

# طبقه‌بندی داده‌های نامتوازن با استفاده از ترکیب طبقه‌بندها و توصیف‌گرهای بردار پشتیبان

پایان‌نامه کارشناسی ارشد

عباس پیرمحمدی

استاد راهنما: دکتر علی امیری

استاد مشاور: دکتر محسن افشارچی

۲۴ اردیبهشت ۱۳۹۳



تقدیم به آنهایی که می‌خوانند بیشتر  
بدانند.

## چکیده

مسئله طبقه‌بندی داده‌های نامتوازن به‌عنوان یکی از چالش‌های اصلی در حوزه‌ی داده‌کاوی، مورد توجه بسیاری از محققان و پژوهش‌گران قرار گرفته است. در سال‌های اخیر تحقیقات ارزشمند زیادی برای حل مسئله طبقه‌بندی داده‌های نامتوازن انجام شده است. در بین این تحقیقات، رهیافت‌های مبتنی بر ترکیب طبقه‌بندها از موفقیت قابل توجهی برخوردار بوده‌اند.

علیرغم کارهای بسیار مؤثر انجام شده در ترکیب طبقه‌بندها هنوز برخی چالش‌ها از قبیل؛ عدم توجه به اهمیت نمونه‌ها در متوازن سازی، تعیین تعداد مناسب طبقه‌بندها و عدم بهینه‌سازی وزن طبقه‌بندها در ترکیب طبقه‌بندها به‌صورت حل نشده باقی مانده است. بنابراین در این پایان‌نامه تلاش شده است رهیافتی برای چالش‌های مطرح شده ارائه شود. در این راستا، ترکیبی از الگوریتم‌ها جهت ایجاد ساختار پیشنهادی ارائه شده است. سیستم پیشنهادی شامل ۳ بخش کلی است: در بخش اول توصیف‌گر بردار پشتیبان داده جهت نمونه‌برداری از دو کلاس اقلیت و اکثریت استفاده شده است. این روش برخلاف روش‌های پیش نمونه‌برداری و زیر نمونه‌برداری از داده‌هایی که توصیف بهتری از کل داده‌ها دارد، نمونه‌برداری می‌کند. در بخش دوم تعداد بهینه طبقه‌بندهای پایه در ترکیب طبقه‌بندها برای نمونه‌های منتخب تعیین می‌شود و در بخش پایانی با استفاده از الگوریتم ژنتیک عمل بهینه‌سازی وزن‌ها در ترکیب طبقه‌بند انجام می‌شود. روش پیشنهادی با تعدادی از الگوریتم‌های موجود در این حوزه مقایسه شده و نتایج نشان دهنده مطلوبیت ساختار پیشنهادی در مقایسه با برخی از الگوریتم‌ها از لحاظ دقت است.

**واژه‌های کلیدی:** شبیه‌سازی تبرید تدریجی، توصیف‌گر بردار پشتیبان داده، ترکیب طبقه‌بندها، الگوریتم

ژنتیک

# فهرست

چهار	چکیده	
۱	بیان مسئله	۱
۱	مقدمه	۱.۱
۲	بیان مسئله و چالش‌های پیش‌رو	۲.۱
۳	رهیافت پیشنهادی	۳.۱
۴	ساختار پایان‌نامه	۴.۱
۵	داده‌های نامتوازن	۲
۵	مقدمه	۱.۲
۶	مجموعه داده‌های نامتوازن	۲.۲
۸	روش‌های مواجهه با مسئله داده‌های نامتوازن	۳.۲
۹	روش‌های سطح داده	۱.۳.۲
۱۳	روش‌های سطح الگوریتم	۴.۲
۱۳	روش‌های حساس به هزینه	۵.۲
۱۴	روش‌های مبتنی بر ترکیب طبقه‌بندها	۶.۲
۱۵	الگوریتم Bagging	۱.۶.۲

