



ارائه شده توسط:

سایت ترجمه فا

مرجع جدیدترین مقالات ترجمه شده

از نشریات معتبر

مطالعه ای مختصر در خصوص طبقه بندی توالی

چکیده :

طبقه بندی توالی دارای طیف وسیعی از کاربرد ها نظیر تحلیل ژنوم، بازیابی اطلاعات، انفورماتیک سلامت، امور مالی و تشخیص ناهنجاری ها می باشد. متفاوت از انجام طبقه بندی بر روی بردار های ویژگی، توالی ها معمولاً فاقد ویژگی های صریح می باشند. حتی با فنون انتخاب ویژگی پیشرفته، بعدیت ویژگی های بالقوه بسیار بالا است و ماهیت متوالی ویژگی ها را به سختی می توان درک کرد. این موجب می شود تا طبقه بندی توالی و دنباله به یک ویژگی چالش بر انگیز تر از طبقه بندی بردار های ویژگی تبدیل شود. ما اقدام به طبقه بندی توالی از حیث روش ها و حوزه های کاربرد مختلف می کنیم. ما مرور جامعی را در خصوص انواع مختلف مسائل طبقه بندی توالی نظیر طبقه بندی اولیه توالی ها و یادگیری نیمه نظارت شده در زمینه طبقه بندی ها انجام می دهیم.

1- مقدمه

طبقه بندی توالی دارای طیف وسیعی از کاربرد ها می باشد. در تحقیقات ژنومی، طبقه بندی توالی های پروتین به دسته های موجود بر ای یادگیری وظایف و کارکرد های پروتین استفاده می شود (13). در انفورماتیک سلامت، طبقه بندی سری های زمانی ECG (سری های زمانی ضربان قلب) به ما می گوید که آیا داده ها مربوط به یک فرد سالم هستند یا مربوط به یک بیمار مبتلا به بیماری قلبی است (59). در تشخیص نفوذ و ناهنجاری، توالی فعالیت های دسترسی سیستم بر روی یونیکس بر ای تشخیص رفتار های غیر طبیعی و ناهنجاری پایش می شود (33). در بازیابی اطلاعات، طبقه بندی اسناد به مقوله های موضوعی مختلف، توجهات زیادی را به خود جلب کرده است (51). سایر مثال های جالب شامل طبقه بندی توالی های کوئری بر ای تمایز ربات های اینترنتی از کاربر ان انسان (18، 58) و طبقه بندی داده های توالی ترانسفکشن در یک بانک بر ای مبارزه با پول شویی، می باشد (42).

به طور کلی، یک توالی، فهرست منظمی از رویداد هاست. یک رویداد را می توان به صورت یک ارزش نمادین، یک ارزش واقعی عددی، یک بردار با ارزش واقعی یا یک داده نوع پیچیده در نظر گرفت. در این مقاله، داده های توالی یا دنباله ای به زیر انواع زیر در نظر گرفته می شود

- با توجه به القای علایم و نمادها $fE1; E2; E3; :::; Eng$ ، یک توالی نمادین ساده، فهرست منظمی از نمادها از الفبا می باشد. بر ای مثال، یک توالی DNA متشکل از چهار امینو اسید A-C-G-T و قطعه DNA نظیر ACCCCCCGT می باشد که یک توالی نمادین ساده است.

- یک توالی نمادین ساده، فهرستی از بردارها می باشد. هر بردار یک زیر مجموعه ای از الفبا(34) می باشد. بر ای مثال، بر ای توالی ایتم های خریداری شده توسط یک مشتری در یک سال، در نظر گرفتن هر تر اکشن به صورت یک بردار، یک توالی یا دنباله می تواند به صورت ساعت است (شیر؛ نان) (شیر؛ تخم مرغ) (سیب زمینی؛ پنیر؛ کک) در نظر گرفته شود.

- سری های زمانی ساده دنباله ای از ارزش ها و مقادیر واقعی مرتب شده به ترتیب صعودی می باشد

$$h(t_1; 0:1)(t_2; 0:3) \text{ } \text{ } \text{ } (t_n; 0:3)i$$

یک سری زمانی ساده می باشد که داده های مربوط به مهر زمانی t_1 تا t_2 را نشان می دهد

- سری های زمانی چند متغیره، دنباله و توالی از بردارهای عددی می باشد. بر ای مثال

$$h(t_1; h_0:1; 0:3; 0:5i)(t_2; h_0:3; 0:9; 0:8i) \text{ } \text{ } \text{ } (t_n; h_0:3; 0:9; 0:4i)i$$

یک سری زمانی چند متغیره می باشد.

- در بالا، انواع داده های رویدادها ساده می باشد. در برخی از زمینه ها، نوع داده های رویدادها ممکن است

پیچیده باشند. بر ای مثال، در یک مجموعه داده های پرونده بیمار

<http://www.informsdmcontest2009.org/> هر بیمار با یک توالی طولی از ویزیت های بیمارستانی

نشان داده می شود. هر باز بینی یا ویزیت یک رویداد است و با اندازه گیری های عددی چندگانه توصیف می شود.

یک توالی رویداد پیچیده اشاره به شکل کلی از توالی ها و دنباله ها دارد

یک توالی یا دنباله ممکن است حامل یک برچسب باشد. بر ای مثال، سری های زمانی داده های ECG بر گرفته

از افراد سالم یا بیمار می باشند. یک توالی DNA مربوط به یک منطقه ژن کد کننده یا غیر کد کننده می باشد.

با توجه به این L یک مجموعه برچسب باشد، طبقه بندی توالی ایزاری بر ای یادگیری کلاسیفایر توالی C می باشد

شد، که دنباله S را بر ای یک برچسب کلاس $l \in L$ در نظر می گیرد و به صورت $C: s \rightarrow l, l \in L$ تعیین

می شود.

