

مروری بر روش‌های مواجهه با تغییرات زمانی در دسته‌بندی متن

مرضیه سپهر^۱، مرتضی براری^۲، سمیه کافی^۳

^۱ دانشجوی کارشناسی ارشد مهندسی کامپیوتر، دانشگاه صنعتی مالک اشتر

^۲ استادیار مجتمع فناوری اطلاعات و ارتباطات، دانشگاه صنعتی مالک‌اشتر

^۳ دانشجوی دکترای مهندسی کامپیوتر، دانشگاه صنعتی مالک‌اشتر

مسئول مکاتبات: مرضیه سپهر

چکیده

رشد سریع اطلاعات به خصوص رشد سریع وب، دسته‌بندی متون را به صورت دستی غیرممکن می‌کند. از این جهت دسته‌بندی خودکار متن به عنوان یکی از زمینه‌های تحقیقاتی مهم مورد توجه محققین قرار گرفته است. در دسته بندی متن عموماً داده‌ها به صورت آفلاین جمع‌آوری شده و آموزش بر روی این آن‌ها انجام می‌شود سپس مدل آموزش دیده به این روش، می‌تواند برای پیش‌بینی داده‌های دیده نشده‌ی ورودی مورد استفاده قرار بگیرد. اما در دنیای واقعی متون مدام در حال به روز رسانی هستند و توزیع داده‌ها و مفاهیم متنی به مرور زمان دستخوش تغییر می‌شوند به طوری که موجب اثرگذاری بر روی عملکرد دسته بند می‌شوند. در این تحقیق، ابتدا به دسته بندی انواع تغییرات ممکن در طول زمان، روشهای تشخیص آن و سپس انواع روش‌های مواجهه با تغییرات زمانی در دسته‌بندی متون پرداخته شده است و سپس به مقایسه‌ی روش‌ها با توجه به نوع تغییرات پرداخته شده است.

کلمات کلیدی: دسته بندی متن، زمان در دسته‌بندی متن، اثرات زمانی

۱. مقدمه

رشد سریع اطلاعات به خصوص رشد سریع وب، دسته‌بندی متون موجود را به صورت دستی غیرممکن می‌کند. از این جهت دسته‌بندی خودکار متن به عنوان یکی از زمینه‌های تحقیقاتی مهم مورد توجه محققین قرار گرفته است. هدف از دسته‌بندی خودکار متن اختصاص خودکار متون، به دسته‌های از پیش تعیین‌شده، با توجه به محتوای آن‌هاست (Sebastiani ۲۰۰۲). در دسته بندی متن عموماً داده‌ها ابتدا به صورت آفلاین جمع‌آوری شده و آموزش بر روی آن‌ها انجام می‌شود سپس مدل آموزش دیده به این روش، می‌تواند برای پیش‌بینی داده‌های دیده نشده‌ی ورودی مورد استفاده قرار بگیرد. در شرایطی که داده‌ها مدام در حال به روز رسانی هستند، توزیع آنها می‌تواند در طول زمان تغییر کند به طوری که موجب اثرگذاری بر روی عملکرد دسته بند شود. لذا مفهوم زمان را می‌توان به عنوان یک بعد از اطلاعات در نظر گرفت.

دو تعریفی که مرجع (Alonso, Gertz et al. ۲۰۰۷) از مفهوم زمان در بازیابی اطلاعات دارد بدین شرح است: (۱) یک فاصله که دو نقطه‌ی تاریخ را در یک زنجیره‌ی تاریخی از هم جدا می‌کند. (۲) یک زنجیره که در آن وقایع در یک توالی تغییر ناپذیر از گذشته، حال و تا آینده رخ می‌دهد. در تعریف دوم منظور از واقعه یا رخداد^۱ یک اتفاق است که به اندازه‌ی کفایت رخ می‌دهد به نحوی که در داده‌ها تاثیر گذار باشد، مثل یک گردهمایی یا فعالیت اجتماعی مانند انتخابات. بر این اساس، زمان، یک پارامتر تاثیر گذار در کاربردهای بازیابی اطلاعات به حساب آمده و تغییرات و نحوه‌ی مواجهه با آن در تحقیقات به عناوین مختلف مورد بررسی قرار گرفته است.

دسته‌ای از تحقیقات از این تغییرات با عنوان رانش مفهوم^۲ (Forman ۲۰۰۶, Nishida and Yamauchi ۲۰۰۹, Šilić and Bašić ۲۰۱۲, D'hondt, Verberne et al. ۲۰۱۴) و دسته ای دیگر با تغییر زمانی^۳ (Mourão, Rocha et al. ۲۰۰۸, Rocha, Mourão et al. ۲۰۰۸, Rocha, Mourão et al. ۲۰۱۰, Salles, Rocha et al. ۲۰۱۳) یاد کرده اند. بر اساس تحقیقات انجام شده در این مقاله، تفاوت اصلی بین دسته‌ی اول و دسته‌ی دوم مشاهده شد. فرض اصلی دسته‌ی اول این است که داده‌ها به ترتیب زمانی ذخیره نشده‌اند و داده‌های آزمایشی همیشه جدیدتر از داده‌های آموزشی هستند (Rocha, Mourão et al. ۲۰۱۳) بر خلاف دسته‌ی دوم که فرض می‌کند داده‌های آزمایشی می‌توانند از هر بازه‌ی زمانی باشند. در تحقیقات دسته‌ی اول، داده‌ی مورد بررسی اغلب جریان داده‌ی متنی^۴ هستند. در حالی که تحقیقات دسته‌ی دوم اغلب بر روی داده‌های متنی معمولی کار

^۱ Event

^۲ Concept drift

^۳ Temporal variation

^۴ جریان داده دنباله ای پیوسته از داده‌ها است که به همان ترتیب دریافت مورد دسترسی قرار می‌گیرند.

