

بررسی روش های مختلف بیش نمونه گیری در رده بندی داده های نامتوازن

هادی مهدوی نیا

دانشجوی دکتری مهندسی نرم افزار دانشگاه آزاد اسلامی واحد خوراسگان

Email:hadi.mahdavinia1365@gmail.com

چکیده

در استخراج دانش از داده ها، از نظر توازن، به داده های متوازن و نامتوازن برخورد می شود. برای رده بندی داده های نامتوازن، اگر از روش های معمول رده بندی مانند ماشین های بردار و غیره استفاده شود، مشکلاتی مانند: مدل جانبدارانه، رده بندی اشتباه رده ی اقلیت، صرفه نظر کردن از داده های رده ی اقلیت و بیش پوشش به وجود خواهد آمد. برای اجتناب از مشکلات، باید از روش های خاصی برای رده بندی داده های نامتوازن استفاده شود. برای رده بندی داده های نامتوازن، از یکی از روش های سطح داده، روش های سطح الگوریتم، روش های یادگیری حساس به هزینه و روش های یادگیری ترکیبی استفاده می شود که استفاده از هر یک از روش ها به اهداف استخراج دانش بستگی دارد. در روش های سطح داده، یکی از روش های مرسوم، افزایش نمونه های رده ی اقلیت است که به بیش نمونه گیری معروف است. در این تحقیق انواع روش های بیش نمونه گیری در داده کاوی داده های نامتوازن که یکی از زیرمجموعه های رویکرد سطح داده است، بررسی می شوند.

کلمات کلیدی: داده های نامتوازن، رده اکثریت، رده اقلیت، بیش نمونه گیری

۱-مقدمه

یکی از مباحث معمول و رایج یادگیری ماشین^۱، رده بندی^۲ مجموعه داده های نامتوازن^۳ است [۱]. زمانی که یک برچسب یا گروه از مؤلفه ی هدف در مقابل دیگر برچسب ها، دارای نمونه های کمتری باشد، مجموعه داده مورد نظر، یک مجموعه داده ی نامتوازن است [۲]. رده ای^۴ که نسبت به رده های دیگر مؤلفه ی هدف از میزان کمتری نمونه برخوردار است، رده ی اقلیت^۵ و مابقی نیز رده ی اکثریت^۶ نامیده می شوند [۳]. در استخراج دانش از مجموعه داده های نامتوازن، داده های رده ی اقلیت، از اهمیت ویژه ای برخوردار است، بنابراین سعی می شود نمونه های رده ی اقلیت به صورت صحیح افزایش یابد [۴]. داده های نامتوازن در بسیاری از زمینه ها، از جمله تشخیص خطا و عیب، تشخیص پزشکی، تشخیص نفوذ، تشخیص متون، تشخیص تقلب های بانکی، رده بندی جریان داده ها و از این قبیل کاربرد دارند [۵]. خطاهای غیر منتظره ای مانند صرفه نظر کردن از رده ی اقلیت، استخراج مدل جانبدارانه^۷ و محاسبه ی خطاهای رده های مختلف برای کسب هزینه های مختلف، از مشکلات و چالش های استخراج دانش از درون داده های نامتوازن است [۷].

¹ Machine Learning

² Classification

³ Imbalanced datasets

⁴ Class

⁵ Minority Class

⁶ Majority Class

⁷ Bias Model

سومین همایش ملی مهندسی کامپیوتر، داده کاوی و داده های حجیم

برای حل مشکلات ذکر شده، از دسته‌بندی‌های مختلف استفاده می‌شود که هر کدام مزایا و معایب خود را دارند. بارتوز کوارشکی، برای حل مشکلات داده‌های نامتوازن، سه دسته‌بندی از راه کارهای موجود ارائه داده است که شامل: رویکرد سطح داده^۱، یادگیری سطح الگوریتم^۲ و یادگیری ترکیبی^۳ است [۸]. اما بهترین دسته‌بندی توسط ژن و همکاران ارائه شده است که چهار دسته راه حل برای رده‌بندی داده‌های نامتوازن ارائه داده‌اند که عبارتند از رویکرد سطح داده، رویکرد سطح الگوریتم، یادگیری حساس به هزینه و یادگیری ترکیبی [۹].

در این تحقیق، در فصل دوم، خلاصه‌ای از انواع روش‌های رده‌بندی داده‌های نامتوازن تشریح می‌گردد؛ در فصل سوم انواع روش‌های بیش‌نمونه‌گیری معرفی می‌گردد و در فصل چهارم نتیجه‌گیری و نقاط ضعف و قوت هر یک از روش‌ها به صورت اجمالی بیان می‌گردد.

