

استفاده از الگوریتم‌های ژنتیک برای بهینه‌سازی داده‌کاوی در یک سیستم آموزشی مبتنی بر وب

نویسندگان: بهروز مینایی بیدگلی، ویلیام اف. پانچ III^۱

مترجم: محسن مؤمنی

چکیده

این مقاله راه‌بردی را برای دسته‌بندی^۲ دانشجویان، به منظور پیش‌بینی نمرات نهایی آن‌ها، برپایه‌ی ویژگی‌های استخراج‌شده از داده‌های ثبت‌شده، در یک سیستم مبتنی بر وب آموزشی، ارائه می‌کند. ترکیبی از دسته‌بند^۳‌های چندگانه، راه را برای بهبود کارایی دسته‌بندی، می‌گشاید. به‌وسیله‌ی وزن‌دهی به بردار ویژگی‌ها، با استفاده از یک الگوریتم ژنتیک، توانستیم دقت پیش‌بینی را بهینه‌سازی و بهبود قابل ملاحظه‌ای را نسبت به دسته‌بندی عادی^۴ تنها به دست آوریم. بیش‌ازاین، این تحقیق نشان می‌دهد، هنگامی که تعداد ویژگی‌ها اندک است، وزن‌دهی به ویژگی‌ها، به‌تر از انتخاب ویژگی‌ها به تنهایی نتیجه می‌دهد.

¹ Genetic Algorithms Research and Applications Group (GARAGE)
Department of Computer Science&Engineering
Michigan State University
2340 Engineering Building
East Lansing, MI 48824
{minaeibi, punch}@cse.msu.edu
<http://garage.cse.msu.edu>

² classification

³ classifier

⁴ raw classification

بیان مسئله

بسیاری از مؤسسه^۱های آموزشی در برپایی یک پیشگاه^۲ تعلیم و تعلم اینترنتی^۳، تلاش می کنند. چندین سیستم با قابلیت^۴ و راهبردهای مختلف برای به وجود آوردن امکان تحصیل اینترنتی^۴ در محیطی آکادمیک توسعه یافته اند. دانشگاه ایالت میشیگان^۵ نیز، در ایجاد قسمتی از این سیستمها را برای فراهم آوردن یک فراساختار برای آموزش آنلاین^۶ پیشگام شده است. تحقیقی که پیش روی شماست، در بخشی از آخرین سیستم توسعه داده شده در MSU، با نام LON-CAPA^۷ اجرا گردیده است.

در LON-CAPA^۸، ما با دو مجموعه داده ای بزرگ سروکار داریم:

(۱) منابع آموزشی، از قبیل صفحات وب، نمایش ها^۹، شبیه سازی ها، و مسائل فردی ای که برای استفاده به عنوان تکالیف منزل، کوئیزها و امتحانها طراحی شده اند.

(۲) اطلاعاتی درباره ی کاربران که ایجاد، تغییر، ارزیابی^{۱۰}، این منابع را انجام می دهند و یا از آنها استفاده می کنند. به عبارت دیگر، ما در این جا دو مخزن همواره در حال رشد^{۱۱} از داده داریم.

ما به مطالعه ی روش های داده کاوی ای، که در استخراج دانش مفید، از این مجموعه های داده ای بزرگ (درباره ی دانشجویانی که از منابع آموزشی آنلاین استفاده می کنند و مسیر ثبت شده ی آنها در منابع آموزشی وب) به کار می آمدند، پرداختیم.

در این مطالعه، هدف این بوده است که، به سؤالات تحقیقی زیر پاسخ دهیم:

آیا می توانیم دسته هایی را برای دانشجویان بیابیم؟ به عبارت دیگر، آیا گروه هایی از دانشجویان وجود دارند، که به طرز مشابهی این منابع را مورد استفاده قرار دهند؟ اگر بله، آیا می توانیم کلاسی را برای هر دانشجو تشخیص دهیم؟ با این اطلاعات، آیا می توانیم به دانشجویان در بکارگیری بهتر منابع، برپایه ی طرز استفاده ی آنها به وسیله ی دیگر دانشجویانی که در آن گروه بوده اند، راهنمایی لازم را انجام دهیم؟

امید داریم بتوانیم الگوهای استفاده ی مشابهی، در داده ی جمع شده از طریق lon-capa بیابیم، و نهایتاً قادر باشیم، سودمندترین روش را برای مطالعه ی هر آموزنده برپایه ی استفاده ی کنونی اش، پیش بینی کنیم. از آن پس سیستم می تواند، پیشنهادهایی را برای هر آموزنده، در زمینه ی این که چگونه بیشترین پیشرفت را داشته باشد، به او ارائه دهد.

