده‌ها به اشکالی که برای بکار بردن در داده‌کاوی مناسب باشند، تبدیل یا ترکیب می‌شوند .
  ۵. داده کاوی : در این مرحله ، روش‌های هوشمند به منظور استخراج الگوهای داده‌ای بکار گرفته می‌شوند .
  ۶. ارزیابی الگو<sup>۱۱</sup>: در این مرحله الگوهای داده‌ای مورد ارزیابی قرار می‌گیرند .
  ۷. ارائه دانش<sup>۱۲</sup>: در این مرحله , دانش ارائه می‌شود [8].
- داده‌کاوی از دو مرحله اصلی تشکیل شده است؛ مرحله اول پیش پردازش داده‌ها که در این مرحله ویژگی‌های با تأثیر بالاتر از داده‌های سطح پایین استخراج می‌شود. مرحله دوم تشخیص الگو می‌باشد که به کشف الگوی موجود در داده‌ها به کمک صفات و ویژگی‌های بدست آمده می‌پردازد [9].
- کاوش اطلاعات، حجم عظیمی از داده‌های خام را به فرمی تغییر می‌دهد که انسان بتواند آن‌ها را به راحتی بفهمد و برای تصمیم‌گیری بتواند از این اطلاعات استفاده کند. در قدم اصلی داده کاوی ممکن است از چندین الگوریتم داده کاوی استفاده شود. کار اصلی الگوریتم داده کاوی با توجه به نوع مسئله‌ی کشف دانش تغییر می‌کند اما دو نوع اصلی الگوریتم‌های داده کاوی، دسته‌بندی و خوشه‌بندی است.

---

<sup>۷</sup> Data cleaning

<sup>۸</sup> Data Integration

<sup>۹</sup> Data selection

<sup>۱۰</sup> Data transformation

<sup>۱۱</sup> Pattern evaluation

<sup>۱۲</sup> Knowledge presentation

یکی از دلایلی که باعث شد داده کاوی مورد توجه قرار بگیرد، مسأله در دسترس بودن حجم وسیعی از داده‌ها و نیاز شدید به اینکه از این داده‌ها، اطلاعات و دانش استخراج شود. داده‌کاوی عبارت است از استخراج دانش از مجموعه‌ای از داده‌ها.

### ۳-۲- دسته‌بندی

هرگاه داده‌ها دارای خصیصه‌ای خاص باشند که مستقیماً از دیگر خصایص به وجود نیامده باشد اما بین آن مشخصه و دیگر ابعاد رابطه وابستگی وجود داشته باشد، در این صورت می‌توان با کشف مدلی بر اساس دیگر مشخصه‌ها، آن بعد مذکور (که نشان دهنده دسته خاصی از داده‌ها است) را شناسایی نمود. این یک کار تحلیل داده و فرآیند پیدا کردن مدلی است که کلاس‌ها و مفاهیم داده‌ها را توصیف و متمایز می‌کند. دسته‌بندی مسأله شناسایی مجموعه‌ای از طبقات (زیر جمعیت‌ها) و یک مشاهده جدید براساس مجموعه آموزشی داده‌ها است [10].

فرض کنید که قسمتی از نوشته‌های نویسندگان در پایگاه داده‌ای وجود دارد که قبلاً با استفاده از آزمایش خاص سه نویسنده مشخص شده که هر کدام از این نوشته‌ها به کدام نویسنده تعلق دارد، در این جا هیچ نوشته‌ای حق ندارد به دو نویسنده تعلق داشته باشد، به این معنی که دسته‌ها فضای مسئله را افزای می‌کند. در چنین پایگاه داده‌هایی برای هر نوشته یک رکورد خاص وجود دارد و در نهایت نام یا برچسب نویسنده‌ای که آن نوشته به آن تعلق دارد، قرار گرفته است. یک داده کاو تصمیم می‌گیرد سیستمی را ابداع کند که طی آن بدون آزمایش و فقط از روی برخی ویژگی‌های نوشته‌ها بتوان نام نویسنده مربوطه را تشخیص داد. این تصمیم ممکن است به هر دلیلی صورت گرفته باشد. آنچه باید انجام شود عملیات دسته‌بندی نامیده می‌شود.

دو مرحله مهم دسته‌بندی عبارتند از: اولین مرحله ساخت مدل و دومین مرحله استفاده از مدل و پیش‌بینی کلاس از طریق مدل ساخته شده است [11]. بدین منظور باید مجموعه داده‌ها را به سه دسته داده‌های آموزش<sup>۱۳</sup> و داده‌های اعتبار سنجی<sup>۱۴</sup> و داده‌های تست<sup>۱۵</sup> تقسیم کنیم. با استفاده از داده‌هایی که برچسب آموزش خورده‌اند یک دسته‌بند ایجاد می‌شود که بر اساس آن بتوان داده‌های فاقد برچسب را در دسته‌های مربوط به خودشان قرار

---

<sup>۱۳</sup> Train data

<sup>۱۴</sup> Development data

<sup>۱۵</sup> Test data

داد. کارایی دسته‌بند ساخته شده با داده‌های تست و داده‌های اعتبار سنجی (که به صورت تصادفی از میان داده‌ها انتخاب شده‌اند) مورد سنجش قرار می‌گیرد و مدل روی آن‌ها اجرا می‌شود تا دقت پیش‌بینی دسته‌بند بررسی گردد، چنان‌که مدل دارای دقت مناسبی باشد برای دسته‌بندی داده‌ها به کار می‌رود.

## ۲-۴- الگوریتم‌های رایج دسته‌بندی

روش‌های زیادی برای دسته‌بندی وجود دارد که از جمله مواردی که ما در این پژوهش از آنها استفاده کرده ایم می‌توان به موارد زیر اشاره کرد:

۱. دسته‌بندی کننده بردار پشتیبان خطی<sup>۱۶</sup>

۲. شبکه‌های بیزین<sup>۱۷</sup>

۳. رگرسیون منطقی<sup>۱۸</sup>

۴. درخت‌های اضافی<sup>۱۹</sup>

در هر چهار الگوریتم استفاده شده ابتدا  $\text{Binary BOW}^{20}$ ، سپس  $\text{BOW}^{21}$  و در نهایت  $\text{TF-IDF}^{22}$  را محاسبه نموده و با استفاده از  $\text{TF-IDF}$  صد نمونه داده ساخته ایم.

## ۲-۴-۱- دسته‌بندی کننده بردار پشتیبان خطی

به صورت کلی ماشین‌های بردار پشتیبان ( $\text{SVM}^{23}$ ) یکی از الگوریتم‌های طبقه‌بندی است که فاقد پایه نظری

---

<sup>۱۶</sup> Linear SVC : Linear Support Vector Classifier

<sup>۱۷</sup> Naïve Base

<sup>۱۸</sup> Logistic regression

<sup>۱۹</sup> Extra trees

<sup>۲۰</sup> Binary Bag Of Word

<sup>۲۱</sup> Bag Of Word : کیسه کلمات

<sup>۲۲</sup> Term Frequency\_ Inverse Document Frequency

<sup>۲۳</sup> Support vector machine

میباشد و نسبت به تعداد ابعاد در مجموعه داده‌ها حساس است. زمانی که مجموعه داده‌ها بصورت خطی قابل تفکیک باشند، ماشین‌های بردار پشتیبان صفحه ابر را پیدا می‌کنند که به بهترین شکل از هم جدا می‌شوند. این الگوریتم یک جداکننده است که با معیار قرار دادن بردارهای پشتیبان، بهترین دسته بندی و تفکیک بین داده‌ها را برای ما مشخص می‌کند. بردارهای پشتیبان به زبان ساده، مجموعه‌ای از نقاط در فضای  $n$  بعدی داده‌ها هستند که مرز دسته‌ها را مشخص می‌کند و مرزبندی و دسته بندی داده‌ها براساس آنها انجام می‌شود و با جابجایی یکی از آنها، خروجی دسته بندی ممکن است تغییر کند. در فضای دو بعدی بردارهای پشتیبان، یک خط، در فضای سه بعدی یک صفحه و در فضای  $n$  بعدی یک ابر صفحه را شکل خواهند داد. در این الگوریتم فقط داده‌های قرار گرفته در بردارهای پشتیبان مبنای یادگیری ماشین و ساخت مدل قرار می‌گیرند و این الگوریتم به سایر نقاط داده حساس نیست و هدف آن هم یافتن بهترین مرز در بین داده‌هاست به گونه‌ای که بیشترین فاصله ممکن را از تمام دسته‌ها (بردارهای پشتیبان آنها) داشته باشد.

حال در SVC ها با کرنل پارامتر "خطی"، انعطاف‌پذیری بیشتری دارد و باید نسبت به تعداد زیادی از نمونه‌ها مقیاس داشته باشد.

این کلاس هم از ورودی متراکم و هم پراکنده پشتیبانی می‌کند و پشتیبانی چندکلاسه<sup>۲۴</sup> با توجه به یک برنامه تک مرحله‌ای انجام می‌شود.

در SVM ها، مدل چند کلاسی، با استفاده از یک طرح در مقابل یک برنامه اجرا می‌شود در حالی که LinearSVC از یکی در مقابل بقیه استفاده می‌کند که در نهایت اگر ورودی پیوسته باشد یک svc می‌تواند با داده‌های متراکم هماهنگ باشد.

در حالت کلی کد استفاده شده در الگوریتم linear svc اینگونه عمل مینماید که ابتدا تابعی را تعریف نموده سپس الگوریتم را بر اساس پارامترهای خاص خود تشکیل می‌دهد. ما با استفاده از دستور pickle کلیه‌ی مدل‌های ساخته شده را ذخیره کرده و در مرحله بعد محاسبات از قبیل دقت (precision\_recall\_fscore<sup>۲۵</sup>) انجام شده و نتیجه و اینکه میانگیری به چه صورتی باید باشد را نمایش می‌دهد.

---

<sup>۲۴</sup> Multiclass

<sup>۲۵</sup> معیارهای ارزیابی امتیازدهی و دقت و صحت

#### ۱-۱-۴-۲- مزایا:

- در فضاهایی با بعد زیاد، خوب جواب می دهد. یعنی در جاهایی که تعداد ویژگی هایمان زیاد باشد این الگوریتم بسیار مؤثر است.
- در زمانی که تعداد بعد بیشتر از تعداد نمونه ها باشد مؤثر عمل می کند.
- در بعدهای زیاد و همچنین در زمانی که تعداد داده های آموزشی مان نسبت به بعد کم باشد بسیار مؤثر است.
- توابع تصمیم گیری<sup>۲۶</sup> مختلفی دارد که در زیرمجموعه آموزشی<sup>۲۷</sup> در تابع تصمیم گیری استفاده می شود.
- حافظه قوی دارد.
- کرنل، تابع هایی که دارد در انتخاب تصمیم مفید است.

#### ۲-۱-۴-۲- معایب:

- اگر تعداد feature ها و مشخصه هایمان، تعداد ابعاد xtrain و ytrain بیشتر از تعداد sample ها باشد برای جلوگیری از overfitting انتخاب نوع function و c بسیار حساس و مهم هستند که احتمال دارد به مشکل برخورد کنیم. پس زمانی که تعداد مشخصه هایمان بسیار بیشتر از تعداد نمونه هایمان باشد یا اینکه ابعادمان بیشتر از نمونه ها باشد، مفید است.
- سه تا کرنل مهم دارد:
  ۱. کرنل های خطی: که با خط راست دیتاها را جدا میکند.
  ۲. کرنل های قوسی یا rbf: که با تابع های شعاعی (قوسی) دیتاها را از هم جدا می کند. یعنی یک سری موارد مانند نمودارهای قوسی باعث جداکردن دیتا می شود.
  ۳. کرنل های چندجمله ای: که با خطوط چندجمله ای دیتاها را از هم جدا می کند.
- ایراد دیگری که دارد این است که تخمین احتمالات برایش مهم نیست و با استفاده از five-fold-

---

<sup>۲۶</sup> Decision function

<sup>۲۷</sup> Train Set

<sup>۲۸</sup> cross-validation کار می کند و بر اساس تابع تخمین اطلاعات کار نمی کند.

۳-۱-۲- پارامترهای مهم *linear SVC*

• `penalty:string, 'L1' or 'L2' (default='L2')`

نوعی پناستی است که به دو نوع  $L1$  و  $L2$  می باشد. حالت پیش فرض آن  $L2$  است و از آن بهترین جواب بدست می آید. اگر  $L2$  را جایگزین کنیم بهترین جواب را ( $ACCURACY^{۲۹}=0.777$ ) خواهیم داشت.

• `C:float, optional (default=1.0)`

فرمول محاسبه  $C$  به صورت زیر می باشد:

$$C \sum_{i=1}^n \mathcal{L}(f(x_i), y_i) + \Omega(\omega) \quad (۲-۱)$$

مقدار ذاتی  $C$  همان یک است. هرچقدر  $C$  بیشتر باشد، حول همسایه ها می چرخد و اگر سطحی که در می آید کمتر باشد، سطح تخت می شود (هرچه بیشتر باشد سطح خروجی پستی و بلندی بیشتری دارد). مقدار  $C$  هر چقدر کمتر باشد میزان همواری بیشتر است. مقدار این پارامتر بسته به نوع دیتاستی است که ما در اختیار داریم.

• **dual**: در این الگوریتم مسئله ما مقدار **dual** نیست یعنی این معیار برای پارامترهای دوجوهری است اما در اینجا ما یک مجهول داریم و نیازی به تغییر این پارامتر نمی باشد.

• **tol** چیست؟ **tol** یعنی چقدر تابع **optimization** مان یعنی **svc** برای بهینه سازی است. این تابع بر اساس بهینه سازی عمل می کند. حال باید دید اگر چقدر تغییر کند ما آن را قبول می کنیم؟ یعنی از آن حد کمتر تغییر کند سیستم را **break** می کند. اگر عدد به اندازه  $tol=0.0025$  کوچک باشد دقت بیشتر می شود و جواب بهتری خواهیم گرفت اما زمان محاسبات بالاتر خواهد رفت.

یکی دیگر از پارامترهایی که می تواند **svc** بگیرد به شکل زیر است:

• `loss:string, 'hinge' or 'squared_hinge' (default='squared_hinge')`

که تعیین کننده تابع **loss** است. همانطور که در فرمول  $C$  آمده،  $L$  یا همان **loss** پارامتری است که نسبت  $f(x)$  به  $y$  را تعیین می کند. دیفالت تابع **loss** بر اساس میانگین مربعات است. اگر بر اساس میانگین مربعات یا توان ۲ مربعات نخواستیم می توانیم از **hinge** استفاده کنیم که تنها استاندارد را انجام بدهد.

<sup>۲۸</sup> اعتبارسنجی متقابل پنج برابر

<sup>۲۹</sup> صحت



- `dual:bool,(default=true)`

زمانی که بیش از یک پارامتر بهینه سازی داریم یعنی دو یا بیشتر پارامتر بهینه سازی داریم.

- `tol:float,optinal(default=1e-4)`

که همان `tolerance` است که به عنوان یک مدرک `stopping` مورد استفاده قرار می گیرد.

- `c:float,optional(default=1.0)`

`Float c` که یک پارامتر پینالتی برای `error` است.

- `multi_class:string 'ovr' or 'crammer_singer'(default='ovr')`

چند کلاسه برای وقتی است که کلاس دو دسته نباشد و کلاس ها از چند دسته متفاوت باشند . مثلاً دسته یک، دسته دو و .... از چند کلاسه بر اساس `one-vs-res` (یکی و بقیه) استفاده می شود . وقتی از `multi class` `system` استفاده می کند یعنی اینگونه عمل می کند که مثلاً آنهایی که پارامتر یک دارند را یک می کند و بقیه را صفر بعد می آید آنهایی که پارامتر دو دارند را یک می کند و بقیه را صفر و به همین ترتیب ادامه می دهد تا آخر و بر اساس همین روش، داده ها را آموزش می دهد.

- `fit_intercept:Boolean,optional(default=true)`

که می توان گفت بر اساس بولین باشد یا نه که معمولاً بر اساس بولین عمل می کند.

- `max_iter:int,(default=1000)`

این جزو مهمترین پارامترها است . بیشترین میزان تکرار را نشان می دهد . اگر پارامتر `tol` با مقدار مشخص شده فوق به نتیجه نرسد (یعنی میزان تخت بیشتر از حد باشد) در اینجا `max-iter` می تواند آن را بشکند . هم برای اینکه میزان تخت بودن کم باشد و هم برای جوابهای بهتر استفاده می شود و زمان محاسباتی بالا است . باید هم `max-iter` عدد بالایی باشد و هم `desineflprance` عدد کمتری باشد ، زمان محاسباتی بالاتر اما جواب بهتر و درست تری به ما می دهد . پارامترهای دیگر تأثیر چندانی در جواب بهینه ندارند.

حال ما با تنظیم پارامترها به صورت (`C=2.5, penalty="l2", dual=False, tol=0.0025`) توانستیم صد مدل با

استفاده `train` با استفاده از `tfidf` بسازیم که خروجی آنها همانند جدول ذیل می باشد [13][12]:

Bow binary=f-clf0=0.7676200204

Bow=f-clf1=0.6755771195

Tfidf=fclf2=0.7946884576

f-clf	مقدار
۳	۰,۷۷۳۲۳۷۹۹۷۹

۴	۰,۶۷۸۱۳۰۷۴۵۶
۵	۰,۷۶۳۵۳۴۲۱۸۵
۶	۰,۷۷۴۲۵۹۴۴۸۴
۷	۰,۷۶۷۶۲۰۰۲۰۴
۸	۰,۷۶۵۰۶۶۳۹۴۲
۹	۰,۷۶۶۵۹۸۵۶۹۹
۱۰	۰,۷۷۲۲۱۶۵۴۷۴
۱۱	۰,۷۷۱۷۰۵۸۲۲۲
۱۲	۰,۷۷۴۷۷۰۱۷۳۶
۱۳	۰,۷۷۱۷۰۵۸۲۲۲

## ۲-۴-۲- شبکه بیزین

دومین الگوریتمی که در این پژوهش از آن جهت دسته بندی استفاده شده است شبکه بیزین می باشد . در روش های دسته بندی آماری برخلاف سایر دسته بندها میزان عضویت یک نمونه به هر کلاس را با یک احتمال نشان می دهند. روش شبکه های بیزین رایج ترین روش دسته بندی آماری و از روش های ساده و موثر محسوب می شود . در این روش احتمال شرطی هر صفت داده شده را توسط برچسب دسته مربوطه از داده های آموزشی یاد می گیرید . سپس عمل دسته بندی توسط بکار بردن قوانین بیز برای محاسبه مقدار احتمالی دسته نتیجه نمونه داده شده با دقت بالایی انجام می شود . در حالت معمولی این کار با تخمین احتمالاتی هر ترکیب ممکن از صفات صورت می گیرد ولی هنگامی که تعداد صفات خیلی زیاد باشد، این امر امکان پذیر نیست. بنابراین یک فرض مستقل قوی اتخاذ می شود که همه صفات با مشخص بودن مقدار صفت دسته مستقل می باشند. با در نظر گرفتن این فرض لازم است که فقط احتمالات حاشیه ای هر صفت دسته محاسبه گردد. با این حال این فرض به صورت غیرواقعی می باشد و شبکه های بیزین با مدل کردن صریح، وابستگی بین صفات آن را در نظر نمی گیرند [14].

