

دسته‌بندی روش‌های شناسایی رانش مفهوم

نگارنده: امیر کتول

دانشجوی کارشناسی ارشد دانشگاه گلستان گرگان (گرایش نرم افزار) amirkatool@gmail.com

چکیده

بخش عظیمی از دنیای تکنولوژی امروز را سیستم‌هایی تشکیل داده‌اند که با جریانی پیوسته و سریع از داده کار می‌کنند (جریان داده). بحث یادگیری در جریان داده با چالش‌های مختلفی از جمله دسترسی محدود به داده‌ها و تغییر توزیع داده‌ها در گذر زمان (رانش مفهوم) روبه‌رو است. در طول دهه اخیر توجه به بحث رانش مفهوم بیشتر شده و شناسایی آن به مسئله محوری در یادگیری آنلاین و داده کاوی از جریان داده تبدیل شده است. بسیاری از روش‌های شناسایی رانش مفهوم بر اساس کنترل عملکرد کلاس‌بندی‌کننده و پیش‌بینی کلاس حقیقی داده‌ها کار می‌کنند. دسترسی به کلاس حقیقی داده‌ها و میزان تأخیر این دسترسی عاملی تعیین‌کننده در میزان کارایی روش‌های شناسایی رانش مفهوم در شرایط مختلف است. بر همین اساس، دسته‌بندی برای روش‌های شناسایی رانش مفهوم در سال ۲۰۱۵ صورت گرفته است که به سه دسته روش‌های نظارتی، بدون نظارت و نیمه نظارتی تقسیم می‌شود. با توجه به اینکه این دسته‌بندی شامل روش‌های مطرح شده بعد از سال ۲۰۱۵ نمی‌باشد؛ در این گزارش سعی شده تا با پر کردن این جای خالی و تکمیل این دسته‌بندی با روش‌های شناسایی رانش مفهوم مطرح شده بعد از سال ۲۰۱۵ و معرفی اجمالی هر کدام از آنها، قدمی مثبت در تفکیک این روش‌ها از هم برداشته شود.

کلمات کلیدی: جریان داده، یادگیری ماشین، رانش مفهوم، روش‌های شناسایی نظارتی، بدون نظارت، نیمه نظارتی

۱- مقدمه

در دنیای امروز سرعت رشد و توسعه تکنولوژی فراتر از تصورات و پیش‌بینی‌های گذشته است. تخمین حجم اطلاعات تولیدی روزانه در دنیای امروز کاری تقریباً نشدنی است. داده، در هر لحظه، در بسترهای مختلف با سرعت بالا تولید و منتشر شده و توانایی تحلیل و یادگیری از داده در این شرایط را به چالشی جدید برای دنیای تکنولوژی تبدیل می‌کند. یادگیری از این جریان بی‌اندازه داده که به صورت پیوسته و پی‌درپی تولید و ارسال می‌شود و آن را با نام جریان داده می‌شناسیم (Data Stream)، چالش‌های خاص خود را دارد. از چالش‌های اصلی برای یادگیری در جریان داده، حجم و اندازه غیرقابل اندازه‌گیری و بیش از اندازه داده‌ها؛ سرعت انتشار بسیار بالا و دسترسی محدود به داده‌ها که در بسیاری از شرایط تنها یکبار قابل دسترسی هستند و تغییر داده‌ها در گذر زمان (رانش مفهوم)؛ را می‌توان نام برد. اصلی‌ترین چالش در یادگیری جریان داده، شناسایی رخداد رانش مفهوم است. رانش مفهوم که به تغییر ویژگی‌های آماری متغیر مورد مطالعه با گذر زمان در این محیط‌های پویا گفته می‌شود، اصلی‌ترین دلیل افت کیفی سیستم‌های مبتنی بر استخراج دانش در گذر زمان شناخته شده است. هرچند توجه به رانش مفهوم و تلاش برای شناسایی آن در سال‌های اخیر بیشتر شده است، اما هنوز این مبحث در قدم‌های اولیه خود قرار داشته و ناشناخته‌های زیادی پیرامون آن وجود دارد. بیشتر روش‌های شناسایی رانش مفهوم با توجه به بررسی عملکرد کلاس‌بندی‌کننده کار می‌کنند و ضعف اصلی آن‌ها این است که فرض می‌کنند که کلاس حقیقی داده‌ها در هر لحظه قابل دسترسی بوده که در بسیاری از شرایط دنیای حقیقی امری ناممکن است. در سال‌های اخیر روش‌های جدیدی با در نظرگیری تغییرات ابعاد ویژگی‌های داده، برای شناسایی رانش مفهوم مطرح شده است. مقالات متعددی نیز بر اساس معیارهای مختلف، دسته‌بندی‌هایی روی این روش‌ها انجام داده اند اما مقاله‌ای که تقسیم بندی خود را براساس نظارتی یا بدون نظارت بودن، با توجه به جدیدترین روش‌ها (۲۰۱۵ به بعد) انجام دهد، تا آنجا که در دانش این گزارش می‌گنجد صورت نپذیرفته است. در این گزارش بعد از تعریف و بررسی رانش مفهوم، سعی شده تا جدیدترین روش‌های شناسایی رانش مفهوم در یک دسته‌بندی بر اساس نظارتی یا نیمه نظارتی و بدون نظارت، دسته‌بندی شود. ضمناً این روش‌ها به صورت اجمالی بررسی و معرفی شوند.

۲- رانش مفهوم

۲.۱- تعریف

به پدیده تغییر توزیع داده در گذر زمان رانش مفهوم (Concept Drift) می‌گویند. در جریان داده، بسیاری مواقع، داده‌های جدید، از الگوها جدیدی پیروی می‌کنند که در داده‌های پیشین مشاهده نشده است. این تغییرات به دلایل پیش‌بینی نشده‌ای ممکن است رخ دهد. به عنوان مثال در یک سیستم تشخیص کلاهبرداری، یک گونه جدید از کلاهبرداری‌های کارت اعتباری ممکن است مشاهده شود که پیشتر وجود نداشته است. با این تغییرات داده، الگوها و کلاس‌بندی‌کننده‌هایی که قبل‌تر و با داده‌های قبل از این تغییرات ایجاد شده اند، دچار افت عملکرد خواهند شد. در تعریف رانش مفهوم داریم:

در بازه زمانی صفر تا t ، یک مجموعه از نمونه‌های $S_{0:t} = \{d_0, \dots, d_t\}$ داریم که $d_i = (X_i, y_i)$ است. بردار ویژگی‌های داده و y_i کلاس مربوط به آن است. رانش مفهوم بدین صورت تعریف می‌شود که در طول گذر زمان توزیع داده دچار تغییراتی می‌شود که در اثر این تغییرات خواهیم داشت $\exists t: p_t(X, y) \neq p_{t+1}(X, y)$ احتمال مشترک بین X و y است. رانش مفهوم پدیده‌ای پیچیده است که به شکل‌های مختلفی می‌تواند رخ دهد که بررسی آن می‌پردازیم.

۲.۲- انواع رانش مفهوم

اگر رابطه احتمال مشترک را براساس رابطه بیز بنویسیم خواهیم داشت: $p_t(X, y) = p_t(X) \times p_t(y|X)$

براساس این رابطه تغییرات احتمال مشترک می‌تواند ناشی از تغییرات:

الف) $p_t(X)$ باشد. در اثر این تغییرات توزیع موجود درون کلاس‌های تعیین شده دچار تغییر شده اما این تغییرات اثری روی مرزهای کلاس‌بندی نداشته و عملکرد کلاس‌بندی‌کننده را تضعیف نمی‌کند و به آن، رانش مجازی نیز گفته می‌شود.

ب) $P_t(y|X)$ باشد. در اثر این تغییرات، کلاس تعیین شده برای مجموعه ویژگی‌های معین دچار تغییر شده و به دنبال آن مرزهای کلاس-بندی هم دچار تغییر می‌شود. به آن رانش واقعی گفته شده و در اثر آن عملکرد کلاس‌بندی کننده افت خواهد کرد. لازم به ذکر است که در بسیاری از پدیده‌های واقعی هر دو حالت الف و ب می‌تواند به صورت همزمان رخ دهد. همچنین دو پدیده پدیداری کلاس جدید (Novel Class Emergence) و ادغام کلاس‌های موجود (Existing Class Fusion) که در گذشته به عنوان گونه‌های مجزا از رانش مفهوم در نظر گرفته نمی‌شد را می‌توان به عنوان دسته سوم انواع رانش که در اثر تغییرات احتمال پیشین $P_t(y)$ رخ میدهد، در نظر گرفت.