۲-۲- دسته‌بندی روش‌های رده‌بندی داده‌های نامتوازن

برای رده‌بندی داده‌های نامتوازن، می‌توان از روش‌های مختلفی استفاده کرد که هر کدام مزایا و معایب خود را دارند. در رده‌بندی داده‌های نامتوازن، روش‌های اکتشافی در چهار دسته قرار می‌گیرند. حتی می‌توان هر روش در یک رویکرد را با روشی دیگر در رویکرد دیگر ترکیب کرد تا به کارایی مناسب و هدف مدنظر در رده‌بندی مجموعه داده‌ی نامتوازن رسید.

۲-۱- روش‌های سطح الگوریتم

در این روش‌ها، ساختار پایه‌ای الگوریتم‌ها تغییر پیدا می‌کند تا وابستگی و جانبداری مدل استخراجی از داده‌های نامتوازن، کاهش پیدا کند. یکی از تکنیک‌های مرسوم در این روش‌ها استفاده و اختصاص وزن به رده‌های پراهمیت است [۱۰].

۲-۲- روش‌های سطح داده

با دیدگاه پیش‌پردازش داده‌ها، این روش، سعی در ایجاد توازن در رده‌های مؤلفه‌ی هدف دارد که این روش‌ها با تغییر در تعداد نمونه‌های رده‌های مختلف به این هدف می‌رسد [۶]. به عبارتی توزیع داده‌ها در مرحله یادگیری، باعث حل مشکل استخراج دانش و به خصوص رده‌بندی داده‌های نامتوازن می‌شود [۱۱].

از زیرروش‌های این رویکرد که بسیار استفاده می‌شوند، می‌توان به روش بیش‌نمونه‌گیری و کم‌نمونه‌گیری اشاره کرد [۱۲].

۲-۳- روش‌های حساس به هزینه

در این رویکرد، هزینه‌های خطاهای هر رده یکسان نیست؛ بنابراین برای افزایش دقت و صحت، به نمونه‌های هر رده، وزن اختصاص داده می‌شود. در واقع در این رویکرد، برای رده‌های اقلیتی که به صورت اشتباهی رده‌بندی می‌شوند، وزن بیشتری نسبت به رده‌ی اکثریت در نظر گرفته می‌شود؛ بنابراین در فرآیند یادگیری، هدف، کاهش خطاهای وزنی، به جای افزایش نرخ صحت است. به عبارتی هدف اصلی این رویکرد کاهش هزینه رده‌های اشتباه رده‌بندی شده به جای کاهش خطای رده‌های اشتباه رده‌بندی شده است [۹] [۶].

از روش‌های این رویکرد می‌توان به درخت‌های تصمیم حساس به هزینه، شبکه‌های عصبی حساس به هزینه، BEE-Miner و MEPAR اشاره کرد [۶].

۲-۴- روش‌های ترکیبی

¹ Data level approach

² Algorithm level approach

³ Ensemble learning

در این روش، تکنیک‌های مختلف از سه دسته بندی فوق، با یکدیگر ترکیب شده و روش ترکیبی را برای ارائه یک مدل در رده‌بندی داده‌های نامتوازن مورد استفاده قرار می‌دهد. در این نوع از روش‌های داده‌کاوی نامتوازن، دقت بالایی از مدل ارائه می‌شود [۶][۱۴].

از روش‌های مرسوم این رویکرد می‌توان به AdaBoost و Boosting اشاره کرد [۱۵].

۳- روش‌های بیش‌نمونه‌گیری^۱

در فصل ۲، انواع روش‌های استخراج دانش از داده‌های نامتوازن در رده‌بندی داده‌های نامتوازن بیان گردید. یکی از روش‌های ذکر شده، روش سطح داده است که دارای زیر رویکردهایی همچون: کم‌نمونه‌گیری^۲، بیش‌نمونه‌گیری و روش‌های تلفیقی نمونه‌گیری است.

از آنجائی که رده‌ی اقلیت یک رده‌ی مهم است و داده‌های موجود در آن نیز از اهمیت ویژه‌ای برخوردار است، روش‌های ارائه شده در بیش‌نمونه‌گیری از حساسیت بیشتری برخوردار است. بنابراین در این فصل، انواع روش‌های موجود در زیررویکرد بیش‌نمونه‌گیری بررسی می‌گردد.