در طبقه بندی توالی، هر توالی مربوط به تنها یک برچسب کلاسی بوده و توالی کل بر ای دسته بند قل از دسته بندی وجود دارد. هم چنین سناریو ها ی مختلف دیگری بر ای طبقه بندی توالی وجود دارد. بر ای مثال، بر ای توالی یا دنباله علایم یک بیمار در بلند مدت، شرایط سلامتی بیمار می تواند تغییر کند. هم چنین سناریو ها ی دیگری بر ای دسته بندی توالی وجود دارد. بر ای مثال، بر ای یک دنباله و توالی ای از علایم بیمار در یک دوره بلند مدت، شرایط سلامتی بیمار می تواند تغییر کند. بر ای توالی استریمینگ، که به صورت توالی نامحدود در نظر گرفته می شود، به جای پیش بینی یک برچسب کلاس، امکان پیش بینی یک توالی از رویداد ها وجود دارد. این مسئله نیز به صورت یک فرایند دسته بندی توالی قوی در نظر گرفته شده است. در این مقاله، ما در مورد انواع مختلف طبقه بندی توالی در بخش 3 صحبت می کنیم.

سه چالش مهم دیگر در طبقه بندی توالی وجود دارد. اولاً، بیشتر دسته بند ها نظیر درختان تصمیم گیری و شبکه ها ی عصبی، تنها داده ها ی ورودی را به صورت بردار ویژگی ها در نظر می گیرد. دوماً، حتی با روش ها ی مختلف انتخاب ویژگی، می تواند یک توالی را به مجموعه ای از ویژگی ها تبدیل می کند. این ویژگی ها بسیار مهم هستند. بعدیت فضای ویژگی بر ای داده ها ی توالی می تواند بسیار بالا باشد و محاسبه می تواند پرهزینه باشد. سوماً، علاوه بر نتایج طبقه بندی صحیح در برخی از موارد، امکان دست یابی به دسته بند قابل تفسیر وجود دارد. ایجاد ک دسته بند توالی قابل تفسیر سخت است زیرا ویژگی ها ی صریح وجود ندارد.

در این مقاله ما مرور مختصری بر روش ها ی طبقه بندی و دسته بندی توالی موجود ارائه می کنیم. چون بیشتر مطالعات فعلی بر دسته بندی توالی سنتی متمرکز هستند، بخش دوم خلاصه ای از روش ها ی اصلی را بر ای این کار ارائه می کند. در بخش سوم، ما در مورد انواع فعالیت ها ی دسته بندی توالی نظیر استریمینگ دسته بندی توالی و دسته بندی اولیه توالی ها صحبت می کنیم. در بخش چهارم، ما خلاصه ای از دسته بندی توالی را از دیدگاه کاربردی نظیر داده ها ی سری زمانی، داده ها ی متنی و داده ها ی ژنومی را ارائه می کنیم. در بخش 5 نیز نتیجه گیری ارائه شده است.

2- روش های طبقه بندی توالی

روش ها ی طبقه بندی توالی را می توان به سه نوع اصلی تقسیم کرد:

- اولین دسته، طبقه بندی مبتنی بر ویژگی است که توالی را به یک بردار ویژگی تبدیل کرده و سپس از روش های طبقه بندی سنتی استفاده می کند. انتخاب ویژگی نقش مهمی در این نوع روش ها ایفا می کند.
- دومین دسته، طبقه بندی مبتنی بر فاصله توالی است. تابع فاصله ای تشابه بین توالی ها را اندازه گیری کرده و کیفیت طبقه بندی را تعیین می کند
- سومین مقوله طبقه بندی مبتنی بر مدل نظیر استفاده از مدل مارکوف پنهان و سایر مدل های آماری بر ای طبقه بندی توالی هاست

در ادامه این بخش ما یک سری روش های معرف را در سه دسته ارائه می کنیم. برخی از روش ها از چندین مقوله استفاده کرده اند. بر ای مثال، می توان از SVM با استخراج ویژگی ها یا تعریف شاخص های فاصله ای استفاده کرد. طبقه بندی توالی با استفاده از SVM در بخش 2-3 خلاصه شده است. همه روش های بحث شده در این بخش، طبقه بندی توالی سنتی می باشند.

1-2 طبقه بندی مبتنی بر ویژگی

روش های طبقه بندی سنتی نظیر درختان تصمیم گیری و شبکه های عصبی، بر ای طبقه بندی بردار های ویژگی طراحی می شوند. یک شیوه بر ای حل مسئله طبقه بندی توالی، تبدیل توالی به بردار ویژگی ها از طریق انتخاب ویژگی است

بر ای یک توالی نمادین، ساده ترین راه، در نظر گرفتن هر عنصر به صورت یک ویژگی است. بر ای مثال، یک توالی CACG را می توان به صورت یک بردار ACCG در نظر گرفت. با این حال، ماهیت متوالی توالی ها را نمی توان با این تبدیل توجیه کرد. بر ای حفظ ترتیب عناصر در یک توالی، یک قطعه کوچک از نماد های متوالی K موسوم به K-گرم معمولاً به عنوان یک ویژگی انتخاب می شود. با توجه به مجموعه ای از این K-گرم ها، یک توالی را می توان به صورت بردار حضور یا غیاب K-گرم تعیین کرد. گاهی اوقات امکان تطبیق غیر دقیق با K-گرم وجود دارد. با استفاده از ویژگی های K-گرم، توالی ها را می توان با روش طبقه بندی سنتی نظیر SVM تعیین کرد. یک خلاصه ای از روش های انتخاب ویژگی مبتنی بر K-گرم بر ای طبقه بندی توالی را می توان در (16) یافت.

اندازه ویژگی‌ها ی کاندید که به صورت K - گرم هستند $1 \leq k \leq l$ is $2^l - 1$ می باشد. در صورتی که K یک عدد بزرگ باشد، اندازه ویژگی‌ها می تواند بزرگ تر باشد. چون همه ویژگی‌ها به طور یکسان بر ای طبقه بندی اهمیت دارند، چانزولو (12) از تست گاما بر ای انتخاب یک زیر مجموعه ای از ویژگی‌ها با K - گرم استفاده کرده است. یک الگوریتم ژنتیکی بر ای یافتن زیر مجموعه بهینه از ویژگی‌ها ی محلی استفاده می شود.