مسئله یادگیری ساختار شبکه بیزین به این صورت بیان می شود که با داشتن یک مجموعه آموزشی  $A = \{u_1, u_2, \dots, u_n\}$  از  $n$  نمونه  $u$ ؛ یک شبکه پیدا کنیم که بهترین تطبیق را با  $A$  داشته باشد . معمول ترین روش

برای این مسأله معرفی یک تابع هدف است که هر شبکه با توجه به داده‌های آموزشی و جستجوی بهترین شبکه بر اساس این تابع ارزیابی شود. چالش‌های بهینه‌سازی کلیدی انتخاب تابع هدف و تعیین روال جستجو برای بهترین شبکه می‌باشد.

تئوری بیز یکی از روش‌های آماری برای رده‌بندی به شمار می‌آید. در این روش کلاس‌های مختلف، هر کدام به شکل یک فرضیه دارای احتمال در نظر گرفته می‌شوند. هر رکورد آموزشی جدید، احتمال درست بودن فرضیه‌های پیشین را افزایش و یا کاهش می‌دهد و در نهایت، فرضیاتی که دارای بالاترین احتمال شوند، به عنوان یک کلاس در نظر گرفته شده و برچسبی بر آن‌ها زده می‌شود. این تکنیک با ترکیب تئوری بیز و رابطه سببی بین داده‌ها، به دسته‌بندی می‌پردازد. روشی برای دسته‌بندی پدیده‌ها، بر پایه احتمال وقوع یا عدم وقوع یک پدیده است و در نظریه احتمالات با اهمیت و پرکاربرد است. اگر برای فضای نمونه‌ای مفروضی بتوانیم چنان افزایی انتخاب کنیم که با دانستن اینکه کدامیک از پیشامدهای افزاز شده رخ داده است، بخش مهمی از عدم قطعیت تقلیل می‌یابد. این قضیه از آن جهت مفید است که می‌توان از طریق آن، احتمال یک پیشامد را با مشروط کردن نسبت به وقوع یا عدم وقوع یک پیشامد دیگر محاسبه کرد. در بسیاری از حالت‌ها، محاسبه احتمال یک پیشامد به صورت مستقیم کاری دشوار است. با استفاده از این قضیه و مشروط کردن پیشامد مورد نظر نسبت به پیشامد دیگر، می‌توان احتمال مورد نظر را محاسبه کرد.

این رابطه به خاطر بزرگداشت توماس بیز فیلسوف انگلیسی به نام فرمول بیز معروف است.

به طور خلاصه قضیه بیز به صورت زیر عمل می‌نماید:

زمانی که از قبل وقوع یک پیشامد تصادفی را بدانیم، به کمک فرمول‌های احتمال شرطی می‌توانیم مقدار احتمال برای هر پیشامد دیگر را محاسبه کنیم.

۴-۲-۲- پارامترهای  $NB$ :

• `priors:array_lik,shape(n_class)`

ورودی‌هایی که دارد ورودی‌های خاصی هستند. ورودی‌هایی با عنوان توزیع احتمالات کلاس‌های قبلی است.

• `var_smoothing:float.optional(default=1e-9)`

پارامتر هموار سازی است به نوعی پارامتر هموار سازی هر چقدر بیشتر باشد، بهتر است. طبق دیتاست جواب بهینه به ما خواهد داد. اما نسبت بیشترین `variance` به عنوان پارامتر پایداری در نظر می‌گیرد. دقت کنید که

هرچقدر پارامتر پایداری بیشتر باشد سیستم smoot تر می شود (همان مقدار  $1e-9$ ) [1][15][16]  
 ما با استفاده از این الگوریتم صد مدل داده train با استفاده از tfidf ساخته ایم که یک نمونه از آن به شرح جدول ذیل میباشد:

Bow binary=f-clf103=0.6603677221

Bow=f-clf104=0.6613891726

Tfidf=f-clf105=0.6736465781

f-clf	مقدار
106	0.6068028600
107	0.6869254341
108	0.6675178753
109	0.6741573033
110	0.6603677221
111	0.6680286006
112	0.6527068437
113	0.6690500510
114	0.6721144024
115	0.6629213483
116	0.6690500510

## ۲-۴-۳- رگرسیون منطقی

یک مدل آماری رگرسیون برای متغیرهای وابسته دوسویی مانند بیماری یا سلامت، مرگ یا زندگی است. این مدل را می توان به عنوان مدل خطی تعمیم یافته ای که از تابع لججیت به عنوان تابع پیوند استفاده می کند و خطایش از توزیع چندجمله ای پیروی می کند، به حساب آورد. منظور از دوسویی بودن، رخ داد یک واقعه تصادفی در دو موقعیت ممکنه است. به عنوان مثال خرید یا عدم خرید، ثبت نام یا عدم ثبت نام، ورشکسته شدن یا ورشکسته نشدن و ... متغیرهایی هستند که فقط دارای دو موقعیت هستند و مجموع احتمال هر یک آن ها در نهایت یک خواهد شد.

این الگوریتم برای برآورد مقادیر گسسته (مقادیر باینری مانند ۰/۱، آری/خیر، درست/نادرست) بر اساس مجموعه ای از متغیر (های) وابسته استفاده می شود. به بیان ساده تر این الگوریتم احتمال رخداد یک رویداد را بر حسب گنجانیدن داده ها در یک تابع منطقی پیش بینی می کند. از این رو به نام رگرسیون logit نیز شناخته می شود. از آنجا که این الگوریتم، احتمال را پیش بینی می کند، مقادیر خروجی آن بین ۰ تا ۱ هستند

اگر از دید ریاضیاتی به این الگوریتم نگاه کنیم خروجی به صورت یک ترکیب خطی از متغیرهای پیش‌بین مدل‌سازی شده است.

روش‌های سنتی شناسایی نویسنده نیازمند یک روش یادگیری دسته‌بندی ضعیف و پیاده‌سازی بسیار مقیاس پذیر هستند. در مقابل، می‌توان به دلیل ماهیت احتمالی آن به رگرسیون منطقی چند جمله‌ای رفت. از آنجا که این مدل تخمینی از احتمال تعلق ورودی به هر یک از کلاس‌های ممکن را نشان می‌دهد، می‌توان به راحتی هزینه‌های نسبی واحدهای مختلف را هنگام تصمیم‌گیری طبقه‌بندی در نظر گرفت. اگر این هزینه‌ها تغییر کنند، طبقه‌بندی می‌تواند بدون آموزش مجدد مدل تغییر کند.

به علاوه، دیدگاه Bayesian درباره آموزش یک مدل رگرسیون چند جمله‌ای باعث می‌شود اطلاعات آموزشی و دانش حوزه به راحتی ترکیب شوند. رگرسیون منطقی به طور گسترده در طبقه‌بندی متن استفاده شده است. در حالی که طرفدار ارزش‌های نزدیک است، موافق نیست که دقیقاً برابر با صفر باشد. از آنجا که مدل‌های رگرسیون چند جمله‌ای برای شناسایی مولف می‌تواند به راحتی میلیون‌ها پارامتر داشته باشد، چنین تخمین پارامتر متراکم می‌تواند منجر به طبقه‌بندی کننده‌های ناکارآمد شود. با این حال برآوردهای پارامترهای پراکنده را می‌توان به سادگی در چارچوب Bayesian به دست آورد [15][17].

ما با استفاده از این الگوریتم صد مدل داده train با استفاده از tfidf ساخته ایم که یک نمونه از آن به شرح جدول ذیل می‌باشد:

$$\text{Bow binary} = f\text{-clf}_{206} = 0.7967313585$$

$$\text{Bow} = f\text{-clf}_{207} = 0.7941777323$$

$$\text{Tfidf} = f\text{-clf}_{208} = 0.8105209397$$

f-clf	مقدار
209	0.7860061287
210	0.7773237997
211	0.7808988764
212	0.7900919305
213	0.7911133810
214	0.7839417773
215	0.7819203268
216	0.7911133810
217	0.7921348314

218	0.7951991828
219	0.7870275791

## ۲-۴-۴- درخت های اضافی

یک طبقه بندی کننده درخت اضافی است . این کلاس یک  $\text{estimator}^{30}$  متا را پیاده می کند که با تعدادی از درخت های تصمیم گیری تصادفی متناسب است . که در زیر مجموعه داده ها چند نمونه از مجموعه داده ها استفاده شده است و از میانگین گیری برای بهبود دقت پیش بینی و کنترل بیش از حد استفاده می کند .

این الگوریتم یک طبقه بندی کننده درخت بسیار تصادفی است . درخت های اضافی از درخت های تصمیم گیری کلاسیک تفاوت دارند . هنگامی که به دنبال بهترین شکاف برای جدا کردن نمونه های یک نود به دو گروه بگردید ، splits تصادفی برای هر کدام از ویژگی ها انتخاب می شوند که به طور تصادفی ویژگی ها را انتخاب می کنند و بهترین شکاف بین آن ها انتخاب می شود. حداکثر ویژگی ها ۱ تنظیم می شود ، این مقدار برای ساخت یک درخت تصمیم گیری کاملاً تصادفی است .

هشدار : درختان اضافی تنها باید در روش های گروهی مورد استفاده قرار گیرند .

مقادیر پیش فرض برای پارامترهایی که اندازه درختان را کنترل می کنند ( به عنوان مثال عمق ، نمونه \_ نمونه های ، برگ و غیره ) منجر به درخت های کاملاً بالغ و unpruned می شوند که می تواند به طور بالقوه در برخی از مجموعه های داده بسیار بزرگ باشد . برای کاهش مصرف حافظه ، پیچیدگی و اندازه درختان باید با تنظیم مقادیر پارامتر کنترل شوند [18][19].

ما با استفاده از این الگوریتم صد مدل داده train با استفاده از tfidf ساخته ایم که به شرح جدول ذیل میباشد:

Bow binary=f-clf308= 0.65985699

Bow=f-clf309=0.6813074565

Tfidf=f-clf310=0.6884576098

f-clf	مقدار
311	0.6716033677
312	0.6894790602

<sup>30</sup> تخمین زننده ، برآورد گر

313	0.6782431052
314	0.6710929519
315	0.6705822676
316	0.6797752808
317	0.6782431052
318	0.6889683350
319	0.6613891726
320	0.6552604698
321	0.6552604698

## ۲-۵- یادگیری تجمعی

یکی از اهداف اولیه در داده کاوی، پیش بینی یک مقدار نامعلوم از یک نمونه جدید براساس نمونه های مشاهده شده قبلی است. دستیابی به چنین نتیجه ای با دو گام حاصل می شود:

الف) مرحله آموزش:

ایجاد یک مدل پیش بینی از نمونه های آموزشی با استفاده از یکی از الگوریتم های یادگیری باناظر ب) مرحله آزمون:

در گام بعدی مدل ایجاد شده با استفاده از نمونه های آزمایش، مدل مورد ارزیابی قرار میگیرد. تجربیات گوناگون از بکارگیری روشهای آموزش نشان میدهد که هیچ تک الگوریتم آموزشی مشخصی وجود ندارد که بتواند برای تمامی کاربردها بهترین و دقیقترین باشد. در واقع هر الگوریتم، مدل خاصی است که بر پایه مفروضات معینی شکل گرفته است. برخی مواقع این مفروضات برقرار بوده و گاهی نیز نقض میشوند. بنابراین، هیچ الگوریتمی به تنهایی نمیتواند در تمامی شرایط و برای همه بطور موفق عمل نماید. برای رفع این مشکل روش آموزش دسته جمعی معرفی شده است.

انگیزه اصلی بر توسعه چنین روشی، کاهش نرخ خطا می باشد. فرض مبنایی این متدولوژی آن است که در حالت دسته جمعی احتمال اشتباه در تشخیص دسته یا جایگاه یک نمونه جدید خیلی کمتر از حالت پیش بینی با تنها یک مدل میباشد.

این ایده توسط هنسن بصورت نظری به ترتیب زیر اثبات شده است:

اگر  $N$  دسته بند مستقل با احتمال خطای  $e > 0.5$  داشته باشیم، در این صورت میتوان نشان داد که خطای دسته جمعی  $E$  بر حسب  $N$  بصورت یکنواخت کاهش پیدا میکند. همچنین اگر از دسته بندهای وابسته استفاده کنیم، عملکرد کلی شدت کاهش می یابد.

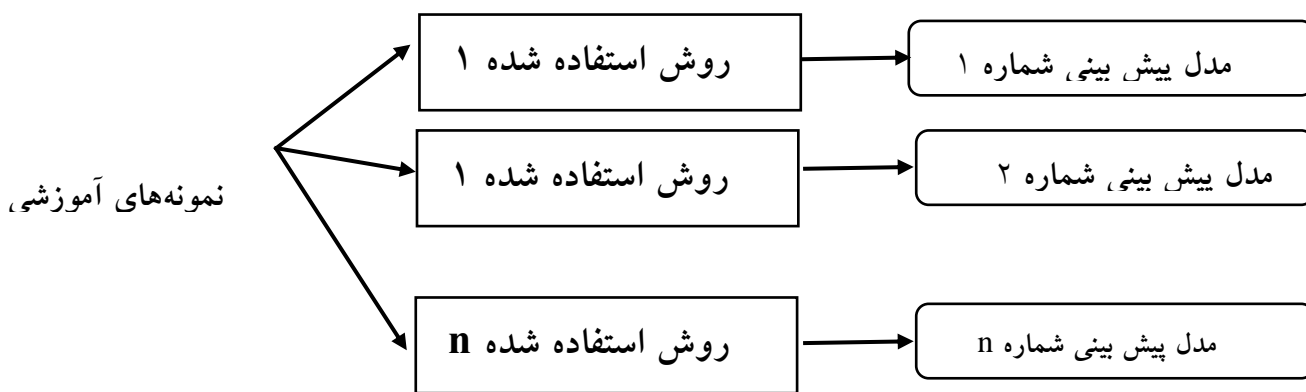
#### ۱-۵-۲- متدولوژی آموزش دسته جمعی

متدولوژی آموزش دسته جمعی از دو مرحله متوالی تشکیل شده است :

الف) مرحله آموزش

ب) مرحله آزمایش

در مرحله آموزش روش دسته جمعی، چندین مدل پیش بینی با استفاده از نمونه های آموزش تولید میشوند .



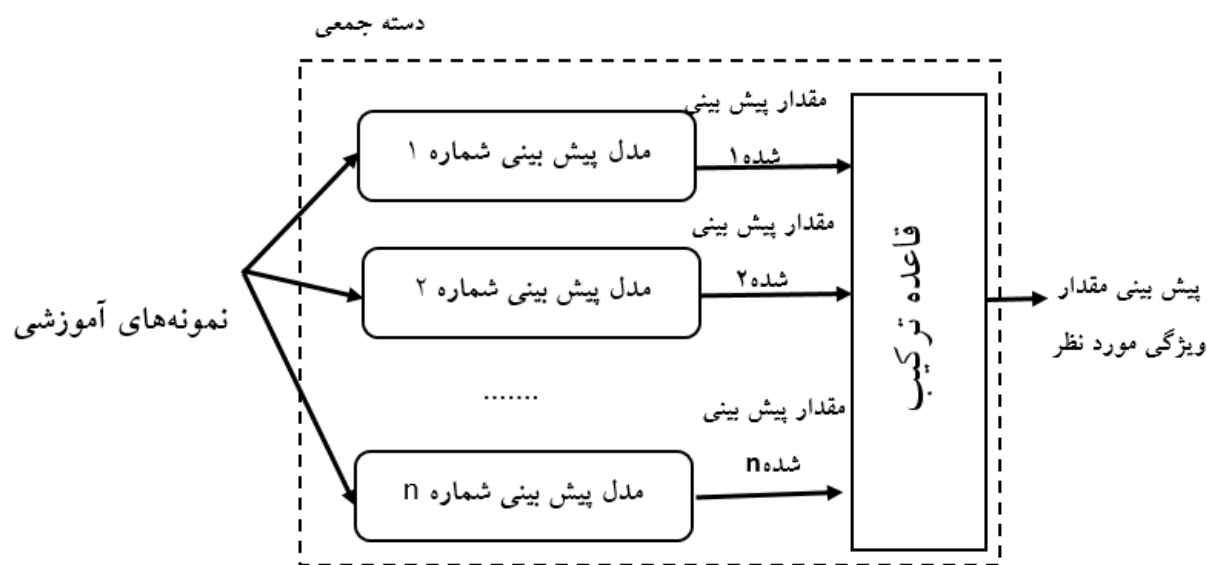
(الف)

۲- ۱- مرحله آموزش روش دسته جمعی

در مرحله بعد، برای پیش بینی دسته یا جایگاه نمونه آزمایشی، روش دسته جمعی، خروجی تک تک مدل های

پیش بینی را محاسبه و با یکدیگر تجمیع می کند .





(ب)

۲-۲- مرحله آزمایشی روش دسته جمعی

برای اینکه مدل دسته جمعی عملکرد بهتری نسبت به مدل تکی داشته باشد، باید هر کدام از مدل های تشکیل دهنده آن مستقل از بقیه بوده و دارای دقتی بزرگتر از ۰,۵ باشد .

شرط استقلال مدل ها مانع از همبسته شدن خطای مدل ها خواهد شد . در نهایت برای تشخیص دسته یا جایگاه نمونه آزمایشی، خروجی همه مدل ها با یکدیگر تجمیع می شوند .مدل دسته جمعی، تصمیم گیری در مورد دسته نمونه آزمایشی را با توجه به بیشترین تکرار در میان خروجی دسته بندها، انجام میدهد[23].

## ۲-۶- انواع روش ترکیب دسته بندها

برای ترکیب دسته بندها روش های متفاوتی وجود دارد، در این پژوهش روش های ترکیب زیر مورد بررسی قرار داده شد.

### ۲-۶-۱- ترکیب همه دسته بندها

در این نوع ترکیب برای هر نمونه با توجه به صورت مسئله تعدادی دسته بند وجود دارد که با توجه به خروجی هر دسته بند، خروجی که بیشترین تکرار را داشته باشد به عنوان خروجی نهایی در نظر گرفته می شود.

