

تکامل تفاضلی با ترکیب شبکه عصبی برای استخراج قوانین از پایگاه داده‌های پزشکی

مهدی جبه داری^۱، حمید پایگذار^۲

^۱ گروه کامپیوتر، دانشکده فنی و مهندسی، واحد خمین، دانشگاه آزاد اسلامی، خمین، ایران

^۲ گروه کامپیوتر، دانشکده فنی و مهندسی، واحد خمین، دانشگاه آزاد اسلامی، خمین، ایران (نویسنده مسئول)

چکیده

امروزه در علم پزشکی سیستم‌های پشتیبانی برای تصمیمات بالینی تعریف شده‌اند که با استفاده از دو یا چند آیتیم از ویژگی‌های بیمار، پیشنهاد‌های مفیدی ارائه می‌دهند. این پیشنهادها به متخصص در تشخیص بیماری یا روند درمان کمک می‌کنند. طبقه‌بندی داده‌ها در این حوزه نیازمند شفافیت و ارائه‌ی علت برای تخصیص داده‌ها می‌باشد. از این رو لازم است علاوه بر طبقه‌بندی درست داده‌های پزشکی، علت تخصیص هر داده شرح داده شود. در این تحقیق، یک روش جدید بر پایه‌ی تکامل تفاضلی برای طبقه‌بندی خودکار آیتیم‌ها در پایگاه داده‌های پزشکی موردنظر است. بر اساس آن، یک ابزار به نام DERE_X معرفی شده است که به‌طور خودکار، دانش واضح را از پایگاه داده، به شکلی از قوانین IF-THEN (اگر-سپس) شامل شرط‌های متصل به AND (و) بر روی متغیرهای پایگاه داده، استخراج می‌کند. هر فرد از تکامل تفاضلی یک قانون برای طبقه‌بندی محسوب می‌شود. مجموعه افراد تکامل تفاضلی مجموعه‌ای از قوانین می‌باشند که از تمام آن‌ها برای طبقه‌بندی داده‌های موجود در پایگاه داده استفاده می‌شود. این قوانین می‌توانند به‌صورت منطقی متصل به OR (یا) نیز در نظر گرفته شوند؛ بنابراین تمام قوانین طبقه‌بندی برای تمام دسته‌ها یک‌باره در یک مرحله پیدا می‌شوند. ابزار اولین ابزار طبقه‌بندی است که بر پایه‌ی تکامل تفاضلی می‌باشد و به‌طور خودکار و بدون دخالت سازوکار دیگری دسته‌هایی از قوانین IF-THEN را استخراج می‌کند.

واژه‌های کلیدی: طبقه بندی، تکامل تفاضلی، پایگاه داده، استخراج قوانین.

۱- مقدمه

در حال حاضر داده‌کاوی مهم‌ترین فناوری جهت بهره‌برداری مؤثر از داده‌های حجیم است و اهمیت آن رو به افزایش است. چند سال پیش تخمین زده شده بود که مقدار داده‌ها در جهان هر ۲۰ ماه به حدود دو برابری می‌رسد. در یک تحقیق مربوط به سال ۲۰۰۴ که بر روی گروه‌های تجاری بسیار بزرگ در جمع‌آوری داده‌ها صورت گرفت مشخص گردید که ۱۹ درصد از این گروه‌ها دارای پایگاه داده‌هایی با سطح بیشتر از ۵۰ گیگابایت می‌باشند و ۵۹ درصد از آن‌ها انتظار دارند که در آینده‌ای نزدیک در چنین سطحی قرار گیرند (هلت^۱، ۱۹۹۳). لذا در چنین شرایطی که پیشرفت‌های فناوری امکان ذخیره‌سازی این حجم عظیم از داده‌ها را برای ما فراهم آورده دیگر بعید است که بتوان با روش‌های معمول و سنتی پردازش اطلاعات به تمام واقعیات نهفته در داده‌ها دست یافت. لذا داده‌کاوی برای یاری در این عرصه با روش‌ها و تکنیک‌های جدید و رو به بهبود به میدان آمده است. بهادرا^۲ و همکاران (۲۰۱۲) یک الگوریتم توسعه‌یافته از بهینه‌سازی پارامترهای الگوریتم ماشین بردار پشتیبان^۳ (SVM) با استفاده از تکامل تفاضلی ارائه داده‌اند که دقت طبقه‌بندی را نسبت به استفاده انفرادی از SVM افزایش می‌دهد. اخیراً SVM به‌عنوان تکنیکی مؤثر برای حل مسائل طبقه‌بندی‌های واقعی، توجه زیادی را به خود جلب نموده است. به‌هرحال عملکرد SVM حساس به انتخاب تابع اساسی مناسب و مقادیر پارامترها به کمک تابع اساسی است. جیانث^۴ و همکاران (۲۰۱۵)، داده‌های گرفته شده توسط ماهواره‌های فضایی را با استفاده از یک الگوریتم کلونی زنبور مصنوعی انجام داده‌اند. علت استفاده از آن‌ها این روش جدید این است که روش‌های آماری قدیمی موفقیت محدودی در طبقه‌بندی داده‌ها ارائه می‌دهند و بسیاری از جنبه‌های اطلاعات مانند همبستگی را در نظر نمی‌گیرند. الگوریتم کلونی زنبور مصنوعی^۵ (ABC) یک الگوریتم مبتنی بر هوش جمعی است. این الگوریتم برای داده‌های حس شده از راه دور توسط ماهواره‌ها استفاده شده است و نتایج نشان می‌دهند که ۵٪ افزایش دقت نسبت به روش سنتی آماری یعنی طبقه‌بند درستی ماکزیمم^۶ (MLC) و شبکه‌های عصبی مصنوعی و همچنین ۳٪ افزایش دقت نسبت به الگوریتم SVM دارد.