۲ تصویر کردن مسئله بر الگوریتم های ژنتیک

¹ institution

² presence

³ online

⁴ Online education

⁵ Michigan State university (MSU)

⁶ Online instruction

⁷ Learning Online Network with Computer-Assisted Personalized Approach

⁸ See <http://www.lon-capa.org>

⁹ demonstrations

¹⁰ assess

¹¹ ever-growing pools

الگوریتم‌های ژنتیک نشان‌داده‌اند که ابزارهای مؤثری برای استفاده در داده‌کاوی و تشخیص الگو هستند. [7][10][6][16][5][13][7]. یک جنبه‌ی مهم از الگوریتم‌های ژنتیک در یک بستر آموزشی، کارکرد آن‌ها در

تشخیص الگوست. دو راهبرد متفاوت برای به‌کارگیری الگوریتم ژنتیک در تشخیص الگو وجود دارد:

(۱) اعمال مستقیم الگوریتم ژنتیک به عنوان یک دسته‌بند^۱. بندی‌پدهایای و مورتی در [3] الگوریتم ژنتیک را برای یافتن محدوده‌ی تصمیم در یک فضای ویژگی n -بُعدی به‌کارگرفتند.

(۲) استفاده از یک الگوریتم ژنتیک به عنوان یک ابزار بهینه‌سازی برای بازمقداردهی^۲ پارامترها در یک دسته‌بند دیگر. بیش‌تر کاربردهای الگوریتم‌های ژنتیک در تشخیص الگوها در بهینه‌سازی برخی پارامترها در فرآیند دسته‌بندی است. تاکنون، بسیاری از محققان الگوریتم ژنتیک را در انتخاب ویژگی‌ها به‌کارگرفته‌اند. [18], [12], [9], [2]. الگوریتم‌های ژنتیک در یافتن مجموعه‌ی بهینه از ویژگی‌های وزن‌دار نیز به‌کارگرفته شده‌اند، که دقت دسته‌بندی را فزونی می‌بخشد. در این روش، در ابتدا یک روش استخراج ویژگی سنتی نظیر PCA^3 به‌کارگرفته می‌شود و سپس یک ابزار دسته‌بندی نظیر $k-NN^4$ برای محاسبه‌ی تابع برازندگی برای الگوریتم ژنتیک [19], [17] استفاده می‌شود. ترکیب دسته‌بندها عرصه‌ی دیگری است که از الگوریتم‌های ژنتیک برای بهینه‌سازی آن استفاده می‌گردد. کانچیوا و جین در [11] از الگوریتم ژنتیک برای انتخاب ویژگی، همانند انتخاب نوعی از دسته‌بندها در طراحی خود از یک سیستم همجوشی دسته‌بندها، بهره‌برده‌اند. الگوریتم‌های ژنتیک همچنین در انتخاب پیش‌گونه^۵ در دسته‌بندی مبتنی برمورد^۶ استفاده شده‌اند [20].

در این مقاله ما بر راهبرد دوم متمرکز می‌شویم و از الگوریتم ژنتیک برای بهینه‌سازی ترکیبی از دسته‌بندها بهره‌می‌جویم. مأموریت ما، پیش‌بینی نمرات نهایی دانشجویان، بر اساس ویژگی‌های بهره‌گیری از وب آن‌هاست، که از میان داده‌های تکلیفی استخراج می‌گردد. ما مجموعه‌ای از دسته‌بندهای الگو را با پارامترهای مختلف طراحی، پیاده‌سازی و ارزیابی نمودیم، تا کارایی آن‌ها را در مجموعه‌ی داده‌ای بدست‌آمده از لن-کاپا، مقایسه کنیم. نرخ خطاها برای دسته‌بندهای منفرد، ترکیب‌های آن‌ها و ترکیب بهینه‌شده توسط الگوریتم ژنتیک ارائه شد.