رانش مفهوم پدیده پیچیده‌ای هست که از لحاظ سرعت، شدت و قابل پیش‌بینی بودن نیز به انواع مختلفی تقسیم شده و می‌تواند به شکل‌های مختلفی رخ دهد. این تنوع تشخیص رانش مفهوم را سخت‌تر کرده و باعث می‌شود بسیاری از روش‌ها تنها در شناسایی گونه‌هایی خاص ممتاز باشند. (Khamassi, Sayed-Mouchaweh et al. (2018)

۳- شناسایی رانش مفهوم

۳.۱- مطالعات پیشین

در سال‌های اخیر مقالات متعددی در رابطه با بحث روش‌های شناسایی رانش مفهوم مطرح شده است. در این قسمت سعی شده تا آخرین مقالات جامع و مروری پیرامون روش‌های شناسایی رانش مفهوم بررسی و به صورت مختصر بیان شوند. در مقاله Lu, Liu et al. (2018) چارچوبی کلی برای شناسایی رانش مفهوم ارائه شده است که متشکل از چهار مرحله مختلف است: بازیابی داده (Data Retrieval)، مدل‌سازی داده (Data Modeling)، محاسبه آزمون‌های آماری (Test Statistics Calculation) و آزمون فرضیه (Hypothesis Test). با توجه به چارچوب مطرح شده تقسیم روش‌های مختلف شناسایی رانش مفهوم در این مقاله براساس تکنیک‌های مختلفی است که در مرحله سوم یعنی محاسبه آزمون آماری استفاده می‌شود. این مقاله روش‌های شناسایی رانش مفهوم را به سه دسته کلی تقسیم می‌کند. دسته اول روش‌هایی است که براساس کنترل نرخ خطا کلاس‌بندی کننده کار می‌کنند. دسته دوم شامل روش‌هایی است که از اندازه‌گیری‌های فاصله بین داده‌های جدید و قدیم برای تخمین میزان تغییرات استفاده می‌کنند و مبتنی بر کنترل توزیع داده هستند. دسته سوم شامل روش‌هایی است که از چندین آزمون فرضیه با هم برای شناسایی رانش مفهوم استفاده می‌کنند که خود به دو دسته موازی و سری تقسیم می‌شود.

در مقاله Hu, Kantardzic et al. (2020) دسته‌بندی روش‌های شناسایی رانش مفهوم به دو دسته کلی مبتنی بر عملکرد (performance-based) و مبتنی بر توزیع داده (data distribution-based) تقسیم شده است. روش‌های مبتنی بر عملکرد به صورت پیوسته، معیاری مربوط به میزان خطا مانند درستی (accuracy)، میزان اف (F-measure)، دقت (precision) و بازیابی (recall) را کنترل می‌کنند. به صورت کلی در صورتی که افت محسوس در این معیارها رخ دهد خبر از رخداد رانش مفهوم می‌دهد. روش‌های مبتنی بر توزیع داده بر اساس کنترل تغییر معیارهای توزیع داده مانند مکان (location)، چگالی (density) و دامنه (range) کار می‌کنند. لازم به ذکر است روش‌های موجود در این دسته با روش‌های دو دسته آخر مقاله Lu, Liu et al. (2018) همپوشانی دارند.

در مقاله Wares, Isaacs et al. (2019) دسته‌بندی مختص روش‌های نظارتی مطرح می‌شود. این دسته‌بندی به چهار گروه روش‌های آماری (Statistical methods)، روش‌های مبتنی بر پنجره (window-based methods)، روش‌های دسته جمعی مبتنی بر بلوک (block-based ensembles) و روش‌های دسته جمعی افزایشی (incremental-based ensembles) تقسیم می‌شوند. دسته اول، شامل روش‌هایی است که از آزمون‌های آماری نظیر Page-Hinckley برای شناسایی رانش استفاده می‌کنند. دسته دوم شامل روش‌هایی است که براساس درستی کلاس‌بندی کننده برای نمونه‌های درون یک پنجره کار می‌کنند. دو دسته آخر براساس درستی کلاس‌بندی کننده دسته جمعی کار می‌کنند و تفاوت آن‌ها بر اساس نحوه واکنش به رانش مفهوم است. دسته سوم نمونه‌ها را به صورت بلوک داده برای آموزش مجدد استفاده می‌کنند در حالی که دسته آخر با رسیدن هر داده جدید به صورت افزایشی مدل را بازآموزی می‌کند.

در مقاله (Gemaque, Costa et al. (2020) دسته‌بندی براساس نحوه عملکرد روش شناسایی رانش مفهوم مختص روش‌های نیمه نظارتی و بدون نظارت انجام شده است. این دسته‌بندی روش‌های نیمه نظارتی و بدون نظارت را بر اساس اینکه داده‌ها را با چه قالبی بررسی می‌کنند به دو دسته روش‌های شناسایی رانش مفهوم بسته‌ای (batch drift detection methods) و روش‌های شناسایی رانش مفهوم آنلاین (online drift detection methods) تقسیم می‌کند. روش‌های بسته‌ای داده‌ها را در قالب یک بسته از جریان داده جمع-آوری کرده و سپس رخداد رانش آن را بررسی می‌کنند که خود بر اساس اینکه از تمام داده‌های یک بسته استفاده می‌کنند یا تنها از بخشی از آن به دو دسته تشخیص با کل بسته و تشخیص با جزئی از بسته تقسیم می‌شوند. در مقابل روش‌های شناسایی آنلاین با رسیدن هر داده رخداد رانش را بررسی می‌کنند و بر اساس اینکه از پنجره‌ای با اندازه ثابت یا متغیر برای مقایسه داده‌های جدید و قدیم استفاده می‌کنند به دو دسته روش‌های آنلاین با پنجره ثابت و متغیر تقسیم می‌شوند.

در مقاله (Khamassi, Sayed-Mouchaweh et al. (2018) دسته‌بندی کلی بر اساس پاسخ به پنج سوال صورت گرفته است: داده چگونه پردازش می‌شود؟ یادگیری چگونه انجام می‌شود؟ رانش مفهوم چگونه نظارت می‌شود؟ رانش مفهوم چگونه مدیریت می‌شود؟ و معیارهای عملکرد چه هستند؟ در پاسخ به سوال سوم مطرح شده در این مقاله که در این گزارش مد نظر ماست، انواع روش‌های شناسایی رانش مفهوم به سه دسته نظارتی، نیمه نظارتی و بدون نظارت براساس نیاز به کلاس حقیقی داده، تقسیم شده‌اند. با توجه به اینکه دسته‌بندی که در مقاله (Khamassi, Sayed-Mouchaweh et al. (2018) صورت گرفته است، شامل بسیاری از روش‌های نیمه نظارتی و بدون نظارت که پس از آن ارائه شده‌اند نیست، در این گزارش سعی شده تا براساس همان روش دسته‌بندی روش‌هایی که بعد از سال ۲۰۱۵ مطرح شده‌اند بررسی و دسته‌بندی شوند. با توجه به این دسته‌بندی، روش‌های دسته اول مقاله‌های Lu, Liu et al. (2018) و Hu, Kantardzic et al. (2020) در دسته نظارتی قرار می‌گیرند. باقی روش‌های مطرح شده در این مقالات، برخی در دسته نظارتی و برخی در دسته بدون نظارت قرار می‌گیرند که در ادامه به بررسی آن‌ها می‌پردازیم.