۱-۳- روش بیش‌نمونه‌گیری اقلیت مصنوعی^۳

معروف‌ترین روش افزایش رده‌ی اقلیت، این روش است. SMOTE جایگزین روش تولید رده‌ی اقلیت مصنوعی با جایگزینی است. اولین بار این روش از تشخیص موفقیت‌آمیز دست نوشته‌ها الهام گرفت. روش افزایش مصنوعی رده‌ی اقلیت به جای عمل در فضای داده‌ای در فضای مؤلفه‌ای یا صفت خاصه عمل می‌کند که هر سطر به عنوان نقطه مرکزیت و سپس تعداد k سطر به عنوان نزدیکترین همسایه نقطه مرکزیت به عنوان نمونه مصنوعی رده‌ی اقلیت تولید می‌شود [۱۶].

یکی از ورودی‌های الگوریتم، عدد T است که بیانگر تعداد نمونه‌های رده‌ی اقلیت است؛ عدد N نیز مقدار درصد تولید نمونه تصادفی از رده‌ی اقلیت است و عدد k نیز تعداد نزدیک‌ترین همسایه به نقطه مرکزیت است. آرایه‌ی Synthetic نیز تمامی نمونه‌های مصنوعی تولید شده را در خود نگه می‌دارد و مقدار بازگشتی الگوریتم SMOTE آرایه‌ای از نمونه‌های مصنوعی تولید شده از رده‌ی اقلیت است [۱۶].

در سال‌های اخیر تلاش‌هایی برای بهبود این روش صورت گرفته است و روش‌هایی مانند Safe-Level-SMOTE و همچنین Borderline-SMOTE ارائه شده است [۱۷][۱۸].

۳-۲- بیش‌نمونه‌گیری خوشه‌بندی^۴

در این روش، الگوریتم مورد نظر، ابتدا تعداد سرخوشه‌های مشخصی را از بین داده‌های رده‌ی اقلیت پیدا می‌کند و سپس k نزدیک‌ترین همسایه به هر کدام از سرخوشه‌ها را محاسبه کرده و ارائه می‌دهد. بدین ترتیب داده‌های مصنوعی نزدیک به رده‌ی اقلیت داده‌ها را پیدا کرده و به مجموعه داده‌ی رده‌ی اقلیت اضافه می‌کند. قابل ذکر است روش k -means یک روش خوشه‌بندی است که با ارزش محاسباتی کم برای داده‌هایی با بُعد بالا می‌تواند مورد استفاده قرار گیرد [۱۹]. این روش می‌تواند با انواع روش‌های خوشه‌بندی مانند k -Medoids یا DBSCAN و CFDP^۵ و روش‌های دیگر خوشه‌بندی صورت گیرد.

¹ Oversampling

² Undersampling

³ Synthetic Minority Oversampling Technique (SMOTE)

⁴ Cluster-based oversampling

⁵ Clustering by Fast search and find of Density Peaks

۳-۳- روش CTD^۱

در ردهی اقلیت، داده‌های مرزی بسیار حیاتی‌تر نسبت به مراکز رده‌های اقلیت هستند. بنابراین روش‌های موجود، بیشتر، نمونه‌های مرزی ردهی اقلیت را در نظر می‌گیرند. در زمان انتخاب نمونه‌های رده اقلیت برای افزایش رده اقلیت با استفاده تولید نمونه تصادفی و مصنوعی، نمونه‌ها می‌توانند در سه دسته قرار بگیرند [۲۰]:

- ۱- نمونه‌هایی با تأثیر تجمعی آشکار در ردهی اقلیت، که می‌توانند تأثیر کمی در کارایی مدل داشته باشند.
- ۲- نمونه‌های دارای نویز یا نمونه‌های جدا شده از رده اقلیت که از نمونه‌های ردهی اقلیت دور هستند و تأثیر تجمعی آشکار دارند. این نمونه‌ها نه تنها تأثیری مثبتی در کارایی مدل ندارند بلکه تأثیر منفی بر روی کارایی مدل نیز هم دارد.
- ۳- نمونه‌هایی که به دو دسته‌ی بالا مرتبط نیستند و نقش بسزایی در بهبود تشخیص رده‌های اقلیت دارد. در واقع این نمونه‌ها تأثیر زیادی در بهبود رده‌بندی‌های پایین رده‌ای دارد. شکل (۱) نمای شماتیک از سه دسته داده را نشان می‌دهد. این روش هر سه نمونه را در تولید نمونه‌های رده اقلیت استفاده می‌کند. از داده‌ها با دامنه‌ی مثبت (pos)، نمونه‌ها با دامنه‌ی منفی (neg) و از نمونه‌ها با دامنه مرزی (bnd) استفاده می‌شود. به زبانی ساده، استفاده از روش‌هایی برای تشخیص این سه دسته و استفاده از آنها برای تولید رده‌های اقلیت، از عملیات‌های روش CTD است. در این روش همچنین برای نرخ نمونه‌گیری از روش SMOTE استفاده می‌شود [۲۰].