۲-۱ مجموعه‌های داده و برجسب‌های دسته

به عنوان داده‌ی آزمایشی ما داده‌های دانشجویان و دروس، درس فیزیک ۱۸۳ (فیزیک رشته‌های علوم پایه و فنی مهندسی I)، لن-کاپا را انتخاب کردیم، که در ترم بهاره سال ۲۰۰۲ در دانشگاه میشیگان، برقرار بوده است. در این درس ۱۲ مجموعه تکلیف شامل ۱۸۴ مسئله، جمع‌گردید، که همه‌ی آن‌ها آنلاین بودند. حدود ۲۶۱ دانشجو از لن-کاپا برای این درس استفاده کرده بودند. بعضی از دانشجویان پس از انجام دو مجموعه تمرین این درس را حذف کردند، که بنابراین آن‌ها هیچ نمره‌ی پایانی‌ای نداشتند. پس از حذف این دانشجویان، ۲۲۷ نمونه‌ی معتبر باقی ماند. توزیع نمرات دانشجویان در شکل ۱ نشان‌داده شده است.

¹ classifier

² resetting

³ Principal Component Analysis

⁴ K nearest neighbor

⁵ prototype

⁶ case-based classification

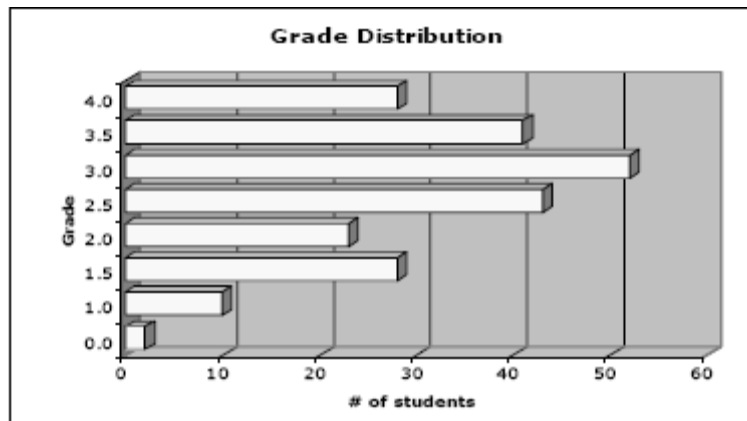


Fig. 1. Graph of distribution of grades in course PHY183 SS02

می‌توانیم دانشجویان را براساس نمرات نهایی آن‌ها به چندین روش گروه‌بندی کنیم، سه تا از این روش‌ها عبارت‌اند از:

۱. اجازه‌دهیم هر ۹ برچسب دسته‌ی ممکن، همان نمرات دانشجویان باشد، آن‌گونه که در جدول ۱ نشان داده شده‌است.
۲. می‌توانیم برچسب دانشجویان را مرتبط با نمره‌ی آن‌ها و گروهی که آن‌ها در سه کلاس، نمرات «بالا»، که نشان‌دهنده‌ی نمره‌ی بین ۳/۵ تا ۴ است، «متوسط» که ارائه‌دهنده‌ی نمره‌ی بین ۲/۵ تا ۳ است و «پایین» که نشان‌گر نمرات کمتر از ۲/۵ است، دارند، تعیین کنیم. جدول ۲ را ببینید.
۳. هم‌چنین می‌توانیم دانشجویان را با یک یا دو برچسب کلاس طبقه‌بندی کنیم. «قبول» برای نمرات بالاتر از ۲ و «رد» برای نمرات کم‌تر یا برابر ۲، آن‌گونه که در جدول ۳ نشان داده شده‌است.