۱۶	.....	الگوریتم Boosting	۲.۶.۲
۱۷	.....	الگوریتم AdaBoost	۳.۶.۲
۲۱	.....	الگوریتم RUSBoost	۴.۶.۲
۲۳	.....	ارزیابی کارایی در دامنه طبقه‌بندی داده‌های نامتوازن	۷.۲
۲۳	.....	ماتریس در هم‌ریختگی	۱.۷.۲
۲۵	.....	نمودار ROC	۲.۷.۲
۲۷	.....	خلاصه فصل	۸.۲
۲۸	.....	بررسی توصیف‌گر بردار پشتیبان داده و الگوریتم تکاملی	۳
۲۸	.....	مقدمه	۱.۳
۳۰	.....	توصیف‌گر بردار پشتیبان داده	۲.۳
۳۲	.....	تعمیم به سایر کرنل‌ها	۳.۳
۳۶	.....	تنظیم پارامترهای SVDD	۴.۳
۳۶	.....	اعتبار سنجی $k$ بخشی	۱.۴.۳
۳۷	.....	جستجوی شبکه‌ای	۵.۳
۳۸	.....	شبیه‌سازی تبرید تدریجی	۶.۳
۴۱	.....	الگوریتم ژنتیک	۷.۳
۴۲	.....	تولید راه‌حل‌های اولیه	۱.۷.۳
۴۲	.....	انتخاب تابع ارزیابی مناسب	۲.۷.۳
۴۲	.....	عملگر انتخاب	۳.۷.۳
۴۴	.....	عملگر تقاطع	۴.۷.۳
۴۶	.....	عملگر جهش	۵.۷.۳
۴۶	.....	خلاصه فصل	۸.۳

۴۷	رهیافت پیشنهادی در طبقه‌بندی داده‌های نامتوازن	۴
۴۷	مقدمه	۱.۴
۴۸	چالش‌ها و انگیزه‌ها	۲.۴
۴۹	چهار چوب کلی رهیافت پیشنهادی	۳.۴
۵۰	نمونه‌برداری با استفاده از SVDD	۴.۴
۵۳	تعیین تعداد طبقه‌بندی‌های پایه در ترکیب طبقه‌بندها	۵.۴
۵۶	طبقه‌بندی با استفاده از ترکیب طبقه‌بندها	۶.۴
۵۷	بهینه‌سازی وزن‌ها با استفاده از الگوریتم ژنتیک	۷.۴
۶۱	خلاصه فصل	۸.۴
۶۲	بررسی نتایج	۵
۶۲	مقدمه	۱.۵
۶۳	مجموعه داده‌های مورد استفاده	۲.۵
۶۴	تعداد نمونه‌ها و طبقه‌بندها در هر بخش از مجموعه داده‌ها	۳.۵
۶۶	ارزیابی ساختار پیشنهادی	۴.۵
۶۸	مقایسه روش پیشنهادی با سایر الگوریتم‌ها	۵.۵
۷۲	خلاصه فصل	۶.۵
۷۳	نتیجه‌گیری و کارهای آتی	۶
۷۳	مقدمه	۱.۶
۷۳	نتیجه‌گیری	۲.۶
۷۴	پیشنهاداتی برای مطالعات آتی	۳.۶
۸۲	واژه‌نامه فارسی به انگلیسی	

# فهرست تصاویر

۸	در هم تداخل کلاسی . . . . .	۱.۲
۸	نمودار جدا کننده‌های کوچک . . . . .	۲.۲
۱۰	حذف تصادفی نمونه‌های کلاس اکثریت . . . . .	۳.۲
۱۰	جایگزینی تصادفی نمونه‌های کلاس اقلیت (نمودار سمت راستی اصلاح شده است). . . . .	۴.۲
۵.۲	نمودار (a) مثالی از ۶ همسایه نزدیک نمونه $x_i$ و نمودار (b) ایجاد داده براساس فاصله . . . . .	
۱۱	اقلیدسی . . . . .	
۱۹	نمودار $\alpha$ به‌عنوان تابعی از خطای آموزشی . . . . .	۶.۲
۲۶	منحنی‌های ROC برای دو طبقه‌بند مختلف . . . . .	۷.۲
۳۲	نمایش دو بعدی SVDD . . . . .	۱.۳
۳۳	توصیف داده‌ها با ۳ مقدار مختلف . . . . .	۲.۳
۳۴	توصیف داده‌ها با ۳ مقدار مختل . . . . .	۳.۳
۳۸	جستجوی شبکه‌ای برای تعیین بهینه . . . . .	۴.۳
۴۳	نمونه‌ای از انتخاب به روش چرخه گردان . . . . .	۵.۳
۴۴	تقاطع تک نقطه‌ای برای تولید فرزند در الگوریتم ژنتیک . . . . .	۶.۳
۴۵	تقاطع دو نقطه‌ای برای تولید فرزند در الگوریتم ژنتیک . . . . .	۷.۳
۴۵	تقاطع یکنواخت برای تولید فرزند در الگوریتم ژنتیک . . . . .	۸.۳