وینودهینی^۷ و همکاران (۲۰۱۴) یک ارزیابی عملکرد مقایسه‌ای برای روش‌های مبتنی بر شبکه‌های عصبی مصنوعی برای طبقه‌بندی عواطف انسانی در فضاهای آنلاین ارائه داده‌اند. کمک اصلی طبقه‌بندی عواطف تعیین هویت عواطفی است که در یک پیام متنی نهفته است. روش‌های یادگیری ماشین برای طبقه‌بندی عواطف انسانی به‌طور گسترده‌ای مورد مطالعه قرار گرفته‌اند ولی روش‌هایی که بر اساس شبکه‌های عصبی مصنوعی باشند به ندرت در طبقه‌بندی عواطف مورد استفاده قرار گرفته‌اند. سانینو^۸ و همکاران (۲۰۱۴)، یک روش استخراج خودکار قوانین با استفاده از تکامل تفاضلی برای نظارت بر روی آپنه انسدادی خواب^۹ (OSA) توسط یک سیستم موبایل ارائه داده‌اند. چن^{۱۰} و همکاران (۲۰۱۳) در مقاله‌شان روشی برای بهبود عملکرد ماشین بردار پشتیبان دوقلوی لاپلاسی^{۱۱} (Lap-TSVM) با استفاده از تکامل تفاضلی ارائه داده‌اند. ماشین بردار پشتیبان دوقلوی لاپلاسی اخیر تولید خوبی با حل یک جفت مسئله‌ی برنامه‌ریزی خطی (QPP) حاصل می‌کند. به‌هرحال فرایند آموزش Lap-TSVM وابسته به زمان است. علاوه بر آن درمقایسه با SVM، Lap-TSVM پارامترهای بیشتری برای تعدیل نیاز دارد که بر روی برنامه‌های عملی آن تأثیر می‌گذارد.

¹ Holte

² Bhadra

³ Support Vector Machine

⁴ Jayanth

⁵ Artificial Bee Clony

⁶ Maximum Legitimacy Classifier

⁷ Vinodhini

⁸ Sannino

⁹ occlusion Sleep Apne

¹⁰ - Chen

¹¹ - Laplasian Twin Support Vector Machine

۲-۱ فرایند کار DERE_x

با ترکیب تمام اطلاعات بالا یک پایگاه داده داده شده با N_V متغیر، DERE_x به عنوان شکل کلی نشان داده شده در شبه کد الگوریتم DE کار می کند. قبل از پرداختن به جزئیات DERE_x لازم است توجه شود که این ابزار می تواند فوراً خصیصه های پایگاه داده را که با مقادیر حقیقی نشان داده شده اند، بررسی کند. در این فرایند یک کدگذاری ساده برای خصیصه های بولی، بی قید و صحیح باید بررسی شود. یک خصیصه صحیح در یک محدوده صحیح برای مثال [10,35]، با یک مقدار حقیقی در محدوده ی مقادیر حقیقی [10.0,35.0] کدگذاری می شود. یک خصیصه بی قاعده می تواند به عنوان یک مورد صحیح کنترل شود. یک خصیصه بی قاعده می تواند مثل یک خصیصه صحیح که یک رابطه ی سفارشی را تأمین کرده، کنترل شود و می تواند توسط مقادیر احتمالی اش تعریف شود. ممکن است نمونه ای از یک خصیصه مثل سن = {جوان، میان سال، پیر} به صورت سن = {۱ و ۲ و ۳} کنترل شود. یک خصیصه بولی در [0,1] نیز به صورت عددی حقیقی در [0.0,1.0] کدگذاری می شود. سپس مقادیر پارامترهای DERE_x تعیین می شوند. این در میان دیگران اشاره دارد به بیشترین تعداد قوانین مطلوب یعنی N_R . در اینجا با دادن پایگاه داده ای با N_V متغیر، DERE_x به طور تصادفی جمعیتی از N_{pop} راه حل ممکن مقداردهی اولیه می کند. هر مقدار از یک توزیع تصادفی یکنواخت که دارای کران بالا و پایین تعیین شده ایست انتخاب می شود. در پایان اجرا قوانین پیداشده با استفاده از مکانیسم اصلی شرح داده شده به طور اتوماتیک کدگذاری می شوند. آن مکانیسم اصلی کدگذاری خصیصه های پایگاه داده مقدار حقیقی را شرح می دهد. برای مقادیر صحیح، مقادیر بی قاعده و خصیصه های بولی پایگاه داده در قوانین باید یک فرایند کدگذاری ثانویه و سراسر با شروع از مقدار حقیقی تأمین شده توسط DERE_x اجرا شود. برای تمام این گونه متغیرها این عمل بدین معناست که مقدار حقیقی را به نزدیک ترین مقدار صحیح گرد کنیم. برای مثال برای یک خصیصه صحیح، یک مقدار ۳۱.۵۷ که توسط DERE_x تولید شده به شکل ۳۲ کدگذاری می شود. در حالی که مقدار حقیقی ۳۱.۲۴ به شکل ۳۱ کدگذاری می شود. به طور مشابه برای خصیصه بی قاعده ی سن، یک مقدار ۲.۱۴ به عدد ۲ گرد می شود که به مقدار بی قاعده ی میان سال اشاره دارد. در نهایت برای یک خصیصه بولی مقدار حقیقی ۰.۱۴ به صورت صفر کدگذاری می شود در حالی که مقدار ۰.۷۴ به عنوان ۱ کدگذاری می شود. خلاصه ی DERE_x در این مراحل قابل ذکر است:

۱. ایجاد مجموعه ای از قوانین به صورت تصادفی نام مجموعه را $pop = \{X_1, X_2, \dots, X_{N_R}\}$ می گذاریم.
۲. برای هر قانون (X_i) از مجموعه ی pop روند زیر اجرا می شود:
 - یک قانون جدید (X'_i) با یکی از استراتژی های تکامل تفاضلی (به عنوان مثال، تصادفی) ایجاد می شود.
 - قانون X' را به جای X وارد مجموعه ی pop می کنیم و دقت طبقه بندی با استفاده از کل مجموعه را اندازه می گیریم. اگر دقت طبقه بندی افزایش یافت، قانون X'_i را به جای X_i قرار داده و قانون X_i را از مجموعه حذف می کنیم. در غیر این صورت قانون X_i در مجموعه باقی می ماند.
۳. مرحله ی دوم برای چندین بار (مثلاً ۱۰۰ دور) تکرار می شود. به عبارتی مجموعه ی pop چندین بار اصلاح می شود تا بهترین قوانین برای طبقه بندی به وجود بیایند.
۴. قوانین استخراجی از مرحله ی سوم به عنوان قوانین نهایی معرفی شده و ازین پس از این قوانین برای طبقه بندی داده های جدید استفاده خواهد شد.

۲-۲ آموزش شبکه عصبی مصنوعی با تکامل تفاضلی

همان طور که پیش تر گفته شد از شبکه های عصبی مصنوعی نیز می توان برای طبقه بندی داده ها استفاده نمود. شبکه های عصبی مصنوعی همانند دیگر روش های هوش مصنوعی نیازمند آموزش می باشند. داده هایی که برای ارزیابی صلاحیت در روش پیشنهادی با تکامل تفاضلی داده می شود، به شبکه عصبی مصنوعی نیز داده می شود و شبکه با تنظیم وزن ها و بایاس ها آموزش می بیند. در این قسمت برخلاف روش معمول، قصد داریم شبکه های عصبی مصنوعی را توسط تکامل تفاضلی

پیشنهادی آموزش داده و طبقه‌بندی را انجام دهیم. برای ارائه روش پیشنهادی یک شبکه عصبی چندلایه پرسپترون^۱ پیشخور^۲ در نظر گرفته شده است.

خروجی شبکه‌های عصبی پیشخور یک تابع از وزن‌های پیوندگاهی W و مقادیر ورودی X می‌باشد؛ یعنی $y=f(X,W)$. در پردازش‌های استاندارد آموزش، هم بردار ورودی X و هم بردار خروجی Y شناخته شده هستند و بردار وزن W برای دستیابی به یک مسیریابی درست برای رسیدن ورودی X به خروجی Y تنظیم می‌شود. به صورت عمومی وزن‌ها می‌توانند با حداقل-سازی تابع خطای شبکه E به دست آیند. تابع خطای شبکه به صورت زیر قابل تعریف است (ایلونن و همکاران^۳، ۲۰۰۳):

$$E(y, f(x, W)): (y^{D1}, x^{D2}, W^{D3}, f) \rightarrow R. \quad (1)$$

هدف بهینه‌سازی حداقل‌سازی تابع هدف E توسط بهینه‌سازی مقادیر وزن‌های شبکه است.

$$W = (w_1, \dots, w_D) \quad (2)$$

مشابه دیگر الگوریتم‌های تکاملی، تکامل تفاضلی مبتنی بر یک جمعیت (P_G) از راه‌حل‌های کاندیدا کار می‌کند. این راه‌حل‌های کاندیدا افراد جمعیت می‌باشند. تکامل تفاضلی یک جمعیت با اندازه ثابت NP ، بردارهایی با مقادیر حقیقی، $W_{i,G}$ که i اندیس جمعیت و G شماره تولیدی را نشان می‌دهد که جمعیت به آن متعلق است.