کردند. تحقیقات دسته‌ی اول اغلب سعی در شناسایی دقیق زمان تغییر و نوع تغییر داشته و سپس برای مواجهه با تغییر روشی را پیشنهاد داده اند در حالی که در تحقیقات دسته‌ی دوم عموماً از روش‌های ابتکاری و یا پیشفرض‌ها برای تعیین زمان مواجهه با تغییر استفاده کرده‌اند.

در این تحقیق سعی شده تا حد امکان یک دید کلی و همه جانبه در این زمینه‌ی اثر زمان در دسته‌بندی خودکار متن و بررسی روش‌های مواجهه با آن ارائه شود. باید توجه داشته که پیش از بررسی روش‌های مواجهه با تغییر، ابتدا به سه سوال مهم باید پاسخ داده شود:

- اثرگذاری تغییرات زمانی به چه نحوی قابل مشاهده هستند؟
- تغییرات زمانی به چه تناوبی رخ می‌دهند؟
- چه زمانی باید نسبت به تغییرات عکس‌العمل نشان داد؟

بر این اساس، در تحقیق حاضر ابتدا به سه سوال مذکور به ترتیب در بخش‌های دوم، سوم و چهارم پاسخ داده می‌شود و سپس در بخش پنجم به بررسی روش‌های مواجهه با تغییر پرداخته خواهد شد. در بخش پنجم به مقایسه‌ی روش‌ها با توجه به نوع تغییرات پرداخته شده است و در نهایت نتیجه‌گیری ارائه می‌گردد.

۲. اثرات زمانی

از یک دیدگاه تغییرات زمانی را می‌توان بر اساس نحوه‌ی اثرگذاری آن به ۳ دسته تقسیم کرد (Rocha, Mourão et al. ۲۰۱۳, D'hondt, Verberne et al).

- تغییر در توزیع کلاسها^۱
- تغییر در توزیع کلمات^۲
- تغییر در شباهت کلاسها^۳

در دسته‌ی اول از این نوع تغییر، گذر زمان موجب تغییر در توزیع کلاس‌ها می‌شود. یعنی ممکن است دسته‌هایی از بین بروند یا دسته‌های جدیدی به وجود آیند. علت این مسئله تقسیم شدن یک دسته به دسته‌های کوچکتر و یا ادغام چند دسته با یکدیگر است. برای مثال در ساختار سلسله‌مراتبی موضوعات زیرمجموعه‌ی علوم کامپیوتر در پایگاه ACM، در سال ۱۹۶۴ دو موضوع بازیابی اطلاعات و هوش مصنوعی متعلق به یک کلاس برنامه‌های کاربردی بودند ولی در حال حاضر هر کدام به ترتیب متعلق به کلاس‌های سیستم‌های اطلاعاتی^۴ و روش‌های محاسباتی^۵ هستند.

گاهی اوقات با گذر زمان، کلمات ممکن است در دسته‌ی خاصی از متون وارد شوند، از یک دسته خارج شوند، از دسته‌ای به دسته‌ی دیگر مهاجرت کرده و یا اینکه وزن آن‌ها در یک دسته تغییر کند. در مرجع (Wenerstrom and Giraud-Carrier ۲۰۰۶) از این نوع تغییر با نام تغییر متنی^۶ نام برده شده است. انتقال کلماتی مانند انرژی هسته‌ای، غنی‌سازی، سانتیفیوژها و غیره از دامنه‌ی فیزیک به دامنه‌ی اخبار سیاسی در چند سال اخیر مثالی از مهاجرت کلمات از دسته‌ای به دسته‌ی دیگر است و ورود کلماتی مانند یارانه در دامنه‌ی اخبار سیاسی در زبان فارسی را می‌توان مثالی از ورود کلمات جدید به یک کلاس دانست.

تغییر در شباهت کلاس‌ها اشاره دارد به تکامل تدریجی شباهت بین کلاس‌ها در طول زمان. برای مثال در گذشته، دسته‌ی جرم‌شناسی و بیولوژی به هم شباهتی نداشتند. اما اخیراً به علت استفاده از علم بیولوژی در شناسایی مجرمین این دو دسته به هم ارتباط بیشتری پیدا کرده اند. لازم به ذکر است که تغییر در شباهت کلاس‌ها به معنی تغییر در توزیع کلمات بین دو کلاس متفاوت است در حالی که تغییر در توزیع کلمات به تغییر کلمات درون یک کلاس اشاره دارد که می‌تواند منجر به افزایش شباهت بین دو کلاس نیز بشود.