## ۲-۶-۲- انتخاب رو به جلو<sup>۳۱</sup>

انتخاب رو به جلو یک نوع رگرسیون گام به گام است که با یک مدل خالی شروع می شود و دسته بندهایی را که به صورت نزولی بر مبنای یکی از معیارهای ارزیابی مرتب شده اند را یک به یک اضافه می کند و بهترین متغیر، با برخی از معیارهای از پیش تعیین شده تعیین می شود و به مدل اضافه می گردد. معیار استفاده شده برای تعیین اینکه کدام مدل اضافه شود، متفاوت است. در هر مرحله رو به جلو، شما یک دسته بند را اضافه می کنید که بهترین نتیجه را برای شما به ارمغان می آورد این یکی از دو روش معمول استفاده از رگرسیون گام به گام است. تغییر بیش از حد زمانی اتفاق می افتد که ما متغیرهای بیشتری را در اختیارمان قرار دهیم تا اینکه برای مدل مناسب باشد. به طور معمول یک اطلاعات بسیار دقیق از داده های مورد استفاده در رگرسیون را نشان می دهد، اما این مدل از نقاط داده های اضافی دور خواهد ماند و برای ایتترلوسیون خوب نیست.

در این روش مدل ها بر مبنای  $F\_score$  به صورت نزولی مرتب شد، در ابتدا اولین مدل مورد ارزیابی قرار گرفت و سپس مدل دوم با مدل اول ترکیب شد و ترکیب این دو، مورد ارزیابی قرار گرفت اگر  $f\_score$  ترکیب این دو مدل بهتر از مدل اول به تنهایی باشد که ترکیب دو مدل را ذخیره می کنیم و مدل بعد را اضافه می کنیم در غیر این صورت مدل دوم حذف می شود و سومین مدل ترکیب می شود، به همین ترتیب تا آخرین مدل (۱۰۰ مدل) انتخاب گردید.

## ۲-۶-۳- حذف رو به عقب<sup>۳۲</sup>

در این روش ابتدا مدل ها را بر مبنای  $F\_score$  به صورت صعودی مرتب کردیم سپس همه ی مدل ها با هم ترکیب شد. ترکیب مدل ها ارزیابی شد و اولین مدل را حذف کردیم و نتیجه را مورد ارزیابی قرار دادیم به همین ترتیب تا آخرین مدل ادامه دادیم [24].

---

<sup>۳۱</sup> Forward selection

<sup>۳۲</sup> Backward elimination

## ۷-۲- معیارهای ارزیابی دسته بندها

معیارهای زیادی برای ارزیابی کارایی الگوریتم‌های دسته‌بندی ارائه می‌شود که مهم‌ترین این معیارها عبارتند از: نرخ صحت<sup>۳۳</sup>، دسته‌بندی، سرعت یادگیری در مرحله آموزش و دسته‌بندی، سادگی و شفافیت مدل، پایداری (توانایی مدل در مواجهه با داده‌های غیر معمول یا مقادیر مفقوده)، نحوه برخورد الگوریتم با صفت‌ها با دامنه مقادیر مختلف (پیوسته گسسته و دودویی) و قابلیت تفسیر.

در این مرحله از مدل‌های ساخته شده در بخش قبلی روی داده‌های توسعه ارزیابی انجام گرفت، برای هر مدل<sup>۳۴</sup> F\_score، precision<sup>۳۵</sup> و recall<sup>۳۶</sup> محاسبه شده و مدل‌ها در یک آرایه بر مبنای F\_score به صورت نزولی مرتب و بهترین مدل انتخاب شد که در ذیل به روشهای ارزیابی اشاره ای خواهیم داشت.

دقت و صحت عبارت‌هایی هستند که برای بیان نتیجه اندازه‌گیری به وفور مورد استفاده قرار می‌گیرند. در گفتگوهای روزانه، بسیار پیش می‌آید که این دو عبارت را به اشتباه به جای هم به کار می‌بریم. در حالی که هریک مفهوم مشخصی دارند و باید در موقعیت مناسب خود استفاده شود.. برای شروع، بهتر است هریک از این اصطلاحات را جداگانه تعریف کنیم.

### ۱-۷-۲- نگاهی دقیق‌تر به مسائل دسته‌بندی

در مسائل دسته‌بندی به دنبال پیش‌بینی دسته یا گروهی هستیم که یک رکورد یا داده به آن تعلق دارد. مثلاً می‌خواهیم با داشتن مشخصات مختلف یک دانشجو مانند رشته، سهمیه قبولی، وضعیت مالی خانواده، تعداد هم‌اتاقی‌ها، وضعیت تحصیلی دوران دبیرستان، تعداد واحدهای گذرانده، معدل ترم‌های پیش، تعداد ترم‌های مشروطی و مانند آن تعیین کنیم که این دانشجو مشروط خواهد شد یا نه. یا به عبارتی به دسته مشروطی‌ها تعلق دارد یا دسته دانشجویان غیر مشروط.

اینکه احتمال دیابت داشتن یک بیمار با توجه به آزمایش‌های انجام شده چقدر است یعنی این بیمار جزء دسته دیابتی‌ها قرار می‌گیرد یا افراد غیر دیابتی، باز هم به یک مدل یا الگوریتم دسته‌بندی نیاز داریم.

---

<sup>۳۳</sup> Accuracy

<sup>۳۴</sup> نرخ امتیاز

<sup>۳۵</sup> دقت

<sup>۳۶</sup> یادآوری

همانطور که از دو مثال فوق مشاهده می‌کنید، دسته‌بندی‌های دوگانه بیشترین کاربرد را در دنیای واقعی دارند و ما هم برای ساده کردن مطلب، تمرکز را بر این نوع دسته‌بندی خواهیم گذاشت اما تمامی موارد بیان شده، به راحتی قابل تعمیم به دسته‌بندی‌های چندگانه هم خواهند بود.

بعد از ساخت یک مدل دسته‌بندی یعنی یافتن الگوریتمی که با مشاهده یک داده جدید، دسته یا گروه (Class) آنرا مشخص کند عمل (Classification<sup>۳۷</sup>)، برای سنجش میزان کارایی و دقت مدل پیشنهادی، آنرا بر روی داده‌های آموزشی یعنی داده‌هایی که از قبل دسته‌بندی آنها را می‌دانیم، اعمال می‌کنیم. خروجی این مدل را با خروجی واقعی به صورت زیر مقایسه خواهیم کرد:

		مقایسه واقعی	
		مثبت	منفی
مقادیر پیش‌بینی شده	مثبت	درست مثبت TP	نادرست مثبت FP خطای نوع یک
	منفی	نادرست منفی FN خطای نوع دو	درست منفی TN

## ۲-۳- ارزیابی کارایی یک مدل

به این جدول مهم که اساس تحلیل و ارزیابی کارایی یک مدل در مباحث دسته‌بندی است، ماتریس درهم‌ریختگی یا اغتشاش گفته می‌شود. اما با توجه به نامفهوم بودن این ترجمه، از عبارت ماتریس پراکنش که پراکندگی توزیع دسته‌ها را از لحاظ درستی یا نادرستی نمایش می‌دهد، در ادامه این سلسله آموزش‌ها برای اشاره به این ماتریس استفاده خواهد شد.

تعداد سطرها و ستون‌های این ماتریس به تعداد دسته‌ها وابسته است. جدول فوق یک ماتریس دو در دو را برای یک مساله دسته‌بندی دوگانه مثلاً دیابت داشتن یا نداشتن، مشروط شدن یا نشدن، ترک کار یا عدم ترک کار نشان می‌دهد. در حالت کلی برای یک مساله دسته‌بندی  $n$  حالت، یک ماتریس پراکنش  $n \times n$  باید رسم شود. در حالت کلی به دنبال این هستیم که بخشهای قرمز رنگ این ماتریس یعنی پیش‌بینی‌های نادرست، به حداقل برسد.

<sup>۳۷</sup> دسته بندی

با فرض اینکه هدف ما پیش‌بینی دیابت یک بیمار باشد یعنی اگر پیش‌بینی مثبت باشد یعنی بیمار، مبتلا به دیابت است و اگر پیش‌بینی منفی باشد، یعنی بیمار به دیابت مبتلا نیست، به تحلیل سلول‌های این ماتریس می‌پردازیم :

- درست مثبت<sup>۳۸</sup> (True Positive-TP) اگر بیمار واقعا دیابت داشته باشد و مقدار پیش‌بینی شده هم دیابت را نشان دهد.

- نادرست مثبت<sup>۳۹</sup> (FP) اگر بیمار دیابت نداشته باشد اما نتیجه پیش‌بینی ما، نشانگر دیابت بیمار باشد.

- نادرست منفی<sup>۴۰</sup> (FN) اگر بیمار دیابت داشته باشد اما پیش‌بینی ما، دیابت را منفی نشان دهد.

- درست منفی<sup>۴۱</sup> (TN) اگر بیمار دیابت نداشته باشد و پیش‌بینی ما هم همین را نشان بدهد.

همانطور که مشخص است ایده آل ما این است که موارد نادرست (نادرست مثبت و نادرست منفی) صفر باشند اما در عمل این اتفاق نمی‌افتد و نیازمند مکانیزم‌ها و معیارهایی برای بررسی دقت و صحت و کارایی مدل ایجاد شده از داده‌ها هستیم.

## ۲-۷-۲- دقت - صحت - بازخوانی

اولین معیار یا سنجه‌ای که به ذهن مان می‌رسد، معیار دقت یا میزان تشخیص درست مدل است. یعنی نسبت تشخیص‌های درست (TP+TN) به کل داده‌ها

$$\text{دقت (Accuracy)} = \frac{TP+TN}{N} = \frac{\text{تشخیص‌های درست}}{\text{کل داده‌ها}}$$

---

<sup>۳۸</sup> True Positive

<sup>۳۹</sup> False Positive

<sup>۴۰</sup> False Negative

<sup>۴۱</sup> True Negative

برای بسیاری از مسائل دسته‌بندی دنیای واقعی این معیار، بسیار کارآمد است چون هم داده‌های در نظر نگرفته شده را لحاظ کرده است (مخرج کسر) و هم داده‌های شناسایی شده را (صورت کسر). هدف ما هم رسیدن این عدد به مقدار یک یا همان صد در صد است. اما اگر با ادبیات یادگیری ماشین آشنا باشید این معیار چندان برایتان آشنا نخواهد بود و امروزه کمتر مورد استفاده قرار می‌گیرد.

این معیار، برای داده‌های نامتعادل یعنی داده‌هایی که تعداد برجسب‌های مثبت و منفی آن در دنیای واقعی از لحاظ عددی اختلاف بسیار زیادی دارند، معیار مناسبی نیست. بسیاری از مسائل دنیای واقعی هم دقیقاً جزء این گروه قرار می‌گیرند. اگر قرار باشد ابتلا به ایدز را از روی آزمایشات مختلف برای هزار نفر تشخیص دهید، شاید یک یا دو نفر از این بین، به ایدز مبتلا باشند که اختلاف زیادی بین دسته مثبت (افراد دارای ایدز) و دسته منفی (افراد سالم) حاکم است. مسایلی مانند تشخیص هرزنامه بودن یک ایمیل، تروریست بودن یک مسافر هواپیما، شناسایی دانشجویان مشروط، پیش‌بینی خروج یک کارمند از شرکت و مانند آن، نمونه‌های دیگری از مسایل دسته‌بندی با دسته‌های نامتعادل در دنیای واقعی هستند. این اختلاف معنی دار بین دسته‌های مختلف داده‌ها، باعث عدم کارایی معیار دقت می‌شود.

مشکل اصلی هم نامتعادل بودن داده‌ها و تفاوت معنی دار تعداد نمونه‌های هر دسته است که باعث می‌شود یک مدل متمایل به دسته پرتعداد، دقت کلی را بالا نشان دهد. بنابراین نیاز به معیاری دقیق‌تر برای سنجش دقت و کارایی الگوریتم‌های پیشنهادی دسته‌بندی هستیم.

در این گونه مسایل بهتر است بر تعداد نمونه‌های مثبت شناسایی شده به کل نمونه‌های مثبت تمرکز کنیم. یعنی ببینیم از دو نفر بیمار مبتلا به ایدز، چند نفر شناسایی شده‌اند. معیار صحت را برای این منظور به صورت زیر تعریف می‌کنیم:

$$\text{بازخوانی (Recall)} = \frac{\text{تعداد های نمونه تشخیصی درست مثبت}}{\text{کل های نمونه واقعاً مثبت}} = \frac{TP}{TP+FN}$$

توضیح اینکه کل نمونه‌های واقعاً مثبت شامل نمونه‌هایی است که درست، مثبت شناسایی شده‌اند (TP) و نمونه‌هایی که مثبت بوده‌اند اما نادرست، منفی شناسایی شده‌اند (FN).

سنجش بازخوانی برای روش پیشنهادی تشخیص ایدز، صفر است (چون هیچ نمونه مثبتی را شناسایی نکرده ایم – صورت کسر برابر صفر است) که نشانگر ضعیف بودن مدل پیشنهادی است و بنابراین آنرا می‌توانیم به راحتی رد کنیم.

## ۲-۷-۳- الگوریتم (مدل) پیشنهادی

تمام نمونه ها را مثبت اعلام کن! در این صورت تمام دو نفر بیمار ایدز را تشخیص داده ایم. یعنی بازخوانی ما برابر حداکثر ممکن یعنی ۱ شده است. توضیح اینکه تعداد داده های درست تشخیص داده شده برابر ۲ و تعداد داده های نادرست منفی اعلام شده (کسانی که ایدز دارند اما نتیجه آزمایش ایدز آنها را منفی اعلام کرده ایم - FN)، برابر صفر است و تقسیم دو بر دو هم که یک می شود. یعنی گاهی بازخوانی ما به خاطر ضعیف بودن مدل پیشنهادی، بالاست. این ضعیف بودن را با معیار دیگری باید اندازه بگیریم.

برای حل این مشکل، در کنار معیار بازخوانی معیار دیگری را به نام صحت (Precision)، برابر تعداد نمونه های تشخیصی درست مثبت به کل نمونه های مثبت اعلام شده به صورت زیر تعریف می کنیم تا میزان مثبت های اشتباه را هم در نظر گرفته باشیم:

$$\text{صحت (Precision)} = \frac{\text{تعداد های نمونه تشخیصی درست مثبت}}{\text{تعداد کل های نمونه تشخیصی مثبت}} = \frac{TP}{TP+FP}$$

در این فرمول، وجود FP در مخرج باعث می شود که اگر تعداد تشخیص های اشتباه مان بالا باشد، صحت الگوریتم عددی نزدیک به صفر نشان دهد و بنابراین کارایی مدل، زیر سوال برود. با این توضیحات، معیارهای بازخوانی و صحت به جای معیار اولیه دقت، کاربرد وسیع تری در دنیای امروز یادگیری ماشین پیدا کرده است. در اغلب موارد، این دو معیار با هم رشد و حرکت نمی کنند. گاهی ما صحت مدل را با الگوریتم های دقیق تر بالا می بریم، یعنی آنهایی را که مثبت اعلام می کنیم، اکثراً درست هستند و موارد نادرست مثبت ما بسیار کم هستند یعنی صحت الگوریتم ما بسیار بالاست اما ممکن است جنبه یا ویژگی خاصی از داده ها را در نظر نگرفته باشیم و تعداد کل نمونه های مثبت، بسیار بیشتر از نمونه های اعلام شده ما باشد یعنی بازخوانی بسیار پایینی داشته باشیم.

از طرفی ممکن است کمی الگوریتم تشخیصی خود را ساده تر بگیریم تا تعداد مثبت های تشخیصی خود را بالا ببریم، در این صورت میزان اشتباهات ما زیادتر شده، صحت الگوریتم عدد پایین تر و بازخوانی آن، عدد بالاتری را نشان می دهد. اگر بتوانیم معیاری ترکیبی از این دو معیار برای سنجش الگوریتم های دسته بندی به دست آوریم، تمرکز بر آن معیار به جای بررسی همزمان این دو، مناسب تر خواهد بود مثلاً از میانگین این دو به

عنوان یک معیار جدید استفاده کنیم و سعی در بالا بردن میانگین حسابی این دو داشته باشیم. اگر بخواهیم میانگین معمولی دو معیار بازخوانی و صحت را ملاک کار در نظر بگیریم، برای حالت هایی که صحت بالا و بازخوانی پایینی داریم (و یا بالعکس)، میانگین معمولی عددی قابل قبول خواهد بود در صورتی که الگوریتم پیشنهادی نباید نمره قبولی بگیرد.

برای رفع این نقیصه و تولید یک معیار واحد که متمایل به عدد کوچکتر باشد، از میانگین هارمونیک استفاده می کنیم. این میانگین هارمونی برای دو مقدار بازخوانی و صحت را با نام F1-Score می شناسیم که طبق فرمول فوق برابر است با :

$$F1-Score = \frac{2}{\frac{1}{Recall} + \frac{1}{Precision}} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

در این فرمول همانطور که مشاهده می کنید اگر یکی از دو مقدار عددی کوچک باشد، یا حتی صفر باشد، نتیجه نهایی عددی کوچک و یا صفر خواهد بود. توضیح این امر هم ساده است چون دو معیار بازخوانی و صحت اعدادی بین صفر تا یک هستند و در صورت کسر در هم دیگر ضرب شده اند بنابراین نتیجه نهایی به سمت عدد کوچکتر، متمایل خواهد بود و اگر هر دو با هم، عددی بزرگ (نزدیک ۱) باشند، نتیجه نهایی به سمت یک حرکت خواهد کرد .

در دنیای واقعی یک حد آستانه پذیرش هم می توانیم برای F1-Score تعیین کنیم مثلاً ۰٫۹ و اعلام کنیم که مدل هایی با نمره بالاتر از این آستانه، مورد تایید نهایی قرار خواهند گرفت.

در فرمول میانگین هارمونیک ، وزنی مساوی به هر دو پارامتر داده ایم و بسته به نیاز می توانیم میانگین هارمونیک مرتبه های بالاتر یعنی F2 ، F3 و غیره را هم به کار ببریم

آخرین مطلب درباره این معیارهای اصلی دسته بندی این است که این معیارها کاملاً بستگی به بستر و حوزه دسته بندی دارند. مثلاً در تشخیص ایدز یا تشخیص کلاه برداری در تراکنش های بانکی، ما نیاز به شناسایی تمامی موارد ایدز و کلاه برداری داریم یعنی نیاز داریم که بازخوانی ما بسیار بالا باشد و اگر خطایی هم تولید شد مثلاً بیماری به اشتباه ایدزی تشخیص داده شد و یا یک تراکنش سالم، متهم به کلاه برداری شد، کافی است با کمی آزمایش بیشتر، نتایج را بهبود خواهیم بخشید و موارد خطا را از لیست تشخیص داده شده ها حذف خواهیم کرد. اما



گاهی اوقات به دنبال صحت بیشتر هستیم مثلاً با خواندن توئیت های روزانه، قرار است آنالیز احساسات روی آنها انجام دهیم. در این حالت، صحت الگوریتم یعنی تشخیص درست و دقیق احساسات هر توئیت خوانده شده و نه همه توئیت ها (با فرض اینکه تعداد توئیت های بررسی شده زیاد باشد) برای ما اهمیت زیادیتری از بررسی تمامی توئیت ها دارد. بنابراین همیشه و در همه موارد، ما از F1-Score استفاده نمی کنیم، بلکه با بررسی نیازمندیها، بهترین معیار را برای کار خود انتخاب خواهیم کرد.