Table 1. Selecting 9 class labels regarding to students' grades in course PHY183 SS02

Class	Grade	Student #	Percentage
1	0.0	2	0.9%
2	0.5	0	0.0%
3	1.0	10	4.4%
4	1.5	28	12.4%
5	2.0	23	10.1%
6	2.5	43	18.9%
7	3.0	52	22.9%
8	3.5	41	18.0%
9	4.0	28	12.4%

Table 2. Selecting 3 class labels regarding to students' grades in course PHY183 SS02

Class	Grade	Student #	Percentage
High	Grade \geq 3.5	69	30.40%
Middle	2.0 < Grade < 3.5	95	41.80%
Low	Grade \leq 2.0	63	27.80%

Table 3. Selecting 2 class labels regarding to students' grades in course PHY183 SS02

Class	Grade	Student #	Percentage
Passed	Grade > 2.0	164	72.2%
Failed	Grade \leq 2.0	63	27.80%

می‌توانیم پیش‌بینی کنیم که نرخ خطا در گروه‌بندی رده‌ی اول^۱، باید بیش از دیگران باشد، زیرا اندازه‌ی نمونه‌ی نمرات بر ۹ کلاس تفاوت نسبتاً زیادی دارد. واضح است که ما داده‌ی کم‌تری برای ۳ کلاس نخست در فاز آموزش خواهیم‌داشت، و بنابراین نرخ خطا در فاز ارزیابی یحتمل بالاتر خواهد بود.

۲-۲ استخراج ویژگی‌ها

یکی از مراحل حیاتی در انجام دسته‌بندی، انتخاب ویژگی‌هایی است که برای دسته‌بندی به‌کار می‌رود. در زیر ما درباره‌ی ویژگی‌های بدست‌آمده از لُن-کاپا که مورد استفاده قرار گرفته‌اند، بحث خواهیم‌نمود، که چگونه می‌توان آن‌ها را متجسم‌نمود (تا به کار انتخاب کمک‌شود). و چرا ما داده‌ها را پیش از دسته‌بندی به‌هنجار^۱ می‌کنیم. ویژگی‌هایی که در پی می‌آید، به‌وسیله‌ی سیستم لُن-کاپا ذخیره‌شده‌اند:

۱. نمره‌ی کلّ پاسخ‌های صحیح (نرخ کامیابی)
۲. فهم درست مسئله در تلاش اولیه، در برابر آن‌ها که با تعدادی تلاش نمره‌ی بالا کسب‌کرده‌اند. (کامیابی در نخستین تلاش)
۳. تعداد کل تلاش‌ها برای انجام تکلیف. (تعداد تلاش‌ها پیش از به‌دست‌آمدن پاسخ صحیح)
۴. زمان صرف‌شده برای مسئله پیش از آن‌که حل‌شود (به‌عبارت دقیق‌تر، تعداد ساعات تا رسیدن به پاسخ صحیح. تفاوت زمان رخداد کامیابی و نخستین زمانی که مسئله آزمایش می‌گردد). هم‌چنین، زمانی که هر دانشجو به پاسخ صحیح مسئله می‌رسد، با توجه به روز محاسبه می‌گردد. معمولاً دانشجویان به‌تر زودتر تکالیف را تکمیل می‌کنند.
۵. کل زمانی که بر مسئله صرف می‌شود، جدای از این‌که دانشجو به پاسخ صحیح برسد یا نه. (تفاوت زمان میان آخرین ارائه‌ی پاسخ و شروع زمان ارائه‌ی مسئله برای آزمون).
۶. شرکت در مکانیسم ارتباط، به جای کار تک نفره. لُن-کاپا امکان تعامل آنلاین را هم با دانشجویان دیگر و هم با مربّی، فراهم می‌سازد. کجا باید استفاده‌گردد؟
۷. مطالعه‌ی مواد پشتیبانی پیش از تلاش برای حل تکلیف، به‌جای تلاش ابتدایی برای حل و سپس بازگشت به خواندن آن.
۸. ارسال تلاش‌های فراوان در زمان کوتاه، بدون توجه به موارد درون آن، به‌جای دادن آن با یک تلاش، بازخوانی، ارسال دوباره و به‌همین‌صورت.

۹. دست‌برداشتن از یک مسئله به جای دانشجویانی که تلاش را بر یک مسئله تا آخرین مهلت صرف می‌کنند.