۳. انواع تغییر

از دیدگاهی دیگر می‌توان تغییرات زمانی را بر اساس دوره‌ی تناوب آن به ۳ دسته‌ی زیر، تقسیم کرد (D'hondt, Verberne et al. 2014):

- تغییر ناگهانی^۷
- تغییر تدریجی^۸
- تغییر بازگشت کننده^۹

¹ Class distribution

² Term distribution

³ Class similarity

⁴ Information system

⁵ Computing Methodologies

⁶ Contextual drift

⁷ Sudden drift

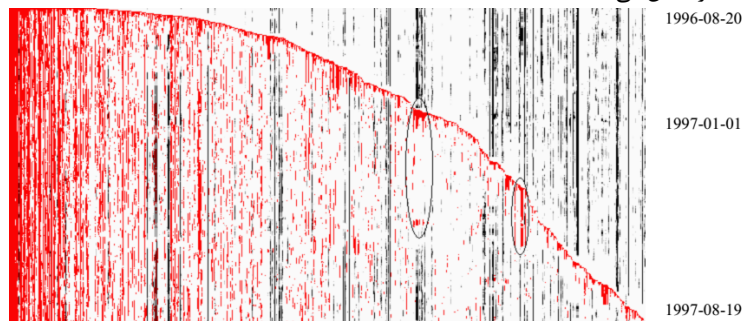
⁸ Gradual

⁹ Recurrent

در تغییرات ناگهانی لحظاتی وجود دارد که در آن توزیع داده در پیکره (توزیع کلاس یا کلمات) به صورت قابل ملاحظه‌ای تغییر می‌کند. این نوع تغییرات در سیستم‌هایی که اطلاعات را متناسب با علاقه‌ی کاربر فیلتر می‌کنند، رایج است. در چنین سیستم‌هایی تغییر ناگهانی علاقه‌ی کاربر موجب کاهش عملکرد سیستم می‌شود. در بعضی از تحقیقات از این نوع تغییر با نام تغییر جهت مفهوم^۱ یاد شده است (Nishida and Yamauchi ۲۰۰۹).

در تغییرات تدریجی نرخ تغییر پایین است. در ساده‌ترین حالت تغییر تدریجی، نمونه‌های آموزشی از ۲ توزیع P_j و P_{j+1} تشکیل شده‌اند و به مرور زمان، احتمال مشاهده‌ی مثال‌هایی با توزیع P_j کاهش و احتمال مشاهده‌ی مثال‌هایی با توزیع P_{j+1} افزایش می‌یابد. نوع دیگری از تغییر تدریجی، تغییر افزایشی^۲ است که در آن داده‌ها از بیش از دو نوع توزیع تشکیل می‌شوند. در این نوع از تغییر تفاوت بین توزیع‌ها بسیار کم است و تغییرات در دوره‌ی طولانی‌تری از زمان دیده می‌شوند (Brzezinski and Stefanowski ۲۰۱۴). بعضی از مراجع تغییرات افزایشی را دسته‌ای جدا از تغییرات تدریجی در نظر گرفته‌اند (Ditzler ۲۰۱۲). تغییرات بازگشت‌کننده، هر چند وقت یکبار رخ داده و موجب تغییر در مفاهیم ویا توزیع آن‌ها می‌شوند ولی پس از چندی مجدداً مفاهیم قدیمی تکرار می‌شوند و این سیر ادامه پیدا می‌کند (Gama, Žliobaitė et al. ۲۰۱۴). به عقیده‌ی (Forman ۲۰۰۶) این نوع تغییر بیشتر در متون خبری رخ می‌دهد. هر چند بررسی‌های حاصل از تحقیق (Šilić and Bašić ۲۰۱۲) نشان از تغییرات تدریجی در متون خبری است. لذا تشخیص نوع تغییر به پیکره‌ی مورد بررسی وابسته است و نیاز به انجام آزمایشات بر روی پیکره‌ی مورد نظر دارد.

برای تشخیص این نوع تغییرات معمولاً از روش‌های بصری استفاده می‌شود. برای مثال در مرجع (Forman ۲۰۰۶) برای این منظور یک نمایش بصری جدید از فضای ویژگی ارائه شده است. روش مذکور به این صورت است که با داشتن مجموعه‌ی کل کلمات در تمام روزها (مثلاً برای مدت یک سال)، کلمات بر اساس تاریخی که اولین بار در ۱۰۰ ویژگی برتر دیده شدند مرتب می‌شوند، سپس تصویری تشکیل داده می‌شود که در آن ستون‌ها نماینده‌ی کلمات و سطرها نماینده‌ی روزها هستند و ترتیب روزها به این نحو است که از بالای تصویر به سمت پایین تصویر، به تاریخ روزها افزوده می‌شود. یک پیکسل به شرطی به رنگ قرمز درمی‌آید که کلمه‌ی مربوط به ستون آن پیکسل در آن روز، جزء ۱۰۰ کلمه‌ی برتر باشد. وگرنه بر اساس میزان اهمیت، به رنگی در طیف سیاه و سفید درمی‌آید به طوری که هر چقدر مشکلی‌تر باشد بی‌اهمیت‌تر و هر چه سفیدتر باشد با اهمیت‌تر است. شکل ۱ یکی از تصاویر حاصل از روش مذکور در مرجع (Forman ۲۰۰۶) را نشان می‌دهد.