۳.۲- روش‌های نظارتی

روش‌های نظارتی روش‌هایی هستند که نتیجه پیش‌بینی کلاس هر داده را در هر زمان آماده فرض می‌کنند. این روش‌ها براساس کنترل معیارهای میزان خطا کلاس‌بندی کننده مانند دقت، درستی، بازیابی و... کار می‌کنند و کنترل تغییرات $P(Y|X)$ را بر عهده دارند. ساختار کار این دسته از روش‌های شناسایی رانش مفهوم، شرایط استقلال از روش یادگیری و درستی بالا در شناسایی همراه با امکان شناسایی انواع رانش مفهوم را ممکن می‌سازد. در مقابل، فرض وجود کلاس حقیقی هر داده در هر لحظه در بسیاری از شرایط ناممکن و در صورت وجود بسیار هزینه بر است.

روش (Fuzzy Windowing Drift Detection Method) FW-DDM در (Liu, Zhang et al. (2017) یکی از روش‌های نظارتی است که برخلاف روش‌های شناسایی رانش مفهوم نظارتی مرسوم که از پنجره استفاده می‌کنند، با استفاده از پنجره زمانی فازی و اجازه به اینکه دوره همپوشان از مفهوم جدید و قدیم در کنار هم در پنجره لغزنده قرار گیرند، توانست نسبت به روش‌های مرسوم نظارتی دقت و درستی بهتری کسب کند.

روش (Dynamic Extreme Learning Machine) DELM که در (Xu and Wang (2017) مطرح شد، این روش از ELM که یک روش شبکه عصبی با یک لایه پنهان است استفاده می‌کند. این روش از DDM که یکی از قدیمی‌ترین روش‌های شناسایی رانش مفهوم نظارتی است (Gama, Medas et al. (2004) استفاده می‌کند تا رانش را شناسایی کند. بعد از اخطار رخداد رانش تعداد گره‌های لایه مخفی افزایش پیدا می‌کند تا کلاس‌بندی کننده عمومی‌تری یادگیری شود.

روش (Linear Four Rate) LFR که در (Wang and Abraham (2015) مطرح شد، از روش‌های چند آزمون فرضیه موازی است که با کنترل همزمان چهار معیار FN, TN, FP, TP و رخداد رانش را کنترل می‌کند. با توجه به اینکه این روش به طور همزمان چهار معیار مختلف طبقه‌بندکننده را کنترل می‌کند حساسیت به تغییر آن بسیار زیاد است و به همین خاطر رخداد رانش مفهوم را با نرخ بالا و همین طور بسیار سریع شناسایی می‌کند و در عین حال از مشکلات آن نرخ بالای False Alarm به خاطر همین حساسیت زیاد است.

یکی دیگر از روش‌های نظارتی چند آزمون فرضیه روش سه لایه (three-layer) مطرح شده در (Zhang, Chu et al. (2017) است. این روش که برای شناسایی رانش مفهوم در متن در ابتدا مطرح شد، از سه لایه تشکیل شده است که براساس ارزش اطلاعات (Information Value) و شباهت ژاکارد (Jaccard similarity) رانش مفهوم را شناسایی می‌کند. در لایه یک تغییرات کلاس حقیقی داده‌ها، $P(y)$ ؛ در لایه دوم تغییرات فضای ویژگی‌ها، $P(X)$ ؛ و در لایه سوم تغییرات رابطه بین ویژگی‌ها و کلاس حقیقی داده، $P(y|X)$ بررسی می‌شود. در این روش با محاسبه ارزش اطلاعات و وزن شواهد (Weight of Evidence) بین دو بسته داده جدید و قدیمی میزان تغییرات محاسبه و مقایسه می‌شود.

روش Hierarchical Change-Detection Tests (HCDDTs) که در (Alippi, Boracchi et al. (2016) مطرح شد ساختار کلی روش‌های چند آزمون فرضیه سری را بیان می‌کند. این روش از دو لایه شناسایی (detection layer) و لایه تصدیق (validation layer) تشکیل شده است. در لایه شناسایی هر کدام از روش‌های موجود شناسایی رانش که تاخیر شناسایی پایین و نیاز محاسباتی کمتری دارند قابل استفاده است. بعد از شناسایی رانش در این لایه، لایه تصدیق که دارای محاسبات سنگین‌تر و پیچیده‌تری است و می‌تواند از آزمون‌های فرضیه معروف در آن استفاده کرد، فعال شده و رخداد رانش یا عدم آن را تایید می‌کند. با وجود اینکه می‌تواند از روش‌های مختلفی در دو لایه شناسایی و تصدیق استفاده کرد و این روش به شکل‌های مختلف مطرح می‌شود. چون نگارنده مقاله از تکنیک‌های آماری برای لایه تشخیص استفاده کرده است که نظارتی هستند؛ این روش در دسته نظارتی در این گزارش مطرح شده است.

روش Hierarchical Linear Four Rate (HLFR) مطرح شده در (Yu and Abraham (2017) یکی دیگر از روش‌های چند آزمون فرضیه سری است. در این روش برای لایه شناسایی، روش LFR مورد استفاده قرار گرفته است. نرخ بالای امکان تشخیص اشتباه روش LFR، حال با وجود لایه تصدیق و وجود روشی برای تایید یا رد رخداد رانش، بهتر شده و امکان تشخیص انواع شکل‌های رانش را برای این روش ممکن می‌سازد. در لایه تصدیق از یک آزمون جایگشتی (Permutation Test) برای تایید رخداد رانش استفاده شده است. در مقاله (de Lima Cabral and de Barros (2018) روش Nishida Statistical Test of Equal Proportions (STEPD) and Yamauchi (2007) که یکی از روش‌های معروف برای شناسایی رانش مفهوم نظارتی است؛ با استفاده از آزمون فیشر بهبود یافته است. روش مطرح شده از نظر اینکه دو پنجره از داده را با هم مقایسه کرده و در صورت تجاوز از حد آستانه اخطار رانش مفهوم را دهد مانند STEPDP عمل می‌کند، با این تفاوت که با بکارگیری آزمون دقیق فیشر (Fisher's Exact test) به جای اینکه براساس نرخ پیش-بینی درست کار کند، براساس تغییرات میزان خطا کار می‌کند. استفاده از این روش آزمون آماری، این امکان را بوجود می‌آورد که برخلاف روش STEPDP که بیشتر برای شرایطی نظیر تعداد نمونه داده محدود یا داده نامتعادل (imbalanced)، که جوابگو نبود؛ اکنون جوابگو باشد. در بررسی‌های انجام شده، این روش توانست گونه‌های بیشتری از رانش را بهتر از روش STEPDP شناسایی کند.

روش McDiarmid Drift Detection Method (MDDM) مطرح شده در (Pesaranghader, Viktor et al. (2018) با استفاده از نابرابری مکدیاریمد (McDiarmid's inequality) و وزن‌دهی به پیش‌بینی داده در پنجره لغزان، توانست تاخیر شناسایی رانش مفهوم را نسبت برخی روش‌ها نظیر Page-Hinckley بهبود ببخشد. در این روش به داده‌های جدید وزن بیشتری داده می‌شود تا اهمیت آنها در محاسبات بیشتر شود. در پنجره لغزان میانگین وزن داده‌ها محاسبه شده و با ماکسیمم میانگین وزن موجود مقایسه می‌شود، در صورتی که میزان تفاوت از حدود نابرابری مکدیاریمد بیشتر باشد، رخداد رانش تایید می‌شود.

روش Hoeffding Drift Detection Method (FHDDM) مطرح شده در (Pesaranghader and Viktor (2016) نیز برای کاهش تاخیر شناسایی رانش مفهوم روشی مبتنی بر نابرابری هافدینگ (Hoeffding's inequality) در پنجره لغزان ارائه می‌کند. این روش براساس مقایسه احتمال پیش‌بینی درست کلی با احتمال پیش‌بینی درست داده‌های جدید، رخداد رانش را شناسایی می‌کند. در صورتی که رانش رخ نداده باشد تفاوت محسوس نخواهد بود.

روش Multiscale Drift Detection Test (MDDT) مطرح شده در (Wang, Kang et al. (2018) با نگاه کاهش حجم محاسباتی برای شناسایی رانش مفهوم، روشی مبتنی بر نمونه‌گیری مجدد (resampling) و آزمون تی (paired student's t test) ارائه

می‌کند. در این روش داده‌های قدیمی با نمونه‌گیری و به صورت تصادفی با داده‌های جدید برای تایید رخداد رانش مقایسه می‌شوند. کاهش حجم محاسبات برای شناسایی رانش باعث شده این روش برای شناسایی رانش بسیار سریع عمل کند.