بر عکس مجموعه ویژگی‌ها ی مبتنی بر K - گرم، لیش و همکاران (30-34) یک روش انتخاب ویژگی مبتنی بر مدل را ارائه کرده اند. این ویژگی‌ها به صورت قطعات توالی کوتاهی هستند که معیارهای زیر را دارند 1- حداقل در یک دسته دیده می شوند 2- حداقل در یک کلاس یا دسته متمایز هستند 3- اضافی نباشند. معیار (2) به معنی این است که ویژگی بایستی ارتباط و همبستگی معنی داری با یک کلاس داشته باشد. افزونگی در معیار 3 را می توان بر ای طبقه بندی ویژگی و تعمیم ویژگی استفاده کرد. یک الگوریتم ویژگی کلوی کار آمد، بر ای کاوش ویژگی‌ها بر طبق معیارها پیشنهاد می شود. نتایج آزمایشی در (30) نشان می دهد که مقایسه با روش در نظر گرفتن هر عنصر به عنوان ویژگی، انتخاب ویژگی مبتنی بر الگو موجب بهبود صحت تا 10 تا 15 درصد می شود.

چالش انتخاب ویژگی مبتنی بر الگو در توالی‌های نمادین ابزاری بر ای جست و جوی ویژگی‌ها مطابق با معیار هاست. جی و همکاران یک الگوریتم را بر ای کاوش توالی‌های با محدودیت فاصله پیشنهاد کرده اند. این الگوریتم که از عملیات بولین و چارچوب رشد پیشوند استفاده می کند، حتی با استانه فرکانس پایین کار آمد است سری‌های زمانی به صورت عددی هستند. روش‌های انتخاب ویژگی بر ای توالی‌های نمادین را نمی توان به آسانی به داده‌های سری‌های زمانی به کار برد. گسسته سازی می تواند موجب از بین رفتن اطلاعات شود. یی و همکاران (65) روش انتخاب ویژگی ای را پیشنهاد کردند که قابل کاربرد مستقیم به سری‌های زمانی عددی است. اشکال مختلف سری‌های زمانی نشان دهنده یک دسته خاص بوده و به عنوان یک ویژگی بر ای طبقه بندی سری‌ها مورد استفاده قرار می گیرد. بر ای یک فرایند دو دسته ای، یک واحد به صورت بخشی از سری‌های زمانی است که بر ای تفکیک داده‌های آموزشی به دو بخش بر اساس فاصله از واحد استفاده شده و موجب بیشینه سازی بهره وری اطلاعات می شود. استانه فاصله و شکل از داده‌های آموزشی بر ای بهینه سازی اسفاده می شود. بر ای ایجاد یک دسته بند، فرایند انتخاب شکل با ساخت درخت تصمیم گیری تلفیق می شود

اگرچه زیر توالی ها جزو ویژگی های مفید می باشند، با این حال آنها قادر به توصیف ویژگی های محلی با توالی بلند هستند. اگر اول و همکاران (5) یک روشی را برای پوشش دادن ویژگی های جهانی و محلی توالی ها برای طبقه بندی توسعه داده اند. اگر اول و همکاران (5) تجزیه موجک را برای توصیف توالی نمادین بر روی وضوح های مختلف اصلاح کرده است. با ضرایب تجزیه متفاوت، موجک ها بیانگر تغییرات در دامنه ها و بازه های مختلف جهانی تا محلی است.

با استفاده از تجزیه موجک و دسته بندی مبتنی بر قاعده، روش تجزیه موجک عملکرد بهتری نسبت به دسته بندی نزدیک ترین همسایه بر روی مجموعه داده های توالی و مجموعه داده های توالی ژنومی دارد.

به طور خلاصه، روش های موجود از دیدگاه های زیر متفاوت می باشند:

- کدام معیار ها بایستی برای انتخاب ویژگی های نظیر منحصر به فرد، فرکانس و طول استفاده شوند؟
- در کدام زمینه انتخاب ویژگی منعکس کننده ماهیت متوالی توالی، محلی و جهانی می باشند؟
- آیا انطباق ها بایستی دقیق یا غیر دقیق با فاصله ها باشند؟
- آیا انتخاب ویژگی بایستی در فرایند ایجاد یک دسته بندی و یا مرحله پیش پردازش مجزا تلفیق شوند؟

2-2 دسته بندی مبتنی بر فاصله توالی

روش های مبتنی بر فاصله توالی، یک تابع فاصله ای را برای اندازه گیری تشابه بین جفت توالی ها ارائه می کنند. وقتی این توابع فاصله ای انتخاب شدند، می توان از روش های دسته بندی موجود نظیر دسته بندی نزدیک ترین همسایه K و SVM با هسته همسوی محلی برای دسته بندی توالی استفاده کرد.

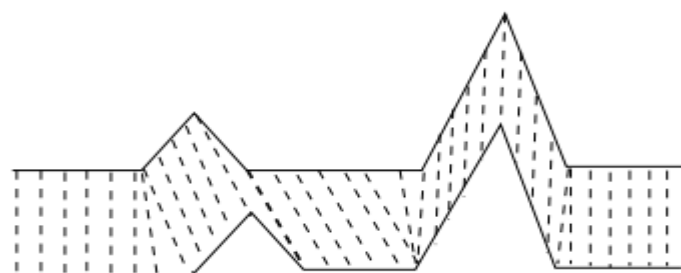
KNN یک روش یادگیری کند بوده و مدل دسته بندی را از قبل محاسبه نمی کند. با توجه به مجموعه داده های توالی T، یک عدد صحیح مثبت K و یک توالی جدید S، دسته بندی KNN، نزدیک ترین همسایه K از S را در T پیدا می کند و سپس برچسب دسته غالب را در KNN به صورت برچسب S قرار می دهد.