در مواردی که دسته ها، متعادل هستند، مثلاً تعیین جنسیت ارسال کننده یک توئیت، می توانیم همان معیار دقت که معیار اول مورد بحث بود را هم به کار ببریم. در هر صورت باید بدانیم که دنبال چه هستیم. هدف اصلی ما در یافتن یک مدل دسته بندی افزایش F1-Score آن خواهد بود [21].

پس به طور خلاصه می توان اینگونه بیان کرد :

صحت (Accuracy) : به طور کلی، دقت به این معناست که مدل تا چه اندازه خروجی را درست پیش بینی می کند.

$$Accuracy = \frac{true\ positives + true\ negatives}{total\ examples}$$

با نگاه کردن به دقت ، بلافاصله می توان دریافت که آیا مدل درست آموزش دیده است یا خیر و کارایی آن به طور کلی چگونه است. اما این معیار اطلاعات جزئی در مورد کارایی مدل ارائه نمی دهد.

دقت (Precision) : وقتی که مدل نتیجه را مثبت (positive) پیش بینی می کند، این نتیجه تا چه اندازه درست است؟

$$precision = \frac{true\ positives}{true\ positives + false\ positives}$$

زمانی که ارزش false positives بالا باشد، معیار صحت، معیار مناسبی خواهد بود. فرض کنید، مدلی برای تشخیص سرطان داشته باشیم و این مدل Precision پایینی داشته باشد. نتیجه این امر این است که این مدل، بیماری بسیاری از افراد را به اشتباه سرطان تشخیص می دهد. نتیجه این امر استرس زیاد، آزمایش های فراوان و هزینه های گزافی را برای بیمار به دنبال خواهد داشت.

$$Recall = \frac{true\ positives}{true\ positives + false\ negatives}$$

یادآوری (Recall) :

زمانی که ارزش false negatives بالا باشد، معیار Recall، معیار مناسبی خواهد بود. فرض کنیم مدلی برای تشخیص بیماری کشنده ابولا داشته باشیم. اگر این مدل Recall پایینی داشته باشد چه اتفاقی خواهد افتاد؟ این

مدل افراد زیادی که آلوده به این بیماری کشنده هستند را سالم در نظر می‌گیرد و این فاجعه است.

نرخ امتیاز (F1-Score) :

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

معیار F1، یک معیار مناسب برای ارزیابی دقت یک آزمایش است. این معیار Precision و Recall را با هم در نظر می‌گیرد. معیار F1 در بهترین حالت، یک و در بدترین حالت صفر است [22].

فصل سوم :

روش تحقیق

### ۱-۳- مقدمه

همان طور که بیان شد روش‌های مختلفی برای دسته‌بندی نمونه‌های ورودی ارائه شده‌اند که با توجه به معیارهای ارائه شده برای برازش کارایی این روش‌ها، هر کدام از این روش‌ها دارای مزایا و معایبی هستند. در این فصل قصد داریم عملیات اصلی را بر روی دادگان<sup>۴۲</sup> انجام داده و نتایج حاصله را بررسی نماییم. روال کار بر اساس روش‌های منطقی انجام شده که هر کدام در نحوه انجام و انتخاب انواع روش‌ها و کتابخانه‌ها دارای جزییات منحصر به فرد می‌باشند. در مورد هر کدام از کتابخانه‌هایی که استفاده شده به طور مختصر توضیح داده خواهد شد.

هدف آن است که پس از انتخاب روش مناسب برای خواندن داده، ابتدا در مورد دیتاست اطلاعات دقیقی بدست آورده و پس از تحلیل دقیق آن، برای ادامه کار با استفاده از روش‌های معمول تقسیم‌بندی و پاکسازی که در هر بخش بطور مجزا بررسی خواهد شد سپس دادگان آموزشی را آماده خواهد شد تا بتوان انواع مختلفی از دسته‌بند<sup>۴۳</sup>ها را بر روی آن‌ها پیاده‌سازی کرده و نتایج را بررسی کرد. همچنین در قدم بعدی، برخی از روش‌های ترکیب دسته‌بندهای نیز استفاده شده تا امکان بهبود نتایج توسط ترکیب دسته‌بندها نیز بررسی گردد.

### ۳-۲ روش کار

#### ۱-۲-۳ خواندن مجموعه داده

اولین قدم برای آغاز کار هوش مصنوعی، بارگذاری داده است [25]. برای این کار کتابخانه‌ها و روش‌های متفاوتی وجود دارد. برخی از این روش‌ها عبارتند از:

- بارگذاری فایل‌های CSV با کتابخانه استاندارد پایتون

- بارگذاری فایل‌های CSV با کتابخانه Numpy

- بارگذاری فایل‌های CSV با کتابخانه Pandas

---

<sup>۴۲</sup> Dataset

<sup>۴۳</sup> Classifier

در این پژوهش از کتابخانه pandas برای بارگذاری داده استفاده شده است. کتابخانه Pandas بر مبنای NumPy ساخته شده است و ساختارها و ابزارهای تحلیل داده با کاربری آسان را برای زبان برنامه‌نویسی پایتون ارائه می‌کند. این کتابخانه قدرتمند و انعطاف‌پذیر به طور مکرر از سوی علاقه‌مندان علم داده برای دریافت داده‌ها در ساختارهای داده‌ای که برای تحلیل‌هایشان کاملاً گویا هستند مورد استفاده قرار می‌گیرد [26].

در کدهای نوشته شده در این پژوهش برای شناخت بیشتر کتابخانه‌ها و مشاهده آن‌ها به صورت یکجا، تمام کتابخانه‌های مورد استفاده در ابتدای کد‌های برنامه فراخوانی شده‌اند. اما در زمان توضیح روال برنامه، در هر زمان که از کتابخانه‌ای استفاده شده باشد، دستور مورد نیاز برای فراخوانی آن نیز توضیح داده خواهد شد.

```

1 import re
2 import nltk
3 import pickle
4 import numpy as np
5 import pandas as pd
6 import seaborn as sns
7 from scipy import stats
8 from sklearn import tree
9 from sklearn.svm import SVC
10 from random import randrange
11 from bs4 import BeautifulSoup
12 from nltk import word_tokenize
13 import matplotlib.pyplot as plt
14 from operator import itemgetter
15 from nltk.corpus import stopwords
16 from sklearn.svm import LinearSVC
17 from nltk.stem.porter import PorterStemmer
18 from sklearn.naive_bayes import GaussianNB
19 from sklearn.ensemble import VotingClassifier
20 from sklearn.ensemble import ExtraTreesClassifier
21 from sklearn.ensemble import RandomForestClassifier
22 from sklearn.linear_model import LogisticRegression
23 from sklearn.model_selection import cross_val_score
24 from sklearn.model_selection import train_test_split
25 from sklearn.metrics import precision_recall_fscore_support
26 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
27 NB = GaussianNB()
28 %matplotlib inline

```

### فراخوانی کتابخانه‌های مورد نیاز

در تصویر زیر فراخوانی داده توسط کتابخانه pandas را مشاهده می‌فرمایید. برای فراخوانی کتابخانه pandas از کد `import pandas as pd` استفاده شده است.

```
In [3]: dataset = pd.read_csv("train.csv")
print("Number of rows in train dataset : ",len(dataset))
dataset.head()
```

Number of rows in train dataset : 19579

Out[3]:

	id	text	author
0	id26305	This process, however, afforded me no means of...	EAP
1	id17569	It never once occurred to me that the fumbling...	HPL
2	id11008	In his left hand was a gold snuff box, from wh...	EAP
3	id27763	How lovely is spring As we looked from Windsor...	MWS
4	id12958	Finding nothing else, not even gold, the Super...	HPL

## خواندن مجموعه داده

### ۳-۲-۲- تجزیه و تحلیل داده

از مهم ترین قسمت های کار با داده ، شناخت داده ورودی است . تسلط بر اینکه داده ورودی، شامل چه آیتم هایی است ، چه ویژگی ۴۴ و چه شیء ۴۵ هایی در داده وجود دارد، شامل چند شیء است و نوع ویژگی و اشیاء چیست . همچنین بررسی برچسب ۴۶ ها ، نوع آن ها و تعداد برچسب تاثیر بسیاری در نحوه کار با داده و انتخاب روش های مختلف کار با داده دارد .

داده استفاده شده در این پژوهش ، شامل ۱۹۵۷۹ شیء و سه ویژگی می باشد .  
به عنوان ویژگی ، داده هایی نظیر نام نویسنده، ID و بخشی از متن کتاب قرار گرفته است . کل داده شامل ۱۹۵۷۹ جمله از کتاب های این سه نویسنده است .

```
In [4]: 1 dataset = pd.read_csv("train.csv")
2 print("Number of rows in train dataset : ",len(dataset))
3 dataset.head()
4 dataset.shape
```

Number of rows in train dataset : 19579

Out[4]: (19579, 3)

## Shape داده ورودی

<sup>۴۴</sup> Attribute

<sup>۴۵</sup> Object

<sup>۴۶</sup> Label

از سه نویسنده که با برچسب های HPL , MWS , EAP شناسایی میشوند ، تعداد ۷۹۰۰ متن از EAP و 6044 متن از MWS و ۵۶۳۵ متن از HPL در داده قرار داده شده است . و برای هر متن یک ID خاص در نظر گرفته شده است .

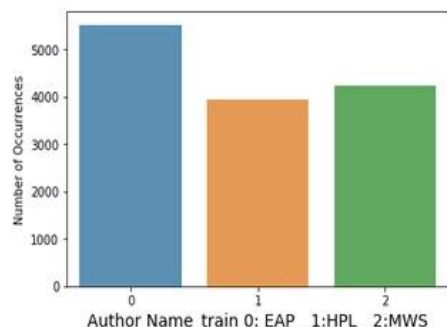
```
In [5]: 1 dataset['author'].value_counts()

Out[5]: EAP    7900
        MWS    6044
        HPL    5635
        Name: author, dtype: int64
```

### تعداد متن های مربوط به نویسنده

داده از لحاظ توازن، تقریباً متوازن محسوب می شود زیرا حدود تعداد متن های در نظر گرفته شده از هر نویسنده به هم نزدیک می باشد. در شکل زیر می توان تعداد متن هایی که برای هر نویسنده در داده ورودی در نظر گرفته شده را با استفاده از کتابخانه seaborn که با کد `import seaborn as sns` در ابتدای برنامه فراخوانی شده است ، به صورت بصری مشاهده فرمایید .

```
In [41]: 1 cnt_srs=pd.DataFrame(ytrain)[0].value_counts()
        2 plt.figure(figsize=(6,4))
        3 sns.barplot ( cnt_srs.index , pd.DataFrame(ytrain)[0].value_counts(),alpha=0.8)
        4 plt.ylabel('Number of Occurrences', fontsize=10)
        5 plt.xlabel('Author Name_train 0: EAP_1:HPL_2:MWS', fontsize=14)
        6
        7 plt.show()
```



### نمودار مربوط به تعداد متن های هر نویسنده در داده ورودی

کتابخانه Seaborn یک کتابخانه بسیار مفید مصورسازی داده در پایتون محسوب می گردد. کتابخانه Matplotlib ساخته می شود و قابلیت های پیشرفته بسیاری در رابطه با مصورسازی داده ارائه می کند. کتابخانه Seaborn را جهت رسم انواع چارت هایی چون نمودارهای ماتریسی، نمودارهای شبکه ای (Grid)،

نمودارهای رگرسیونی ، برای رسم نمودارهای توزیعی و غیره می توان مورد استفاده قرار داد [27].

### ۳-۲-۳- تقسیم بندی دیتاست

در این پژوهش تقسیم بندی دیتاست بر اساس سه دسته انجام گرفته است .  
دسته اول داده آموزش ۲۸ است که ۷۰ درصد از داده اصلی را شامل میشود  
دسته دوم داده آزمون ۲۹ است که ۲۰ درصد از داده اصلی را شامل می شود  
و دسته سوم به عنوان داده اعتبارسنجی ۳۰ شناخته می شود که با نام اختصاری Dev تعریف شده و شامل ۱۰ درصد از داده اصلی می باشد .  
در واقع در این روش پس از یادگیری از روی داده آموزش و انجام آزمون روی داده آزمون، ارزیابی بر روی داده اعتبارسنجی انجام میگردد.  
مهم ترین فاکتور در ساخت این سه دسته، رعایت توزیع کلاس هاست . به نحوی که وقتی از داده اصلی، ۷۰ درصد را به عنوان داده آموزشی استفاده می کنیم، توزیع کلاس ها یا برچسب ها باید به همان میزانی باشد که در داده اصلی موجود بوده است .  
به عنوان مثال با فرض اینکه در داده اصلی، ۲۰ درصد از داده مربوط به کلاس A ، ۳۰ درصد مربوط به کلاس B و ۵۰ درصد مربوط به کلاس C باشد، این توزیع کلاس باید در داده های آموزش، آزمون و اعتبار سنجی نیز به همین صورت باشد.  
در شکل زیر کدهای مربوط به اختصاص داده را مشاهده می فرمایید .  
برای انجام عملیات بر روی داده ، با توجه به اینکه بعضی از الگوریتم های یادگیری ماشین صرفاً توانایی کار با عدد را دارند، برچسب ها را برای جلوگیری از خطاهای احتمالی به عدد تبدیل کرده ایم (خطوط ۳ و ۴) .



```

In [ ]: 1 from sklearn import preprocessing, decomposition, model_selection, metrics, pipeline
2
3 lbl_enc = preprocessing.LabelEncoder()
4 y = lbl_enc.fit_transform(dataset.author.values)
5
6 xrest, xdev, yrest, ydev = train_test_split(dataset.text.values, y,
7                                           stratify=y,
8                                           random_state=42,
9                                           test_size=0.1, shuffle=True)
10 Xrest=pd.DataFrame(xrest)
11 Yrest=pd.DataFrame(yrest)
12 Xrest["author"]=Yrest[0]
13 Xrest.columns = ['text', 'author']
14
15 y2 = lbl_enc.fit_transform(Xrest.author.values)
16
17 xtrain, xtest, ytrain, ytest = train_test_split(Xrest.text.values, y2,
18                                                stratify=y2,
19                                                random_state=42,
20                                                test_size=0.2222, shuffle=True)
21
22 Ydev=pd.DataFrame(ydev)
23 z0=pickle.dump(Ydev, open('Ydev.pkl', 'wb'))
24
25 Ytrain=pd.DataFrame(ytrain)
26 z0=pickle.dump(Ytrain, open('Ytrain.pkl', 'wb'))
27
28 Ytest=pd.DataFrame(ytest)
29 zc=pickle.dump(Ytest, open('Ytest.pkl', 'wb'))
30

```

### تقسیم بندی داده

در هر قسمت پس از انجام تقسیم بندی ، داده های مربوطه جهت تسهیل در فراخوانی به عنوان یک دیتافریم با استفاده از تابع pickle ذخیره شده اند .  
نمایش مقادیر عددی برجسب ها بعد از تبدیل به عدد

```

In [6]: dataset['author'].value_counts(),pd.DataFrame(y)[0].value_counts()
print(0,'\t:EAP\n',1,'\t:HPL\n',2,'\t:MWS\n')

```

```

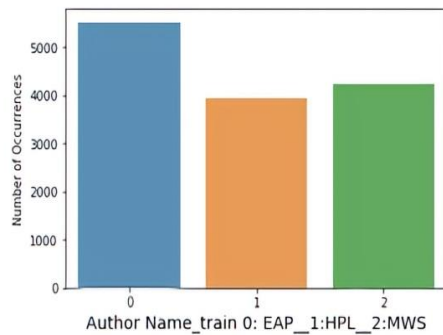
0      :EAP
1      :HPL
2      :MWS

```

### نمایش برجسب ها مقادیر عددی

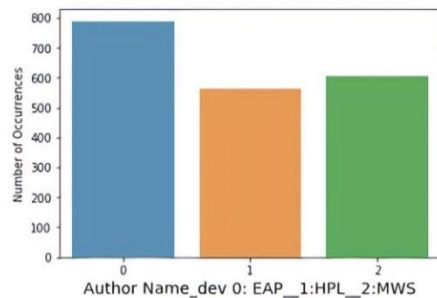
همانطور که در شکل های زیر مشاهده میفرمایید ، توازن کلاس ها در تمام تقسیم بندی های داده رعایت شده است .

```
In [41]: 1 cnt_srs=pd.DataFrame(ytrain)[0].value_counts()
2 plt.figure(figsize=(6,4))
3 sns.barplot ( cnt_srs.index , pd.DataFrame(ytrain)[0].value_counts(),alpha=0.8)
4 plt.ylabel('Number of Occurrences', fontsize=10)
5 plt.xlabel('Author Name_train 0: EAP__1:HPL__2:MWS', fontsize=14)
6
7 plt.show()
```



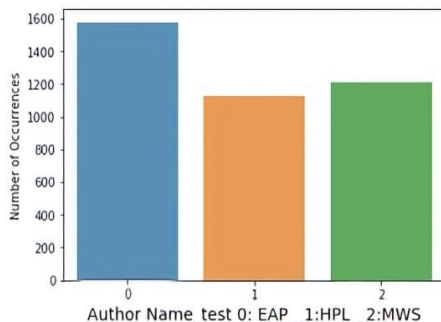
توازن کلاس ها در داده آزمایش

```
In [42]: 1 cnt_srs=pd.DataFrame(ydev)[0].value_counts()
2 plt.figure(figsize=(6,4))
3 sns.barplot ( cnt_srs.index , pd.DataFrame(ydev)[0].value_counts(),alpha=0.8)
4 plt.ylabel('Number of Occurrences', fontsize=10)
5 plt.xlabel('Author Name_dev 0: EAP__1:HPL__2:MWS', fontsize=14)
6
7 plt.show()
8
```



توازن کلاس ها در داده اعتبارسنجی

```
In [43]: 1 cnt_srs=pd.DataFrame(ytest)[0].value_counts()
2 plt.figure(figsize=(6,4))
3 sns.barplot ( cnt_srs.index , pd.DataFrame(ytest)[0].value_counts(),alpha=0.8)
4
5 plt.ylabel('Number of Occurrences', fontsize=10)
6 plt.xlabel('Author Name_test 0: EAP_1:HPL_2:MWS', fontsize=14)
7
8 plt.show()
```



توازن کلاس ها در داده آزمایش

#### ۴-۲-۳- پاک سازی داده

پاک سازی داده ها یا تمیز کردن داده ها فرایند پیدا کردن، اصلاح کردن (یا حتی حذف کردن) داده های بی ارزش و اشتباه از دادگان یا پایگاه داده ۳۲ است. فرایند تمیز کردن داده ها ممکن است که از طریق ابزارهای داده کاوی یا پردازش دسته ای از طریق اسکریپت ها انجام شود. بعد از پاک سازی، مجموعه داده باید با سایر مجموعه داده های مشابه در سیستم سازگار باشد. ناسازگاری داده ها شناسایی و حذف (اصلاح) شده ممکن است بر اثر اشتباه انسانی هنگام ورود اطلاعات، انحراف در هنگام انتقال و ذخیره سازی اطلاعات یا به دلیل واژه نامه های داده مختلف باشد.