بسیاری از روش‌های شناسایی رانش مفهوم در مدل‌های دست‌جمعی کلاس‌بندی‌کننده‌ها کار می‌کنند. روش Number and Distance of Errors (NDE) در (Dehghan, Beigy et al. (2016) مطرح شد، روشی است که هر یک از نمونه‌ها را بررسی کرده و توزیع خطا مدل دست‌جمعی کلاس‌بندی را کنترل می‌کند. NDE تخمینی از میزان خطا کلاس‌بندی هر نمونه ارائه می‌دهد، در صورتی که NDE داده جدید با NDE کلی مجموعه تفاوت محسوسی داشته باشد، رخداد رانش اعلام می‌شود.

روش diversified dynamic-weighted majority (DDWM) در (Sidhu and Bhatia (2018) مطرح شده از دو مدل دسته‌جمعی کلاس‌بندی‌کننده استفاده می‌کند، یکی با تنوع پایین و دیگری با تنوع بیشتر، با اینکار روش پیشنهاد شده دقت بهتری برای کلاس‌بندی و شناسایی انواع مختلف رانش مفهوم پیدا می‌کند. به طور کلی اگر هر کدام در شناسایی رانش گونه‌ای از رانش مفهوم ضعیف عمل کنند، دیگری آنرا جبران خواهد کرد. هر مدل موجود در هر دو کلاس‌بندی‌کننده دسته‌جمعی (تنوع پایین و بالا) برای پیش‌بینی کلاس هر داده جدید تاثیر گذار خواهند بود. این روش با وزن‌دهی برای پیش‌بینی کلی هر دو مدل دسته‌جمعی و کلاس‌بندی‌کننده‌های پایه درون آنها کار می‌کند. با پیش‌بینی نادرست وزن‌ها کاهش می‌یابد و اگر از میزان آستانه‌ای عبور کند، آن کلاس‌بندی‌کننده پایه حذف می‌شود. در صورتی که پیش‌بینی کلی مجموعه اشتباه باشد، هر دو مدل دسته‌جمعی متنوع و کم تنوع کنترل شده و کلاس‌بندی‌کننده ضعیف حذف شده و کلاس‌بندی‌کننده جدیدی به مجموعه اضافه می‌شود.

در مدل کلاس‌بندی‌کننده‌های دسته‌جمعی، تکنیک‌های مختلفی برای ایجاد تنوع کلاس‌بندی‌کننده‌های پایه وجود دارد که هر کدام نقاط ضعف و قوتی دارند. در روش Ensemble EDIST2 مطرح شده در (Khamassi, Sayed-Mouchaweh et al. (2019) با ترکیب سه روش مختلف تنوع مبتنی بر بلوک داده (block-based data)، وزن داده (weighting-data) و فیلتر کردن داده (filtering-data) سعی شده تا عملکرد هر یک را در کنار دیگری بهبود داده و ضعف‌های آنرا جبران کند. این روش از EDIST2 که یک روش شناسایی رانش مفهوم است که با نظارت بر عملکرد مدل دسته‌جمعی (ensemble) رانش را شناسایی می‌کند؛ استفاده می‌کند. (Khamassi, Sayed-Mouchaweh et al. (2015) با ترکیب این سه روش ایجاد تنوع، این روش قادر است تا انواع مختلف رانش مفهوم را شناسایی کند.

روش Entropy Based Ensemble (EBE) که در (Mahdi, Pardede et al. (2018) مطرح شده است از آنتروپی (entropy) برای شناسایی رانش مفهوم در کلاس‌بندی‌کننده دسته‌جمعی استفاده می‌کند. در این روش آنتروپی اطلاعات (information entropy) برای هر یک از کلاس‌بندی‌کننده‌های پایه از یک بلوک داده با اندازه ثابت محاسبه می‌شود. کلاس‌بندی‌کننده‌ها براساس اینکه سطح آنتروپی آنها مورد تایید است یا خیر در مدل دسته‌جمعی وارد یا از آن حذف می‌شوند.

بیشتر روش‌های دسته‌جمعی از روش رای‌گیری سنتی برای کلاس‌بندی کردن استفاده می‌کنند. در (Krawczyk and Cano (2018) تغییر در روش رای‌گیری صورت گرفته است. در این روش امکان اینکه کلاس‌بندی‌کننده‌ها بتوانند از شرکت در رای‌گیری نهایی پرهیز کنند فراهم شده است. اطمینان (confidence) تمام کلاس‌بندی‌کننده‌ها برای هر یک از نمونه‌ها کنترل شده و تنها کلاس‌بندی‌هایی که میزان اطمینان آن از میزان آستانه بیشتر باشد قادر به شرکت در رای‌گیری نهایی هستند. این میزان آستانه به صورت پویا از جریان داده بدست می‌آید. امکان خودداری از شرکت در رای‌گیری نهایی در این روش مقاوت این روش دست‌جمعی نسبت به نویز را افزایش داده است. نویسندگان مقاله بیان می‌کنند که تکنیک اعمال شده، تکنیکی کم هزینه و ساده است که می‌توان برای بهبود همه روش‌های دسته‌جمعی آنلاین مورد استفاده قرار گیرد.

برخی از روش‌ها از ابزار متفاوتی برای شناسایی رانش مفهوم استفاده می‌کنند. در (Jaworski, Duda et al. (2017) Restricted Boltzmann Machine (RBM) برای شناسایی رانش مفهوم استفاده می‌شود که خود گونه‌ای از روش‌های شبکه عصبی است. RBM قادر است تا مدلی بنیادی از داده را یادگیری کند. بعد از مرحله آموزش RBM حاوی اطلاعاتی فشرده شده از توزیع داده‌ای است که از آن یادگیری انجام داده است. نویسندگان مقاله بیان می‌کنند که می‌توان از این اطلاعات فشرده شده برای مقایسه یکسانی توزیع داده‌های

جدید و قدیم استفاده شود. دو شاخص انرژی آزاد (free energy) و خطای بازسازی (reconstruction error) برای محاسبه و مقایسه پیشنهاد شده است. روش دیگری که از شبکه عصبی برای شناسایی رانش مفهوم استفاده می‌کند در (Lobo, Del Ser et al. (2018) مطرح شده است که از Evolving Spiking Neural Networks استفاده می‌کند که براساس تغییرات میزان ادغام نورون‌های خود می‌تواند رخداد رانش را شناسایی کند. این دو روش از جهت اینکه از ابتدا به شبکه عصبی مورد استفاده خود برای شناسایی رانش نیازمند بوده و در مرحله یادگیری خود در واقع شناسایی رانش مفهوم را شروع کرده‌اند که در آن‌ها کلاس حقیقی داده مورد نیاز است. به همین خاطر در این گزارش در دسته نظارتی قرار گرفته‌اند.

روش Online Histogram-based Naïve Bayes Classifier (OHNBC) که در (Astudillo, González et al. (2016) مطرح شده است با استفاده از هیستوگرام‌های تک بعدی (uni-dimensional histograms) پیش‌بینی کلاس نمونه‌ها را انجام می‌دهد و با کنترل تغییرات عملکرد کلاس‌بندی رخداد رانش مفهوم را شناسایی می‌کند. در این روش آموزش و آزمون و استقرار کلاس‌بندی‌ها به صورت متناوب انجام شده و امکان اینکه این روش در جریان داده‌ای که در آن داده‌هایی با کلاس معین و نامعین در کنار هم وجود دارد جوابگو باشد؛ را ممکن می‌کند.

۳،۳- روش‌های نیمه نظارتی

اهمیت توانایی فعالیت به صورت بدون نظارت باعث شده است تا بسیاری از روش‌ها با ترکیب روش‌های نظارتی با روش‌های بدون نظارت رانش مفهوم را شناسایی کنند. این دسته از روش‌ها در دسته نیمه نظارتی در این مطالعه قرار گرفته‌اند. روش‌های نیمه نظارتی روش‌هایی هستند که به تعداد کمتری نمونه با کلاس حقیقی نیازمند هستند. این روش‌ها در شرایطی که کلاس حقیقی همه داده‌ها در دسترس نیست جوابگو خواهند بود.