انتخاب شاخص های فاصله ای برای عملکرد دسته بندی KNN اهمیت زیادی دارند. در ادامه این بخش، ما بر خلاصه سازی شاخص های فاصله ای مختلف برای داده های توالی تاکید داریم. برای دسته بندی سری های زمانی ساده، فاصله اقلیدسی استفاده می شود. برای دو سری زمانی فاصله اقلیدسی به صورت زیر است

$$dist(s, s') = \sqrt{\sum_{i=1}^L (s[i] - s'[i])^2}.$$

فاصله اقلیدسی معمولاً مستلزم دو سری زمانی با طول یکسان است. کوگ و همکاران (26) نشان داده اند که هنگام کاربرد دسته بند 1NN بر روی سری های زمانی، فاصله اقلیدسی از حیث صحت در مقایسه با سایر شاخص های تشابه، رقابتی خواهد بود. فاصله اقلیدسی به شدت به اعوجاج و انحراف در بعد زمانی حساس است. فاصله تابیدگی زمانی پویا (DTW) برای حل این مسئله استفاده شده و نیازمند دو سری زمانی با طول یکسان نیست. ایده DTW همسوسازی سری های زمانی و فاصله ای را ارائه شده است. یک مثال از DTW در شکل 1 نشان داده شده است. زای و همکاران (61) نشان می دهد که در مجموعه ای از داده های کوچک، شاخص های الاستیک نظیر تابیدگی زمانی پویا صحیح تر از فاصله اقلیدسی است. با این حال نتایج تجربی اخیر نشان می دهد که بر روی مجموعه داده های بزرگ، صحت شاخص های الاستیک با فاصله اقلیدسی هم پوشانی دارد

تابیدگی زمانی پویا با برنامه ریزی پویا محاسبه شده و دارای پیچیدگی زمانی چند جمله ای است. از این روی، در مجموعه داده های بزرگ هزینه است. راتامانتا و همکاران (48) روشی را برای تسریع فرایند جست و جوی تشابه DTW با استفاده از محاسبات ارائه کردند. زای و همکاران (61) از کاهش عددی برای تسریع محاسبه DTW استفاده کرده است. ایده اصلی کاهش تعداد مثال های آموزشی مورد استفاده توسط دسته بند 1NN می باشد و در عین حال امکان تنظیم پویای پنجره تابیدگی وجود دارد



شکل 1: تابیدگی زمانی پویا

برای توالی های نما دین نظیر توالی های پروتینی و DNA، فواصل مبتنی بر همسوسازی معمولاً استفاده می شوند. با توجه به ماتریس تشابه و نیز جریمه فاصله ای، الگوریتم نیدل من یک امتیاز همسوسازی جهانی بهینه را بین دو توالی و دنباله از طریق برنامه نویسی پویا محاسبه می کند. برعکس الگوریتم های همسوسازی جهانی،

الگوریتم های همسو سازی محلی نظیر اسمیت واترمن و بلاست، تشابه را بین دو توالی با در نظر گرفتن شبیه ترین مناطق اندازه گیری می کنند

2-3 ماشین بردار پشتیبان

SVM یک روش موثر بر ای دسته بندی توالی [13; 52; 55; 54; 35; 39; 43] می باشد. ایده اصلی استفاده از SVM بر روی داده های توالی، تهیه نقشه توالی بر روی فضای ویژگی و یافتن ابر صفحه با حاشیه حداکثر بر ای تفکیک دو کلاس است. گاهی اوقات، ما نیازی به انتخاب ویژگی صریح نداریم. یک تابع هسته متناظر با یک فضای ویژگی با بعد بالا است. با توجه به دو توالی، X, Y برخی از توابع هسته ای، $K(x; y)$ را می توان به صورت تشابه بین دو توالی در نظر گرفت. چالش های استفاده از SVM بر ای طبقه بندی توالی شامل شیوه تعریف فضای ویژگی و توابع هسته و شیوه تسریع محاسبه ماتریس های هسته است

یکی از رایج ترین هسته ها بر ای طبقه بندی توالی، هسته طیف K یا هسته رشته می باشد که توالی را به یک بردار ویژگی تبدیل می کند. لزی و همکاران (35) یک هسته طیف K را بر ای طبقه بندی پروتین پیشنهاد کرده اند. با توجه به الفبای آمینو اسیدی 20 مولفه ای، این طیف شامل همه دنباله های احتمالی با طول K می باشد که متشکل از عناصر و مولفه های القبا است. بر ای مثال اگر $K=3$ باشد، طیف K دارای ARN, AND, DCN می باشد. با توجه به الفبا، A یک دنباله X می باشد که به فضای ویژگی از طریق تابع ترانسفورماسیون تبدیل می شود

$$\Phi_k(x) = (\phi_a(x))_{a \in A^k}$$

که $\phi_a(x)$ تعداد دفعاتی است که a در x اتفاق می افتد. تابع هسته، محصول نقطه ای بردار های ویژگی است

$$K(x, y) = \Phi_k(x) \cdot \Phi_k(y)$$

با استفاده از الگوریتم درخت پسوند، $K(x, y)$ را می توان در زمان $O(kn)$ محاسبه کرد

لودهی و همکاران (43) یک هسته رشته ای را بر ای طبقه بندی متنی پیشنهاد کرده است. مشابه با هسته طیف k در (35)، هسته رشته ای از زیر توالی طول k استفاده کرده است. با استفاده از ضریب تجزیه با طول توالی ها در متن، فاصله به صورت جریمه کد گذاری می شود. تابع هسته یک حاصل نقطه ای از بردار های ویژگی بوده و به