Category GCAT (government & social issues, 30%): 729 top predictive words

شکل ۱: نمایش تغییر توزیع کلمات در مرجع (Forman ۲۰۰۶)

در این شکل، مرز قرمز رنگی که از پیکسل‌ها به وجود آمده زمانی را نشان می‌دهد که هر یک از کلمات برای اولین بار در ۱۰۰ کلمه‌ی برتر رویت شده‌اند. این نمایش بصری کمک می‌کند تا نرخ تغییر و نوع تغییر به وضوح مشاهده شود. برای مثال قسمت‌هایی از تصویر که در آن، در یک ستون بعد از چند پیکسل قرمز پیکسل‌های سیاه رویت شده‌اند نشان دهنده‌ی کاهش اهمیت کلمات مهم بعد از چند روز است و بر عکس، قسمت‌هایی که چندین پیکسل عمودی پیوسته از رنگ قرمز دیده می‌شود یعنی کلمه‌ی مربوط به ستون مذکور در طی چند هفته مدام استفاده شده است (بعضی سمت راست در شکل). همینطور برای تشخیص نوع تغییر نیز با مشاهده‌ی رنگ قرمز در یک ستون به صورتی که در بعضی سمت چپ در تصویر نشان داده شده است می‌توان به یک تغییر بازگشت‌کننده پی برد.

در مرجع (D'hondt, Verberne et al ۲۰۱۴). برای تشخیص تغییر بازگشت‌کننده و تدریجی از دو آزمایش متفاوت استفاده شده است. برای تشخیص تغییر تدریجی، توزیع دسته‌های مختلف در طول سال‌های متوالی بررسی می‌شود. منطق نویسنده این است که اگر بین دو سال متوالی اندازه‌ی یک دسته تغییر چندانی نداشته باشد، ولی در طول یک دوره‌ی طولانی از زمان، اندازه‌ی آن دسته به طور قابل ملاحظه‌ای رشد داشته باشد، تغییر داده‌های مربوط به آن دسته تدریجی است. برای تشخیص تغییر بازگشت‌کننده، توزیع کلمات در سال‌های مختلف و برای هر کلاس به صورت جداگانه بررسی می‌شود.

در مقاله‌ی (Šilić and Bašić ۲۰۱۲) برای شناسایی انواع تغییرات ابتدا انواع روش‌های نمونه‌برداری داده‌های آموزشی و آزمایشی از کل داده‌ها ارائه شده است و سپس با بررسی میزان عملکرد دسته‌بندی‌های ساخته شده بر روی انواع نمونه برداری‌ها، نوع تغییر موجود در پیکره‌ی مورد بررسی تشخیص

¹ Concept shift

² Incremental

داده شده است. بیشتر روش های مواجهه با تغییرات محدود به نوع خاصی از تغییر می شوند، در حالی که در واقعیت اغلب ترکیبی از چنین تغییراتی رخ می دهد.

۴. انواع روش های تشخیص تغییر

بسیاری از الگوریتم هایی که برای تطبیق با تغییرات زمانی ارائه شده اند، از یک مدل به روز رسانی ساده استفاده می کنند (Katakis, Tsoumakas et al. ۲۰۰۶). به این ترتیب که به محض ورود یک داده به سیستم، به روز رسانی مدل انجام می شود. در حالی که این روش ها هزینه ی بالایی دارند. راه دیگر، تشخیص زمان تغییرات و به روز رسانی مدل در صورتی است که این تغییرات قابل توجه باشند (Tsymbal ۲۰۰۴). در حقیقت هدف از تشخیص تغییر، شناسایی نقاط و یا بازه هایی است که در طول آنها تغییر رخ می دهد. در حالت کلی تشخیص تغییر می تواند با نظارت و بررسی داده ی خام صورت گیرد یا اینکه بر اساس عملکرد حاصل از دسته بند انجام شود. بر اساس دسته بندی انجام شده در (Gama, Žliobaitė et al. ۲۰۱۴) به طور کلی تشخیص دهنده های تغییر به سه دسته تقسیم می شوند:

۴.۱. تشخیص دهنده های مبتنی بر آنالیز ترتیبی

در این روش ابتدا فرض می شود که X_1^n یک توالی از نمونه هاست که در آن زیرمجموعه ای از نمونه ها مثل X_1^w ($1 < w < n$) از یک توزیع مانند P_0 و زیرمجموعه ی دیگر از نمونه ها X_w^n (یعنی سایر نمونه ها) از توزیعی مانند P_1 تولید شده اند. زمانی که توزیع مثال ها در نقطه های مانند W از P_0 به P_1 تغییر می کند، انتظار می رود احتمال مشاهده ی زیردنباله ای با توزیع P_1 به شدت بالاتر از زیردنباله ای با توزیع P_0 شود. میزان شدت باید طوری باشد که نسبت دو احتمال از یک حد آستانه کمتر نباشد. با فرض اینکه نمونه های X_i از هم مستقل باشند، یک شاخص آماری برای آزمودن این فرضیه که آیا W نقطه ی تغییر است یا نه بر اساس رابطه ی (۱) محاسبه می شود:

$$T_w^n = \log \frac{P(x_w \dots x_n | P_1)}{P(x_w \dots x_n | P_0)} = \sum_{i=w}^n \log \frac{P_1[x_i]}{P_0[x_i]} = T_w^{n-1} + \log \frac{P_1[x_n]}{P_0[x_n]} \quad (1)$$

در این رابطه اگر $T_w^n > L$ باشد آنگاه تغییر تشخیص داده می شود. (L یک حد آستانه تعریف شده توسط کاربر است.)