روش knowledge-maximized ensemble (KME) که در (Ren, Liao et al. (2018) مطرح شده است از همین تکنیک استفاده می‌کند. در این روش از داده‌هایی با کلاس حقیقی معین و داده‌هایی بدون کلاس به صورت همزمان استفاده می‌شود تا دانش مربوط به مفهوم کنونی بیشتر شود. در این روش هم عملکرد کلاس‌بندی‌کننده و هم تغییرات توزیع نمونه‌ها کنترل می‌شود. از اطلاعات بدست آمده از داده‌های با کلاس حقیقی و بدون کلاس حقیقی، وزن اعضای مدل دسته‌جمعی محاسبه و کنترل می‌شود. با اینکه این روش از داده‌های بدون کلاس حقیقی نیز استفاده می‌کند اما بیشتر مبتنی بر داده‌های با کلاس معین است و می‌توان آن را بیشتر روشی نظارتی تا بدون نظارت در نظر گرفت.

در مقاله (Lughofer, Weigl et al. (2016) روشی با ترکیب دو تکنیک نیمه نظارتی و بدون نظارت برای شناسایی رانش مفهوم ارائه شده است. روش اول که یک روش نیمه نظارتی است با استفاده از فیلترهای یادگیری فعال داده‌هایی که برای کلاس‌بندی نظارتی مناسب‌تر هستند را شناسایی می‌کند. این داده‌ها، داده‌هایی هستند که کلاس‌بندی‌کننده کمترین میزان اطمینان را برای کلاس‌بندی آنها دارد و معمولا داده‌های نزدیک مرزهای کلاس‌بندی هستند. روش دوم براساس همپوشانی توزیع قطعیت دو کلاس مختلف پیش‌بینی کلاس‌بندی‌کننده کار می‌کند. در هر دو روش برای شناسایی رانش از آزمون PHT آنلاین و بهبود یافته (Page-Hinkley test) استفاده میشود. در روش اول براساس معنی دارترین داده‌ها خطای پیش‌بینی کلاس‌بندی‌کننده با این آزمون تخمین زده می‌شود. در روش دوم فرض بر این است که توزیع قطعیت کلاس‌های مختلف جدا از هم هستند و به صورت توزیع‌های گوسی مجزا فرض می‌شوند. تغییرات همپوشانی این کلاس‌ها به عنوان متغیر مورد بررسی رانش در آزمون PHT استفاده می‌گردد. با اینکه روش بدون نظارت تنها قادر به شناسایی رانش‌های درون فضای ورودی است، با استفاده در کنار روش نیمه نظارتی مطرح شده، تفکیک رانش‌های فضای ورودی و رانش‌های هدف از هم را فراهم می‌کند.

در (Kim and Park (2017) نیز روشی برای شناسایی رانش مفهوم با تعداد محدود نمونه با کلاس حقیقی و بدون کلاس حقیقی ارائه کرده است. این روش براساس تخمین احتمالی نتیجه کلاس‌بندی رانش مفهوم را شناسایی می‌کند. در این روش نمونه‌های تصادفی انتخاب شده و بردار اطمینان کلاس‌بندی‌کننده محاسبه می‌شود. با استفاده از پنجره، تغییرات اطمینان بررسی می‌شود و از این طریق

رخداد رانش شناسایی می‌شود. برای حالت بدون نمونه با کلاس معین، نگارنده پیشنهاد می‌دهد که از میانگین کی (k-means) برای ایجاد ساختار خوشه‌ای استفاده شود تا به کلاس‌بندی‌کننده مجازی دست پیدا کنیم. ادامه روند مشابه حالت نیمه نظارتی مطرح شده است. در (Yu, Wang et al. (2018) روشی نیمه نظارتی بر اساس چند آزمون فرضیه سری ارائه شده است. در این مقاله دو روش Hierarchical Hypothesis Testing و Hierarchical Hypothesis Testing with Classification Uncertainty (HHT-CU) و HHT-AG (HHT-AG) "Goodness-of-fit" with Attribute-wise برای شناسایی رانش مفهوم پیشنهاد شده است. HHT-CU در لایه تشخیص از نابرابری هافدینگ که تغییرات عدم قطعیت کلاس‌بندی‌کننده را بررسی می‌کند، استفاده می‌کند و برای لایه تصدیق هم از آزمون جایگشتی برای تایید رخداد رانش، استفاده می‌شود. HHT-AG در لایه تشخیص از آزمون Kolmogorov-Smirnov روی توزیع هر یک از ویژگی‌ها رخداد رانش را بررسی می‌کند. در لایه تصدیق کلاس حقیقی داده‌ها برای رد یا تایید رخداد رانش درخواست می‌شود. در (Hosseini, Gholipour et al. (2016) یک کلاس‌بندی‌کننده نیمه نظارتی دسته‌جمعی برای شناسایی رانش مفهوم ارائه شده است. در این روش، هر کلاس‌بندی‌کننده یک مفهوم منفرد را در خود دارد. وزن کلاس‌بندی‌کننده‌ها بر اساس بسته داده‌هایی که فقط تعدادی از آنها کلاس حقیقی دارند تعیین شده و از همین داده‌ها برای آپدیت کردن کلاس‌بندی‌کننده‌ها استفاده می‌شود. روش (Haque, Khan et al. (2015) با نام SAND (Semi-supervised Adaptive Novel Class Detection and Classification over Data Stream) از دو ماژول شناسایی رانش مفهوم و شناسایی داده پرت (Outlier detection) تشکیل شده است. این روش هم توانایی شناسایی رانش مفهوم را داشته و هم می‌تواند با جمع‌آوری داده‌هایی که در مدل کلاس‌بندی‌کننده دسته‌جمعی کلاس آن‌ها تعیین نمی‌شود؛ کلاس جدید موجود را به مجموعه اضافه کند. این روش بر اساس کلاسیک ترین (k-nearest neighbor) کلاس‌بندی‌کننده دسته‌جمعی، کلاس نمونه‌ها را پیش‌بینی کرده و میزان اطمینان هر کلاس‌بندی‌کننده را تخمین می‌زند. با کنترل تغییرات اطمینان، این روش قادر است رانش مفهوم را شناسایی کند. بعد از شناسایی رانش SAND فقط کلاس حقیقی داده‌هایی را درخواست می‌کند که اطمینان خیلی پایینی دارند و باقی داده‌ها بر اساس همان مدل پیش‌بینی کلاس‌بندی می‌شوند. با وجود اینکه روش SAND از تعداد نمونه با کلاس حقیقی کمتری استفاده می‌کند و درستی پیش‌بینی خوبی هم دارد، اما از نظر هزینه محاسباتی و زمان محاسبه به علت حجم بالای محاسبات کاملاً کاربردی نبوده است. از همین رو، روش SAND در (Haque, Khan et al. (2016) بهبود پیدا کرده است. روش جدید که با نام Efficient Concept Drift and Concept Evolution Handling over Stream Data (ECHO) شناخته می‌شود. این روش عملکردی مانند روش SAND دارد با این تفاوت که ماژول شناسایی رانش مفهوم حال دیگر در هر بار محاسبه میزان آستانه اطمینان مدل دسته‌جمعی فعال نشده و فراخوانی آن وابسته به میزان اطمینان مجموعه تست خواهد بود.