$$P_G = (w_{1,G}, \dots, w_{NP,G}), \quad G = 0, \dots, G_{max} \quad (3)$$

علاوه بر این در آموزش شبکه‌ی عصبی مصنوعی هر بردار شامل D وزن شبکه است (کروموزوم‌های افراد):

$$W_{i,G} = (w_{1,i,G}, \dots, w_{D,i,G}), \quad i = 1, \dots, NP, \quad G = 0, \dots, G_{max} \quad (4)$$

طرح تولید مجدد جمعیت تکامل تفاضلی با دیگر روش‌های تکاملی متفاوت است. پس از مقداردهی اولیه‌ی جمعیت نخست، بردارها در جمعیت جاری یعنی P_G به طور تصادفی نمونه‌برداری و ترکیب می‌شوند تا بردارهای کاندیدا را برای نسل بعدی یعنی P_{G+1} تولید کنند. جمعیت کاندیداها یا بردارهای آزمایشی $P'_{G+1} = U_{i,G+1} = u_{j,i,G+1}$ به صورت زیر تولید می‌شوند:

$$\begin{aligned} v_{j,i,G+1} &= w_{j,r3,G} + F \cdot (w_{j,r1,G} - w_{j,r2,G}) \\ u_{j,i,G+1} &= v_{j,i,G+1}, \text{ if } \text{rand}_j[0,1] \leq CR \\ &w_{j,i,G}, \text{ otherwise} \end{aligned} \quad (5)$$

به طوریکه:

$$\begin{aligned} i &= 1, \dots, NP, \quad j = 1, \dots, D \\ r_1, r_2, r_3 &\in [1, \dots, NP], \text{ randomly}, r_1 \neq r_2 \neq r_3 \neq i \\ CR &\in [0,1], \quad F \in (0,1). \end{aligned} \quad (6)$$

و CR عبارت است از مقدار حقیقی فاکتور گذردهی که احتمال پارامترهای بردار آزمایشی را کنترل می‌کند. عموماً مقادیر F و CR سرعت همگرایی و قدرت فرایند جستجو را تحت تأثیر قرار می‌دهند. مقادیر بهینه‌ی آن‌ها بستگی به ویژگی‌های تابع هدف و اندازه‌ی جمعیت، NP دارد.

طرح انتخاب تکامل تفاضلی نیز متفاوت از دیگر روش‌هاست. جمعیت برای تولید بعدی، P_{G+1} از جمعیت جاری P_G به صورت زیر انتخاب می‌شود:

$$\begin{aligned} W_{i,G+1} &= U_{i,G+1}, \text{ if } E(y, f(x, W_{i,G+1})) \leq E(y, f(x, W_{i,G})) \\ &W_{i,G}, \text{ otherwise} \end{aligned} \quad (7)$$

¹- Multi Layer Perceptron Neural Network

²- Feed Forward

³ Ilonen

بنابراین هر فرد از جمعیت موقت، با همتای خود در جمعیت جاری مقایسه می‌شود. با فرض حداقل‌سازی تابع هدف، برداری که مقدار کمتری برای تابع هدف ایجاد کند در مقایسه پیروز شده و در جمعیت قرار می‌گیرد؛ بنابراین افراد در جمعیت بعدی یا به خوبی جمعیت جاری هستند، یا از آن‌ها بهتر می‌باشند. نکته جالب در طرح جایگزینی DE این است که یک بردار آزمایشی فقط با یک فرد مقایسه می‌شود نه تمام افراد. این نکته از این بابت ارزشمند است که طرح جایگزینی اطمینان می‌دهد که جمعیت از بهترین پاسخ‌ها واگرا نمی‌شود (ایلون و همکاران، ۲۰۰۳).

۴- نتایج

۴-۱ پایگاه داده‌های مورد استفاده

در این تحقیق از سه پایگاه داده استاندارد پزشکی برای بررسی روش پیشنهادی استفاده شده است. این هشت پایگاه داده عبارت‌اند از:

۱. دیتابیس هابرمین شامل مواردی روی بقای بیماری است که برای سرطان سینه تحت جراحی است. این داده‌ها توسط دانشگاه بیمارستان بیلینگ شیکاگو جمع‌آوری شده است. هر آیتمی توسط سه مقدار روشن صحیح نشان داده شده که به ترتیب شامل: سن بیمار هنگام عمل، سال عمل بیمار (به صورت: سال منهای ۱۹۰۰) و تعداد گره‌های گوشه‌ای مثبت پیداشده، می‌شوند. هر آیتم به کلاس صفر اختصاص می‌یابد به شرطی که بیمار ۵ سال یا بیشتر زنده بماند و اگر بیمار در کمتر از ۵ سال فوت شود به کلاس یک اختصاص می‌یابد (هابرمین^۱، ۱۹۷۶).

۲. دیتابیس بریست^۲ موارد سرطان سینه را که توسط دکتر ویلیام اچ وولبرگ از بیمارستان دانشگاه ویسکانسین^۳ جمع‌آوری شده، نشان می‌دهد. کلاس صفر شامل موارد خوش‌خیم و کلاس یک شامل موارد بدخیم است (ولبرگ^۴، ۱۹۹۰).