۴۹	فلوچارت چارچوب کَلّشی رهیافت پیشنهادی	۱.۴
۵۲	نمایش دو بعدی مجموعه داده ۵۰ عضوی	۲.۴
۵۲	نمایش مرز بسته روی مجموعه داده بعد استفاده از روش SVDD	۳.۴
۵۳	نمونه‌های انتخاب شده به‌وسیله SVDD	۴.۴
۵۸	نمایش کروموزوم	۵.۴
۶۰	عملگر تقاطع	۶.۴
۶۰	استفاده از عملگر جهش روی کروموزوم منتخب	۷.۴
	نمودار ستونی میانگین دقت هر یک از معیارهای ارزیابی در روش پیشنهادی و	۱.۵
۶۸	الگوریتم RUSBoost	۲.۵
۶۸	نمودار ستونی میانگین زمان اجرایی روش پیشنهادی و الگوریتم RUSBoost	۳.۵
	نمودار ستونی میانگین دقت هر یک از معیارهای ارزیابی در روش پیشنهادی و	۳.۵
۷۱	الگوریتم‌های Easyensemble و SMOTEBOOST	۴.۵
	نمودار ستونی میانگین زمان اجرایی روش پیشنهادی و الگوریتم‌های Easyensem-	۴.۵
۷۱	SMOTEBOOST و ble	۷.۱

# فصل اول

## بیان مسئله

### ۱.۱ مقدمه

امروزه پردازش داده‌های با توزیع احتمال نامتوازن<sup>۱</sup>، به دلیل گستردگی آنها در بسیاری از مسائل دنیای واقعی، از قبیل تشخیص تقلب [۱، ۲]، تشخیص ناهنجاری [۳]، تشخیص پزشکی [۴] و تشخیص نشت نفت [۵]، مورد توجه بسیاری از محققان و پژوهش‌گران قرار گرفته است. در داده‌های نامتوازن، معمولاً تعداد نمونه‌های یکی از کلاس‌ها خیلی بیشتر از نمونه‌های کلاس دیگر است. کلاس با تعداد داده‌های بیشتر را کلاس اکثریت<sup>۲</sup> و کلاس با داده‌های کمتر را کلاس اقلیت<sup>۳</sup> می‌گوییم. در این مسائل معمولاً نسبت داده‌های کلاس اقلیت به اکثریت، اغلب ۱ به ۱۰۰ و یا بیشتر است [۶]. در داده‌های نامتوازن چالش اصلی شناسایی صحیح نمونه‌های کلاس اقلیت است. به‌عنوان مثال در حوزه پزشکی، تعداد نمونه‌های مثبت از یک بیماری در مقابل تعداد نمونه‌های منفی بسیار کمتر است. در حالی که

---

<sup>۱</sup> Imbalanced

<sup>۲</sup> Majority

<sup>۳</sup> Minority

اهمیت شناسایی نمونه‌های مربوط به دسته مثبت، بسیار زیاد است. توزیع کلاس، نسبت نمونه‌های متعلق به هر کلاس در یک مجموعه داده، نقش کلیدی در طراحی طبقه‌بند ایفا می‌کند، در الگوریتم‌های استاندارد طبقه‌بندی داده‌ها، توزیع کلاس‌ها متوازن در نظر گرفته می‌شود. از این رو در صورت استفاده از این الگوریتم‌ها در طبقه‌بندی داده‌های نامتوازن، نمی‌توان به نتایج قابل قبولی دست یافت؛ زیرا در این الگوریتم‌ها طبقه‌بند به سمت نمونه‌های آموزشی کلاس بزرگ‌تر متمایل می‌شود که این موضوع سبب افزایش تعداد خطاها در شناسایی نمونه‌های مثبت می‌شود. بنابراین در سال‌های اخیر مسئله‌ی طبقه‌بندی داده‌های نامتوازن به عنوان یکی از چالش‌های اصلی در حوزه‌ی داده‌کاوی، توجه بسیاری از محققین را به خود جلب نموده است [۷]. در این پایان‌نامه، مسئله طبقه‌بندی داده‌های نامتوازن مورد توجه قرار گرفته است. بیان دقیق مسئله و چالش‌های اصلی پایان‌نامه، طرح کلی رهیافت پیشنهادی و ساختار کلی فصل‌های آتی در ادامه این فصل آمده است.