طور کار آمد از طریق برنامه نویسی پویا محاسبه می شود. لسلوی و همکاران (36) هسته طیف k را برای مدیریت عدم انطباق توسعه داده اند. سانبرگ و همکاران (55) یک هسته با طیف k با عدم انطباق را پیشنهاد کرده اند. یکی از معایب روش های مبتنی بر هسته این است که تفسیر آن سخت است و کاربر آن به سختی می توان دانش آن را در کنار نتایج طبقه بندی حفظ کنند. سانبرگ یک روشی را برای یادگیری SVM های قابل تفسیر با استفاده از مجموعه ای از هسته های رشته پیشنهاد کرده اند. ایده اصلی، استفاده از یک ترکیب خطی از اوزان می باشد. هر هسته از یک مجموعه ویژگی های منحصر به فرد استفاده می کند. اوزان بیانگر اهمیت ویژگی هاست. پس از یادگیری SVM، کاربر آن می تواند اطلاعات ارزشمندی را در زمینه اهمیت ویژگی ها ارائه کند هسته رشته یا هسته با طیف k را می توان به عنوان یک روش مبتنی بر ویژگی در نظر گرفت. سیگو و همکاران (49) یک هسته همسوسازی محلی را برای دسته بندی توالی پروتین ارائه کرده است که به عنوان یک روش مبتنی بر فاصله در نظر گرفته می شود. اگرچه فاصله همسوسازی محلی به طور موثر قادر به توصیف تشابه بین توالی هاست، با این حال به طور مستقیم به عنوان تابع هسته ای در نظر گرفته نمی شود زیرا فاقد ویژگی معین بودن مثبت است. سیگو و همکاران (49)، فاصله همسویی محلی را اصلاح کرده و یک هسته معتبر را موسوم به هسته همسوسازی محلی ایجاد کردند که از رفتار همسوسازی محلی پیروی می کند. ارتباط نظری بین هسته همسوسازی محلی و فاصله همسوسازی محلی اثبات می شود. با توجه به دو توالی X و Y ، هسته همسوسازی محلی را می توان با برنامه نویسی پویا محاسبه کرد.

هسته های دیگر مورد استفاده بر ای طبقه بندی توالی شامل هسته های چند جمله ای، هسته های برگرفته از مدل احتمال گرایی و هسته های انتشار می باشند.

4-2 دسته بندی مبتنی بر مدل

یک مقوله دیگر از روش های دسته بندی توالی، بر مبنای مدل های تولید کننده است که فرض می کند توالی ها در یک دسته از طریق مدل m تولید می شود. با توجه به یک دسته ای از توالی ها، M توزیع احتمال اولی ها را در یک دسته مدل سازی می کند. این مدل بر اساس فرضیات خاص تعریف می شود. در مرحله آموزش، پارامترهای M یاد گرفته می شوند. در مرحله دسته بندی، توالی جدید به یک کلاس با بیشترین احتمال نسبت داده می شود.

ساده ترین مدل، دسته بند توالی بیز ساده می باشد. فرض بر این است که با توجه به یک کلاس یا دسته، ویژگی ها در دنباله ها و توالی ها مستقل از یک دیگر می باشند. احتمالات شرطی ویژگی ها در یک دسته در مرحله آموزشی یاد گرفته می شوند. به دلیل سادگی خود، بیز ساده به طور گسترده توسط دسته بندی تصویر و دسته بندی توالی ژنومی استفاده شده اند

با این حال، فرض مستقل بودن در بیز ساده در عمل نقض می شود. مدل مارکوف و مدل مارکوف پنهان قادر به مدل سازی وابستگی میان عناصر در توالی هاست

یاک شنوف و همکاران (64) از مدل مارکوف رتبه K بر ای دسته بندی پروتین و داده های توالی متنی استفاده کردند. در فرایند آموزش این مدل در شرایط متفاوت و افتراقی بر ای افزایش قدرت دسته بندی روش های مبتنی بر مدل آموزش داده می شود. متفاوت از مدل مارکوف، مدل مارکوف پنهان فرض می کند که سیستم مدل سازی شده یک فرایند مارکوف با حالت های مشاهده نشده است. سریوستاوا و همکاران (56) از پروفیل HMM بر ای دسته بندی توالی های زیستی استفاده کرده است. یک پروفیل HMM دارای سه نوع حالت، افزایش، انطباق و حذف می باشد. مثال های آموزشی بر ای یادگیری احتمال تغییر بین حالت ها و احتمالات حذف استفاده می شود. HMM یاد گرفته شده بیانگر پروفیل مجموعه داده های آموزشی است. یک پروفیل HMM را می توان از توالی های غیر همسو با همسو سازی تدریجی هر مثال با پروفیل موجود یاد گرفت. بر ای هر دسته، یک پروفیل HMM یاد گرفته می شود. در مرحله دسته بندی یک توالی با HMM در هر دسته با برنامه نویسی پویا همسو سازی می شود یک توالی مجهول به یک دسته بالاترین امتیاز همسویی دسته بندی می شود.

3- نسخه های مختلف دسته بندی توالی

در این بخش، ما به مرور مسائل مربوط به طبقه بندی و دسته بندی توالی می پردازیم. این مسائل بر ای حل چالش ها در زمان استفاده از طبقه بندی توالی در زمینه های مختلف استفاده می شوند. از این روی از پیشوند های بر ای دست یابی به دسته بندی اولیه، طبقه بندی توالی ها با استفاده از هر دو داده های دارای برچسب و بدون برچسب و پیش بینی توالی برچسب ها به جای یک برچسب بر ای توالی های استریمینگ استفاده می کنیم