۳-۴- روش FRO^۲

این روش به طور خاص به مشکل رده‌بندی داده‌های نامتوازن با خصوصیات عددی می‌پردازد. برای تعیین وزن قوانین در این روش از رابطه (۱) استفاده می‌شود.

$$rw_j = \frac{\sum_{x_i \in class C_j} f_{A_j}(x_i)}{\sum_{i=1}^m f_{A_j}(x_i)} \quad \text{رابطه (۱)}$$

در این رابطه C_j یک برچسب از یک رده است؛ m تعداد قوانین فازی است؛ rw_j وزن محاسبه شده قانون j است و A_j یک ارزش زبانی است که توسط یک تابع عضویت به نام f تعریف شده است [۲۱].
از این رابطه می‌توان نتیجه گرفت که هر چه وزن قوانین بیشتر باشد، مساحت پوشش داده شده توسط قوانین امن‌تر خواهد بود. این الگوریتم در دو نوع ردهی دودویی و ردهی چند متغیره استفاده می‌شود. همچنین در این روش مقادیر گمشده^۳ نیز پوشش داده می‌شوند [۲۱].

۳-۵- رویکرد نمونه‌گیری مصنوعی سازگار^۴

یکی دیگر از روش‌های افزایش ردهی اقلیت، استفاده از روش ADASYN است. تحقیقات زیادی بر روی این روش و تفاوت آن با روش SMOTE شده است. ایده‌ی اصلی این روش استفاده از یک توزیع وزن دار برای نمونه‌های متفاوت ردهی اقلیت مطابق با سطوح مختلف یادگیری است که این داده‌ها، نسبت به داده‌های واقعی رده اقلیت، یادگیری سخت‌تری دارند. این روش یادگیری از داده‌های نامتوازن، سعی در بهبود دو جنبه از داده‌های نامتوازن را دارد [۲۲]:

- ۱- کاهش معرفی یک مدل جانبدارانه

¹ CCA based three-way decision model

² Fuzzy Rule-based Oversampling

³ Missing values

⁴ ADASYN: Adaptive Synthetic Sampling Approach

سومین همایش ملی مهندسی کامپیوتر، داده کاوی و داده های حجیم

۲- تغییر مرز رده بندی در مورد داده های سخت به طور سازگار

ورودی ADASYN یک مجموعه داده ی آموزشی به تعداد m نمونه است. رویه ی اصلی ADASYN به صورت مراحل ذیل است [۲۲]:

۱- ابتدا درجه ی عدم توازن محاسبه می شود. در رابطه (۲) تعداد نمونه های رده ی اقلیت با m_s و تعداد نمونه های رده ی اکثریت با m_l نشان داده شده است.

$$d = m_s/m_l, d \in (0, 1] \quad \text{رابطه (۲)}$$

۲- اگر $d < d_{th}$ باشد:

I. محاسبه تعداد نمونه های مصنوعی مورد نیاز برای رده ی اقلیت از رابطه (۳) حاصل می شود.

$$G = (m_l - m_s) \times \beta, \beta \in [0, 1] \quad \text{رابطه (۳)}$$

β مشخص کننده ی سطح مطلوب تولید داده های مصنوعی بعد از تولید داده های مصنوعی است که $\beta = 1$ به معنای تولید یک مجموعه داده ی متعادل کامل بعد از فرآیند تولید نمونه ی مصنوعی از رده ی اقلیت است. همچنین G کل تعداد نمونه های داده ی مصنوعی مورد نیاز برای رده ی اقلیت است.