## ۲.۱ بیان مسئله و چالش‌های پیش‌رو

مسئله اصلی این پایان‌نامه طراحی یک طبقه‌بند برای داده‌های نامتوازن است. همان‌طور که در بخش قبل اشاره شد، استفاده از طبقه‌بندهای استاندارد نمی‌تواند برای حل این مسئله مؤثر باشد. در سال‌های اخیر تحقیقات ارزشمند زیادی برای حل این مسئله ارائه شده است. در بین روش‌های ارائه شده، رهیافت‌های مبتنی بر ترکیب طبقه‌بندها از موفقیت قابل توجهی برخوردار بوده‌اند. از جمله این الگوریتم‌ها می‌توان *AdaBoost* [۸]، *SMOTEBoost* [۹]، *RUSBoost* [۱۰] و غیره اشاره کرد. علیرغم کارهای بسیار موثر انجام شده در ترکیب طبقه‌بندها هنوز چالش‌های زیر به صورت حل نشده باقی مانده است:

۱. الگوریتم‌های ترکیب طبقه‌بندها عمدتاً از روش‌های زیر نمونه‌برداری<sup>۱</sup> و بیش نمونه‌برداری<sup>۲</sup> برای متوازن نمودن داده‌ها استفاده می‌کنند. این کار با افزایش تصادفی تعداد نمونه‌های اقلیت یا کاهش تعداد نمونه‌های اکثریت صورت می‌گیرد. بدیهی است افزایش بی‌رویه و کورکورانه نمونه‌ها می‌تواند منجر به تولید و گسترش برخی نمونه‌ها که در توصیف نقش موثری را ندارند، گردد. بدیهی است داده‌هایی که در مرز داده‌ها قرار می‌گیرد نسبت به داده‌هایی که در مرکز داده‌ها قرار دارند از اهمیت بالایی برخوردار هستند. عدم توجه به اهمیت نمونه‌ها در متوازن‌سازی داده‌ها یک نقطه ضعف در کارهای گذشته بوده و از این رو به‌عنوان یکی از چالش‌های اصلی در کار ما قرار گرفته است.

۲. در الگوریتم‌های ترکیب طبقه‌بند ارائه شده در کارهای قبل، تعداد طبقه‌بندهای پایه و اهمیت آنها مورد توجه قرار نگرفته است. معمولاً بیشتر کارهای انجام شده با تعداد طبقه‌بند پایه ثابت به حل مسئله پرداخته‌اند و هیچ توجهی به انتخاب بهینه تعداد آنها نشده است.

۳. در بیشتر کارهای گذشته، رأی‌گیری طبقه‌بندها معمولاً به‌صورت وزن‌دار و بر اساس قدرت تصمیم آن روی کل داده‌ها تعیین می‌شد. تنظیم این وزن‌ها جهت افزایش قدرت تصمیم کل طبقه‌بند مرکب به‌دست آمده تاکنون مورد توجه قرار نگرفته است.

## ۳.۱ رهیافت پیشنهادی

در این پایان‌نامه تلاش شده است رهیافتی برای چالش‌های مطرح شده در بخش قبل ارائه شود. در این راستا، ترکیبی از الگوریتم‌ها جهت ایجاد ساختار پیشنهادی ارائه شده است. سیستم پیشنهادی شامل ۳ بخش کلی است: در بخش اول جهت نمونه‌برداری از هر دو کلاس اقلیت و

---

<sup>۱</sup> Under sampling

<sup>۲</sup> Over sampling

اکثریت از روش SVDD استفاده شده است. این روش برخلاف سایر روش‌های نمونه‌برداری که از کل مجموعه داده‌ها نمونه‌برداری می‌کنند، تنها از نمونه‌هایی که در مرز داده‌ها قرار دارند نمونه‌برداری می‌کند. در بخش دوم روش پیشنهاد شده توسط هرناندز-لباتو<sup>۱</sup> و همکارانش روی نمونه‌های منتخب اعمال می‌شود که منجر به تعیین حداقل تعداد طبقه‌بندی‌های پایه در ترکیب طبقه‌بندی می‌شود. در بخش پایانی با استفاده از الگوریتم ژنتیک عمل بهینه‌سازی وزن‌ها در ترکیب طبقه‌بندی انجام می‌شود.

## ۴.۱ ساختار پایان‌نامه

این پایان‌نامه شامل شش فصل است. در فصل اول، مسئله اصلی و چالش‌های پیش‌رو بیان شد. در فصل دوم، داده‌های نامتوازن و الگوریتم‌های طبقه‌بندی داده‌های نامتوازن مطالعه شده است. مفهوم توصیف‌گر بردار پشتیبان داده و مفهوم شبیه‌سازی تبرید تدریجی<sup>۲</sup> و سایر مفاهیم مورد نیاز در فصل سوم آمده است. در فصل چهارم رهیافت پیشنهادی برای طبقه‌بندی داده‌های نامتوازن بیان شده است. در فصل پنجم نتایج آزمایشات آمده است و فصل ششم شامل نتیجه‌گیری و کارهای آتی است.