3-1 دسته بندی اولیه

برای توالی‌های نمادین موقت و سری‌های زمانی، مقادیر توالی به ترتیب صعودی مهر زمانی در نظر گرفته می‌شود. گاهی اوقات پایش و دسته‌بندی توالی‌ها در اسرع وقت مطلوب است. برای مثال، در مطالعه کودکان پذیرفته شده در بخش مراقبت‌های ویژه، نتایج نشان داد که کودکان دارای الکوی سری زمانی ضربان قلب غیر عادی بودند. به عنوان مثال دیگر برنالی و همکاران (9) نشان داده‌اند که تنها مشاهده پنج بسته اول TCP را می‌توان دسته‌بندی کرد. استفاده از ترافیک آنلاین را می‌توان بدون انتظار برای پایان یافتن جریان TCP شناسایی کرد. تا آنجا که ما می‌دانیم، دیز و همکاران (14) اولین بار مفهوم دسته‌بندی اولیه سری‌های زمانی را ارائه کرده‌اند. آن‌ها به توصیف سری‌های زمانی از طریق عباراتی نظیر افزایش و ماندگاری و معمولاً و همیشه پرداختند. به این ترتیب ادا بوست (19) برای شبیه‌سازی دسته‌بندی اصلی استفاده می‌شوند. این دسته‌بندی قادر به انجام پایش بینی بر روی داده‌های ناقص با مشاهده پسوند‌های غیر قابل دسترس توالی‌ها می‌باشد

انیبال و همکاران (8) از روش استدلال مبتنی بر مورد برای طبقه‌بندی سری‌های زمانی جهت پایش خرابی سیستم در سیستم پویا استفاده کردند. دسته‌بندی KNN بر روی دسته‌بندی سری‌های زمانی ناقص با فواصل مختلف استفاده می‌شود نظیر فاصله اقلیدسی و تائیدگی زمانی پویا. مطالعات شبیه‌سازی با استفاده از استدلال مبتنی بر مورد نشان داده‌اند که مهم‌ترین افزایش صحت طبقه‌بندی بر روی پیشوند‌ها از طریق سی تا چهل درصد طول کل اتفاق می‌افتد

اگرچه در (8-14)، اهمیت دسته‌بندی اولیه بر روی سری‌های زمانی شناسایی می‌شود و برخی از نتایج مفید حاصله شده است، این مطالعه دسته‌بندی اولیه را به صورت دسته‌بندی پایش‌وند‌های توالی‌ها در نظر می‌گیرد. زینگ و همکاران (62) بر چالش‌های طبقه‌بندی اولیه تأکید کردند که آن‌هم مطالعه تعادل بین صحت طبقه‌بندی و زمان است. روش‌های پیشنهادی تنها بر انجام پایش بینی بر اساس اطلاعات جزئی متمرکز است، با این حال به بررسی شیوه انتخاب کوتاه‌ترین پیشوند برای ارائه پایش بینی مطمئن تأکید دارد. این موجب می‌شود تا نتایج طبقه‌بندی اولیه را نتوان به خوبی مورد استفاده قرار داد.

زینگ و همکاران (62) مسئله طبقه‌بندی اولیه را به صورت طبقه‌بندی توالی ضمن حفظ صحت مورد انتظار فرموله کردند. یک روش مبتنی بر ویژگی‌های بر روی طبقه‌بندی اولیه در توالی‌های نمادین پیشنهاد می‌شود. ایده اصلی، انتخاب مجموعه‌ای از ویژگی‌های منحصر به فرد و سپس ایجاد یک رابطه بین دسته‌بندی قاعده و یک

دسته بند درخت تصمیم گیری با استفاده از این ویژگی ها است. در مرحله دسته بندی یک توالی با همه قواعد و شاخه ها به طور هم زمان منطبق می شود تا زمانی که یک پیشوند پیدا شده و توالی طبقه بندی شود. به این ترتیب، یک توالی فوراً زمانی طبقه بندی می شود که صحت و دقت مورد انتظار کاربر حاصل شده باشد. روش های پیشنهادی در 62 برخی موفقیت ها را در مدیریت توالی های نمادین با دست یابی به صحت رقابتی با استفاده از نیمی از طول توالی ها نشان داده اند.

یکی از معایب روش ها در (62) این است که قادر به مدیریت سری های زمانی عددی نیست. چون سری های زمانی عددی بایستی به طور آنلاین طبقه بندی شوند، از دست رفت اطلاعات موجب می شود تا نتوان برخی از ویژگی های منحصر به فرد را پوشش داد. زینگ و همکاران (63) یک دسته بند اولیه را برای سری های زمانی عددی با استفاده از یادگیری مبتنی بر فاصله ارائه کردند. این روش، طول پیش بینی حداقل را برای هر سری زمانی در مجموعه داده های آموزشی از طریق خوشه بندی یاد گرفته و از MLP برای طبقه بندی استفاده می کند. همان طور که در بخش 2 نشان داده شده است، دسته بند 1NN با فاصله اقلیدسی، یک دسته بند صحیح برای دسته بندی سری های زمانی است. یک ویژگی جالب این روش این است که بدون نیاز به صحت مورد انتظار کاربر، دسته بند موجب فعال سازی دسته بندی شده ضمن این که صحت را نیز حفظ می کند.

2-3 دسته بندی توالی نیمه نظارت شده

تعداد داده های بدون شناسه و برچسب بیش از داده های دارای شناسه است. برخی از داده های بدون شناسه دارای ویژگی های مشترکی با داده های با شناسه هستند و در عین حال حاوی ویژگی های اضافی ای هستند که توصیف جامع تری را از یک دسته در اختیار می گذارند. از این روی، با استفاده از داده های بدون برچسب و شناسه، یک دسته بند صحیح تر را می توان تولید کرد.

برای طبقه بندی و دسته بندی متن، طیف وسیعی از داده های بدون برچسب و بدون شناسه وجود دارند. نیگام و همکاران (46) یک روش دسته بندی نیمه نظارت شده را برای نشانه گذاری و برچسب گذاری اسناد پیشنهاد کرده اند. در ابتدا، دسته بند ساده بیز برای دسته بندی نمونه های بدون برچسب پیشنهاد شد. سپس، فرآیند پیشینه سازی انتظار برای تعدیل پارامترهای دسته بند بیز ساده و دسته بندی مجدد داده های بدون برچسب

در یک تکرار پیشنهاد شد. این فرآیند زمانی به پایان می‌رسد که نتایج دسته‌بندی پایدار باشد. یک سند ممکن است متعلق به چندین دسته باشد و یا دارای برچسب‌های متعددی باشد.