۳. دیتابیس کلورلند^۵ ویژگی‌های بیماران با مرض قلبی را از مرکز پزشکی Long Beach, V.A و موسسه‌ی بالینی کلورلند لیست می‌کند. آیتم‌های آن هم می‌تواند به صورت ۲ کلاس و هم ۵ کلاس طبقه‌بندی شوند. در این پژوهش مورد پنج کلاس مورد استفاده قرار گرفته است. هدف، یافتن وجود بیماری قلبی در بیمار است. کلاس‌ها به صورت افزایش خطر از صفر (بدون بیماری) تا چهار (بیشترین اعلام خطر) طبقه‌بندی شده‌اند (دترانو و همکاران^۶، ۱۹۸۹).

۴-۲ نتایج اجرای DERE_x

الگوریتم برای هر سه پایگاه داده اجرا شده و نتایج با الگوریتم‌های طبقه‌بندی KNN و NB مقایسه گردیده که به شرح زیر می‌باشند:

۴-۲-۱ پایگاه داده Breast

برای اجرای DERE_x روی این پایگاه داده مقادیر RT، CT و CR هر سه را برابر ۰.۵ در نظر گرفته و تعداد قوانین تولیدی را ۱۰ در نظر گرفته‌ایم. همچنین تعداد دوره‌های تولید را ۱۰۰ دور قرار داده‌ایم. قوانین استخراج شده پس از اجرا به صورت زیر می‌باشند.

$if\ 1 \ \&\&\ 1 \ \&\&\ 1 \ \&\&\ 1 \ \&\&\ A5 \geq 0.29804 \ \&\&\ in(A6, 1.5961, 2.6723) \ \&\&\ 1 \ \&\&$
 $in(A8, 2.1922, 4.3447) \ \&\&\ 1 \ \&\&\ then\ CLASS = 1;$

¹ Haberman

² Breast

³ Wisconsin

⁴ Wolberg

⁵ Cleveland

⁶ Detrano et al.

if A1<15.5653 && A2==0.82261 && A3>-0.43475 && out(A4,-0.12868,13.7069) && A5== -0.088695 && A6<=0.82261 && A7>-0.088695 && 1 && A9== -0.088695 then CLASS= 1;

if 1 && 1 && in(A3,24.5059,77.5088) && 1 && A5<=0.41012 && A6<=1.8202 && out(A7,0.41012,1.4702) && 1 && A9>0.41012 then CLASS= 1;

if out(A1,32.275,67.9236) && 1 && 1 && A4>7.892 && 1 && A6>1.491 && out(A7,0.2455,0.95847) && in(A8,1.982,4.8339) && 1 then CLASS= 1;

if 1 && A2>1.3633 && A3>13.0828 && in(A4,6.3597,24.7892) && 1 && A6>1.3633 && 1 && A8>1.7266 && 1 then CLASS= 1;
else CLASS=0;

این پایگاه داده ۹ متغیر و دو کلاس (۰ و ۱) دارد. هر متغیر یکه در شرط لحاظ نشده به جای آن ۱ قرار گرفته است که در کل شرط بی تأثیر است.

۲-۲-۴ پایگاه داده هابرمین

برای اجرای DEREX روی این پایگاه داده مقادیر CT، RT و CR هر سه را برابر ۰.۵ در نظر گرفته و تعداد قوانین تولیدی را ۴ در نظر گرفته ایم. همچنین تعداد دوره های تولید را ۱۰۰ دور قرار داده ایم.
قوانین استخراج شده:

if 1 && in(A2,61.8457,66.2149) && A3>18.1797 then CLASS= 1;
else CLASS=0;

این پایگاه داده شامل سه متغیر و دو کلاس (۰ و ۱) است.