II. پیدا کردن k نزدیکترین همسایه برای هر نمونه ای که عضو مجموعه داده ی رده ی اقلیت است. در روش ADASYN از فاصله اقلیدسی در فضای مؤلفه ای (صفات خاصه) برای پیدا کردن k همسایه نزدیک استفاده می شود. نام مجموعه k نزدیکترین همسایه، r_i است که با رابطه ی (۴) بدست می آید که همان تعداد نمونه ها در k نزدیکترین همسایه ی X_i که وابسته به رده ی اکثریت است. بنابراین $r_i \in [0, 1]$ است.

$$r_i = \Delta_i/K, i = 1, \dots, m_s \quad \text{رابطه (۴)}$$

III. با استفاده از رابطه (۵)، نرمال سازی r_i با توزیع چگالی $(\sum_i \hat{r}_i = 1)$ صورت می گیرد.

$$\hat{r}_i = r_i / \sum_{i=1}^{m_s} r_i \quad \text{رابطه (۵)}$$

IV. محاسبه تعداد نمونه های رده ی اقلیت مورد نیاز که عضو رده ی اقلیت است. این کار به وسیله ی رابطه ی (۶) صورت میگیرد.

$$g_i = \hat{r}_i \times G \quad \text{رابطه (۶)}$$

V. برای هر نمونه X_i رده ی اقلیت تعداد g_i نمونه ی مصنوعی تولید می شود که این تولید نمونه با انتخاب تصادفی یک نمونه ی رده ی اقلیت در یک حلقه با رابطه (۷) صورت می گیرد که $(X_{zi} - X_i)$ همان تفاوت برداری صفات خاصه یا همان مؤلفه های رده ی اقلیت است و یک عدد تصادفی بین صفر و یک است.

$$S_i = X_i + (X_{zi} - X_i) \times \lambda, \lambda \in [0, 1]$$

رابطه (۷)

۴- نتایج

هر یک از روش های بیش نمونه گیری با توجه به نوع داده و یادگیرنده انتخاب می شوند. یکی از کاربردهای مهم SMOTE در بازبایی اطلاعات در ساخت کوله های کلمات است. SMOTE و دیگر انشعابات آن مانند SMOTE و Borderline- SMOTE با یکدیگر تفاوت دارد که این تفاوت باعث تغییر در به کارگیری و همچنین نوع هدف دارد. هر یک از روش های مطرح شده در بیش نمونه گیری برای مؤلفه های هدف دودویی و چند رده ای می توانند مورد استفاده قرار بگیرند.

استفاده انفرادی از روش های بیش نمونه گیری بازدهی کامل در بسیاری از روش های داده کاوی و به خصوص رده بندی داده های نامتوازن ندارد و سعی می شود برای افزایش کارایی سیستم، روش های بیش نمونه گیری به همراه روش های کم نمونه گیری مورد استفاده قرار می گیرند.

در روش ADASYN برای نمونه های رده ای اقلیت وزن در نظر گرفته می شود و براساس وزن هر کدام از نمونه های رده ای اقلیت، نمونه ای مصنوعی تولید می شود و این در حالی است که در روش SMOTE، تولید نمونه های مصنوعی براساس K نزدیکترین همسایه صورت می گیرد که کلیه نمونه های رده ای اقلیت براساس یک شانس برابر، برای تولید نمونه ای تصادفی از نزدیکترین همسایه استفاده می کنند. در روش ADASYN صحت و دقت مدل ایجاد شده و یا دانش استخراج شده، بالاتر است زیرا نمونه های تولید شده از رده ای اقلیت، بیشترین شباهت را به اصل نمونه های وزن دار رده ای اقلیت دارند. در روش FRO، حتی اگر مؤلفه ها بین صفر و یک نباشند، خروجی تولید شده بین صفر و یک تولید می شوند، بنابراین نیاز به نرمال سازی نیست.

در روش CTD نمونه های مرزی بسیار پر اهمیت هستند، اما در روش بیش نمونه گیری خوشه بندی، آن نمونه هایی از اهمیت بیشتری برخوردارند که در سرخوشه باشند.

در روش های بیش نمونه گیری نباید داده های مصنوعی تولید شده از رده ای اقلیت، باعث یادگیری سخت شوند و یا مدل غیرمناسب و ناکارا ارائه دهند.