در این پژوهش برای پاک سازی داده به انجام عملیات نشانه گذاری ۳۳، حذف علائم، کوچک کردن تمام حروف، حذف کلمات کلیدی ۳۴ و یافتن ریشه لغات ۳۵ بسنده کرده ایم. یک تابع کلی برای انجام پاک سازی در ابتدا تعریف شده است تا بعدا بتوان به سادگی با فراخوانی این تابع، تمامی عملیات مربوطه را بر روی هر سه دیتا انجام داد.

```
In [54]: 1 def cleaning_and_stemming(text, non_alpha=True, normalization=True, stemming=True, stopword=True):
2     BeautifulSoup_text = BeautifulSoup(text).get_text()
3     words = word_tokenize(BeautifulSoup_text)
4     if non_alpha:
5         words = [word for word in words if word.isalpha()]
6     if normalization:
7         words = [w.lower() for w in words]
8     if stopword:
9         stop_words = set(stopwords.words('english'))
10        words = [w for w in words if w not in stop_words]
11    if stemming:
12        porter = PorterStemmer()
13        words = [porter.stem(word) for word in words]
14    return(" ".join(words))
15
```

حالا این تابع را بر روی داده های آموزش ، آزمون و اعتبارسنجی پیاده سازی می کنیم .

```
In [34]: 1 clean_train=list(range(len(xtrain)))
2 clean_dev = list(range(len(xdev)))
3 clean_test = list(range(len(xtest)))
4 xTrain= pd.DataFrame(xtrain)
5 xDev=pd.DataFrame(xdev)
6 xTest=pd.DataFrame(xtest)
7 clean_train=xTrain[0].apply(lambda x: cleaning_and_stemming(x)).tolist()
8 clean_dev=xDev[0].apply ( lambda x: cleaning_and_stemming(x)).tolist()
9 clean_test=xTest[0].apply ( lambda x: cleaning_and_stemming(x)).tolist()
```

### پاکسازی داده

#### ۵-۲-۳- تهیه ورودی برای الگوریتم های پایه با استفاده از TF-IDF و BOW(TF) و BOW(binary)

یک نمایش ساده است که در پردازش زبان های طبیعی و بازیابی اطلاعات ۳۶ استفاده می شود. همچنین به عنوان مدل فضای بردار شناخته می شود. در این مدل، یک متن (مانند یک جمله یا سند) به صورت یک بسته چند مجموعه از کلمات آن، بی توجه به دستور زبان و حتی نظم کلمات نمایش داده می شود (McTear et al 2016, p. 167) مدل بسته کلمات برای بینایی کامپیوتر ۳۷ استفاده شده است (Sivic, Josef (April 2009)) مدل بسته کلمات معمولاً در روش های دسته بندی اسناد مورد استفاده قرار می گیرد که در آن وقوع هر کلمه به عنوان یک ویژگی برای آموزش طبقه بندی آماری استفاده می شود (Harris, Zellig (1954))

برای این منظور از ماژول های TfidfVectorizer, CountVectorizer که به صورت زیر فراخوانی می شوند ، استفاده شده است.

```
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
```

tfidfVectorizer, CountVectorizer های کتابخانه های

برای ساخت بردار کلمات در متن از ماژول CountVectorizer پارامترهای زیر استفاده شده است:

**min\_df=3**: در زمان ساختن دایره لغات، کلماتی که حداقل تکرار خاصی داشته باشند (در این پژوهش عدد ۳ انتخاب شده است) را حذف می کند .

**max\_features=None**: مشخص کننده حداکثر تعداد کلماتی است که از ویژگی ها انتخاب می شود. می توان این محدودیت را با وارد کردن یک عدد (integer) مشخص کرد . در این پژوهش با انتخاب پارامتر None ، هیچ حداکثری برای این مورد در نظر نگرفته ایم .

**strip\_accents='unicode'**: این پارامتر برای حذف لهجه ها و نرمال سازی کلمات به کار می رود . به سه صورت قابل تنظیم است .

**ASCII**: در این حالت که یک روش سریع است، صرفاً بر روی کاراکترهایی تمرکز می شود که به کد اسکی خاصی نگاشت می شوند.

**UNICODE**: در این روش که نسبت به روش قبل کندتر است، عملیات حذف لهجه و نرمال سازی بر روی تمام کاراکترها انجام می شود .

**NONE**: در این روش هیچ کاری بر روی کاراکترها انجام نمی شود .

**analyzer='word'**: با این پارامتر مشخص میکنیم که ویژگی های ما بر اساس کلمات، کاراکترها یا n-grams ساخته شده اند. با توجه به feature های دادگان و نوع عملیاتی که قصد داریم انجام دهیم، این پارامتر را word در نظر گرفته ایم .

**ngram\_range= (1, 3)**: مشخص کننده حداقل و حداکثر تعداد کلمات چندتایی که باید استخراج شوند استفاده می شود.

```
vectorizer = CountVectorizer(min_df=3, max_features=None, strip_accents='unicode', analyzer='word', ngram_range=(1, 3))
```

تعیین پارامترهای تابع vectorizer

پس با تابع fit\_transform و تعیین کردن داده آموزش پاک سازی شده به عنوان پارامتر ورودی مدل بردار کلمات ساخته شده است.

```
train_BOW_bi = vectorizer.fit_transform(clean_train).toarray()
```

تابع fit\_transform

در ادامه بردار داده‌های اعتبارسنجی و آزمایش با تابع fit ساخته شده است.

```
6 dev_BOW_bi = vectorizer.transform(clean_dev).toarray()
7 test_BOW_bi = vectorizer.transform(clean_test).toarray()
```

ساخت بردار دادگان اعتبارسنجی و آزمایش

در انتها با توجه به اینکه در جدول بردار کلمات، تعداد تکرار کلمات موجود می باشد و درکیسه کلمات دو دویی، معیار فقط وجود و یا عدم وجود کلمه است، با استفاده از دو حلقه for تمام مولفه‌هایی که کلمه در آن ها موجود بوده و به عبارتی مخالف صفر بوده اند، با عدد یک جایگزین شده و بدین صورت تمام داده های آموزش، اعتبارسنجی و آزمون به کیسه کلمات دودویی تبدیل شده‌است.

```
In [35]: 1 print ("Creating the bag of words binary...\n")
2 vectorizer = CountVectorizer(min_df=3, max_features=None, strip_accents='unicode',
3 analyzer='word', ngram_range=(1, 3))
4
5 train_BOW_bi = vectorizer.fit_transform(clean_train).toarray()
6 dev_BOW_bi = vectorizer.transform(clean_dev).toarray()
7 test_BOW_bi = vectorizer.transform(clean_test).toarray()
8
9 for i in range (len(train_BOW_bi)):
10     for j in range (len(train_BOW_bi[i])):
11         if train_BOW_bi[i][j]!=0 :
12             train_BOW_bi[i][j]=1
13
14 for i in range (len(dev_BOW_bi)):
15     for j in range (len(dev_BOW_bi[i])):
16         if dev_BOW_bi[i][j]!=0 :
17             dev_BOW_bi[i][j]=1
18
19 for i in range (len(test_BOW_bi)):
20     for j in range (len(test_BOW_bi[i])):
21         if test_BOW_bi[i][j]!=0 :
22             test_BOW_bi[i][j]=1
```

ساخت Binary BOW

ساخت کیسه کلمات با فرکانس کلمات

```
In [ ]: 1 print ("Creating the bag of words...\n")
2 train_BOW = vectorizer.fit_transform(clean_train).toarray()
3 dev_BOW = vectorizer.transform(clean_dev).toarray()
4 test_BOW = vectorizer.transform(clean_test).toarray()
```

ساخت BOW

ساخت فراوانی کلمه کلیدی (tf-idff)

```
In [36]: 1 print ("Creating the Tfidf ...\n")
2 vectorizer = TfidfVectorizer(min_df=3, max_features=None,
3 strip_accents='unicode', analyzer='word',
4 ngram_range=(1, 3), use_idf=1, smooth_idf=1, sublinear_tf=1)
5
6
7 train_TFIDF = vectorizer.fit_transform(clean_train).toarray()
8 dev_TFIDF = vectorizer.transform(clean_dev).toarray()
9 test_TFIDF = vectorizer.transform(clean_test).toarray()
```

### ساخت TFIDF BOW

ابتدا تمام برچسب های دادگان آموزش ، اعتبار سنجی و آزمون را در لیست های جدیدی که تعریف می کنیم قرار داده ایم

```
In [37]: 1 Ytrain=pd.DataFrame(ytrain)
2 Ydev=pd.DataFrame(ydev)
3 Ytest=pd.DataFrame(ytest)
```

### فراخوانی برچسب های دادگان

با توجه به اینکه تعداد مدل هایی که قصد ساخت آن ها را داریم ، بالاست و در حال حاضر بصورت دقیق قابل پیش بینی نیست، لیست هایی به طول ۱۰۰۰ برای مدل ها ، نتایج و مقادیر پیش بینی شده در نظر گرفته شده است . در آینده با توجه به تعداد مدل هایی که ساخته شده اند، میتوان برای ذخیره سازی تعداد مورد نظر را انتخاب کرد.

```
In [38]: 1 allModels_clf=list (range (1000)) # [(model,fscore,predict),...,(model,fscore,predict)]
2 Fscore_clf=list (range (1000))
3 Pred_clf=list (range (1000))
```

### تعریف لیست های مورد نیاز

## ۶-۲-۳- ساخت ۱۰۰ نمونه داده آموزش با رویکرد جایگشتی

برای آموزش هوش مصنوعی ، بهتر است تعداد داده های بیشتری به عنوان ورودی به دسته بند های متفاوت داده شود .

با توجه به اینکه داده های دیتاست دارای محدودیت هستند، می توان با استفاده از همان داده ها، تعداد بیشتری داده آموزش ساخت . روش کار به اینصورت است که به با استفاده از تابع rand در هر بار برداشت، یکی از داده



های دیتاست برداشته می شود و در جدول جدیدی قرار میگیرد. از تابع `rand` به این دلیل استفاده شده که داده های دیتاست اصلی به ترتیب برداشته نشده و در نتیجه ترتیب داده ها یکسان نباشند. این کار باعث می شود دسته بند، اگر چه داده های تکراری را آموزش میبیند اما با تغییر جایگاه داده ها بتواند بصورت بهتر عمل یادگیری را انجام دهد.

توجه داشته باشید که این روش بصورت جایگشتی بوده و با توجه به اینکه داده ها بصورت رندوم انتخاب می شوند ممکن است داده تکراری داشته باشیم. این روش را `sampling with replacement` میگویند. همانطور که در شکل زیر ملاحظه میفرمایید، ابتدا باید داده های آموزش، آزمون و اعتبارسنجی که پاکسازی شده بصورت دادگان جدید ذخیره شده اند را فراخوانی کرده و در لیست های جداگانه قرار دهیم. سپس با تعریف دو تابع، عملیات ساخت ۱۰۰ عدد داده آموزش را انجام می دهیم.

برای انجام این کار دو تابع `Subsample` و `cerat_multi_dataset` به صورت زیر تبدیل شده است. در تابع `subsample` یک لیست به اندازه طول اشیاء دادگان ساخته شده و با روش جایگشتی، هر یک از اشیاء به صورت رندوم انتخاب شده است تا بتوانیم از روی داده آموزش، یک نمونه جدید بسازیم. در تابع `cerat_multi_dataset` عملیات فوق به تعداد ۱۰۰ بار انجام شده تا در نتیجه از روی اشیاء دادگان اصلی ۱۰۰ داده آموزش ساخته باشیم.

حال این داده ها آماده هستند تا در دسته بندهای مختلف به عنوان داده ورودی استفاده شوند.

```
In [40]: 1 train_TFIDF_D=pd.DataFrame(train_TFIDF)
2 train_TFIDF_D['lable']=ytrain
3 train_TFIDF_D=train_TFIDF_D.as_matrix()
4
5 dev_TFIDF_D=pd.DataFrame(dev_TFIDF)
6 dev_TFIDF_D['lable']=ydev
7 dev_TFIDF_D=dev_TFIDF_D.as_matrix()
8
9 def subsample(dataset):
10     sample = list()
11     for i in range (dataset.shape[0]):
12         index = randrange(dataset.shape[0])
13         sample.append(dataset[index])
14     return sample
15
16 def cerat_muliti_dataset(dataset):
17     COM_LIST=list(range(100))
18     for i in range (len(COM_LIST)):
19         COM_LIST[i]=subsample(dataset)
20     return COM_LIST
21 COM_LIST=cerat_muliti_dataset(train_TFIDF_D)
```

ساخت ۱۰۰ داده train و Dev



### ۷-۲-۳- انتخاب الگوریتم های پایه و ساخت ۱۰۰ دسته بند با استفاده از هر الگوریتم

الگوریتم، روشی که برای جستجوی الگو در داده‌ها مورد استفاده قرار می‌گیرد را تعیین می‌کند و در واقع مانند یک روال ریاضی برای حل یک مساله خاص است. الگوریتم‌های گوناگونی برای تحلیل داده<sup>۴۷</sup> موجود هستند و لذا انتخاب الگوریتم داده‌کاوی<sup>۴۸</sup> مناسب یک مساله، برای پژوهشگران و تحلیلگران کاری دشوار است. برخی از سازمان‌ها به دلیل دشواری انتخاب الگوریتم داده‌کاوی مناسب، به طور مکرر از برخی الگوریتم‌های داده‌کاوی استفاده می‌کنند. هنگامی که یک الگوریتم نامناسب پیاده‌سازی می‌شود، دانش کشف شده اغلب برای سازمان کاربرد است زیرا از اطلاعات صحیحی استخراج نکرده و این امر می‌تواند منجر به تصمیم‌گیری‌های غلط در کسب‌وکار شود.

هیچ راهنمای مشخصی برای آنکه پژوهشگران یا تحلیلگران چگونه الگوریتم انتخاب کنند وجود ندارد. انتخاب یک الگوریتم مشخص امری بسیار پیچیده است، لذا در این پژوهش برای ارتقای نتایج داده‌کاوی از چندین الگوریتم استفاده شده و پردازش‌ها را با الگوریتم‌های مختلف تکرار و دست آخر از تکنیک رای گیری اکثریت استفاده می‌شود.

گاه نیاز به استفاده از چندین الگوریتم برای حل یک مساله واحد جهت حل فازهای مختلف مساله است. در مجموع می‌توان گفت هدف مساله، ساختار مجموعه داده موجود، نتایج مورد انتظار در خروجی، شناخت داده‌کاو از یک الگوریتم و مولفه‌های پیکربندی پایگاه داده در انتخاب الگوریتم مناسب داده‌کاوی تاثیرگذار هستند [1]. الگوریتم‌های پایه مورد استفاده در این مقاله شامل SVC خطی، بیز ساده، رگرسیون منطقی و درختان اضافی می‌باشد. همه این الگوریتم‌ها با ۳ رویکرد BOW\_Binary، BOW\_TF و TF-IDF مدل‌سازی شده‌اند که به شرح و تفصیل هر یک می‌پردازیم.

کیسه کلمات فرکانس: کیسه کلمات (BOW) یک مدل در پردازش زبان‌های طبیعی است که با هدف دسته‌بندی مستندات و متون استفاده می‌شود. ایده اصلی آن، به این صورت هست که به هر کدام از کلمات یک عدد Unique نسبت می‌دهیم و Feature بدست آمده بر اساس فرکانس تکرار هر کدام از کلمات به دست خواهد

آمد. به طور متعارف، مجموعه‌ای از کلمات با استفاده از فرکانس کلمه یا اهمیت به عنوان ویژگی اغلب برای طبقه‌بندی یا بازیابی سند استفاده می‌شود. یکی از مشکلات این روش، انفجار ابعاد است که با افزایش کلمات، بعد هم افزایش می‌یابد [28].

کیسه کلمات دودویی: تنها بودن و یا نبودن یک کلمه در آن داده یا همان رکورد را بررسی می‌کنه که آن را با صفر و یک نشون میده.

فراوانی کلمه کلیدی<sup>۴۹</sup>:

جدول تکرار کلمات (TF): روش اصلی یافتن صفحات مرتبط با یک جستجو، روش سنجش تعداد تکرار (TF) یک کلمه است. هر چه که یک کلمه در یک صفحه بیشتر تکرار شده باشد، آن صفحه ارتباط بیشتری با آن کلمه دارد بنابراین اگر کاربر کلمه‌ای را جستجو کرد، صفحاتی را نمایش خواهیم داد که آن کلمه در آنها بیشتر تکرار شده باشد.

جدول معکوس تکرار در صفحات (IDF): به ازای هر لغت، باید فرمولی را استفاده کنیم که هر چه تکرار یک لغت در یک کتاب کمتر باشد، به آن امتیاز بیشتری بدهد (رابطه معکوس). مثلاً می‌توانیم از فرمول  $N/DF$  استفاده کنیم که  $N$  تعداد کل کتابها و  $DF$  تعداد کتابهای حاوی آن لغت است. با این فرمول هر چه یک لغت کمتر تکرار شده باشد، عدد بزرگتری تولید می‌شود.

جدول TF-IDF: با داشتن جدول TF و IDF، می‌توانیم مرتبط بودن یک لغت با یک صفحه را با ضرب این دو در هم نمایش دهیم:

$$TF * IDF = \text{میزان ارتباط لغت با یک صفحه}$$

یعنی هر چه یک لغت در یک متن بیشتر به کار رفته باشد و در سایر متن‌ها خیلی کم به کار رفته باشد، امتیاز آن صفحه برای آن لغت بیشتر می‌شود که منطقی هم به نظر می‌رسد.