۱۰. زمانی که نخستین بار کار را شروع می‌کنند. (شروع ثبت نام، میان هفته، یا آخرین دقایق!) بر مبنای تعداد تلاش‌ها یا تعداد مسائل حل‌شده. یک دانشجو که همه‌ی جواب‌های صحیح را می‌دهد، لزوماً در گروه کامیابان دسته‌بندی نمی‌شود، اگر آن‌ها به‌طور متوسط برای هر مسئله ۵ تلاش داشته‌باشند، امّا این مسئله باید در این تحقیق مورد بازخوانی قرار گیرد.

در این مقاله ما بر ۶ ویژگی نخست در مجموعه‌ی داده‌ای درس فیزیک ۱۸۳، که برای آزمون دسته‌بندی برگزیده شده‌است، بهره می‌گیریم.

۲-۳ دسته‌بندها

¹ normalized

تشخیص الگوها کاربردهای متنوعی در عرصه‌های مختلف دارد، بدین دلیل امکان‌ندارد ابزار دسته‌بندی‌ای تهیه‌کنیم، که در همه‌ی موارد نتایج خوبی را عرضه‌کند. دسته‌بند بهینه در هر مورد به‌شدت به محدوده‌ی مسئله بستگی دارد. در عمل، ممکن است به موردی برخوردیم که هیچ ابزار دسته‌بندی‌ای به‌تنهایی نتواند، دسته‌بندی‌ای با دقت قابل قبول ارائه‌دهد. در چنین مواردی شاید بهتر باشد که از ترکیب نتایج دسته‌بندهای متفاوت برای دستیابی به دقت بهینه استفاده‌کنیم. هر دسته‌بند با جنبه‌های متفاوتی از بردار ویژگی در پرورش و آزمون کارکند. نتیجتاً، به‌فرض وجود شرایط مناسب، ترکیب چندین دسته‌بند، شاید به بهبود کارایی دسته‌بندی، در قیاس با یک دسته‌بند به‌تنهایی، بیانجامد. [4]

تمرکز این مرور، به مقایسه‌ی چند دسته‌بند الگوی غیرپارامتریک مشهور و یک دسته‌بند پارامتریک، بر طبق تخمین خطاها، محدود می‌شود. ۶ دسته‌بند مختلف با استفاده از مجموعه‌ی داده‌ای گن-کاپا، در این مطالعه، مورد مقایسه قرار گرفته‌اند. دسته‌بندهایی که در این مطالعه مورد استفاده قرار گرفته‌اند، عبارت‌اند از:

*Quadratic Bayesian classifier, 1-nearest neighbor (1-NN), k-nearest neighbor (k-NN) Parzen-window, multi-layer perceptron (MLP), Decision Tree.*¹

این دسته‌بندها، برخی از دسته‌بندهای معمول هستند، که در بیش‌تر مسائل کاربردی دسته‌بندی مورد استفاده قرار می‌گیرند. پس از آن که قدری اعمال پیش‌پردازشی روی مجموعه‌ی داده‌ای اعمال شد، نرخ خطا در هر دسته‌بند گزارش می‌شود. در نهایت، برای بهبود کارایی، ترکیبی از دسته‌بندها ارائه می‌شود.

۲-۴ به‌هنجارسازی^۲

در دسته‌بندهای بیز و پنجره‌ی پارزن، فرض بر آن است که ویژگی‌ها دارای توزیع نرمال هستند. به همین دلیل لازم است که همه‌ی ویژگی‌ها به‌هنجار شوند. این امر اطمینان می‌دهد که هر ویژگی وزن یکسانی را در فرآیند تصمیم‌گیری خواهد داشت. با فرض این‌که، داده‌ی داده‌شده، توزیع گوسی دارد، این به‌هنجارسازی با استفاده از میانگین و انحراف معیار داده‌ی آموزشی^۳ به‌جای آورده می‌شود. به منظور به‌هنجارسازی داده‌ی آموزشی، نخست لازم است میانگین نمونه محاسبه‌گردد و انحراف معیار هر ویژگی، یا ستون، در مجموعه‌ی داده‌ای، بر مبنای آن تعیین شود و سپس داده با استفاده از معادله‌ی (۱) به‌هنجارگردد:

$$x_i = \frac{x_i - \mu}{\sigma} \quad (1)$$

این امر اطمینان می‌دهد که هر ویژگی داده‌ی آموزشی توزیع نرمال با میانگین صفر و انحراف معیار ۱، داشته‌باشد. به‌علاوه، روش k-نزدیک‌ترین همسایگی نیاز به به‌هنجارسازی همه‌ی ویژگی‌ها در یک دامنه‌ی یکسان دارد. به‌هرحال، باید در استفاده از به‌هنجارسازی، پیش از در نظر داشتن تأثیر آن بر کارایی دسته‌بند، بر حذر باشیم.