---

<sup>۱</sup> Hernandez-Lobato

<sup>۲</sup> Simulated Annealing (SA)

## فصل دوم

### داده‌های نامتوازن

#### ۱.۲ مقدمه

در سال‌های اخیر مشکل نامتوازن بودن کلاس‌ها مورد توجه محققان در زمینه داده‌کاوی قرار گرفته است. در موارد متعددی کلاسی که از نقطه نظر دامنه کاربردی، اهمیت زیادی دارد (کلاس اصلی) شامل تعداد نمونه‌های کمتری نسبت به کلاس اکثریت است.

طبقه‌بندی داده‌های نامتوازن یک مسئله رایج در بسیاری از دامنه‌ها، از قبیل تشخیص معاملات تقلبی [۲]، تشخیص پزشکی [۴] و تشخیص نشت نفت [۵] است. در همه سناریوها وقتی کلاس اکثریت شامل ۹۸-۹۹ درصد از کل نمونه‌ها باشد، یک طبقه‌بند ساده در مرحله یادگیری نمونه‌های کلاس اقلیت را نادیده می‌گیرد و در مرحله آزمایش هر چیزی را از کلاس اکثریت برچسب‌دهی می‌کند، بنابراین این طبقه‌بند می‌تواند به بالاترین دقت ممکن برسد. معیارهایی از قبیل: آنالیز منحنی مشخصه عملکرد

دریافت کننده<sup>۱</sup>، دقت<sup>۲</sup>، فراخوانی<sup>۳</sup>، اندازه‌گیری-F<sup>۴</sup> و میانگین-G<sup>۵</sup> برای درک کارایی الگوریتم یادگیری روی مجموعه داده‌های نامتوازن استفاده می‌شود. با توجه به اهمیت این نوع داده‌ها، امروزه الگوریتم‌های متعددی برای طبقه‌بندی داده‌های نامتوازن بوجود آمده است. بدین منظور در ادامه فصل، نگاهی بر داده‌های نامتوازن خواهیم داشت. سپس چگونگی مدیریت این نوع داده‌ها و روش‌های مبتنی بر ترکیب طبقه‌بندها را بررسی خواهیم نمود و در آخر معیارهای اندازه‌گیری کارایی در دامنه نامتوازن را بیان خواهیم کرد.

## ۲.۲ مجموعه داده‌های نامتوازن

مشکل توزیع نامتوازن داده در کلاس‌ها در مواقعی مطرح می‌شود که نمونه‌های موجود در برخی کلاس‌ها از سایر کلاس‌ها بیشتر باشد، این مسئله بخصوص در کاربردهای دوکلاسی مطرح است که يك کلاس از نمونه‌های زیادی نسبت به کلاس دیگر برخوردار است. در تعداد زیادی از حوزه‌های داده‌کاوی، مشکل عدم توازن کلاس‌ها یکی از مسائل مهم می‌باشد.

الگوریتم‌های معمول داده‌کاوی در مواجهه با مشکل کلاس‌های نامتوازن معمولاً عمل کرد ضعیفی دارند. زیرا با صرف‌نظر کردن از نمونه‌های داده در کلاس اقلیت، اقدام به افزایش دقت کلی می‌کنند. در صورتی که نمونه‌های موجود در کلاس اقلیت در بسیاری از کاربردها، نسبت به نمونه‌های کلاس دیگر از اهمیت بیشتری برخوردارند. برای مثال، در يك مسئله تشخیص پزشکی با وجود اینکه نمونه‌های بیماری در مقایسه با نمونه‌های معمول در اقلیت قرار دارند؛ ولی هدف از طبقه‌بندی در واقع یافتن

---

<sup>۱</sup> Receiver Operating Characteristic)(ROC)

<sup>۲</sup> Precision

<sup>۳</sup> Recall

<sup>۴</sup> F-measure

<sup>۵</sup> G-mean

نمونه‌های بیماری است. در مجموعه داده‌ای با نرخ عدم توازن ۱:۱۰۰ (یعنی با فرض وجود یک نمونه کلاس اقلیت به ازای ۱۰۰ نمونه از کلاس اکثریت) روش‌های طبقه‌بندی استاندارد با نادیده گرفتن تمامی نمونه‌های اقلیت به یک طبقه‌بندی با دقت بالای ۹۹٪ دست می‌یابد، در صورتی که در طبقه‌بندی داده‌های نامتوازن با توجه به اهمیت بالای شناخت نمونه‌های کلاس اقلیت، طبقه‌بندی صحیح آنها خیلی مهم است [۶، ۹، ۱۱].