به هر کلمه در متن یک وزن اختصاص دهیم. با این کار، می‌توانیم اهمیت یک کلمه را در فرآیند مهندسی ویژگی بهتر شناسایی کرده و در نهایت ویژگی‌های بهتری را برای تزریق به الگوریتم‌های بعدی مانند طبقه‌بندی یا خوشه‌بندی داشته باشیم. مقدار TF-IDF مخفف دو کلمه است: TF به معنی Term Frequency یعنی تعداد تکرار یک کلمه در یک متن و عبارت IDF به معنی Inverse Document Frequency که می‌توان آن را به

---

<sup>۴۹</sup> TF-IDF

برعکس تعداد تکرار در متون ترجمه کرد. وزن دهی TF-IDF طبق فرمول (۱) محاسبه می شود.

$$f(w) = TF(w).IDF(w) = TF(w). \log N/(n(w) + 1) \quad (۱)$$

۱-۷-۲-۳ الگوریتم *Linear SVC* :

این الگوریتم مشابه SVC هست بدین صورت که پارامتر Kernel آن برابر با Linear می باشد و برای داده های در مقیاس بزرگ مناسب است . در ابتدا باید ماژول و کتابخانه این الگوریتم فراخوانی شود.

```
10 from sklearn.svm import SVC
```

### فراخوانی کتابخانه Linear SVC

جهت ساخت مدل، پارامترهای الگوریتم *LinearSVC*، به صورت زیر مقداردهی شده است:

$C=2.5$  : مقدار ذاتی  $C$  همان یک است. هرچقدر  $C$  بیشتر باشد، حول همسایه ها می چرخد و اگر سطحی که در می آید کمتر باشد، سطح تخت می شود. مقدار  $C$  هر چقدر کمتر باشد میزان همواری بیشتر است مقدار این پارامتر بسته به نوع دیتاستی است که ما در اختیار داریم انتخاب شده چون بهترین نتیجه را می دهد

$Penalty="l2"$  : نوعی پناالتی است که به دو نوع  $L1$  و  $L2$  است . پیش فرض آن  $L2$  است و از آن بهترین جواب بدست می آید

$dual=False$  : در این الگوریتم مسئله ما مقدار  $dual$  نیست یعنی این معیار برای پارامترهای دومجهولی است اما در اینجا ما یک مجهول داریم و نیازی به تغییر این پارامتر نمی باشد.

$tol=0.0025$  : این تابع بر اساس بهینه سازی عمل می کند. یعنی از آن حد کمتر تغییر کند سیستم را  $break$  می کند. هرچقدر عدد کوچک باشد دقت بیشتر می شود و جواب بهتری خواهیم گرفت اما زمان محاسبات بالاتر خواهد رفت. بهینه ترین مقدار  $0.0025$  در نظر گرفته شده است.

```
2 clf=LinearSVC(C=2.5, penalty="l2", dual=False ,tol=0.0025).fit(train_TFIDF,Ytrain)
```

### ساخت مدل با Linear SVC

برای ذخیره مدل، پیش بینی ها و ارزیابی آن از ماژول *Pickle* استفاده شده است. *Pickle* یک ماژول از پیش ساخته شده در پایتون است که یک ذخیره سازی قابل حمل از داده های ساخت یافته را فراهم می کند. از تابع  $dump()$  در ماژول *pickle* برای ذخیره کردن و از تابع  $load()$  برای فراخوانی مجدد جریان داده استفاده شده است.

برای استفاده از این ماژول ابتدا باید آن را فراخوانی کنیم .

```
3 import pickle
```

فراخوانی کتابخانه pickle

تابع دامپ دو پارامتر ورودی دارد، در ابتدا باید داده‌ای که می‌خواهیم ذخیره کنیم را به تابع می‌دهیم در این جا با `clf0` که همان مدل ساخته شده می‌باشد، مقداردهی شده و سپس نام فایل مورد نظر و اینکه می‌خواهیم برای چه کاری این فایل را باز کنیم، `'wb'` مشخص می‌کند که فایل جهت نوشتن باز شده است، مدل تولید شده با نام `bay_model_0.pkl` ذخیره شده است.

```
3 z0=pickle.dump(clf0, open('bay_model_0.pkl','wb'))
```

و با استفاده از تابع `load()` امکان استفاده و فراخوانی مدل ذخیره شده فراهم می‌شود. `'rb'` مشخص می‌کند که فایل `bay_model_0.pkl` جهت خواندن باز شده است.

```
4 z0=pickle.load(open('bay_model_0.pkl','rb'))
```

با استفاده از مدل ذخیره شده بر روی داده‌های اعتبار سنجی پیش‌بینی و ارزیابی را انجام شده و به ترتیب در لیست‌های `Pred_clf[]` و `Fscore_clf[]` قرار داده شده است. لیست `allModels_clf[]` نیز محتوی مدل و ارزیابی آن می‌باشد. معیار ارزیابی مورد استفاده `F-Score` می‌باشد.

```
5 p_clf0=z0.predict(dev_TFIDF)
6 F_clf0 = precision_recall_fscore_support(Ydev, p_clf0, average='micro')[2]
```

به صورت نمونه در این قسمت کد مربوط به `tf-idf` ذکر شده است.

```
In [82]: #tfidf
# clf0=LinearSVC(C=2.5, penalty="L2", dual=False ,tol=0.0025).fit(train_TFIDF,Ytrain) ##only for frist run
# z0=pickle.dump(clf0, open('bay_model_2.pkl','wb')) ##only for frist run
z0=pickle.load(open('bay_model_2.pkl','rb'))
p_clf0=np.array(z0.predict(dev_TFIDF).tolist(),dtype=float)
F_clf0 = precision_recall_fscore_support(Ydev, p_clf0, average='micro')[2]
Pred_clf[2]=p_clf0
Fscore_clf[2]=F_clf0
print('F_clf2=',F_clf0)
```

F\_clf2= 0.7946884576098059

## TF-IDF

بعد از ساخت مدل‌ها با داده‌های آموزش کیسه کلمات دودویی، فرکانس کلمات و فراوانی کلمه کلیدی، طبق بررسی‌های انجام شده مدلی که با استفاده از فراوانی کلمه کلیدی ساخته شده بود بهترین `F-Score` را داشت. سپس در ادامه ارزیابی مدل بر روی داده‌های آزمایش نیز انجام شد.

```
In [83]: 1 p_clf1=np.array(z0.predict(test_TFIDF).tolist(),dtype=float)
2 F_clf1 = precision_recall_fscore_support(Ytest, p_clf1, average='micro')[2]
3 Pred_clf_test[2]=p_clf1
4 Fscore_clf_test[2]=F_clf1
5 print('F_clf2=',F_clf1)
```

F\_clf2= 0.8003064351378958

با توجه به بهینه بودن این رویکرد، ۱۰۰ مدل Linear SVC با استفاده از آن ساخته شد. ساخت این مدل‌ها با استفاده از ۱۰۰ نمونه داده آموزشی که در مراحل قبل ساخته شده و در لیست COM\_LIST قرار گرفته، انجام شده است، بدین صورت که ابتدا ویژگی‌ها و برجسب این داده آموزشی با دستورات زیر جدا و به ترتیب در Xcom و Ycom قرار گرفته که به عنوان ورودی به الگوریتم LinearSVC جهت ساخت ۱۰۰ داده شده است. سپس مدل‌های ساخته شده در فایل‌های جداگانه ذخیره شده است. در ضمن پیش بینی و ارزیابی مدل‌ها به ترتیب در لیست‌های Pred\_clf[] و Fscore\_clf[] قرار داده شده است. تمام این ارزیابی‌ها بر روی داده اعتبارسنجی صورت گرفته است.

```
In [84]: 1 #100 train with tfidf
2
3 for i in range(3,103):
4     ##only for frist run
5     # j=i-3
6     # ff=np.array(COM_LIST[j],dtype=float)
7     # fff=pd.DataFrame(ff)
8     # num=fff.shape[1]-1
9     # xcom=fff.iloc[:,0:num].as_matrix()
10    # Ycom=pd.DataFrame(fff.iloc[:,-1].as_matrix())
11    # clf0=LinearSVC(C=2.5, penalty="l2", dual=False ,tol=0.0025).fit(xcom,Ycom)
12    # z0=pickle.dump(clf0, open('bay_model_%d.pkl'%i,'wb'))
13    z0=pickle.load(open('bay_model_%d.pkl'%i,'rb'))
14    p_clf0=z0.predict(dev_TFIDF)
15    F_clf0 = precision_recall_fscore_support(ydev, p_clf0, average='micro')[2]
16    Pred_clf[i]=p_clf0
17    Fscore_clf[i]=F_clf0
18    # print('F_clf%d=%i,F_clf0)
19
20
```

## ساخت ۱۰۰ مدل LinearSvc

### 3-2-7-2- الگوریتم Naïve bayes

تکنیک نایوبیز مبتنی بر قضیه بیزین<sup>۵۰</sup> است و به خصوص زمانی که ابعاد ورودی بالا باشد، مناسب است. برخلاف سادگی، بیز می‌تواند از بسیاری روش‌های طبقه‌بندی پیچیده تر بهتر عمل کند [29]. برای استفاده از این الگوریتم، ماژول آن بدین صورت فراخوانی شده است.

```
16 from sklearn.naive_bayes import GaussianNB
```

### فراخوانی کتابخانه Naïve bayes

<sup>۵۰</sup> Bayesian

جهت ساخت مدل، از تابع `fit()` برای الگوریتم نایویز استفاده شده است. NB ذاتا پارامتر زیادی جهت تغییر ندارد و به همین دلیل پرانتز خالی می باشد.

```
2 clf0=GaussianNB().fit(train_TFIDF,Ytrain)
```

```
In [90]: 1 #tfidf
2 # clf0=GaussianNB().fit(train_TFIDF,Ytrain) ##only for frist run
3 # z0=pickle.dump(clf0, open('bay_model_105.pkl','wb')) ##only for frist run
4 z0=pickle.load(open('bay_model_105.pkl','rb'))
5 p_clf0=np.array(z0.predict(dev_TFIDF).tolist(),dtype=float)
6 F_clf0 = precision_recall_fscore_support(Ydev, p_clf0, average='micro')[2]
7 Pred_clf[105]=p_clf0
8 Fscore_clf[105]=F_clf0
9 print('F_clf105=',F_clf0)
```

```
F_clf105= 0.6736465781409602
```

### ساخت مدل با Naïve bayes

سپس ۱۰۰ مدل با استفاده از ۱۰۰ نمونه داده آموزشی ساخته و بر روی داده های اعتبارسنجی ارزیابی صورت گرفته است. بعد از آن مدل ها ذخیره گردیدند. در ضمن پیش بینی مدل در لیست `Pred_clf`، ارزیابی در لیست `Fscore_clf` قرار گرفته شده است. سپس ارزیابی بر روی داده های آزمایش صورت گرفت.

```
In [91]: 1 p_clf1=np.array(z0.predict(test_TFIDF).tolist(),dtype=float)
2 F_clf1 = precision_recall_fscore_support(Ytest, p_clf1, average='micro')[2]
3 Pred_clf_test[105]=p_clf1
4 Fscore_clf_test[105]=F_clf1
5 print('F_clf105=',F_clf1)
```

```
F_clf105= 0.6557711950970377
```

```
In [92]: 1 for i in range (106,206):
2         ##only for frist run
3         # j=i-106
4         # ff=np.array(COM_LIST[j],dtype=float)
5         # fff=pd.DataFrame(ff)
6         # num=fff.shape[1]-1
7         # xcom=fff.iloc[:,0:num].as_matrix()
8         # Ycom=pd.DataFrame(fff.iloc[:,1].as_matrix())
9         # clf0=GaussianNB().fit(xcom,Ycom)
10        # z0=pickle.dump(clf0, open('bay_model_%d.pkl'%i,'wb'))
11        z0=pickle.load(open('bay_model_%d.pkl'%i,'rb'))
12        p_clf0=z0.predict(dev_TFIDF)
13        F_clf0 = precision_recall_fscore_support(Ydev, p_clf0, average='micro')[2]
14        Pred_clf[i]=p_clf0
15        Fscoer_clf[i]=F_clf0
16        # print('F_clf%d='%i,F_clf0)
17
```

### ۳-۷-۲-۳- الگوریتم *Logistic regression* ساخت ۱۰۰ مدل با Naïve bayes

رگرسیون منطقی یک الگوریتم یادگیری ماشین آماری است که داده‌ها را با در نظر گرفتن متغیرهای خروجی نهایی طبقه‌بندی می‌کند و سعی می‌کند یک خط لگاریتمی برای تمایز بین آن‌ها ایجاد کند [30].

جهت ساخت مدل، پارامتر الگوریتم *LogisticRegression*، به صورت زیر مقداردهی شده است:

$C=8$  : پارامتر پایدارسازی معکوس - یک متغیر کنترلی که اصلاح قدرت را از طریق ارتباط معکوس با تنظیم‌کننده لامبدا حفظ می‌کند .

با توجه به این *Scikit* آن را اینگونه توصیف می‌کند :

$C = 1/\lambda$  : رابطه، این است که کاهش  $C$  - تنظیم‌کننده لامبدا را تقویت خواهد کرد [31].

بعد از ساخت مدل، با ماژول *pickle* ذخیره شده است .پیش بینی و ارزیابی بر روی داده اعتبارسنجی انجام شده و هر کدام به طور جداگانه در لیست نگه‌داری می‌شود.

```
In [99]: 1 #tfidf
2 # clf0=LogisticRegression(C=8).fit(train_TFIDF,Ytrain) ##only for frist run
3 # z0=pickle.dump(clf0, open('bay_model_208.pkl','wb')) ##only for frist run
4 z0=pickle.load(open('bay_model_208.pkl','rb'))
5 p_clf0=np.array(z0.predict(dev_TFIDF).tolist(),dtype=float)
6 F_clf0 = precision_recall_fscore_support(Ydev, p_clf0, average='micro')[2]
7 Pred_clf[208]=p_clf0
8 Fscoer_clf[208]=F_clf0
9 print('F_clf208=',F_clf0)
```

F\_clf208= 0.810520939734423

ساخت مدل با *logistic regression* و ارزیابی بر روی dev

سپس در ادامه ارزیابی بر روی داده‌های آزمایش انجام شده است.

```
In [100]: 1 p_clf1=np.array(z0.predict(test_TFIDF).tolist(),dtype=float)
2 F_clf1 = precision_recall_fscore_support(ytest, p_clf1, average='micro')[2]
3 Pred_clf_test[208]=p_clf1
4 Fscore_clf_test[208]=F_clf1
5 print('F_clf208=',F_clf1)
```

F\_clf208= 0.8158835546475995

۱۰۰ نمونه مدل با استفاده از الگوریتم پایه لوجستیک رگرسیون ساخته شده است. ارزیابی و پیش بینی همه ۱۰۰ مدل بر روی داده اعتبارسنجی انجام شده است و به طور جداگانه در لیست‌ها نگه‌داری شده‌اند.

```
In [101]: 1 for i in range (209,308):
2 #only for frist run
3 # j=i-209
4 # ff=np.array(COM_LIST[j],dtype=float)
5 # fff=pd.DataFrame(ff)
6 # num=fff.shape[1]-1
7 # xcom=fff.iloc[:,0:num].as_matrix()
8 # Ycom=pd.DataFrame(fff.iloc[:,-1].as_matrix())
9 # clf0=LogisticRegression(C=8).fit(xcom,Ycom)
10 # z0=pickle.dump(clf0, open('bay_model_%d.pkl'%i,'wb'))
11 z0=pickle.load(open('bay_model_%d.pkl'%i,'rb'))
12 p_clf0=z0.predict(dev_TFIDF)
13 F_clf0 = precision_recall_fscore_support(ydev, p_clf0, average='micro')[2]
14 Pred_clf[i]=p_clf0
15 Fscore_clf[i]=F_clf0
16 # print('F_clf%d'%i,F_clf0)
17
```

ساخت ۱۰۰ مدل با logistic regression

۴-۷-۲-۳- الگوریتم *ExtraTrees*:

درختان اضافی اصلاح دیگری از کیسه هستند که در آن درخت‌های تصادفی از نمونه‌های مجموعه داده آموزشی ساخته می‌شوند. شما می‌توانید یک مدل درختی اضافی برای طبقه‌بندی با استفاده از کلاس *ExtraTreesClassifier* بسازید.

جهت استفاده از این الگوریتم، به صورت زیر فراخوانی آن انجام گرفته است.

```
from sklearn.ensemble import ExtraTreesClassifier
```

فراخوانی الگوریتم *Extra trees*



برای این الگوریتم، پارامتری تغییر داده نشده و از پیش فرض‌های خود الگوریتم استفاده و مدل ساخته شده است. ارزیابی مدل بر روی داده اعتبارسنجی انجام شده است.

```
In [107]: 1 #tfidf
2 # clf0=ExtraTreesClassifier().fit(train_TFIDF,Ytrain) ##only for frist run
3 # z0=pickle.dump(clf0, open('bay_model_310.pkl','wb')) ##only for frist run
4 z0=pickle.load(open('bay_model_310.pkl','rb'))
5
6 p_clf0=np.array(z0.predict(dev_TFIDF).tolist(),dtype=float)
7 F_clf0 = precision_recall_fscore_support(ydev, p_clf0, average='micro')[2]
8 Pred_clf[310]=p_clf0
9 Fscore_clf[310]=F_clf0
10 print('F_clf310=',F_clf0)
```

F\_clf310= 0.702247191011236

ارزیابی مدل بر روی داده های Dev

در ادامه ارزیابی بر روی داده آزمایش نیز انجام شده است.

```
In [108]: 1 p_clf1=np.array(z0.predict(test_TFIDF).tolist(),dtype=float)
2 F_clf1 = precision_recall_fscore_support(ytest, p_clf1, average='micro')[2]
3 Pred_clf_test[310]=p_clf1
4 Fscore_clf_test[310]=F_clf1
5 print('F_clf310=',F_clf1)
```

F\_clf310= 0.6999489274770173

سپس ۱۰۰ نمونه مدل از این الگوریتم ساخته و ذخیره شده است.

```
In [109]: 1 for i in range (311,411):
2 #only for frist run
3 # j=i-311
4 # ff=np.array(COM_LIST[j],dtype=float)
5 # fff=pd.DataFrame(fff)
6 # num=fff.shape[1]-1
7 # xcom=fff.iloc[:,0:num].as_matrix()
8 # Ycom=pd.DataFrame(fff.iloc[:,-1].as_matrix())
9 # clf0=ExtraTreesClassifier().fit(xcom,Ycom)
10 # z0=pickle.dump(clf0, open('bay_model_%d.pkl'%i,'wb'))
11 z0=pickle.load(open('bay_model_%d.pkl'%i,'rb'))
12 p_clf0=z0.predict(dev_TFIDF)
13 F_clf0 = precision_recall_fscore_support(Ydev, p_clf0, average='micro')[2]
14 Pred_clf[i]=p_clf0
15 Fscoer_clf[i]=F_clf0
16 # print('F_clf%d='%i,F_clf0)
17
```

ساخت ۱۰۰ مدل با Extra trees  
در انتها تمام لیست‌های Pred\_clf، F\_score و allModels\_clf با استفاده از ماژول pickle در فایل‌های جداگانه‌ای ذخیره شده‌است.

```
In [113]: 1 #result dev
2 #z0_Pred_clf=pickle.dump(Pred_clf[:411], open('Pred_clf0.pkl','wb')) ##only for frist run
3 z0_Pred_clf=pickle.load(open('Pred_clf0.pkl','rb'))
4
5 # z0_Fscoer_clf=pickle.dump(Fscoer_clf[:411], open('Fscoer_clf0.pkl','wb'))##only for frist run
6 z0_Fscoer_clf=pickle.load(open('Fscoer_clf0.pkl','rb'))
7
8 z0_allModels_clf=[]
9 for i in range (len(z0_Fscoer_clf)):
10 z0=pickle.load(open('bay_model_%d.pkl'%i,'rb'))
11 temp=[]
12 temp.append(z0)
13 temp.append(z0_Fscoer_clf[i])
14 z0_allModels_clf.append(temp)
```

در ادامه ارزیابی‌هایی که بر روی داده‌های آزمایش نیز انجام شده بود نیز به طور جداگانه ذخیره شد.