۲-۵ ترکیب چندین دسته‌بند^۴ (CMC)

در ترکیب چند دسته‌بند، ما قصد داریم، کارایی دسته‌بند را بهبود بخشیم. راه‌های گوناگونی در ترکیب دسته‌بندها به‌نظر می‌رسد:

¹ The first five classifiers are coded in MATLABM 6.0, and for the decision tree classifiers we have used some available software packages such as C5.0, CART, QUEST, and CRUISE.

² normalization

³ Training Data

⁴ Combination of Multiple Classifiers

❖ آسان‌ترین راه، یافتن نرخ خطای روی هم‌رفته‌ی دسته‌بندها و انتخاب آن ابزار دسته‌بندی‌ای است که حداقل نرخ خطای داده‌شده در مجموعه‌ی داده‌ای را دارد. این راه، یک ترکیب دسته‌بندهای *offline* نامیده می‌شود. این مورد ممکن است فی‌الواقع ترکیبی از دسته‌بندها به نظر نرسد؛ ولی در هر حال، عموماً کارایی به تری نسبت به یک دسته‌بند به تنهایی دارد.

❖ روش دوم، که ترکیب دسته‌بندهای آنلاین نامیده می‌شود، از همه‌ی دسته‌بندها به همراه یک رای‌گیری استفاده می‌کند. دسته‌ای که بیش‌ترین آراء را، از یک دسته‌بند کسب کند، به نمونه‌ی آزمون واگذار می‌شود. از طریق شهودی، به نظر می‌رسد، که این روش به‌تر از روش قبلی باشد. در هر حال، وقتی بر مواردی از مجموعه‌ی داده‌ای خود این تلاش را انجام دادیم، نتایج به‌تر از نتایج روش قبلی نبودند. بنابراین، ما معیار رای اکثریت را از «کسب بیش از ۵۰٪ آراء» به «کسب بیش از ۷۵٪ آراء» تغییر دادیم. این امر منجر به آن شد که بهبود قابل ملاحظه‌ای نسبت به ترکیب دسته‌بندهای افلاین به دست آید.

با استفاده از روش دوم، ما در جدول ۴ نشان داده‌ایم که CMC می‌تواند، به بهبود دقت قابل ملاحظه‌ای، در هر سه مورد دسته‌های ۲، ۳ و ۹، دست یابد. اینک قصد داریم از الگوریتم ژنتیک، برای کشف این‌که، آیا می‌توانیم کارایی CMC را پیشینه کنیم، یا نه؟، استفاده کنیم.

۳ بهینه‌سازی CMC با استفاده از الگوریتم ژنتیک

3 Optimizing the CMC Using a GA

We used GAToolBox3 for MATLAB to implement a GA to optimize classification performance. Our goal is to find a population of best weights for every feature vector, which minimize the classification error rate.

The feature vector for our predictors are the set of six variables for every student: Success rate, Success at the first try, Number of attempts before correct answer is derived, the time at which the student got the problem correct relative to the due date, total time spent on the problem, and the number of online interactions of the student both with other students and with the instructor.