۵- منابع

- [1] Q. Zou, S. Xie, Z. Lin, M. Wu, and Y. Ju, "Finding the best classification threshold in imbalanced classification," *Big Data Res.*, vol. 5, pp. 2–8, 2016.
- [2] D. Devi and B. Purkayastha, "Redundancy-driven modified Tomek-link based undersampling: a solution to class imbalance," *Pattern Recognit. Lett.*, vol. 93, pp. 3–12, 2017.
- [3] S. Abdellatif, M. A. Ben Hassine, S. Ben Yahia, and A. Bouzeghoub, "ARCID: A new approach to deal with imbalanced datasets classification," in *International Conference on Current Trends in Theory and Practice of Informatics*, 2018, pp. 569–580.
- [4] M. A. H. Farquand and I. Bose, "Preprocessing unbalanced data using support vector machine," *Decis. Support Syst.*, vol. 53, no. 1, pp. 226–233, 2012.
- [5] Y. Qian, Y. Liang, M. Li, G. Feng, and X. Shi, "A resampling ensemble algorithm for classification of imbalance problems," *Neurocomputing*, vol. 143, pp. 57–67, 2014.
- [6] J. F. Díez-Pastor, J. J. Rodríguez, C. García-Osorio, and L. I. Kuncheva, "Random Balance: Ensembles of variable priors classifiers for imbalanced data," *Knowledge-Based Syst.*, vol. 85, pp. 96–111, 2015.
- [7] S. Wang, W. Liu, J. Wu, L. Cao, Q. Meng, and P. J. Kennedy, "Training deep neural networks on imbalanced data sets," in *Neural Networks (IJCNN), 2016 International Joint Conference on*, 2016, pp. 4368–4374.

- [8] B. Krawczyk, "Learning from imbalanced data: open challenges and future directions," *Prog. Artif. Intell.*, vol. 5, no. 4, pp. 221–232, 2016.
- [9] Z. Zhang, B. Krawczyk, S. García, A. Rosales-Pérez, and F. Herrera, "Empowering one-vs-one decomposition with ensemble learning for multi-class imbalanced data," *Knowledge-Based Syst.*, 2016.
- [10] J. M. Johnson and T. M. Khoshgoftaar, "Survey on deep learning with class imbalance," *J. Big Data*, vol. 6, no. 1, p. 27, 2019.
- [11] I. Nekooimehr and S. K. Lai-Yuen, "Cluster-based weighted oversampling for ordinal regression (CWOS-Ord)," *Neurocomputing*, vol. 218, pp. 51–60, 2016.
- [12] T. Elhassan and M. Aljurf, "Classification of Imbalance Data using Tomek Link (T-Link) Combined with Random Under-sampling (RUS) as a Data Reduction Method.," 2016.
- [13] L. Zhao et al., "A cost-sensitive meta-learning classifier: SPFCNN-Miner," *Futur. Gener. Comput. Syst.*, 2019.
- [14] R. Barandela, R. M. Valdovinos, and J. S. Sánchez, "New applications of ensembles of classifiers," *Pattern Anal. Appl.*, vol. 6, no. 3, pp. 245–256, 2003.
- [15] A. Zhukov, N. Tomin, V. Kurbatsky, D. Sidorov, D. Panasetsky, and A. Foley, "Ensemble methods of classification for power systems security assessment," *Appl. Comput. Informatics*, 2017.
- [16] N. V Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: synthetic minority over-sampling technique," *J. Artif. Intell. Res.*, vol. 16, pp. 321–357, 2002.
- [17] C. Bunkhumpornpat, K. Sinapiromsaran, and C. Lursinsap, "Safe-level-smote: Safe-level-synthetic minority over-sampling technique for handling the class imbalanced problem," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 2009, pp. 475–482.
- [18] H. Han, W.-Y. Wang, and B.-H. Mao, "Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning," in *International Conference on Intelligent Computing*, 2005, pp. 878–887.
- [19] M. S. Santos, P. H. Abreu, P. J. García-Laencina, A. Simão, and A. Carvalho, "A new cluster-based oversampling method for improving survival prediction of hepatocellular carcinoma patients," *J. Biomed. Inform.*, vol. 58, pp. 49–59, 2015.
- [20] Y. T. Yan, Z. B. Wu, X. Q. Du, J. Chen, S. Zhao, and Y. P. Zhang, "A three-way decision ensemble method for imbalanced data oversampling," *Int. J. Approx. Reason.*, vol. 107, pp. 1–16, 2019.
- [21] G. Liu, Y. Yang, and B. Li, "Fuzzy rule-based oversampling technique for imbalanced and incomplete data learning," *Knowledge-Based Syst.*, vol. 158, pp. 154–174, 2018.
- [22] H. He, Y. Bai, E. A. Garcia, and S. Li, "ADASYN: Adaptive synthetic sampling approach for imbalanced learning," in *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, 2008, pp. 1322–1328.