### ۳-۲-۸- گرفتن میانگین F-score برای هر مدل و اعمال فیلترینگ

برای اینکه مدل دسته جمعی عملکرد بهتری نسبت به مدل تکی داشته باشد، باید هرکدام از مدل‌های تشکیل

دهنده آن مستقل از بقیه بوده و دارای دقتی بیشتر از میانگین امتیاز مدل‌ها داشته باشد.

بدین سه لیست خالی به نام‌های `allModels_clf`، `Fscore_clf` و `Pred_clf` به ترتیب برای مدل‌ها، امتیازها و پیش‌بینی‌ها تعریف شده است. طبق خط شماره ۵ میانگین امتیاز همه مدل‌ها به دست می‌آید. سپس با استفاده از یک حلقه `for` و پیمایش امتیاز تمام مدل‌ها، امتیازی که از میانگین امتیازها بیشتر باشد به لیست `Fscore_clf` اضافه می‌شود. مدل و پیش‌بینی هم به لیست‌های `allModels` و `Pred_clf` اضافه می‌گردد. همان طور که مشاهده می‌کنید از بین ۴۱۱ مدل ساخته شده فقط ۲۰۵ مدل امتیازی بالاتر از میانگین امتیازها داشتند.

```
In [196]: 1 allModels_clf=[]
2 Fscore_clf=[]
3 Pred_clf=[]
4
5 mean_fscore = sum(z0_Fscore_clf) / len(z0_allModels_clf)
6 for i in range(len(z0_allModels_clf)):
7     if z0_Fscore_clf[i] > mean_fscore :
8         Fscore_clf.append(z0_Fscore_clf[i])
9         allModels_clf.append(z0_allModels_clf[i])
10        Pred_clf.append(z0_Pred_clf[i])
11
12 len(allModels_clf),len(z0_allModels_clf)
```

Out[196]: (205, 411)

سپس رویکردها بر روی داده‌های اعتبار سنجی، اجرا و ارزیابی شده‌است.

```
[209]: 1 best_clasfire_test_filter=best_clasfire(allModels_clf_testf,Fscore_clf_testf,Pred_clf_testf,'test')
2 print ('best_clasfire_test_filter=',best_clasfire_test_filter[0])
3
4 voting_test_filter=modle_voting(Ytest,Pred_clf_testf,'test')
5 print ('voting_test_filter=',voting_test_filter[0])
6
7 forward_test_filter=forward (allModels_clf_testf,Pred_clf_testf,Ytest,'test')
8 print ('fscore_Forward_test_filter=',forward_test_filter[0])
9
10 backward_test_filter= backward(allModels_clf_testf,Pred_clf_testf,Ytest,'test')
11 print ('fscore_backward_test_filter=',backward_test_filter[0])
12
13 best_predict_final_test_filter = backward_test_filter[1]+forward_test_filter[1]+[np.array(voting_test_filter[1],dtype=float)]
14 all_comb_test_filter=all_comb (best_predict_final_test filter,Ytest,"test")
15 print('best_score_final_test_filter=',all_comb_test_filter)
```

```
best_clasfire_test_filter= 0.8158835546475995
voting_test_filter= 0.8128192032686414
fscore_Forward_test_filter= 0.813585291113381
fscore_backward_test_filter= 0.813585291113381
best_score_final_test_filter= 0.8133299284984679
```

### ۳-۲-۹- انتخاب بهترین دسته بند

بدین منظور تابعی برای انتخاب بهترین دسته بند با نام `best_clasfire` نوشته شد. ورودی این تابع لیست تمام مدل‌ها، ارزیابی‌ها و پیش‌بینی‌ها می‌باشد. با تبدیل لیست مدل‌ها با استفاده از تابع `DataFrame` کتابخانه پاندا به دیتافریم تبدیل و نام‌های `Modle` و `Fscore` برای ستون‌های آن در نظر گرفته شده است. با استفاده از تابع `max` بیشترین `Fscore` از بین تمام `Fscore`ها انتخاب شده است. سپس طبق خط شماره ۴ تکه کد زیر بیشترین امتیاز بدست آمده با امتیازهای دیتافریم `all_M_D_frame` مقایسه شده و مدل متناسب با آن بدست آمده است. که در `best_Modle` به صورت یک تاپلی از مدل و امتیاز قرار گرفته است. با استفاده از مشخصه `index` اندیس صفر آن که همان نام مدل است به دست آمده است. در انتها نام و امتیاز و پیش‌بینی بهترین مدل برگردانده می‌شود.

```
In [116]: 1 def best_clasfire(list_allModle,list_Fscore,list_pred,name):
2     all_M_D_frame =pd.DataFrame(list_allModle,columns=['Modle','Fscore'])
3     best_Fscore=(max(list_Fscore))
4     best_Modle=all_M_D_frame[all_M_D_frame['Fscore']==best_Fscore]
5     best_pre=list_pred[best_Modle.index[0]]
6     return best_Fscore,best_pre,best_Modle
```

با استفاده از تابع تعریف شده بهترین مدل با امتیاز ۰,۸۱۰ انتخاب شد.

```
1 best_clasfire_dev=best_clasfire(z0_allModels_clf,z0_Fscoer_clf,z0_Pred_clf,'dev')
2 print ('best_clasfire_dev=',best_clasfire_dev[0])
```

### ۳-۲-۱۰- رای گیری اکثریت

هر مدل یک پیش‌بینی (رای‌گیری) برای هر نمونه آزمایش و پیش‌بینی خروجی نهایی همان چیزی است که بیش از نیمی از آرا را دریافت می‌کند. اگر هیچ یک از پیش‌بینی‌ها بیش از نیمی از آرا را به دست نیاورند، می‌توانیم بگوییم که روش گروهی نمی‌تواند یک پیش‌بینی پایدار برای این نمونه ایجاد کند. اگر چه این یک تکنیک بسیار مورد استفاده است، شما ممکن است بیش‌ترین پیش‌بینی را امتحان کنید (حتی اگر آن کم‌تر از نیمی از آرا) به عنوان پیش‌بینی نهایی باشد [32]. که در این پژوهش رویکرد بیش‌ترین پیش‌بینی انتخاب شده است. در ابتدا ماژول رای‌گیری اکثریت به برنامه اضافه شده است.

```
17 from sklearn.ensemble import VotingClassifier
```

فراخوانی ماژول رای‌گیری اکثریت

فایل پیش‌بینی‌های که در فایل ذخیره شده بود با ماژول `pickle` جهت خواندن فراخوانی شده است. جهت

انجام رای گیری اکثریت یک تابع نوشته شده است، که لیست پیش بینی ها و برچسب های داده اعتبارسنجی را دریافت کرده ، سپس یک لیست خالی به نام Columns ایجاد شده با استفاده از حلقه for ای که به اندازه طول یکی از لیست پیش بینی ها می باشد با تابع append لیستی از تمام شماره های پیش بینی ها به آن اضافه می شود . این لیست به اندازه طول تمام مدل هایی است که ساخته شده است.

```
5 Columns=[]
6 for i in range(len(predict_list[0])) :
7     Columns.append('pred%d'%i)
```

Columns= [pred1, pred2, pred3 ... pred411]

سپس لیست پیش هایی با استفاده از تابع DataFrame و با در نظر گرفتن لیست Columns به عنوان سرستون ها به یک دیتافریم تبدیل شده است.

```
8 df=pd.DataFrame(predict_list,columns=Columns)
```

یک لیست خالی به نام pre\_final برای پیش بینی های نهایی ایجاد شده است.

حلقه for دومی به تعداد پیش بینی های انجام شده توسط مدل ها که در دیتافریم موجود می باشد، بررسی می کند که برای نمونه تمام مدل ها چه برچسبی را پیش بینی کرده اند و طبق آن ها با استفاده از max\_freq رای گیری اکثریت انجام داده است. و در انتها امتیاز Fscore و پیش بینی نهایی را برمی گرداند. امتیاز Fscore با این رویکرد ۰٫۸۲۰ به دست آمده است.

```
In [117]: 1 def modle_voting (lable,predict_list,name):
2     Columns=[]
3     for i in range(len(predict_list[0])) :
4         Columns.append('pred%d'%i)
5     df=pd.DataFrame(predict_list,columns=Columns)
6     pre_final=[]
7     for i in range (df.shape[1]):
8         a=df['pred%d'%i].tolist()
9         max_freq=np.bincount(a).tolist()
10        pre_final.append(max_freq.index(max(max_freq)))
11    F_eclf1=precision_recall_fscore_support(lable,pre_final, average='micro')[2]
12    return F_eclf1,pre_final
13
14
```

پیش بینی نهایی

۳-۲-۱۱- انتخاب رو به جلو

انتخاب رو به جلو یک نوع رگرسیون گام به گام است که با یک مدل خالی شروع می شود و طبقه بندی را

که به صورت نزولی بر مبنای یکی از معیارهای ارزیابی مرتب شده اند را یک به یک اضافه می کند و بهترین متغیر ، با برخی از معیارهای از پیش تعیین شده تعیین می شود و به مدل اضافه می گردد. معیار استفاده شده برای تعیین اینکه کدام مدل اضافه شود ، متفاوت است. در هر مرحله رو به جلو، شما یک طبقه بند را اضافه می کنید. در اینجا برای انجام مقایسات، Fscore ها مدنظر قرار داده شده اند [33].

بدین منظور تابعی با نام forward نوشته شده است. این تابع لیست مدل ها، پیش بینی ها، برچسب و نام را به عنوان ورودی دریافت می کند.

```
1 def forward (list_modles,list_pred,lable,name):
```

در این تابع ابتدا با استفاده از یک حلقه for تمام مدل ها، پیش بینی ها و امتیازهایشان در لیست clf\_forward به صورت واحد قرار گرفته اند . این عملیات با استفاده از تابع append لیست انجام گرفته است.

```
2     clf_forward=[]
3     for i in range (len(list_modles)):
4         temp=[]
5         temp.append(list_modles[i][1])
6         temp.append(list_modles[i][0])
7         temp.append(list_pred[i])
8         clf_forward.append(temp)
9
```

قرار دادن تمام مدل ها و پیش بینی ها و امتیازهایشان در یک

سپس لیست به صورت نزولی براساس Fscore مرتب شده است. این مرتب سازی با استفاده از تابع Sorted انجام گرفته است . در این تابع reverse=False لیست داده شده clf\_forward را به صورت نزولی مرتب می کند. اولین مدل و Fscore لیست را به عنوان بهترین مدل و امتیاز در نظر گرفته شده است.

```
10     sorted(clf_forward, key = itemgetter(0),reverse=False)
11     best_score_forward=clf_forward[0][0]# افسکور
12     # print(best_score_forward)
13     best_pre_forward=[list_pred[0]]
```

مرتب کردن لیست به صورت نزولی

همان طور که مشاهده می کنید با استفاده از یک حلقه مدل ها یکی یکی به لیست templist اضافه شده اند و رای گیری اکثریت انجام شده، امتیاز و پیش بینی ها در دو لیست score\_forward و pre\_forward قرار گرفته است.

```

16 for i in range(len(list_pred)):
17     templist=best_pre_forward+[clf_forward[i][2]]
18     p_voting= modle_voting (lable,templist,name)
19     score_forward = p_voting[0]
20     pre_forward =p_voting[1]

```

ذخیره کردن امتیاز و پیش بینی ها در دو لیست

در انتها با یک حلقه for مقایسه بین امتیازها صورت گرفته و تابع forward بهترین امتیاز و پیش‌بینی را برمی‌گرداند. بهترین امتیاز این رویکرد ۰.۸۲۲ بدست آمد.

```

22 if score_forward > best_score_forward:
23     best_score_forward = score_forward
24     best_pre_forward=templist
25     # print(best_score_forward)
26 # print('fscore_Forward%s='%name,best_score_forward)
27 return best_score_forward,best_pre_forward

```

مقایسه بین امتیازها

### ۳-۲-۱۲- حذف رو به عقب

انتخاب ویژگی به عقب مرتبط است ، و ممکن است با انتخاب کل مجموعه مدل‌ها شروع می‌شود و از آنجا به عقب کار می‌کند ، ویژگی‌هایی برای پیدا کردن زیر مجموعه بهینه از یک اندازه از پیش تعریف شده را حذف می‌کند [33].

برای این رویکرد هم یک تابع با نام backward تعریف شده‌است.

```

1 def backward (list_modles,list_pred,lable,name):

```

همانند رویکرد قبلی تمام مدل‌ها به همراه پیش‌بینی و امتیازشان به طور مجتمع در یک لیست به اسم clf\_backward قرارگفته است.

```

2 clf_backward=[]
3 for i in range (len(list_modles)):
4     temp=[]
5     temp.append(list_modles[i][1])
6     temp.append(list_modles[i][0])
7     temp.append(list_pred[i])
8     clf_backward.append(temp)

```

قرار دادن تمام مدل ها و پیش بینی ها و امتیازهایشان در یک لیست

در این رویکرد لیست ایجاد شده، به صورت صعودی مرتب شده است.

```

10 sorted(clf_backward, key = itemgetter(0),reverse=True)

```

مرتب کردن لیست به صورت صعودی

یک لیست خالی به نام predictbl ایجاد شده است که با استفاده از یک حلقه for تمام پیش بینی ها در آن قرار داده شده است.

```

12 predictbl=[]
13 for i in range(len(list_pred)):
14     predictbl.append(clf_backward[i][2])

```

ذخیره کردن پیش بینی ها در لیست

رای گیری اکثریت روی پیش بینی تمام مدل ها و برچسب داده اعتبار سنجی انجام گرفته است. در ابتدا امتیاز و پیش بینی حاصل به عنوان بهترین در نظر گرفته شده است.

```

17 backward= modle_voting(lable,list_pred,name)
18 best_score_back=backward[0]
19 best_pre_back = [backward[1]]

```

```

21 templist=list (range (len(predictbl)))
22 for i in range (1,len(predictbl)):
23     templist[i]=predictbl[:-i]

```

یک حلقه for به اندازه طول تمام پیش بینی ها ایجاد شده است. بر روی templist رای گیری اکثریت انجام شده است در صورتی که این امتیاز از بهترین امتیاز قبلی بیشتر است با آن جایگزین می شود. بهترین امتیاز و



پیش‌بینی برگردانده می‌شود. امتیاز نهایی بدست آمده با این رویکرد ۰,۸۲۳ می باشد.

```
25 for i in range(1,len(predictbl)):#
26     p_voting= modle_voting (lable,templist[i],name)
27     if best_score_back < p_voting[0]:
28         best_score_back = p_voting[0]
29 #         print(best_score_back)
30         best_pre_back=templist[i]
31 return best_score_back,best_pre_back
```

۳-۲-۱۳- ترکیب همه رویکردها:

در این رویکرد پیش‌بینی تمام رویکردها با هم جمع و در یک لیست به نام `best_score_final` قرار داده شده و بر روی این پیش‌بینی‌ها و برچسب داده اعتبارسنجی رای‌گیری اکثریت انجام گرفته است. این `Fscore` رویکرد ۰,۸۲۳ بدست آمده است.

```
In [186]: 1 def all_comb (list_allbestpred,lable,name):
2         best_score_final = modle_voting(lable,list_allbestpred,name)
3         return best_score_final[0]
```

امتیاز نهایی در روش ترکیب همه رویکردها

## فصل چهارم :

محاسبات و یافته های تحقیق

#### ۴-۱- داده های مورد استفاده

برای تشخیص نویسنده از مجموعه داده ی Spooky\_Author\_Identification، از مخزن مسابقه سایت Kaggle استفاده شده است. این مجموعه داده شامل ۲ ویژگی به نامهای id (یک شناسه منحصر به فرد برای هر جمله است) و text (متون نوشته شده توسط نویسندگان) و یک مقدار خروجی به نام Author (نویسندگان جملات) است. متون استفاده شده از آثار نویسندگان رمانهای ترسناک Spooky، می باشد. این مجموعه داده شامل سه کلاس (EAP<sup>۵۱</sup>، HPL<sup>۵۲</sup>، MWS<sup>۵۳</sup>) می باشد.

	id	text	author
0	id26305	This process, however, afforded me no means of...	EAP
1	id17569	It never once occurred to me that the fumbling...	HPL
2	id11008	In his left hand was a gold snuff box, from wh...	EAP
3	id27763	How lovely is spring As we looked from Windsor...	MWS
4	id12958	Finding nothing else, not even gold, the Super...	HPL

کلاس EAP تعداد ۷۹۰۰ متن و کلاس HPL ۶۰۴۴ متن و MWS شامل ۵۶۳۵ متن می باشد. که برای هر متن یک شناسه در نظر گرفته شده است.

#### ۴-۲- تنظیم پارامترها

۴-۲-۱- پارامترهای الگوریتم SVC:

۱. C:float,optional(default=1.0)

فرمول محاسبه C به صورت زیر می باشد:

<sup>۵۱</sup>EAP : Edgar Allan Poe

<sup>۵۲</sup> HPL : HP Lovecraft

<sup>۵۳</sup> MWS : Mary Shelley

$$C \sum_{i=1}^n \mathcal{L}(f(x_i), y_i) + \Omega(\omega)$$

مقدار ذاتی  $C$  همان یک است. هرچقدر  $C$  بیشتر باشد، حول همسایه ها می چرخد و اگر سطحی که در می آید کمتر باشد، سطح تخت می شود (هرچه بیشتر باشد سطح خروجی پستی و بلندی بیشتری دارد) مقدار  $C$  هر چقدر کمتر باشد میزان هواری بیشتر است مقدار این پارامتر بسته به نوع دیتاستی است که ما در اختیار داریم.