علاوه بر طبقه‌بندی متن، زانگ و همکاران (66) یک طبقه‌بندی نیمه نظارت شده مبتنی بر HMM را برای داده‌های سری‌های زمانی پیشنهاد کرده‌اند. این روش از داده‌های دارای برچسب برای آموزش پارامترهای اولیه HMM مرتبه اول استفاده کرده و سپس از داده‌های بدون برچسب برای تعدیل مدل در یک فرآیند EM بهره می‌برد. وی و همکاران از دسته‌بندی نزدیک‌ترین همسایه برای طبقه‌بندی سری‌های زمانی نیمه نظارت شده استفاده کردند. این روش برای مدیریت شرایطی استفاده می‌شود که در آن‌ها تنها یک تعداد کمی از داده‌های دارای برچسب در یک دسته مثبت موجود می‌باشد. در مرحله آموزش، در ابتدا، همه داده‌های بدون برچسب به صورت منفی در نظر گرفته می‌شود. سپس دسته‌بندی 1NN برای دسته‌بندی داده‌های بدون برچسب در یک تکرار استفاده می‌شود.

3-3 دسته‌بندی توالی با یک توالی و دنباله‌ای از برچسب‌ها

همان‌طور که در بخش 1 بحث شد، برای استریمینگ و ساده‌سازی طبقه‌بندی توالیف به جای پیش‌بینی یک برچسب دسته، پیش‌بینی توالی برچسب‌ها مطلوب‌تر است. کادوس این مسئله را به صورت طبقه‌بندی توالی قوی در نظر می‌گیرد ولی یک راه حل را برای این مسئله ارائه نمی‌کند.

یک مسئله مرتبط در پردازش زبان طبیعی، موسوم به توالی‌های برچسب‌زنی است. وظیفه این فرآیند برچسب‌زنی هر عنصر در یک توالی است. برای مثال، با توجه به یک توالی، که در آن هر کلمه به صورت یک عنصر در نظر گرفته می‌شود، برچسب‌زنی توالی یک کلمه را به یک مقوله نسبت می‌دهد نظیر هویت، عبارت، فعل یا اسم. یک راه حل پیشرفته، در نظر گرفتن برچسب‌های عناصر در یک توالی مرتبط با یک دیگر است. مسئله برچسب‌زنی با استفاده از زمینه‌های تصادفی شرطی حل می‌شود. این مسئله با سایر روش‌ها نظیر استفاده از مدل ترکیبی HMM و SVM و شبکه عصبی تکراری حل شده است.

4- کاربرد دسته‌بندی توالی

دسته بندی توالی دار ای طیف وسیعی از کاربرد ها می باشد. بر ای حوزه ها ی کاربردی مختلف فر ایند دسته بندی دار ای ویژگی ها ی متفاوتی است. در این بخش، ما روش ها ی اصلی به کار برده شده در زمینه ها ی مختلف را خلاصه می کنیم.

1-4 داده های ژنومیک

در سال ها ی اخیر، حجم زیادی از توالی ها ی دی ان ای و پروتین در دیتابیس ها ی عمومی نظیر بانک ژنی، دیتابیس توالی نوکلوتیدی EMBL و دیتابیس پروتین انرتز موجود بوده است. بر ای درک وظایف ژن ها و پروتین ها، دسته بندی توالی توجه زیادی را در تحقیقات ژنومی به خود معطوف کرده است.

روش ها ی مبتنی بر ویژگی به طور گسترده ای بر ای طبقه بندی توالی ژنوم استفاده شده اند. بر ای یک توالی نمادین، ساده ترین راه، در نظر گرفتن هر عنصر به صورت یک ویژگی است. بر ای مثال، یک توالی C ACG را می توان به صورت یک بردار ACCG در نظر گرفت. با این حال، ماهیت متوالی توالی ها را نمی توان با این تبدیل توجیه کرد. بر ای حفظ ترتیب عناصر در یک توالی، یک قطعه کوچک از نماد ها ی متوالی K موسوم به K- گرم معمولا به عنوان یک ویژگی انتخاب می شود. با توجه به مجموعه ای از این K- گرم ها، یک توالی را می توان به صورت بردار حضور یا غیاب K- گرم تعیین کرد. گاهی اوقات امکان تطبیق غیر دقیق با K- گرم وجود دارد. با استفاده از ویژگی ها ی K- گرم، توالی ها را می توان با روش طبقه بندی سنتی نظیر SVM تعیین کرد. یک خلاصه ای از روش ها ی انتخاب ویژگی مبتنی بر K- گرم بر ای طبقه بندی توالی را می توان در (16) یافت. اندازه ویژگی ها ی کاندید که به صورت K- گرم هستند $1 \leq k \leq l$ می باشد. در صورتی که K یک عدد بزرگ باشد، اندازه ویژگی ها می تواند بزرگ تر باشد. چون همه ویژگی ها به طور یکسان بر ای طبقه بندی اهمیت دارند، چانزولو (12) از تست گاما بر ای انتخاب یک زیر مجموعه ای از ویژگی ها با K- گرم استفاده کرده است. یک الگوریتم ژنتیکی بر ای یافتن زیر مجموعه بهینه از ویژگی ها ی محلی استفاده می شود. بر عکس مجموعه ویژگی ها ی مبتنی بر K- گرم، لیش و همکاران (30-34) یک روش انتخاب ویژگی مبتنی بر مدل را ارائه کرده اند. این ویژگی ها به صورت قطعات توالی کوتاهی هستند که معیار ها ی زیر را دارند 1- حدقل در

یک دسته دیده می شوند 2- حداقل در یک کلاس یا دسته متمایز هستند 3- اضافی نباشند. معیار (2) به معنی این است که ویژگی بایستی ارتباط و همبستگی معنی داری با یک کلاس داشته باشد. افزونگی در معیار 3 را می توان بر ای طبقه بندی ویژگی و تعمیم ویژگی استفاده کرد. یک الگوریتم ویژگی کاوی کار آمد، بر ای کاوش ویژگی ها بر طبق معیار ها پیشنهاد می شود. نتایج آزمایشی در (30) نشان می دهد که مقایسه با روش در نظر گرفتن هر عنصر به عنوان ویژگی، انتخاب ویژگی مبتنی بر الگو موجب بهبود صحت تا 10 تا 15 درصد می شود.