۴.۲. تشخیص دهنده های مبتنی بر جداول کنترلی

جداول کنترلی، یادگیری را به عنوان یک فرآیند در نظر می گیرند و تکامل این فرآیند را بر اساس میزان رخداد خطا در پیشگویی های انجام شده توسط یادگیرنده نظارت می کنند (Nishida and Yamauchi ۲۰۰۹). بر اساس این نظارت سیستم در یکی از ۳ حالت تحت کنترلی، خارج از کنترلی و وضعیت هشدار قرار خواهد گرفت. وضعیت سیستم تا زمانی تحت کنترلی است که خطای سیستم ثابت باشد؛ و این یعنی نمونه ی آزمایشی ورودی از همان توزیعی آمده است که نمونه های قبلی آمده اند. وضعیت خارج از کنترلی وضعیتی است که خطا به طور قابل ملاحظه ای در مقایسه با نمونه های اخیر افزایش داشته است. لذا به احتمال ۹۹ درصد نمونه ی آخر با توزیع متفاوتی نسبت به نمونه های قبلی آمده است. وضعیت هشدار حالتی بین دو وضعیت قبلی است. خطا در این حالت در حال افزایش است ولی هنوز خارج از کنترلی نشده است. افزایش خطا در این حالت ممکن است به علت نویز، کاهش عملکرد مدل پیشگویانه به صورت مقطعی و یا تغییر واقعی باشد. در این وضعیت تعداد بیشتری مثال نیاز است تا بتوان به یک نتیجه ی قطعی رسید و تشخیص داد افزایش خطا حاصل از نویز است یا تغییر واقعی. با توجه به توضیحات مذکور، با توجه به فاصله ی زمانی بین وضعیت هشدار و وضعیت خارج از کنترلی می توان نرخ تغییر را مشخص کرد. اگر فاصله ی زمانی بین وضعیت هشدار و وضعیت خارج از کنترلی کم باشد نشان دهنده ی سرعت بالای تغییر است و اگر این فاصله زیاد باشد نشان از سرعت پایین تغییر دارد.

۴.۳. تشخیص دهنده های مبتنی بر نظارت بر توزیع پنجره های زمانی

این روش ها عموماً از یک پنجره ی مرجع با اندازه ی ثابت استفاده می کنند که شامل داده های قدیمی می شود و از یک پنجره ی لغزنده استفاده می کنند که حاوی جدیدترین نمونه ها هستند. مقایسه ی بین دو پنجره از مقایسه ی بین شاخص های آماری توزیع دو پنجره انجام می شود. اگر هر دو پنجره دارای شاخص های آماری یکسانی نباشند، شروع پنجره ی لغزنده به عنوان محل تغییر تعیین می شود. شاخص آماری مذکور می تواند از روی داده های خام و یا عملکرد مدل آموزشی تعریف شود. دو پنجره ی مذکور می توانند اندازه ی یکسانی داشته باشند و یا اینکه اندازه ی پنجره ی لغزان می تواند به صورت تصاعدی رشد کند. نویسندگان مرجع (Kifer, Ben-David et al. ۲۰۰۴) از معیار Chernoff bound برای مقایسه ی آماری توزیع دو پنجره استفاده کرده اند. تحقیق (Vorburger and Bernstein ۲۰۰۶) یک معیار مبتنی بر آنتروپی ارائه کرده است تا تفاوت بین توزیع دو پنجره را نشان دهد. در این روش اگر توزیع دو پنجره مساوی باشد، آنگاه آنتروپی برابر ۱ می شود و اگر توزیع ها مساوی نباشند مقدار آنتروپی برابر صفر می شود. در این روش شاخص آنتروپی به صورت مداوم و در طول زمان نظارت می شود و اگر این شاخص از یک حد آستانه ی تعریف شده توسط کاربر پایین تر برود نشان دهنده ی یک تغییر است. در تحقیق (Bach and Maloof ۲۰۰۸) از دو مدل با نام های مدل پایدار و مدل واکنشی استفاده شده است. مدل پایدار بر اساس داده های آموزشی در طول یک بازه ی طولانی از زمان و مدل واکنشی بر اساس یکی بازه ی زمانی کوتاه از داده های آموزشی ساخته شده است. این تکنیک از مدل واکنشی به عنوان

شاخصی برای تشخیص تغییر استفاده می کند. به این صورت که وقتی مدل واکنشی برای یک نمونه نسبت به مدل پایدار عملکرد بهتری از خود نشان می دهد نشان دهنده رخداد یک تغییر در بازه زمانی کوتاه مربوط به مدل واکنشی است. در مقایسه بین این سه روش می توان گفت محدودیت تشخیص دهنده های دسته سوم نسبت به دسته اول، نیاز به حافظه است. تشخیص دهنده های دسته اول نیازی به ذخیره داده ورودی ندارند ولی تشخیص دهنده های دسته دوم باید داده های موجود در هر دو پنجره مرجع و لغزنده را در حافظه نگه دارند. مزیت تشخیص دهنده های دسته اول دقت بالاتر آنها نسبت به سایر تشخیص دهنده ها است.