We randomly initialized a population of six dimensional weight vectors with values between 0 and 1, corresponding to the feature vector and experimented with different number of population sizes. We found good results using a population with 200 individuals. The GA Toolbox supports binary, integer, real-valued and floatingpoint chromosome representations. Real-valued populations may be initialized using

3 Downloaded from <http://www.shef.ac.uk/~gaipp/ga-toolbox/>

the Toolbox function *crtrp*. For example, to create a random population of 6 individuals with 200 variables each: we define boundaries on the variables in *FieldD* which is a matrix containing the boundaries of each variable of an individual. `FieldD = [0 0 0 0 0 0; % lower bound
1 1 1 1 1 1]; % upper bound`

We create an initial population with `Chrom = crtrp(200, FieldD)`, So we have for example:

```
Chrom = 0.23 0.17 0.95 0.38 0.06 0.26  
0.35 0.09 0.43 0.64 0.20 0.54  
0.50 0.10 0.09 0.65 0.68 0.46  
0.21 0.29 0.89 0.48 0.63 0.89  
.....
```

We used the simple genetic algorithm (SGA), which is described by Goldberg in [9]. The SGA uses common GA operators to find a population of solutions which optimize the fitness values.

3.1 Recombination

We used “*Stochastic Universal Sampling*” [1] as our selection method. A form of stochastic universal sampling is implemented by obtaining a cumulative sum of the fitness vector, $FitnV$, and generating N equally spaced numbers between 0 and $\text{sum}(FitnV)$. Thus, only one random number is generated, all the others used being equally spaced from that point. The index of the individuals selected is determined by comparing the generated numbers with the cumulative sum vector. The probability of an individual being selected is then given by

$$F(x_i) = \frac{f(x_i)}{\sum_{i=1}^{N_{\text{ind}}} f(x_i)} \quad (2)$$

where $f(x_i)$ is the fitness of individual x_i and $F(x_i)$ is the probability of that individual being selected.

3.2 Crossover

The crossover operation is not necessarily performed on all strings in the population. Instead, it is applied with a probability P_x when the pairs are chosen for breeding. We selected $P_x = 0.7$. There are several functions to make crossover on real-valued matrices.

One of them is *recint*, which performs intermediate recombination between pairs of individuals in the current population, *OldChrom*, and returns a new population after mating, *NewChrom*. Each row of *OldChrom* corresponds to one individual. *recint* is a function only applicable to populations of real-value variables. Intermediate recombination combines parent values using the following formula [14]:

$$\text{Offspring} = \text{parent1} + \text{Alpha} \times (\text{parent2} - \text{parent1}) \quad (3)$$

Alpha is a Scaling factor chosen uniformly in the interval [-0.25, 1.25]

3.3 Mutation

A further genetic operator, mutation is applied to the new chromosomes, with a set probability P_m . Mutation causes the individual genetic representation to be changed according to some probabilistic rule. Mutation is generally considered to be a background operator that ensures that the probability of searching a particular subspace of the problem space is never zero. This has the effect of tending to inhibit the possibility of converging to a local optimum, rather than the global optimum.

There are several functions to make mutation on real-valued population. We used *mutbga*, which takes the real-valued population, *OldChrom*, mutates each variable with given probability and returns the population after mutation, $\text{NewChrom} = \text{mutbga}(\text{OldChrom}, \text{FieldD}, \text{MutOpt})$ takes the current population, stored in the matrix *OldChrom* and mutates each variable with probability by addition of small random values (size of the mutation step). We considered 1/600 as our mutation rate. The mutation of each variable is calculated as follows:

$$\text{Mutated Var} = \text{Var} + \text{MutMx} \times \text{range} \times \text{MutOpt}(2) \times \text{delta} \quad (4)$$

where delta is an internal matrix which specifies the normalized mutation step size; MutMx is an internal mask table; and MutOpt specifies the mutation rate and its shrinkage during the run. The mutation operator *mutbga* is able to generate most points in the hypercube defined by the variables of the individual and the range of the

mutation. However, it tests more often near the variable, that is, the probability of small step sizes is greater than that of larger step sizes.

3.4 Fitness Function

During the reproduction phase, each individual is assigned a fitness value derived from its raw performance measure given by the objective function. This value is used in the selection to bias towards more fit individuals. Highly fit individuals, relative to the whole population, have a high probability of being selected for mating whereas less fit individuals have a correspondingly low probability of being selected. The error rate is measured in each round of cross validation by dividing “the total number of misclassified examples” into “total number of test examples”. Therefore, our *fitness function* measures the error rate achieved by CMC and our objective would be to maximize this performance (minimize the error rate).