۳-۲-۴ پایگاه داده کلورلند

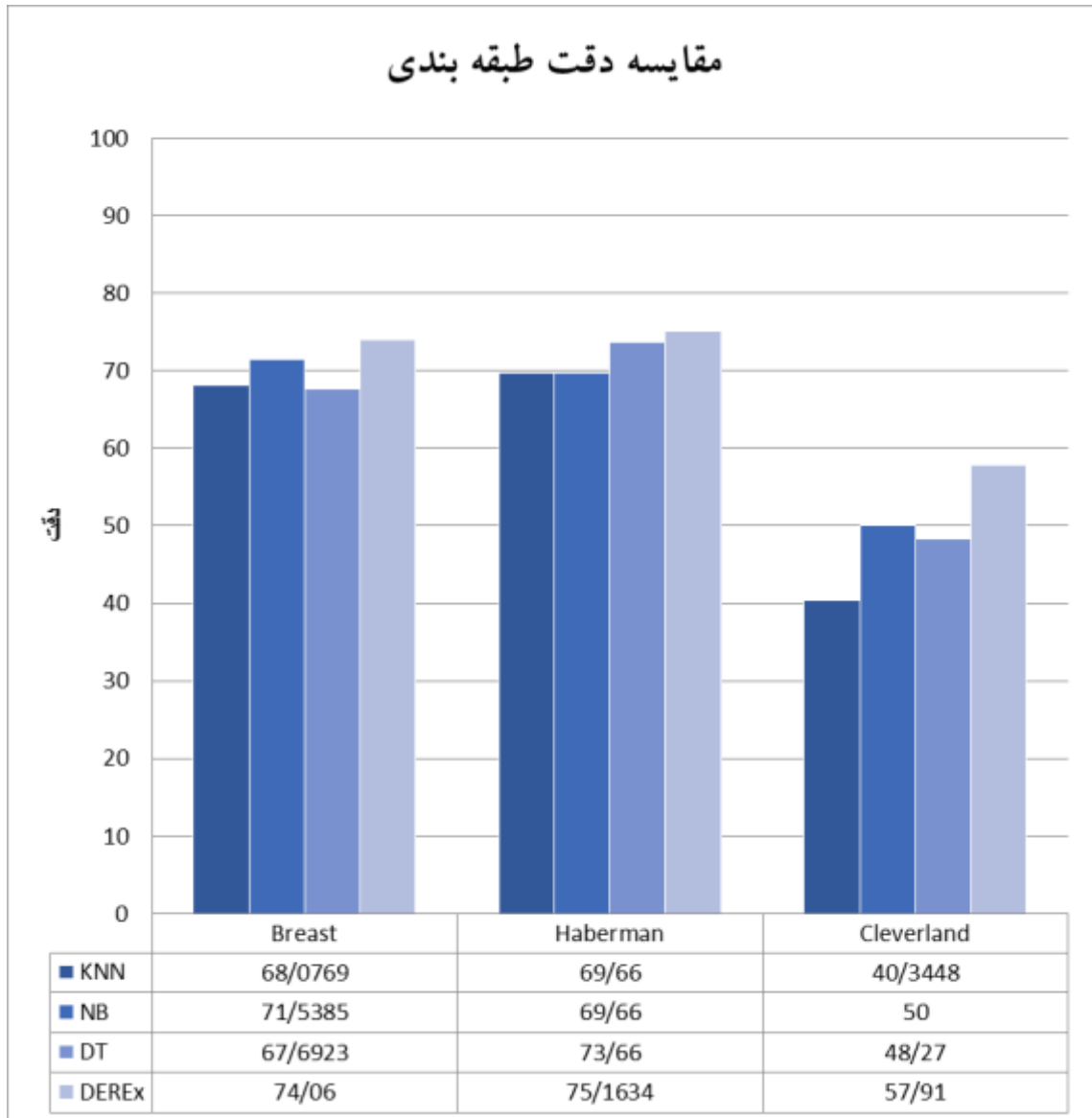
برای اجرای DEREX روی این پایگاه داده مقادیر CT، RT و CR هر سه را برابر ۰.۵ در نظر گرفته و تعداد قوانین تولیدی را ۱۵ در نظر گرفته ایم. همچنین تعداد دوره های تولید را ۱۰۰ دور قرار داده ایم.
قوانین استخراج شده:

if out(A1,35.1223,56.3106) && A2==0.12755 && 1 && 1 && 1 && 1 && 1 && 1 && 1 && 1 && A10>=0.79079 && out(A11,1.2551,2.1379) && 1 && A13<=3.5102 then CLASS= 1;
if A1>=46.3485 && out(A2,0.36143,0.40984) && 1 && out(A4,132.3114,137.4434) && A5>=284.3055 && A6<=0.36143 && 1 && A8<118.3471 && 1 && in(A10,2.2409,2.541) && 1 && out(A12,1.0843,1.2295) && 1 then CLASS= 1;
if 1 && 1 && A3>1.4445 && 1 && A5==190.8943 && 1 && 1 && A8>90.409 && A9>=0.14816 && 1 && A11<=1.2963 && out(A12,0.44448,2.3476) && A13>=3.5926 then CLASS= 1;
if 1 && 1 && 1 && 1 && A5==191.9729 && 1 && A7==0.30125 && 1 && 1 && 1 && 1 && A11<1.3012 && 1 && 1 then CLASS= 2;
if 1 && 1 && A3>1.7791 && 1 && 1 && A6>0.25972 && 1 && in(A8,105.0227,157.0504) && 1 && 1 && 1 && 1 && A13>=4.0389 then CLASS= 3;
if 1 && A2>=0.37964 && 1 && out(A4,134.2421,142.0223) && 1 && 1 && 1 && out(A8,120.7331,130.3483) && A9<0.37964 && out(A10,2.3538,2.8089) && 1 && A12>1.1389 && A13>4.5186 then CLASS= 3;

```

if out(A1,47.4131,53.7349) && 1 && 1 && A4>=134.6623 && A5<=294.0196 &&
out(A6,0.38361,0.51531) && 1 && A8<121.2524 && 1 && 1 && 1 && A12<=1.1508 &&
1 then CLASS= 3;
if A1>=31.2227 && A2>=0.046306 && out(A3,1.1389,3.5945) && 1 && 1 &&
A6<=0.046306 && A7>0.092611 && 1 && out(A9,0.046306,0.86482) && 1 && 1 && 1
&& A13>3.1852 then CLASS= 4;
else CLASS=0;
    
```