۵. روش های مواجهه با تغییرات زمانی

در تحقیقات مطالعه شده، چهار شیوه در مواجهه با اثرات حاصل از تغییرات زمانی در داده ها گزارش شده است. (Nishida, Hoshide et al., ۲۰۱۲, Tsymbal, Fukumoto, Ushiyama et al., Gama, Žliobaitė et al., ۲۰۱۴) انتخاب نمونه (۲) وزن دهی به نمونه ها (۳) یادگیری دسته جمعی (Tsymbol) (۴) انتخاب ویژگی^۱ (Nishida, Hoshide et al., ۲۰۱۲). در ادامه هر یک از روش ها به طور کامل توضیح داده می شوند.

۵.۱. انتخاب نمونه

در روش انتخاب نمونه، هدف انتخاب نمونه های مرتبط با داده هایی است که از سایر داده ها به نمونه های آزمایشی نزدیک ترند. این شیوه که رایج ترین روش برای مواجهه با تغییرات زمانی است، از یک پنجره لغزان تشکیل شده است که بر روی آخرین نمونه ها حرکت کرده و از مدل ساخته شده حاصل از آن ها برای پیشگویی نمونه ها در مثال هایی که در آینده ای بسیار نزدیک وجود دارند استفاده می کند (Widmer and Kubat, ۱۹۹۶, Tsymbal, ۲۰۰۴). در این الگوریتم ها می توان از پنجره هایی با اندازه ی ثابت استفاده کرد یا اینکه با روش های اکتشافی اندازه ی پنجره را متناسب با آخرین بازه ی تغییر در داده ها تنظیم کرد. روش پیشنهادی مرجع (D'hondt, Verberne et al.) سعی در انتخاب پنجره به نحوی دارد که یک تعادل بین جدید بودن داده و اندازه ی پنجره برقرار شود. در مرجع (Šilić and Bašić, ۲۰۱۲) نشان داده شده است که در صورت وجود رانش مفهوم در داده های آموزشی، حتی استفاده از انتخاب نمونه رندم به جای کل داده ها به عنوان داده ی آموزشی، می تواند میزان عملکرد دسته بند را افزایش دهد. در مرجع (Rocha, Mourão et al., ۲۰۱۳) الگوریتمی به نام Chronos پیشنهاد شده است که به ازای هر نمونه ی آزمایشی بخشی از مجموعه ی آموزشی به نام محتوای زمانی^۲ را انتخاب می کند. این الگوریتم کمک می کند تا بزرگترین محتوای زمانی انتخاب شود به طوری که اثرات ناشی از تغییرات زمانی در عملکرد دسته بند به حداقل برسد.

۵.۲. وزن دهی به نمونه

در این روش به بعضی از نمونه ها در طول آموزش وزن بیشتری اختصاص داده می شود. نمونه ها می توانند بر اساس میزان قدمت و میزان ارتباط به مفاهیم جدید وزن دهی شوند (Tsymbal, ۲۰۰۴). روش وزن دهی به نمونه ها نیز، فرض می کند که نمونه های جدیدتر آموزشی برای دسته بندی مثال های آزمایشی جدید بهتر عمل می کنند. (Nishida, Hoshide et al., ۲۰۱۲).

۵.۳. یادگیری دسته جمعی

روش یادگیری دسته جمعی مجموعه ای از مدل های پیشگویی را نگهداری کرده و از رای گیری بین مدل ها یا وزن دهی به آنها در زمان پیشگویی نمونه های آزمایشی استفاده می کند. معمولاً وقتی مدل های جدید اضافه می شوند، مدل های ضعیف یا قدیمی، حذف می شوند. حذف شدن مدل های قدیمی یا به صورت دوره ای است و یا پس از مشاهده کاهش صحت عملکرد دسته بند انجام می شود. هر چند در بعضی از مراجع (Nishida and Yamauchi, ۲۰۰۹)، هیچ دسته بندی تحت هیچ شرایطی از سیستم حذف نمی شود. سیستم پیشنهادی این مرجع (Nishida and Yamauchi, ۲۰۰۹) از چندین دسته بند آفلاین و آنلاین برای مواجهه با تغییرات زمانی استفاده کرده است. با شروع کار سیستم با ورود هر مثال جدید دسته بند آنلاین به روزرسانی می شود، برخلاف دسته بند آفلاین که با ورود مثال های جدید به روزرسانی نمی شوند. در حقیقت دسته بند های آنلاین برای تطبیق با تغییرات تدریجی، مدام در حال یادگیری مثال ها هستند و دسته بند های آفلاین برای پاسخ به تغییرات بازگشت کننده استفاده می شوند. در (Kolter and Maloof, ۲۰۰۳) مجموعه ای از دسته بند ها با قدمت های متفاوت ساخته می شود به طوری که هر یک از دسته بند ها شامل نمونه هایی مربوط به دوره ی خاصی از زمان هستند.