در اغلب الگوریتم‌های طبقه‌بندی، تمایل در جهت کلاسی است که بیشترین تعداد نمونه‌ها را دارند، از این رو توانایی کمی در پیش‌گویی صحیح داده‌های کلاس اقلیت از خود نشان می‌دهند. در چنین الگوریتم‌هایی، نمونه‌های کلاس اقلیت نسبت به نمونه‌های کلاس اکثریت نادرست طبقه‌بندی می‌شوند [۱۲، ۱۳]. مسائل مرتبط در شکل‌گیری داده‌های نامتوازن عبارتند از [۶]:

۱. کم بودن تعداد نمونه‌ها: به‌طور کلی در مجموعه داده‌های نامتوازن، تعداد نمونه‌های اقلیت به اندازه کافی نیست.

۲. تداخل یا کلاس تفکیک‌پذیر<sup>۱</sup>: این حالت زمانی رخ می‌دهد که قوانین متمایز کردن کلاس‌ها سخت باشد. در نتیجه، قوانین کلی‌تر سبب می‌شوند داده‌های کلاس اقلیت نادرست طبقه‌بندی شوند [۱۱]. اگر در هم تداخلی بین کلاس‌ها نباشد، آنگاه هر طبقه‌بند ساده با توجه به توزیع کلاس می‌تواند به‌طور مناسب یاد بگیرد. شکل ۱.۲ بیان‌گر این حالت می‌باشد.

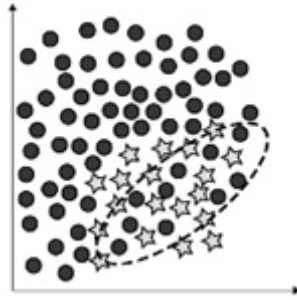
۳. جدا کننده‌های کوچک<sup>۲</sup>: جدا کننده‌های کوچک، جدا کننده‌هایی هستند که در طبقه‌بندی یادگیرنده تعدادی از نمونه‌های آموزشی را پوشش می‌دهند [۱۴]. وجود جدا کننده‌های کوچک در یک مجموعه داده، زمانی روی می‌دهد که مفهوم ارائه شده توسط کلاس اقلیت متشکل از زیرمفهوم‌ها<sup>۳</sup> باشد [۱۵]. بنابراین، در اغلب مسائل شاهد حضور این زیر مفهوم‌ها هستیم.

<sup>۱</sup> Overlapping or class separability

<sup>۲</sup> Small disjuncts

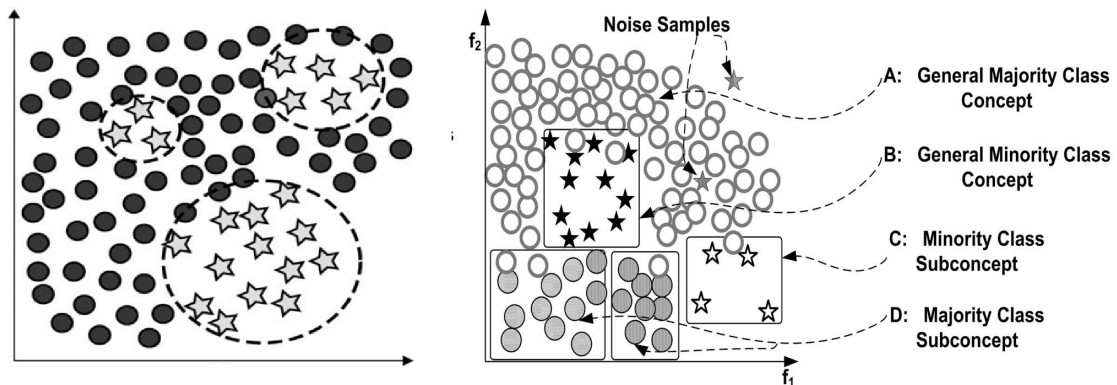
<sup>۳</sup> Subconcepts





شکل ۱.۲: در هم تداخل کلاسی [۶]

وجود زیر مفهوم‌ها، پیچدگی مسئله را افزایش می‌دهد، زیرا تعداد نمونه‌ها معمولاً میان آنها به‌طور متعادل نمی‌باشد. نمودار سمت چپ شکل ۲.۲ از مشکل سه زیر مفهوم از کلاس اقلیت است.