این پایگاه داده نیز دارای ۱۳ متغیر و ۵ کلاس (۰ و ۱ و ۲ و ۳ و ۴) می باشد. توجه شود که: نتایج بالا در هر سه پایگاه داده برای ۱۰۰ دور از تکامل تفاضلی اجرا شده است. با بیشتر کردن تعداد دورها طبیعتاً دقت افزایش خواهد یافت. نمودار دقت طبقه بندی DERE_x در مقایسه با سه الگوریتم دیگر برای هر سه پایگاه داده در شکل ۴-۱ قابل مشاهده است:



شکل ۲. مقایسه دقت طبقه بندی DERE_x با الگوریتم های دیگر

۵- جمع‌بندی

وظایف طبقه‌بندی در حوزه‌ی پزشکی با اشاره به فرایند درمان داری اهمیت هستند. در این حوزه و در بسیاری از حوزه‌های دیگر، یک ویژگی خوش‌آیند طبقه‌بندی کننده این است که برای کاربران دانسته‌های واضحی از پایگاه داده به‌طور اتوماتیک استخراج می‌کند. البته دانش به‌دست‌آمده هرگز نباید جانشین کار کارشناسانه شود، بلکه باید به‌عنوان پشتیبانی برای تصمیم‌گیری در نظر گرفته شود.

در این پژوهش، یک روش جدید مبتنی بر تکامل تفاضلی برای طبقه‌بندی خودکار آیت‌های پایگاه‌داده پیشنهاد شده است. بر اساس آن ابزاری که DERE_x نام دارد معرفی شده است که به‌طور خودکار دانش واضحی را از پایگاه‌داده‌به‌صورت قوانین IF-THEN متصل به شرط‌های AND روی متغیرها استخراج می‌کند. نتیجه برای مجموعه‌ای از قوانین کدگذاری می‌شود و این قوانین می‌توانند به‌صورت متصل به OR منطقی نیز دیده شوند. بعلاوه تمام قوانین طبقه‌بندی‌شده برای تمام کلاس‌ها یک‌باره در یک مرحله به دست می‌آیند.

بحث اصلی این پژوهش این است که DERE_x اولین ابزار طبقه‌بندی است که بر اساس DE است و به‌طور خودکار مجموعه‌ای از قوانین IF-THEN را بدون مداخله‌ی مکانیسم دیگری استخراج می‌کند.

برای آزمایش‌ها سه پایگاه‌داده از حوزه‌ی پزشکی انتخاب شده است. ابزار در سرتاسر این سه پایگاه داده با یک دسته قوانین طبقه‌بندی‌کننده استفاده شده است. نتایج، تأثیر روش موردنظر را ثابت کرده و تصدیق می‌کنند که DERE_x در مقایسه با دیگر الگوریتم‌ها بهترین ابزار طبقه‌بندی است.

مزیت اصلی DERE_x تأمین اطلاعات واضح استخراج شده از پایگاه‌داده برای کاربر است؛ زیرا برخلاف Naïve, RBF, Bayes, Bagging و NBtree می‌تواند به‌طور مستقیم قوانین IF-THEN را برای تشخیص ارائه کند. علاوه بر این DERE_x می‌تواند استخراج ویژگی‌ها را نیز انجام دهد. چون قوانین به‌دست‌آمده شامل تعدادی از متغیرها می‌باشند که می‌توانند برای تشخیص درست بیماری بسیار مهم باشند. در این روش با استفاده از اطلاعات مفید به پزشکان کمک شده است. البته نظر آن‌ها درباره‌ی درستی و مفید بودن این قوانین بیشترین اهمیت را در عمل دارد. وقتی که DERE_x با دیگر ابزار طبقه‌بندی مبتنی بر قانون مقایسه می‌شود، DERE_x نیازمند کمترین میانگین از قوانین برای روبرو شدن با مسئله است که این امتیاز قابل توجهی نسبت به دیگر تکنیک‌هاست. علاوه بر این میانگین تعداد شرط‌ها در قوانین زیاد بزرگ نیست. این ویژگی DERE_x را کاربرپسند می‌کند.

لازم به یادآوری است که این نتیجه باوجود استفاده از AND به‌عنوان تنها متصل‌کننده‌ی محلی در شرط‌ها به‌دست‌آمده است. این نقص با این واقعیت که سیستم قادر است مجموعه‌ای از قوانین و درصورت لزوم برای هر کلاس بیش از یک قانون داشته باشد تا حد زیادی برطرف شده است؛ زیرا قوانین می‌توانند با استفاده از OR در نظر گرفته شوند.