دسفنده و همکاران (13) به مقایسه عملکرد SVM، HMM و KNN بر ای دسته بندی داده های توالی ژنومی پرداختند. آن ها نشان داده اند که SVM عملکرد بهتری از سایرین دارد و انتخاب ویژگی نقش مهمی در تعیین صحت دسته بندی های SVM ایفا می کند. شی و همکاران (52) نیز به این نتیجه رسیده اند که SVM موثر ترین روش بر ای طبقه بندی پروتین است. علاوه بر صحت، چالش های دیگر دسته بندی توالی ژنومی، شامل تسریع طبقه بندی بر ای مدیریت داده های حجیم است.

4-2 داده های سری های زمانی

داده های سری های زمانی نوع مهمی از داده های توالی یا دنباله ای می باشند. در کتابخانه داده های سری های زمانی، داده های سری های زمانی در 22 حوزه از جمله کشاورزی، شیمی، سلامت، مالی، صنعت و غیره جمع اوری شده اند داده های سری های زمانی UCR مجموعه ای از دیتاست ها را به عنوان معیار ارزیابی روش های دسته بندی سری های زمانی ارائه می کند

بر ای داده های سری های زمانی ساده، بر ای استفاده از روش های مبتنی بر ویژگی، انتخاب ویژگی یک فرایند چالش بر انگیز می باشد زیرا ما قادر به شمارش ویژگی ها با داده های عددی نیستیم. از اینرو روش های فاصله ای بر ای طبقه بندی سری های زمانی استفاده می شوند. نشان داده شده است که در مقایسه طیف وسیعی از دسته بندی ها، SVM، HMM و نزدیک ترین همسایه دار ای اهمیت زیادی می باشند. به منظور استفاده از روش های مبتنی بر ویژگی در سری های زمانی ساده، قبل از انتخاب ویژگی، داده های سری های زمانی بایستی به توالی های نمادین از طریق تبدیل نمادین تبدیل شوند. در مقایسه با روش های فاصله ای، روش های مبتنی بر ویژگی موجب تسریع فرایند دسته بندی شده و یک سری نتایج قابل تفسیر را ارائه می کند. روش های مبتنی بر مدل بر ای دسته بندی سری های زمانی ساده نظیر HMM استفاده می شوند.

طبقه بندی سری های زمانی چند متغیره بر ای تشخیص حرکت و اشاره استفاده شده است. داده های چند ماغیره از طریق مجموعه ای از سنسور ها بر ای اندازه گیری حرکت اشیا در مناطق مختلف تولید می شوند. بر ای دسته بندی سری های زمانی چند متغیره، کادوس و همکاران یک دسته بند مبتنی بر ویژگی را پیشنهاد کرده اند. برخی از فرآیندهای تواریخ روندهای افزایشی یا کاهشی، سری های زمانی چند متغیره استفاده می شوند. لی و همکاران روشی را بر ای تبدیل سری های زمانی چند متغیره به یک بردار از طریق تجزیه ارزش عددی و تبدیلات دیگر ارائه کرده اند. SVM بر ای دسته بندی بردار ها استفاده می شود

3-4 داده های متنی

دسته بندی توالی به طور گسترده ای در بازیابی اطلاعات بر ای دسته بندی متن و اسناد استفاده می شود. رایج ترین روش ها بر ای دسته بندی متن شامل بیز ساده، و SVM می باشند. دسته بندی متن دارای انواع مختلفی می باشد که شامل دسته بندی چند برچسبی، دسته بندی سلسله مراتبی، و دسته بندی متن نیمه نظارت شده است. شبستانی و همکاران، یک مطالعه دقیق را بر روی دسته بندی متنی ارائه کرده اند.

5- نتیجه گیری

در این مقاله ما یک مروری بر طبقه بندی و دسته بندی متن داشتیم. داده های توالی به 5 زیر نوع طبقه بندی شد. روش های دسته بندی توالی در روش های مبتنی بر ویژگی، روش های مبتنی بر فاصله توالی و روش های مبتنی بر مدل طبقه بندی شد. سپس دسته بندی توالی متعارف نیز بحث شد. در نهایت، روش های مختلف در حوزه های کاربردی مختلف مقایسه شدند

نتایج نشان می دهد که بیشتر کارها بر دسته بندی توالی های نمادین ساده و سری های زمانی ساده تاکید دارند. اگرچه مطالعات کمی بر روی سری های زمانی چند متغیره وجود دارد مسئله طبقه بندی داده های توالی هنوز حل نشده است. به علاوه بیشتر روش ها به دسته بندی توالی متعارف تخصیص داده شده است. استریمینگ طبقه بندی توالی، اولیه، نیمه نظارت شده و ترکیبی از این مسائل بر روی داده های توالی پیچیده که دارای کاربرد های زیادی می باشند می توانند چالشی بر ای مطالعات آینده باشند.

این مقاله، از سری مقالات ترجمه شده رایگان سایت ترجمه فا میباشد که با فرمت PDF در اختیار شما عزیزان قرار گرفته است. در صورت تمایل میتوانید با کلیک بر روی دکمه های زیر از سایر مقالات نیز استفاده نمایید:

لیست مقالات ترجمه شده ✓

لیست مقالات ترجمه شده رایگان ✓

لیست جدیدترین مقالات انگلیسی ISI ✓

سایت ترجمه فا ؛ مرجع جدیدترین مقالات ترجمه شده از نشریات معتبر خارجی