۳-۴- روش‌های بدون نظارت

روش‌های بدون نظارت شناسایی رانش مفهوم مناسب شرایطی هستند که بازخورد نتایج پیش‌بینی کلاس‌های داده با تاخیر همراه است. این روش‌ها با بررسی تغییرات ابعاد داده رانش مفهوم را شناسایی می‌کنند. در مقاله (Liu, Lu et al. (2018) روشی برای شناسایی رانش‌های منطقه‌ای به صورت بدون نظارت مطرح شده است. طبق گفته نگارنده، بیشتر روش‌های شناسایی رانش مفهوم در تلاش برای شناسایی رانش‌های عمومی هستند و رانش‌هایی که به علت تغییرات چگالی منطقه‌ای به وجود می‌آیند را نمی‌توانند شناسایی کنند. روش Nearest Neighbor-based Density Variation Identification (NN-DVI) یک روش مبتنی بر تخمین چگالی منطقه‌ای است و از سه ماژول جزیندی فضا بر اساس کلاسیک ترین همسایه (kNN-based space partitioning)، تابع فاصله میزان اختلاف چگالی و آزمون معنی‌داری آماری تشکیل شده است. این روش با مقایسه بین پنجره داده‌های جدید و پنجره داده‌های قدیمی تر کار می‌کند. در ماژول اول بر اساس همسایگی هر داده x ، آن داده و همسایه‌های آن تشکیل یک ذره نمونه (instance particle) را می‌دهند. بعد گروه‌های ذره نمونه بر اساس همسایگی مشترک ایجاد شده و در نتیجه هر داده با مجموعه‌ای از ذره نمونه‌های دارای اشتراک نمایش داده می‌شود. این باعث می‌شود که هر ذره نمونه مشترک جزئیات بیشتری از نمونه را معرفی کرده و تغییرات جزئی نیز بهتر مشاهده شوند. در ماژول دوم تفاوت میان داده‌ها بر اساس مجموع

تفاوت تعداد ذره نمونه آنها محاسبه میشود. در ماژول سوم آزمونی برای تایید رانش انجام میگردد. نویسنده مقاله فرض می‌کند که مقادیر تفاوت‌های نمونه‌های مورد مقایسه در یک توزیع نرمال قرار می‌گیرند؛ از این رو از برآورد درست‌نمایی بیشینه (max likelihood) برای تخمین میانگین و واریانس توزیع تفاوت نمونه‌های مورد بررسی استفاده می‌کند. سپس تابع توزیع انباشته (cumulative distribution function) با مقادیر آستانه تعریف شده توسط کاربر مورد بررسی قرار میگیرد و رخدادهای رانش مفهوم تایید می‌گردد.

در (Bashir, Petrovski et al. (2017) روشی دو مرحله‌ای با نام Unsupervised Change Detection for Activity Recognition برای شناسایی رانش مفهوم به صورت بدون نظارت ارائه شده است. این روش مبتنی بر این فرض است که رانش مفهوم به شکل تفاوت زیاد بین داده‌های جدید منسوب به یک کلاس با نمونه‌های آموزشی آن کلاس، ظاهر می‌شود. در مرحله اول این روش یک کلاس بندی کننده تمرین میدهد و ویژگی هر یک از کلاس‌ها را تعیین می‌کند. در مرحله دوم داده‌ها به صورت بسته‌ای دریافت شده و بررسی می‌شوند تا تفاوت آنها با داده‌های پیشین تعیین شود. با اینکه این روش در این مقاله برای شناسایی رفتار فیزیکی نظیر دویدن و نشستن و... مطرح شده است اما تأکیدی بر انحصار آن برای شناسایی حرکت بیان نشده است.

در مقاله (Sethi and Kantardzic (2015) روشی با نام Margin Density Drift Detection مطرح شد که بر اساس کنترل تغییرات در مرز اختلاف کلاس بندی کننده (classifier's margin) رانش مفهوم را شناسایی می‌کند. مرز اختلاف به این صورت تعریف میشود: ناحیه‌ای است در نزدیکی به کرانه‌های مدل و پیش‌بینی‌ها برای این ناحیه‌ها قطعیت بسیار پایینی دارد. فرض بر این است که تغییرات چگالی این ناحیه‌ها بیشتر از آستانه تعریف شده نشان از رخداد رانش است. با استفاده از تعداد نمونه‌هایی که در این ناحیه‌ها قرار دارند چگالی مرز اختلاف محاسبه می‌شود. با استفاده از دیتاست اولیه و روش SVM کلاس بندی کننده اولیه تولید شده و با توجه به اینکه از روش SVM استفاده می‌شود، مرزهای اختلاف طبقه بندی کننده مشخص خواهد شد. در مرحله بعد بیشترین و کمترین مقادیر چگالی بعد از آخرین شناسایی رانش مفهوم محاسبه شده و تحت نظارت قرار میگیرد و برای هر دسته داده جدید که به سیستم می‌رسد چگالی آنها محاسبه شده و با مقایسه آن با میزان آستانه تعریف شده می‌توان رخداد رانش را شناسایی کرد. بعد از اینکه رخداد رانش شناسایی شد مدل بر اساس نمونه‌هایی که کلاس حقیقی آنها مشخص هستند آپدیت میشود. در این مرحله به کلاس حقیقی داده‌ها نیازمند هستیم. یکی از مشکلات اصلی این روش نیازمندی به روش یادگیری SVM است. در سال 2017 در مقاله (Sethi and Kantardzic (2017) روش پیشنهادی برای اینکه مشکلات پیشین را نداشته باشد، بهبود یافت. روش جدید با نام MD3 Ensemble Generic Margin (MD3-EGM) شناخته میشود. در این روش با استفاده از کلاس بندی کننده دسته‌جمعی مرز اختلاف عمومی شناسایی می‌شود که مجموعه‌ای از نمونه‌های ناشناخته است که بیشترین عدم قطعیت را برای پیش‌بینی دارند. این عدم قطعیت، بر اساس میزان اختلاف کلاس بندی کننده‌های مدل دست‌جمعی برای پیش‌بینی کلاس یک داده بدست می‌آید. شناسایی رانش نیز بر اساس همین افزایش میزان اختلاف برای کلاس بندی یک داده ناشناخته صورت می‌گیرد. از آنجایی که نرخ هشدار اشتباه این روش می‌تواند زیاد باشد؛ یک ماژول تایید رانش نظارتی به این روش اضافه میشود و آن را به روشی نیمه نظارتی تبدیل می‌کند.

در (Costa, Albuquerque et al. (2018) روشی برای شناسایی رانش مفهوم با فرض اینکه تغییر چگالی معنی دارترین داده‌ها می‌تواند نشانه رانش مفهوم باشد، مطرح شد. در این روش بر اساس اطمینان کلاس بندی کننده برای داده‌ها حدود مجازی تعیین می‌شود. نمونه‌هایی که اطمینان آنها از میزان آستانه کمتر باشد طبق یادگیری فعال (Active Learning) برای محاسبه چگالی انتخاب می‌شوند. میزان چگالی داده‌های جدید با میزان‌های بیشینه و کمینه موجود چگالی مقایسه شده و در صورتی که از میزان آستانه تایید شده بیشتر باشد؛ رخداد رانش مفهوم شناسایی می‌شود.

در (dos Reis, Flach et al. (2016) روشی بدون نظارت با بکارگیری آزمون کولموگوروف-اسمیرنوف (Incremental Kolmogorov-Smirnov) مطرح شد. ایده اصلی این روش، اعمال این آزمون برای هر یک از ویژگی‌های داده به صورت مجزا است. شناسایی تغییر در یک ویژگی ممکن است که برای شناسایی رانش مفهوم کافی باشد. از ضعف‌های این روش این است که قادر نخواهد بود رانش‌های مفهومی را شناسایی کند که نیازمند بررسی ویژگی‌های نمونه‌ها در کنار هم هست.

در (Mustafa, Ayoade et al. (2017) یک روش برای شناسایی کلاس جدید در داده با نام Nonparametric Multidimensional Drift Detection Method (NM-DDM) ارائه شد که از denoising autoencoders که یکی از الگوریتم‌های شبکه‌های عصبی است؛ استفاده میکند. در این مقاله یک روش برای شناسایی نقطه تغییر و رانش مفهوم ارائه شد که در اینجا مورد نظر ماست. روش شناسایی رانش مفهوم بر اساس محاسبه نسبت درست‌نمایی (log-likelihood ratio) بین دو پنجره، قبل از نقطه احتمالی رانش و بعد از آن، کار میکند. نسبت آزمون درست‌نمایی برای هر یک از ویژگی‌های داده (ابعاد) تخمین زده شده و در صورتی که بیشترین نسبت محاسبه شده میان دو پنجره، بیشتر از مقدار آستانه باشد، رخداد رانش تایید شده و داده‌های پنجره جدیدتر برای یادگیری مورد استفاده قرار می‌گیرد. برای مرحله آپدیت مدل تنها کلاس واقعی داده‌هایی مورد نیاز است که میزان اطمینان آنها در کلاس‌بندی‌کننده پایین بوده است و باقی داده‌ها با همان کلاس پیش‌بینی شده مورد بررسی قرار می‌گیرند.