۵.۴. انتخاب ویژگی

در اکثر تحقیقات تنها سه دسته ی قبلی برای مواجهه با تغییرات زمانی بیان شده است ولی در تحقیق (Nishida, Hoshide et al., ۲۰۱۲, Fukumoto, Ushiyama et al., ۲۰۱۴) دسته روش انتخاب ویژگی^۳ نیز برای مواجهه با تغییرات زمانی ارائه شده است. این رویکرد در چند سال اخیر مورد توجه قرار گرفته و تمرکز را از وزن دهی به نمونه ها به سمت وزن دهی به ویژگی ها سوق می دهد. در این رویکرد برای مواجهه با تغییرات زمانی هر بار به صورت پویا ارزشمندترین ویژگی ها (کلمات) انتخاب می شوند. در تحقیق (Fukumoto, Ushiyama et al., ۲۰۱۴) ویژگی ها یا کلمات مورد استفاده در دسته بندی

¹ Ensemble learning

² Feature selection

³ Temporal context

⁴ Feature selection

متن به دو دسته تقسیم شده‌اند، یکی کلمات مستقل از زمان^۱ و دیگری کلمات وابسته به زمان^۲ است که در مرجع (Nishida, Hoshide et al. ۲۰۱۲) به ترتیب کلمات پایدار^۳ و ناپایدار^۴ خوانده می‌شوند. در (Fukumoto, Ushiyama et al. ۲۰۱۴) از روش TWF (Salles, Rocha et al. ۲۰۱۰) برای وزن دهی استفاده شده است با این تفاوت که به جای وزن دهی در سطح نمونه، وزن دهی در سطح ویژگی انجام می‌شود. به این صورت که اگر یک ویژگی وابسته به زمان بود، از وزن دهی بر پایه‌ی TWF و اگر مستقل از زمان بود وزن دهی معمولی برای آن استفاده می‌گردد.

۶. ارتباط بین روش های مواجهه با تغییر با نوع تغییر

به طور کلی انتخاب روش مواجهه با تغییر، به نوع تغییر موجود در پیکره‌ی مورد بررسی وابسته است (D'hondt, Verberne et al.). نویسندگان تحقیق (Nishida, Hoshide et al. ۲۰۱۲) معتقدند که روش وزن دهی به نمونه‌ها تنها در مواجهه با تغییر تدریجی خوب عمل می‌کند ولی در مواجهه با تغییر ناگهانی شکست می‌خورد. در تحقیق (Carmona-Cejudo, Baena-García et al. ۲۰۱۱) نتایج مناسبی از انتخاب نمونه برای مواجهه با تغییر تدریجی موجود در دسته‌بندی جریانی از ایمیل‌ها به دست آمده است. به عقیده‌ی (D'hondt, Verberne et al.) یادگیری دسته جمعی برای تغییرات تدریجی و بازگشت کننده خوب عمل می‌کند و بهترین راه مواجهه با تغییر ناگهانی، انتخاب داده‌ی آموزشی تنها از نمونه‌هایی است که با توزیع جدید همخوانی دارند و زمانی که تغییر بازگشت کننده وجود ندارد بهترین راه مواجهه با تغییر این است که مدل هر چند وقت یکبار به روز رسانی شود. به این نحو که مدل‌های قدیمی آموزش دیده تخریب و مدل‌های جدید به جای مدل‌های دسته بندی قدیمی جایگزین شود.

در بعضی تحقیقات برای مواجهه با تغییر تنها به یک روش بسنده نمی‌کنند و از ترکیبی از چند روش استفاده می‌کنند. برای مثال در تحقیق (Salles, Rocha et al. ۲۰۱۰) نوعی وزن دهی مبتنی بر زمان به نام TWF برای دسته بندی ارائه شده است. وزن دهی TWF، تغییرات موجود در رابطه ی بین کلاس و واژگان را در دوره‌های زمانی مختلف استخراج کرده و در دسته‌بند آن را به دو روش استفاده می‌کند: TWF بر روی نمونه‌ها و TWF بر روی امتیاز به دسته‌بندها. در TWF بر روی متون، فاصله‌ی زمانی بین متون آموزشی و آزمایشی به منظور اندازه‌گیری میزان شباهت محاسبه شده و بر این اساس یک متن آموزشی که از لحاظ زمانی نزدیک‌تر به متن آزمایشی است، اهمیت بیشتری پیدا می‌کند. در TWF بر روی دسته‌بندها، امتیاز نهایی به یک کلاس از جمع وزن دار امتیازهای تولید شده توسط دسته‌بندهای عادی که در بازه‌های زمانی متفاوت ساخته شده اند به دست می‌آید.

۷. نتیجه گیری

در دسته بندی متن عموماً داده ها ابتدا به صورت آفلاین جمع آوری شده و آموزش بر روی آن‌ها انجام می‌شود سپس مدل آموزش دیده به این روش، می‌تواند برای پیش بینی داده‌های دیده نشده‌ی ورودی مورد استفاده قرار بگیرد. در شرایطی که داده‌ها مدام در حال به روز رسانی هستند، توزیع آنها می‌تواند در طول زمان تغییر کند به طوری که موجب اثرگذاری بر روی عملکرد دسته بند شود. از آنجایی که انتخاب روش مواجهه با تغییر، به نوع تغییر وابسته است، لذا در این مقاله ابتدا به بررسی انواع تغییرات زمانی پرداخته شد. سپس انواع روش‌های مواجهه با تغییر بیان شد و در نهایت بهترین روش‌های مواجهه با تغییر متناسب با نوع تغییر رخ داده مورد بررسی قرار گرفت. در این مقاله سعی شد تا حد امکان یک دید کلی و همه جانبه در زمینه‌ی اثر زمان در دسته‌بندی خودکار متن ارائه شده و مسیر را برای سایر محققین در انتخاب روش مناسب برای مواجهه با تغییرات زمانی هموارتر کند.