یک توضیح ممکن برای کیفیت نتایج به‌دست‌آمده این است که در اصل روش ارائه شده در اینجا تکاملی است. لذا تمام مزایای الگوریتم‌های تکاملی را داراست:

- می‌تواند بدون دانش اولیه‌ی هیچ‌گونه فرضی وارد مسئله شود.
- می‌تواند مقادیر بهینه‌ی محلی را خارج کند.
- نیازی نیست تابع صلاحیت قابل تشخیص باشد

تا اینجا این ویژگی‌ها برای تمام روش‌های مبتنی بر تکامل در روبرو شدن با مسائل متداول ند. مخصوصاً ابزار معرفی شده که بر اساس تکامل تفاضلی است که در نوشته‌ها به‌طور گسترده به‌عنوان روشی سریع‌تر از دیگر روش‌ها در جستجوی محلی شناخته می‌شود.

منابع

1. G. Vinodhini, R.M. Chandrasekaran, A comparative performance evaluation of neural network based approach for sentiment classification of online reviews, *Journal of King Saud University – Computer and Information Sciences* (2016) 28, 2–12
2. Giovanna Sannino, Ivanoe De Falco, Giuseppe De Pietro, Monitoring Obstructive Sleep Apnea by means of a real-time mobile system based on the automatic extraction of sets of rules through Differential Evolution, *Journal of Biomedical Informatics* 49 (2014) 84–100
3. J. Ilonen, J. K. Kamarainen, J. Lampinen, Differential Evolution Training Algorithm for Feed-Forward Neural Networks, *Neural Processing letters* 17 93-105, 2003.
4. J. Ilonen, J. K. Kamarainen, J. Lampinen, Differential Evolution Training Algorithm for Feed-Forward Neural Networks, *Neural Processing letters* 17 93-105, 2003
5. J. Jayanth Shivaprakash Koliwad, Ashok Kumar T. Classification of remote sensed data using Artificial Bee Colony algorithm, *The Egyptian Journal of Remote Sensing and Space Sciences* (2015) 18, 119–126
6. J. Wu, Z. Cai, Attribute weighting via differential evolution algorithm for attribute weighted naive bayes (WNB), *Journal of Computational Information Systems* 7 (5) (2011) 1672–1679.
7. R. Detrano, A. Janosi, W. Steinbrunn, M. Pfisterer, J. Schmid, K. Sandhu, K. Guppy, S. Lee, V. Froelicher, International application of a new probability algorithm for the diagnosis of coronary artery disease, *American Journal of Cardiology* 64 (1989) 304–310.
8. R.C. Holte, Very simple classification rules perform well on most commonly used datasets, *Machine Learning* 11 (1993) 63–91.
9. S. Das, A. Abraham, A. Konar, Automatic clustering using an improved differential evolution algorithm, *IEEE Transactions on Systems, Man and Cybernetics Part A: Systems and Humans* 38 (1) (2008) 218–237.
10. S.J. Haberman, Generalized residuals for log-linear models, in: *Proceedings of the 9th International Biometrics Conference*, 1976, pp. 104–122.
11. Tapas Bhadra, Sanghamitra Bandyopadhyay, Ujjwal Maulik, Differential Evolution Based Optimization of SVM Parameters for Meta Classifier Design, *Procedia Technology* 4 (2012) 50 – 57
12. W.H. Wolberg, O.L. Mangasarian, Multisurface method of pattern separation for medical diagnosis applied to breast cytology, *Proceedings of the National Academy of Sciences of the United States of America* 87 (1990) 9193–9196.
13. Wei-Jie Chena, Yuan-Hai Shaoa, Ya-Fen Ye, Improving Lap-TSVM with successive overrelaxation and differential evolution, *Procedia Computer Science* 17 (2013) 33 – 40
14. Z. Cai, H. Chengyu, Z. Kang, L. Yong, L. Xiaobo, Y. Chao, Differential evolution based band selection in hyperspectral data classification, in: *Proceedings of ISICA 2010 in Lecture Notes in Computer Science*, Springer, 2010, pp. 86–94.

A Combined Differential Evolution and Neural Network Approach to Extract Rules from Medical Databases

Mehdi Jobedari, Hamid Paygozar

Department of Computer Engineering, Islamic Azad University, Khomein Branch, Khomein, Iran

Abstract

Today, clinical decision support systems have been defined in medical science to offer useful suggestions using two or more items of the patient's characteristics. These suggestions help specialists in disease diagnosis or in the treatment process. The classification of data in this area requires transparency and reasoning for the allocation of data. Therefore, in addition to the correct classification of medical data, it is necessary to explain the reason for the data allocation. This research uses a new method based on differential evolution for automatic classification of items in the medical database. This method has introduced a tool called DEREx, which automatically extracts explicit knowledge from the database in the form of an IF-THEN rule, including the conditions connected by AND (and) on the database variables. Each item of differential evolution is considered a rule for classification. A set of differential evolution items is a set of rules, all of which are used to classify data in the database. These rules can also be logically linked by OR (or), so all classification rules are found for all categories at once in one step. The DEREx tool is the first classification tool based on differential evolution which automatically extracts categories of IF-THEN rules without the involvement of another mechanism.

Keywords: Classification, Differential Evolution, Database, Rule Extraction