روش Plover است که در (de Mello, Vaz et al. (2019) مطرح می‌شود که بر اساس معیارهای آماری مختلف نظیر میانگین، واریانس، ... و طیف توانی (power spectrum) رانش مفهوم را شناسایی می‌کند. نگارنده در این مقاله تحلیلی بر مبانی نظری و تضمین‌های یادگیری در شناسایی رانش مفهوم بدون نظارت انجام داده است. آن‌ها نتیجه‌گیری کرده‌اند؛ درحالی که روش‌های نظارتی با نظریه‌های یادگیری آماری پشتیبانی می‌شوند؛ روش‌های بدون نظارت وابسته به معیارهایی درونی و بیرونی هستند که چنین پشتوانه‌ای ندارند. معیارهای درونی نظیر فشردگی خوشه و فاصله گرانیگاه که از خود ساختار داده قابل محاسبه بوده، هستند. در مقابل معیارهای بیرونی وابسته به کلاس حقیقی داده است که در بسیاری از شرایط برای روش‌های بدون نظارت ممکن نخواهد بود.

در (Pinagé, dos Santos et al. (2020) یکی دیگر از روش‌های بدون نظارت معرفی شد. DSDD یا Dynamic Selection Drift Detector روشی است که به صورت پویا از مجموعه کلاس‌بندی‌کننده‌های دسته‌جمعی، کلاس‌بندی‌کننده مناسب هر داده را انتخاب کرده و با کنترل شبه‌خطا (Pseudo Error) می‌تواند رانش مفهوم را شناسایی کند و از سه مازول مختلف تشکیل شده است. مازول‌های ساخت مدل دست‌جمعی (Ensemble Creation)، انتخاب پویا (Dynamic Selection) و شناسایی (Detection). در مازول اول با استفاده از دیتاست اولیه با روش Online Bagging بهبود یافته یک مدل دسته‌جمعی از کلاس‌بندی‌کننده‌ها ایجاد می‌شود و فرض بر این است که هر کدام از این کلاس‌بندی‌کننده‌ها در یک ناحیه از داده‌ها بهترین عملکرد را دارد. مازول دوم بدین صورت کار می‌کند که وقتی یک داده جدید و بدون کلاس معلوم به مازول وارد میشود، براساس اینکه کدام کلاس‌بندی‌کننده بهترین عملکرد را برای داده‌های همسایه آن داده در دیتاست اعتبارسنجی بر اساس (k nearest neighbor) داشته، آن کلاس‌بندی‌کننده برای پیش‌بینی کلاس آن داده انتخاب میشود. کلاسی که کلاس‌بندی‌کننده منتخب در مازول دوم برای داده پیش‌بینی می‌کند به عنوان کلاس شبه درست آن داده تایید میشود (pseudo true). در مازول بعدی برای هر یک از کلاس‌بندی‌کننده‌های مدل دسته‌جمعی یکی از روش‌های تشخیص رانش با نظارت مانند DDM یا EDDM مورد استفاده قرار می‌گیرد. با اینکه این روش‌ها برای کارکرد نیازمند کلاس حقیقی داده هستند اما در اینجا با توجه به کلاس شبه درست داده بررسی‌ها صورت می‌گیرد و در صورتی که تعداد مشخصی از کلاس‌بندی‌کننده‌ها تشخیص رانش دهند؛ رخداد رانش مفهوم تایید شده و دیتاست اعتبارسنجی با توجه به داده‌های جدید آپدیت شده و مقادیر نرخ خطا و انحراف معیار مورد استفاده در DDM ریست می‌شود. در این مقاله از DDM و EDDM استفاده شد اما می‌توان از هر کدام از روش‌های مبتنی بر نرخ خطا برای مازول تشخیص رانش استفاده کرد.

با توجه به اهمیت کارکرد این روش‌ها در حالتی که کلاس حقیقی داده‌ها با تاخیر همراه است؛ بسیاری از روش‌ها از تکنیک‌های مختلفی در کنار هم استفاده می‌کنند که فاصله‌گذاری روش‌ها از یکدیگر را سخت می‌کند. در این قسمت سعی شد تا با توجه به میزان وابستگی به کلاس حقیقی داده‌ها، روش‌ها به این سه گروه تقسیم شوند.

۴- نتیجه‌گیری

در این گزارش سعی شد دسته‌بندی نظارتی، بدون نظارت و نیمه نظارتی که در مقاله Khamassi, Sayed-Mouchaweh et al. (2018) مطرح شده بود با بررسی روش‌های مطرح شده بعد از سال ۲۰۱۵ به روز شود. این تقسیم‌بندی، با توجه به اندازه نیاز این روش-

ها به کلاس حقیقی داده‌ها برای کار کرد، انجام شده است. در بیشتر موقعیت‌های دنیای واقعی که کلاس حقیقی داده‌ها یا در دسترس نیست یا با تاخیر در دسترس است، روش‌های نیمه نظارتی و بدون نظارت توانایی بیشتری برای شناسایی رانش مفهوم دارند. با بررسی روش‌های نظارتی، می‌توان دید که این روش‌ها توانایی شناسایی رانش مفهوم با دقتی بالا را داشته و می‌توان بیشتر گونه‌های رخداد رانش مفهوم را با آن‌ها تشخیص داد، در مقابل فراهم کردن کلاس حقیقی همه داده‌ها هزینه سنگینی است که این روش‌ها برای رسیدن به این دقت باید بپردازند. روش‌های بدون نظارت و نیمه نظارتی، با کار کردن در شرایطی که کلاس حقیقی داده‌ها در دسترس نیست، امکان شناسایی رانش مفهوم در شرایط دنیای واقعی را فراهم می‌کنند اما در مقابل توانایی شناسایی رانش مفهوم آن‌ها پایین‌تر بوده و هر کدام توانایی شناسایی دامنه محدودی از انواع شکل‌های رخداد رانش مفهوم را دارند. بر همین اساس، توازن بین عملکرد و هزینه بین روش‌های مختلف شناسایی رانش مفهوم مشاهده می‌شود که انتخاب بهترین روش را مشکل خواهد کرد. از دیگر مشکلات می‌توان به نبود معیار و مقیاس مطرحی برای انتخاب دیتاست و بررسی عملکرد روش‌های شناسایی رانش مفهوم در این دیتاست‌ها اشاره کرد. با ایجاد این چارچوب، بررسی عملکرد روش‌ها آسان‌تر خواهد بود. معیارهای بررسی هر روش مانند دقت کلاس‌بندی به تنهایی کافی نیستند. نیاز به معیارهای بیشتر برای سنجش میزان عملکرد روش‌های شناسایی از دیگر مشکلات پیرامون این بحث است. در آخر، معرفی و ابداع روش‌هایی که توانایی شناسایی رانش مفهوم با توجه به ناهنجاری‌های داده نظیر وابستگی موقت (temporal dependence)، عدم تعادل کلاس (Class imbalance) و جریان‌های داده با مقادیر از دست رفته (missing value) به شدت می‌تواند به شناسایی رانش مفهوم و مطالعه بهتر داده‌ها کمک کند.

روش‌های شناسایی رانش که در مدل دسته‌جمعی استفاده می‌شود، تکنیکی است که شانس شناسایی انواع رانش مفهوم را بیشتر می‌کند به خصوص برای روش‌های بدون نظارت که ضعف آن‌ها در شناسایی انواع گونه‌های رانش مفهوم است. بررسی مجزا این روش‌ها خود می‌تواند موضوعی برای مقاله‌ای مروری در حوزه شناسایی رانش مفهوم باشد.