۸. منابع مورد استفاده

- Alonso, O., M. Gertz and R. Baeza-Yates (۲۰۰۷). On the value of temporal information in information retrieval. ACM SIGIR Forum, ACM.
- Bach, S. H. and M. Maloof (۲۰۰۸). Paired learners for concept drift. Data Mining, ۲۰۰۸. ICDM'0۸. Eighth IEEE International Conference on, IEEE.
- Brzezinski, D. and J. Stefanowski (۲۰۱۴). "Reacting to different types of concept drift: The accuracy updated ensemble algorithm." Neural Networks and Learning Systems, IEEE Transactions on ۲۵(۱): ۹۴-۸۱
- Carmona-Cejudo, J. M., M. Baena-García, R. M. Bueno, J. Gama and A. Bifet (۲۰۱۱). Using GNUsmail to Compare Data Stream Mining Methods for On-line Email Classification. WAPA.

¹ Temporal independent terms

² Temporal dependent terms

³ Stationary

⁴ Bursty

- D'hondt, E., S. Verberne, N. Oostdijk, J. Beney, C. Koster and L. Boves (۲۰۱۴). "Dealing with temporal variation in patent categorization." Information Retrieval: ۲۵-۱
- Ditzler, G. (۲۰۱۲). Incremental Learning of Concept Drift from Imbalanced Data.
- Forman, G. (۲۰۰۶). Tackling concept drift by temporal inductive transfer. Proceedings of the ۲۹th annual international ACM SIGIR conference on Research and development in information retrieval, ACM.
- Fukumoto, F., S. Ushiyama, Y. Suzuki and S. Matsuyoshi (۲۰۱۴). The Effect of Temporal-based Term Selection for Text Classification. Australasian Language Technology Association Workshop .۲۰۱۴
- Gama, J., I. Žliobaitė, A. Bifet, M. Pechenizkiy and A. Bouchachia (۲۰۱۴). "A survey on concept drift adaptation." ACM Computing Surveys (CSUR) ۴۶(۴): .۴۴
- Katakis, I., G. Tzoumakas and I. Vlahavas (۲۰۰۶). "Dynamic feature space and incremental feature selection for the classification of textual data streams." Knowledge Discovery from Data Streams: ۱۱۶-۱۰۷
- Kifer, D., S. Ben-David and J. Gehrke (۲۰۰۴). Detecting change in data streams. Proceedings of the Thirtieth international conference on Very large data bases-Volume ۳۰, VLDB Endowment.
- Kolter, J. Z. and M. Maloof (۲۰۰۳). Dynamic weighted majority: A new ensemble method for tracking concept drift. Data Mining, ۲۰۰۳. ICDM ۲۰۰۳. Third IEEE International Conference on, IEEE.
- Mourão, F., L. Rocha, R. Araújo, T. Couto, M. Gonçalves and W. Meira Jr (۲۰۰۸). Understanding temporal aspects in document classification. Proceedings of the ۲۰۰۸ International Conference on Web Search and Data Mining, ACM.
- Nishida, K., T. Hoshida and K. Fujimura (۲۰۱۲). Improving tweet stream classification by detecting changes in word probability. Proceedings of the ۳۵th international ACM SIGIR conference on Research and development in information retrieval, ACM.
- Nishida, K. and K. Yamauchi (۲۰۰۹). Learning, detecting, understanding, and predicting concept changes. Neural Networks, ۲۰۰۹. IJCNN ۲۰۰۹. International Joint Conference on, IEEE.
- Rocha, L., F. Mourão, H. Mota, T. Salles, M. A. Gonçalves and W. Meira Jr (۲۰۱۳). "Temporal contexts: Effective text classification in evolving document collections." Information Systems ۳۸(۳): .۴۰۹-۳۸۸
- Rocha, L., F. Mourão, A. Pereira, M. A. Gonçalves and W. Meira Jr (۲۰۰۸). Exploiting temporal contexts in text classification. Proceedings of the ۱۷th ACM conference on Information and knowledge management, ACM.
- Salles, T., L. Rocha, G. L. Pappa, F. Mourão, W. Meira Jr and M. Gonçalves (۲۰۱۰). Temporally-aware algorithms for document classification. Proceedings of the ۳۳rd international ACM SIGIR conference on Research and development in information retrieval, ACM.
- Sebastiani, F. (۲۰۰۲). "Machine learning in automated text categorization." ACM computing surveys (CSUR) ۳۴(۱): .۴۷-۱
- Šilić, A. and B. D. Bašić (۲۰۱۲). Exploring classification concept drift on a large news text corpus. Computational linguistics and intelligent text processing, Springer: .۴۳۷-۴۲۸
- Tsybmal, A. (۲۰۰۴). "The problem of concept drift: definitions and related work".
- Vorburger, P. and A. Bernstein (۲۰۰۶). Entropy-based concept shift detection. Data Mining, ۲۰۰۶ ICDM'۰۶. Sixth International Conference on, IEEE.
- Wenerstrom, B. and C. Giraud-Carrier (۲۰۰۶). Temporal data mining in dynamic feature spaces. Data Mining, ۲۰۰۶. ICDM'۰۶. Sixth International Conference on, IEEE.
- Widmer, G. and M. Kubat (۱۹۹۶). "Learning in the presence of concept drift and hidden contexts." Machine learning ۲۳(۱): .۱۰۱-۶۹