**LinearSVC(C=2.5, penalty="l2", dual=False, tol=0.0025)**

۲. `penalty:string, 'L1' or 'L2' (default='L2')`

نوعی پناالتی است که به دو نوع  $L1$  و  $L2$  است. پیش فرض آن  $L2$  است و از آن بهترین جواب بدست می آید. اگر  $L2$  را جایگزین کنیم همانطور که در تصویر زیر می بینیم بهترین جواب را ( $Fscore=0.794$ ) خواهیم داشت.

۳. `dual`: در این الگوریتم مسئله ما مقدار `dual` نیست یعنی این معیار برای پارامترهای دو مجهولی است اما در اینجا ما یک مجهول داریم و نیازی به تغییر این پارامتر نمی باشد.

۴. `tol`: یعنی چقدر تابع `svc`، برای بهینه سازی است. این تابع بر اساس بهینه سازی عمل می کند. مقدار حد آستانه برای این پارامتر `tol=0.0025` است، یعنی از آن حد کمتر تغییر کند سیستم را `break` می کند. هرچقدر عدد `tol=0.0025` کوچک باشد دقت بیشتر می شود و جواب بهتری خواهیم گرفت اما زمان محاسبات بالاتر خواهد رفت.

لیست پارامترهای تنظیم شده توسط کاربر، برای الگوریتم `LinearSVC` را در جدول (۱-۴) مشاهده می کنید.

جدول ۴-۱: پارامترهای تنظیم شده توسط کاربر

الگوریتم <code>LinearSVC</code>	
نام پارامتر	مقدار
<b>C</b>	2.5
<b>Penalty</b>	L2
<b>Dual</b>	False

0.0025	Tol
--------	-----

۴-۲-۲- پارامترهای الگوریتم Logistic Regression :

این مدل حتی باینکه در نام خود " regression " دارد ، logistic regression یک روش دسته بندی است. در این جا ما هیچ hyperparameters را تنظیم نمی کنیم . ما فقط باید مدل را بسازیم و آموزش دهیم .

LogisticRegression(C=8)

پارامتر های تنظیم شده برای الگوریتم Logistic Regression را در جدول ذیل مشاهده می کنید.

جدول ۴-۲: پارامترهای تنظیم شده توسط کاربر

الگوریتم Logistic Regression	
نام پارامتر	مقدار
c	8
Class Weight	None

۴-۲-۳- Extra trees :

برای این الگوریتم، پارامتری تغییر داده نشده و از پیش فرض های خود الگوریتم استفاده و مدل ساخته شده است. ارزیابی مدل بر روی داده اعتبارسنجی انجام شده است.

ExtraTreesClassifier()

۴-۲-۴- Naïve Bayes :

ذاتا پارامتر زیادی جهت تغییر ندارد و به همین دلیل پراتز خالی می باشد.

GaussianNB()

مقدار این پارامترها ی ذکر شده به وسیله آزمون و خطا بر روی داده های اعتبارسنجی به دست آمده است.

یعنی با آزمایش مقادیر مختلف، مقداری را برای پارامتر ورودی انتخاب می‌کنیم که بهترین کارایی را داشته است.

#### ۴-۳- روش‌های استفاده شده به منظور مقایسه

روش‌هایی که برای مقایسه با روش پیشنهادی در نظر گرفته شده‌اند، عبارتند از:

##### • SVM:

الگوریتم SVM یکی از الگوریتم‌های معروف در زمینه یادگیری با نظارت است که برای دسته‌بندی و رگرسیون استفاده می‌شود. این الگوریتم به طور هم‌زمان حاشیه‌های هندسی را بیشینه کرده و خطای تجربی دسته‌بندی را کمینه می‌کند لذا به عنوان دسته‌بندی حداکثر حاشیه<sup>۵۱</sup> نیز نامیده می‌شود [23].

##### • NB :

یک الگوریتم یادگیری ماشین برای حل مشکلات طبقه‌بندی است. برای ساخت مدل و پیش‌بینی سریع از قضیه احتمال استفاده می‌کند. اساس قضیه Bayes، بر پایه استقلال ویژگی‌ها است. یک روش مبنا برای طبقه‌بندی متن و حل مشکل قضاوت اسناد درباره تعلق یک سند به یک نویسنده است که از روش تعیین فرکانس کلمه به عنوان ویژگی استفاده می‌کند.

##### • ET :

مجموعه‌ای از درخت تصمیم‌گیری یا رگرسیون غیرمعمول را ایجاد می‌کند با توجه به روش کلاسیک بالا به پایین. دو تفاوت عمده نسبت به سایر روش‌های گروهی، بر پایه دسته‌بندی دارد. اولاً گره‌ها را با انتخاب نقطه‌های برش به صورت تصادفی جدا می‌کند و دوماً از کل نمونه یادگیری (بوت استرپ) برای رشد درخت استفاده می‌کند.

##### : Logistic Regression

---

<sup>۵۱</sup> maximum margin

به بیان ساده‌تر این الگوریتم احتمال رخداد یک رویداد را بر حسب گنجاندن داده‌ها در یک تابع منطقی پیش‌بینی می‌کند. از این‌رو به نام رگرسیون **logit** نیز شناخته می‌شود. از آنجا که این الگوریتم، احتمال را پیش‌بینی می‌کند، مقادیر خروجی آن بین ۰ تا ۱ هستند.

- رای گیری اکثریت:

رای گیری اکثریت ما بین پیش‌بینی برای هر نمونه توسط مدل‌های ساخته شده، بدین صورت که هر پیش‌بینی که بیشترین تکرار را داشته باشد به عنوان پیش‌بینی اصلی انتخاب می‌شود.

- Forward :

انتخاب رو به جلو با یک مدل خالی شروع می‌شود و طبقه‌بندی‌هایی را که به صورت نزولی بر مبنای یکی از معیار ارزیابی **F1** مرتب شده‌اند را یک به یک اضافه می‌کند و در هر بار رای گیری اکثریت صورت می‌پذیرد و بهترین متغیر، با برخی از معیارهای از پیش تعیین شده، تعیین می‌شود و به مدل اضافه می‌گردد.

- Back Ward :

انتخاب ویژگی به عقب مرتبط است، و با انتخاب کل مجموعه مدل‌ها شروع می‌شود و از آنجا به عقب کار می‌کند، ویژگی‌هایی برای پیدا کردن زیر مجموعه بهینه از یک اندازه از پیش تعریف شده را حذف می‌کند.

همانند رویکرد قبلی تمام مدل‌ها در هر بار رای گیری اکثریت صورت می‌پذیرد و بهترین متغیر، با برخی از معیارهای از پیش تعیین شده، تعیین می‌شود و به مدل اضافه می‌گردد.

#### ۴-۴- نتایج

روش پیشنهادی در پایتون<sup>۵۵</sup> ۳,۷ با استفاده از Jupyter پیاده‌سازی و در یک کامپیوتر شخصی با پردازشگر پنج هسته‌ای اینتل<sup>۵۶</sup> با سرعت ۵,۲۰۰ گیگا هرتز و حافظه با دست‌یابی ۸ گیگا بایت به اجرا در آمده است. همچنین این روش بر روی مجموعه داده، برای تشخیص نویسنده از مجموعه داده‌ی Spooky\_Author\_Identification، از مخزن مسابقه سایت Kaggle محک زده شده است. مهم‌ترین معیارهایی که برای برازش کارایی الگوریتم پیشنهادی در نظر گرفته شده است، عبارتند از: نرخ افسکور F1. همان‌طور که در فصل دوم بیان شد؛ (میانگین هارمونی برای دو مقدار بازخوانی و صحت)

جدول ۴-۳- بهترین نتیجه بدست آمده از الگوریتم‌های ترکیب بر روی مجموعه داده تست

نام مجموعه داده	نرخ صحت
Test	۰,۸۲۳۵

جدول ۴-۴- مقایسه نتایج بدست آمده برای مجموعه داده Spooky\_Author\_Identification با سایر روش‌ها

نام روش	نرخ سنجش F- BOW_bin	نرخ سنجش F- BOW	نرخ سنجش F- TFIDF	نرخ سنجش F- TFIDF با اعمال فیلتر
LinearSVC	۰,۷۶۵۵	0.7630	0.8003	-----
NB	۰,۶۵۱۵	0.6521	0.6557	-----
LR	۰,۸۰۶۴	0.8082	۰,۸۱۵۸	-----

<sup>۵۵</sup> PYTHON 3.7

<sup>۵۶</sup> Intel Core i3



-----	0.6999	0.6640	0.6695	ER
0.8128	0.8222	-----	-----	رای گیری اکثریت
0.8135	0.8225	-----	-----	Forward
0.8135	0.8235	-----	-----	Backward
0.8133	0.8220	-----	-----	رای گیری کل

ارزش استفاده از الگوریتم با توجه به نتایج به وضوح قابل رویت می باشد.

#### ۴-۵- نتیجه گیری

نتایج بدست آمده نشان می دهد دقت روش پیشنهادی ترکیب مدل ها در دسته بندی مجموعه داده های تشخیص نویسنده بهتر یا حداقل قابل رقابت با الگوریتم های مطرحی مانند NB و LR و ET و SVM می باشد. البته الگوریتم ET نیز در اصل استفاده از گروهی از الگوریتم های درخت تصمیم می باشد. در ارزیابی الگوریتم های ترکیبی نیز دو رویکرد با اعمال فیلتر بر اساس معیار F1 مورد سنجش قرار گرفت. بدین صورت که مدل هایی که دارای معیارهای F1 بالاتر از حد میانگین F1 ها مربوط به مدل ها انتخاب، و در روش های ترکیب مورد استفاده قرار گرفت.

تعداد مدل های انتخاب شده از مجموع ۴۱۱ مدل ساخته شده برابر ۲۰۵ مدل می باشد.

نتایج بدست آمده از این رویکرد بهبودی نسبت به حالت بدون این رویکرد نداشته است و شاید بتوان این نتیجه را استنباط کرد که در روش های ترکیب هرچه تعداد مدل های رای گیری بیشتر باشد نتیجه بهتری می توان گرفت که البته این موضوع را نباید نادیده گرفت که در این پژوهش الگوریتم های استفاده شده اکثرا دارای نتایج بالای ۰,۶۵ می باشند و شاید اگر الگوریتم های دیگری مانند درخت تصمیم که نتایج آن ها حداکثر ۰,۵۵ بوده (بر اساس آزمون های انجام شده) در این پژوهش دخیل بودند شاهد چنین نتایجی نبودیم.

فلذا این فرضیه را باتوجه به ویژگی های مجموعه داده انتخابی و نوع الگوریتم های مورد استفاده باید در نظر گرفت .

یعنی زمانی که الگوریتم ها همگی دارای نتایج نسبتا بالایی می باشند استفاده از رویکرد فیلتر کردن بر اساس حد آستانه میانگین معیار ارزیابی  $F1$  مدل ها ، بهبودی در نتایج ندارد.

لیکن در استفاده از روش های ترکیب ، می توان نتیجه گرفت کاربرد آن ها ، نتایج را بهبود داده و در درنیای واقعی کمتر دچار خطا می شود و یا به عبارتی دارای پایداری نسبی در نتایج می باشند .

## فصل پنجم :

### نتیجه گیری و مشاهدات

## ۵-۱- خلاصه و نتیجه گیری

در این پژوهش روشی را برای واکشی دانش مطرح کردیم. دانش مورد نظر توسط فرایند داده کاوی و به کمک روش دسته بندی به دست آمده است. دسته بند طراحی شده توسط معیارهای **F1** ارزیابی شدند.

که با استفاده از روش های معمول مجموعه داده آموزشی به قسمت های آموزش و اعتبارسنجی و آزمایش تقسیم و پاکسازی و نرمال سازی بر روی آن ها صورت پذیرفت. سپس با استفاده از سه رویکرد کیسه کلمات دودویی و کیسه کلمات با فرکانس کلمات و فراوانی کلمه کلیدی ویژگی های هر متن استخراج و به بردار تبدیل گردید. که با استفاده از الگوریتم های داده کاوی و استفاده از روش های یادگیری تجمعی بر روی مجموعه داده های اعتبارسنجی سنجیده و سپس بر روی داده آزمایش مورد آزمایش قرار گرفت.

نتایج بدست آمده نشان دادند که استفاده از روش قراردادن میانگین **F1** مدل های ساخته شده به عنوان حد آستانه برای انتخاب مدل ها بهبودی در نتایج نسبت به حالت بدون استفاده از این تکنیک نداشته و شاید بتوان این نتیجه را استنباط کرد که در روش های ترکیب هرچقدر تعداد مدل های رای گیری بیشتر باشد نتیجه بهتری می توان گرفت ، که البته این موضوع را نباید نادیده گرفت که در این پژوهش الگوریتم های استفاده شده اکثرا دارای نتایج بالای بودند و شاید اگر الگوریتم های دیگری که دارای نتایج پایین تری بودند در این پژوهش مورد استفاده قرار می گرفت از مقایسه نتایج می توانست نتیجه ای عکس گرفت.

فلذا این فرضیه را می بایست باتوجه به ویژگی های مجموعه داده انتخابی و نوع الگوریتم های مورد استفاده باید در نظر گرفت. یعنی زمانی که الگوریتم ها همگی دارای نتایج نسبتا بالایی می باشند استفاده از رویکرد فیلتر کردن بر اساس حد آستانه میانگین معیار ارزیابی **F1** مدل ها ، بهبودی در نتایج ندارد.

لیکن در مقایسه نتایج روش های ترکیب با سایر روش ها در این پژوهش، می توان نتیجه گرفت کاربرد آن ها ، در بهبود نتایج موثر هستند و در درنایای واقعی به دلیل به کارگیری رای گیری اکثریت مابین مدل های ساخته شده کمتر دچار خطا می شوند و یا به عبارتی دارای پایداری نسبی در نتایج می باشند .

## ۵-۲- پیشنهادات

برای کارهای آتی در این زمینه می توان پیشنهادات زیر را در نظر گرفت:

- با توجه به اینکه در این پژوهش مدل ها اکثرا با استفاده از فراوانی کلمه کلیدی ساخته شده است ،اضافه نمودن مدل های بیشتر با استفاده از کیسه کلمات دودویی و کیسه کلمات با فرکانس می تواند نتایج متفاوت تری را نشان بدهد.
- یادگیری روش های ترکیبی با در نظر گرفتن اینکه در رای گیری اکثریت هریک از پیش بینی ها می بایست بیش از نیمی از آرا را به دست بیاورند.
- استفاده از ویژگی های ظاهری مانند تعداد کل کلمات ، تعداد کلمات اضافه ، میانگین طول کلمات ، تعداد کلمات استپ ورد ، تعداد کلمات شاخص و .... در هرمتن جهت ساخت مدل مربوطه و ترکیب یا روش های یادگیری تجمعی که در این پژوهش آورده شده.
- به کارگیری قواعد انجمنی و کشف روابط با استفاده از ویژگی ها با سطح تکرار مشخص و فراوانی بالا

- [1] O. Maimon and L. Rokach, "Data Mining and Knowledge Discovery Handbook," pp. 1–15, 2005.
- [2] H. Jang, S. Kim, and T. Lam, "Kaggle Competitions : Author Identification & Stotoil / C-CORE Iceberg Classifier Challenge," pp. 1–21, 2017.
- [3] M. Kestemont *et al.*, "Overview of the author identification task at PAN-2018: Cross-domain authorship attribution and style change detection," *CEUR Workshop Proc.*, vol. 2125, 2018.
- [4] <https://searchsqlserver.techtarget.com/definition/data-mining>
- [5] <https://www.geeksforgeeks.org/basic-concept-classification-data-mining/>
- [6] <https://hub.packtpub.com/what-is-ensemble-learning/>
- [7] <https://simplicable.com/new/ensemble-learning>
- [8] [https://www.tutorialspoint.com/data\\_mining/dm\\_knowledge\\_discovery.htm](https://www.tutorialspoint.com/data_mining/dm_knowledge_discovery.htm)
- [9] <https://www.techopedia.com/definition/1181/data-mining>
- [10] <https://www.geeksforgeeks.org/basic-concept-classification-data-mining/>
- [11] <https://www.tutorialride.com/data-mining/classification-in-data-mining.htm>
- [12] <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>
- [13] <https://scikit-learn.org/stable/>
- [14] D. Heckerman and J. S. Breese, "Causal independence for probability assessment and inference using Bayesian networks," *IEEE Trans. Syst. Man, Cybern. Part A Systems Humans.*, vol. 26, no. 6, pp. 826–831, 1996.
- [15] A. Genkin, D. Lewis, "Author Identification on the Large Scale", 2005
- [16] <https://www.edureka.co/blog/naive-bayes-tutorial/>
- [17] <https://www.geeksforgeeks.org/understanding-logistic-regression/>
- [18] [scikit-learn.org/stable/.../sklearn.ensemble.ExtraTreesClassifier.html](https://scikit-learn.org/stable/.../sklearn.ensemble.ExtraTreesClassifier.html)
- [19] <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesClassifier.html>
- [20] <https://blog.faradars.org/>
- [21] <http://www.bigdata.ir/1397/03/>
- [22] <https://howsam.org/>
- [23] <http://research-moghami.ir/1395/05/15/ensemble-learning/>
- [24] M. Brusco, D. Steinley, J. Cradit, "An exact algorithm for hierarchically well-formulated subsets in second-order polynomial regression", 2009
- [25] <https://dataio.ir/how-to-import-dataset-with-pandas-r1gsgvg9xx2>
- [26] <https://blog.faradars.org/python-excel-tutorial/>
- [27] <https://www.datopia.ir/>
- [28] K. Matsumoto, "Classification of Emoji Categories from Tweet Based on Deep Neural Networks," pp. 17–25, 2018

- [29] <http://www.statsoft.com/textbook/naive-bayes-classifier>
- [30] <https://towardsdatascience.com/logistic-regression-b0af09cdb8ad>
- [31] <https://www.quora.com/What-is-the-C-parameter-in-logistic-regression>
- [32] <https://www.toptal.com/machine-learning/ensemble-methods-machine-learning>
- [33] [https://xavierbourretsicotte.github.io/subset\\_selection.html](https://xavierbourretsicotte.github.io/subset_selection.html)
- [34] M.Rastegar, S.Hosseinzadeh, E.Bakhshi, “Application of Logistic Regression with Missclassified Variables in Diabetes Data”, 2018