شکل ۲.۲: نمودار جدا کننده‌های کوچک [۶، ۱۳]

## ۳.۲ روش‌های مواجهه با مسئله داده‌های نامتوازن

با توجه به اهمیت مسئله مجموعه داده‌های نامتوازن در برنامه‌های کاربردی، تعداد زیادی از روش‌ها در برخورد با این مسئله توسعه داده شده‌اند. در این ادبیات، تکنیک‌های متنوعی برای حل مسئله در ارتباط

با کلاس نامتوازن پیشنهاد شده است، که در سه گروه: روش سطح داده<sup>۱</sup>، روش سطح الگوریتم<sup>۲</sup> و روش حساس به هزینه<sup>۳</sup> تقسیم می شوند [۶]. در ادامه به توصیف هر کدام از این روش ها خواهیم پرداخت.

## ۱.۳.۲ روش های سطح داده

در این روش با اضافه کردن مرحله پیش پردازش<sup>۴</sup> قبل از طبقه بندی، موجب متوازن شدن مجموعه داده های نامتوازن می شود. در روش پیش پردازش داده از روش نمونه برداری استفاده می شود. نمونه برداری به دو شکل است، زیر نمونه برداری از کلاس اکثریت، بیش نمونه برداری از کلاس اقلیت، یا ترکیبی از هر دو روش است.

۱. زیر نمونه برداری: یکی از معروفترین روش های نمونه برداری به طور تصادفی است. در زیر نمونه برداری به طور تصادفی نمونه هایی از کلاس اکثریت حذف می شود، تا زمانی که کلاس اقلیت درصدی از کلاس اکثریت شود؛ به این ترتیب توازن در مجموعه آموزشی برقرار می شود. شکل ۳.۲ بیانگر این روش می باشد. از معایب این روش، از دست دادن اطلاعات با ارزش است و موجب زیر برازش<sup>۵</sup> داده های کلاس اکثریت می شود.

۲. بیش نمونه برداری: این روش برای به تعادل رساندن توزیع کلاس، از روش جایگزینی نمونه های کلاس اقلیت استفاده می کند، که نیازی به اطلاعات اضافی ندارد و از داده های موجود دوباره استفاده می کند تا توازن در مجموعه آموزشی برقرار شود. شکل ۴.۲ نمایش روش بیش نمونه برداری

---

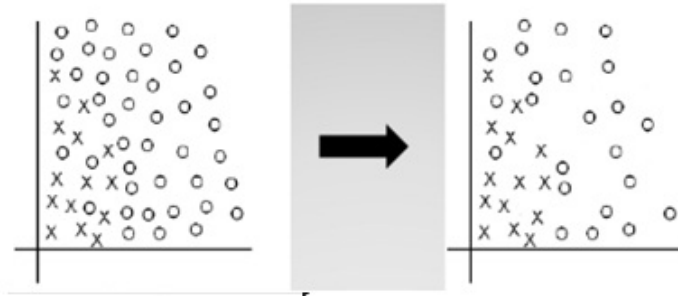
<sup>۱</sup> Data level

<sup>۲</sup> Algorithm level

<sup>۳</sup> Cost-sensitive

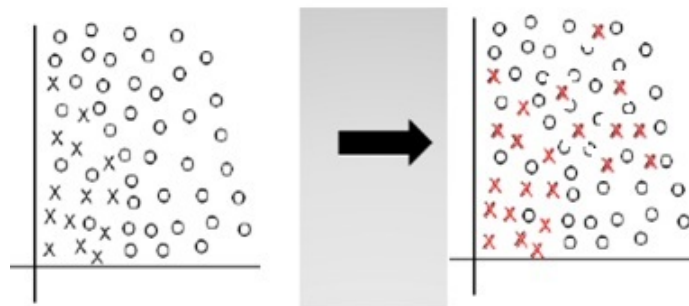
<sup>۴</sup> Preprocessing

<sup>۵</sup> Under fitting



شکل ۳.۲: حذف تصادفی نمونه‌های کلاس اکثریت [۱۶]

است. از معایب این روش باعث بیش برآزش<sup>۱</sup> مدل آموزشی می‌شود و همچنین روش بیش نمونه‌برداری اندازه مجموعه آموزشی را افزایش می‌دهد که این امر منجر به افزایش زمان یادگیری طبقه‌بند می‌شود.