۵-مراجع

- [1]Khamassi, I., M. Sayed-Mouchaweh, M. Hammami and K. Ghédira (2018). "Discussion and review on evolving data streams and concept drift adapting." *Evolving systems* 9(1): 1-23.
- [2]Lu, J., A. Liu, F. Dong, F. Gu, J. Gama and G. Zhang (2018). "Learning under concept drift: A review." *IEEE Transactions on Knowledge and Data Engineering* 31(12): 2346-2363.
- [3]Hu, H., M. Kantardzic and T. S. Sethi (2020). "No Free Lunch Theorem for concept drift detection in streaming data classification: A review." *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 10(2): e1327.
- [4]Wares, S., J. Isaacs and E. Elyan (2019). "Data stream mining: methods and challenges for handling concept drift." *SN Applied Sciences* 1(11): 1412.
- [5]Gemaque, R. N., A. F. J. Costa, R. Giusti and E. M. Dos Santos (2020). "An overview of unsupervised drift detection methods." *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 10(6): e1381.
- [6]Liu, A., G. Zhang and J. Lu (2017). Fuzzy time windowing for gradual concept drift adaptation. 2017 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE), IEEE.
- [7]Xu, S. and J. Wang (2017). "Dynamic extreme learning machine for data stream classification." *Neurocomputing* 238: 433-449.
- [8]Gama, J., P. Medas, G. Castillo and P. Rodrigues (2004). *Learning with drift detection*. Brazilian symposium on artificial intelligence, Springer.
- [9]Wang, H. and Z. Abraham (2015). Concept drift detection for streaming data. 2015 International Joint Conference on Neural Networks (IJCNN), IEEE.
- [10]Zhang, Y., G. Chu, P. Li, X. Hu and X. Wu (2017). "Three-layer concept drifting detection in text data streams." *Neurocomputing* 260: 393-403.
- [11]Alippi, C., G. Boracchi and M. Roveri (2016). "Hierarchical change-detection tests." *IEEE transactions on neural networks and learning systems* 28(2): 246-258.
- [12]Yu, S. and Z. Abraham (2017). Concept drift detection with hierarchical hypothesis testing. Proceedings of the 2017 SIAM International Conference on Data Mining, SIAM.
- [13]de Lima Cabral, D. R. and R. S. M. de Barros (2018). "Concept drift detection based on Fisher's Exact test." *Information Sciences* 442: 220-234.
- [14]Nishida, K. and K. Yamauchi (2007). Detecting concept drift using statistical testing. International conference on discovery science, Springer.

- [15]Pesaranghader, A., H. L. Viktor and E. Paquet (2018). McDiarmid drift detection methods for evolving data streams. 2018 International Joint Conference on Neural Networks (IJCNN), IEEE.
- [16]Pesaranghader, A. and H. L. Viktor (2016). Fast hoeffding drift detection method for evolving data streams. Joint European conference on machine learning and knowledge discovery in databases, Springer.
- [17]Wang, X., Q. Kang, M. Zhou and S. Yao (2018). A multiscale concept drift detection method for learning from data streams. 2018 IEEE 14th International Conference on Automation Science and Engineering (CASE), IEEE.
- [18]Dehghan, M., H. Beigy and P. ZareMoodi (2016). "A novel concept drift detection method in data streams using ensemble classifiers." *Intelligent Data Analysis* **20**(6): 1329-1350.
- [19]Sidhu, P. and M. Bhatia (2018). "A novel online ensemble approach to handle concept drifting data streams: diversified dynamic weighted majority." *International Journal of Machine Learning and Cybernetics* **9**(1): 37-61.
- [20]Khamassi, I., M. Sayed-Mouchaweh, M. Hammami and K. Ghédira (2019). A new combination of diversity techniques in ensemble classifiers for handling complex concept drift. *Learning from data streams in evolving environments*, Springer: 39-61
- [21]Khamassi, I., M. Sayed-Mouchaweh, M. Hammami and K. Ghédira (2015). "Self-adaptive windowing approach for handling complex concept drift." *Cognitive Computation* **7**(6): 772-790.
- [22]Mahdi, O. A., E. Pardede and J. Cao (2018). Combination of information entropy and ensemble classification for detecting concept drift in data stream. *Proceedings of the Australasian Computer Science Week Multiconference*.
- [23]Krawczyk, B. and A. Cano (2018). "Online ensemble learning with abstaining classifiers for drifting and noisy data streams." *Applied Soft Computing* **68**: 677-692.
- [24]Jaworski, M., P. Duda and L. Rutkowski (2017). On applying the restricted Boltzmann machine to active concept drift detection. 2017 IEEE Symposium Series on Computational Intelligence (SSCI), IEEE.
- [25]Lobo, J. L., J. Del Ser, I. Laña, M. N. Bilbao and N. Kasabov (2018). Drift detection over non-stationary data streams using evolving spiking neural networks. *International symposium on intelligent and distributed computing*, Springer.
- [26]Astudillo, C. A., J. I. González, B. J. Oommen and A. Yazidi (2016). Concept drift detection using online histogram-based bayesian classifiers. *Australasian Joint Conference on Artificial Intelligence*, Springer.
- [27]Ren, S., B. Liao, W. Zhu and K. Li (2017). "8Knowledge-maximized ensemble algorithm for different types of concept drift." *Information Sciences* **430**: 261-281.
- [28]Lughofer, E., E. Weigl, W. Heidl, C. Eitzinger and T. Radauer (2016). "Recognizing input space and target concept drifts in data streams with scarcely labeled and unlabelled instances." *Information Sciences* **355**: 127-151.
- [29]Kim, Y. and C. H. Park (2017). "An efficient concept drift detection method for streaming data under limited labeling." *IEICE Transactions on Information and systems* **100**(10): 2537-2546.
- [30]Yu, S., X. Wang and J. C. Principe (2018). "Request-and-reverify: Hierarchical hypothesis testing for concept drift detection with expensive labels." *arXiv preprint arXiv:1806.10131*.
- [31]Hosseini, M. J., A. Gholipour and H. Beigy (2016). "An ensemble of cluster-based classifiers for semi-supervised classification of non-stationary data streams." *Knowledge and information systems* **46**(3): 567-597.
- [32]Haque, A., L. Khan and M. Baron (2015). Semi supervised adaptive framework for classifying evolving data stream. *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Springer.
- [33]Haque, A., L. Khan, M. Baron, B. Thuraisingham and C. Aggarwal (2016). Efficient handling of concept drift and concept evolution over stream data 2016 .IEEE 32nd International Conference on Data Engineering (ICDE), IEEE.
- [34]Liu, A., J. Lu, F. Liu and G. Zhang (2018). "Accumulating regional density dissimilarity for concept drift detection in data streams." *Pattern Recognition* **76**: 256-272.
- [35]Bashir, S. A., A. Petrovski and D. Doolan (2017). "A framework for unsupervised change detection in activity recognition." *International Journal of Pervasive Computing and Communications*.
- [36]Sethi, T. S. and M. M. Kantardzic (2015). Don't Pay for Validation :Detecting Drifts from Unlabeled data Using Margin Density. *INNS Conference on Big Data*.
- [37]Sethi, T. S. and M. Kantardzic (2017). "On the reliable detection of concept drift from streaming unlabeled data." *Expert Systems with Applications* **82**: 77-99.
- [38]Costa, A. F. J., R. A. Albuquerque, n. Saraiva and E. M. dos Santos (2018). A drift detection method based on active learning. 2018 International Joint Conference on Neural Networks (IJCNN), IEEE.
- [39]dos Reis, D. M., P. Flach, S. Matwin and G. Batista .(2016)Fast unsupervised online drift detection using incremental kolmogorov-smirnov test. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [40]Mustafa, A. M., G. Ayoade, K. Al-Naami, L. Khan, K. W. Hamlen, B. Thuraisingham and F. Araujo (2017). Unsupervised deep embedding for novel class detection over data stream. 2017 IEEE International Conference on Big Data (Big Data), IEEE.
- [41]de Mello, R. F., Y. Vaz, C. H. Grossi and A. Bifet (2019). "On learning guarantees to unsupervised concept drift detection on data streams." *Expert Systems with Applications* **117**: 90-102.
- [42]Pinagé, F., E. M. dos Santos and J. Gama (2020). "A drift detection method based on dynamic classifier selection." *Data Mining and Knowledge Discovery* **34**(1): 50-74.