شکل ۴.۲: جایگزینی تصادفی نمونه‌های کلاس اقلیت (نمودار سمت راستی اصلاح شده است.) [۱۷]

۳. روش بیش نمونه‌برداری مصنوعی کلاس اقلیت (SMOTE)<sup>۲</sup> الگوریتم SMOTE داده‌های مصنوعی را بر اساس شباهت فضای ویژگی بین نمونه‌های کلاس اقلیت ایجاد می‌کند. برای زیر مجموعه  $S_{min} \in S$ ، تا از نزدیکترین همسایه نمونه  $x_i \in S_{min}$  را انتخاب می‌کنیم که  $k$  یک عدد صحیح می‌باشد. نحوه انتخاب  $k$  تا از همسایه‌های نزدیک نمونه مورد نظر براساس

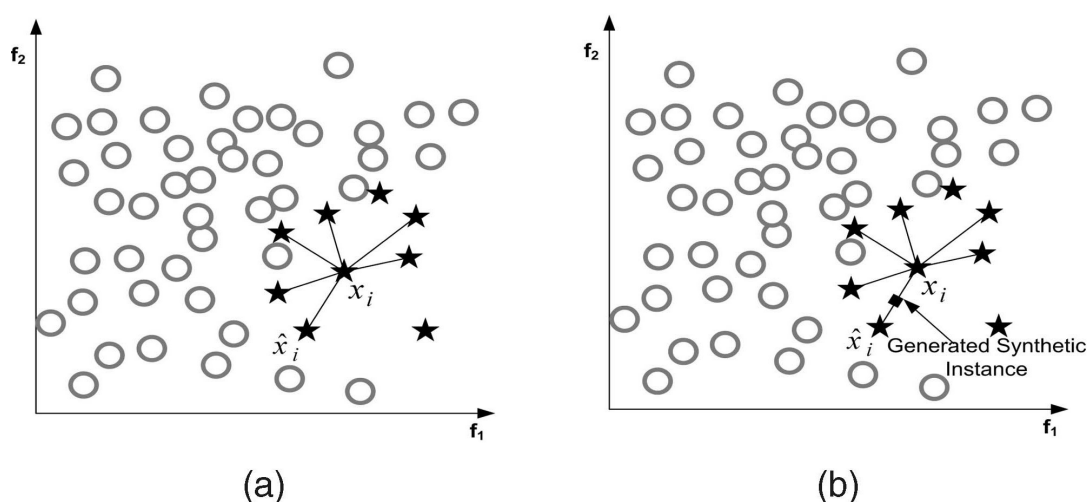
<sup>۱</sup> Overfitting

<sup>۲</sup> Synthetic Minority Oversampling Technique (SMOTE)

فاصله اقلیدسی در فضای  $n$  بعدی می‌باشد. برای ایجاد نمونه مصنوعی، یکی از  $k$  همسایه‌های نزدیک نمونه  $x_i$  به‌طور تصادفی انتخاب می‌کنیم، سپس با ضرب اختلاف این دو نمونه در یک عدد تصادفی بین  $[0, 1]$  و اضافه کردن نتیجه حاصل شده به نمونه  $x_i$ ، نمونه جدید به‌دست می‌آید.

$$x_{new} = x_i + (\hat{x}_i - x_i) \times \delta$$

که  $\hat{x}_i$  یکی از  $k$  همسایه نزدیک نمونه  $x_i$  است و  $\delta \in [0, 1]$  یک عدد تصادفی است. نمونه جدید در بخشی از خط متصل کننده بین دو نمونه قرار می‌گیرد [۱۷].



شکل ۵.۲: (a) نمودار ۶ همسایه نزدیک نمونه  $x_i$  و نمودار (b) ایجاد داده براساس فاصله اقلیدسی [۱۳]

۴. اصلاح شده روش بیش نمونه‌برداری مصنوعی کلاس اقلیت (MSMOTE)<sup>۱</sup> نسخه اصلاح شده روش SMOTE است. این الگوریتم با محاسبه مرکز نمونه‌های کلاس اقلیت و تعیین فاصله اقلیدسی هر یک از نمونه‌ها نسبت به مرکز نمونه‌ها، آن‌ها را در سه گروه، امن، مرزی و نمونه‌های نویزدار نهان طبقه‌بندی می‌کند. روش‌های تولید نمونه‌های مصنوعی به‌وسیله MSMOTE

<sup>۱</sup> Modified Synthetic Minority Oversampling Technique (MSMOTE)