۴ نتایج پژوهش

بدون استفاده از الگوریتم ژنتیک، نتایج کارایی روی هم‌رفته‌ی دسته‌بندها روی مجموعه‌ی داده‌ای ما، درخصوص ۴ دسته‌بند درختی، و پنج دسته بند غیردرختی و حالت ترکیبی در جدول ۴ نشان داده شده است. درخصوص دسته بندهای تکی، برای مورد ۲-دسته‌ای، k-NN به‌ترین کارایی را با ۸۲٫۳٪ دقت داشت. در مورد ۳-دسته و ۹-دسته، CART به‌ترین کارایی را، با حدود ۶۰٪ در حالت ۳-دسته‌ای و ۴۳٪ در حالت ۹-دسته‌ای، داشت. در هر حال، با در نظر گرفتن ترکیب دسته بندهای غیرمبتنی بر درخت، CMC به‌ترین کارایی را در هر سه حالت داشت. که، ۸۶٫۸٪ دقت را در حالت ۲-دسته‌ای، ۷۱٪ را در حالت ۳-دسته‌ای و ۵۱٪ را در حالت ۹-دسته‌ای ارائه داد.

در بهینه‌سازی با الگوریتم ژنتیک، ما از ۲۰۰ فرد در جمعیت خود استفاده کردیم، و الگوریتم را در حدوداً ۵۰۰ نسل اجرا کردیم. برنامه را ۱۰ بار اجرا کردیم و میانگین را به دست آوردیم، آن‌گونه که در جدول ۵ نشان داده شده است. در هر اجرا ۲۰۰×۵۰۰ بار تابع برازندگی صدازده شد که از تقاطع ۱۰-قاج ارزیابی برای اندازه‌گیری کارایی میانگین برای CMC استفاده کردیم. بنابراین هر دسته بند ۱۰۶×۳ بار برای موارد ۲-دسته‌ای و ۳-دسته‌ای و ۹-دسته‌ای صدازده شده است. در این صورت، زمان سربار برای ارزیابی برازندگی، زمان بحرانی تلقی می‌شود. با استفاده از MLP این فرآیند، حدود ۲ دقیقه به طول می‌انجامد و همه‌ی ۴ روش غیردرختی دیگر تنها ۳ ثانیه به طول می‌انجامد. ما MLP را از گروه دسته بندهای خود، حذف کردیم و این نتیجه‌ی معقول از زمان‌های اجرا بود.

Table 5. Comparing the CMC Performance on PHY183 dataset Using GA and without GA in the cases of 2-Classes, 3-Classes, and 9-Classes, 95% confidence interval.

Classifier	Performance %		
	2-Classes	3-Classes	9-Classes
CMC of 4 Classifiers without GA	83.87 ± 1.73	61.86 ± 2.16	49.74 ± 1.86
GA Optimized CMC, Mean individual	94.09 ± 2.84	72.13 ± 0.39	62.25 ± 0.63
Improvement	10.22 ± 1.92	10.26 ± 1.84	12.51 ± 1.75

نتایج در جدول ۵ ارائه شده است، میانگین کارایی با آزمون t دوطرفه با ۹۵٪ فاصله اطمینان. برای بهبود الگوریتم ژنتیک، در برابر نتایج بدون استفاده از این الگوریتم، یک مقدار P که مشخص‌کننده‌ی احتمال فرضیه‌های صفر (عدم بهبود کامل) نیز داده شده است. که بهینه‌سازی قابل ملاحظه‌ی الگوریتم ژنتیک را نشان می‌دهد. در همه‌ی حالت $p < 0.000$ مشخص‌کننده‌ی بهبود قابل ملاحظه است. بنابراین، با استفاده از الگوریتم ژنتیک، در همه‌ی موارد، ما

بیش از ۱۰٪ میانگین بهبود کارائی فردی را داریم و حدود ۱۲ تا ۱۵٪ میانگین بهبود کارائی داریم. شکل ۲، به ترین نتیجه را در ۱۰ اجرا روی مجموعه‌ی داده‌ای ما نشان می‌دهد. این نمودارها میانگین جمعیت را ارائه می‌کنند، و به ترین فرد در هر نسل و به ترین مقدار حاصل شده به وسیله‌ی اجرا.

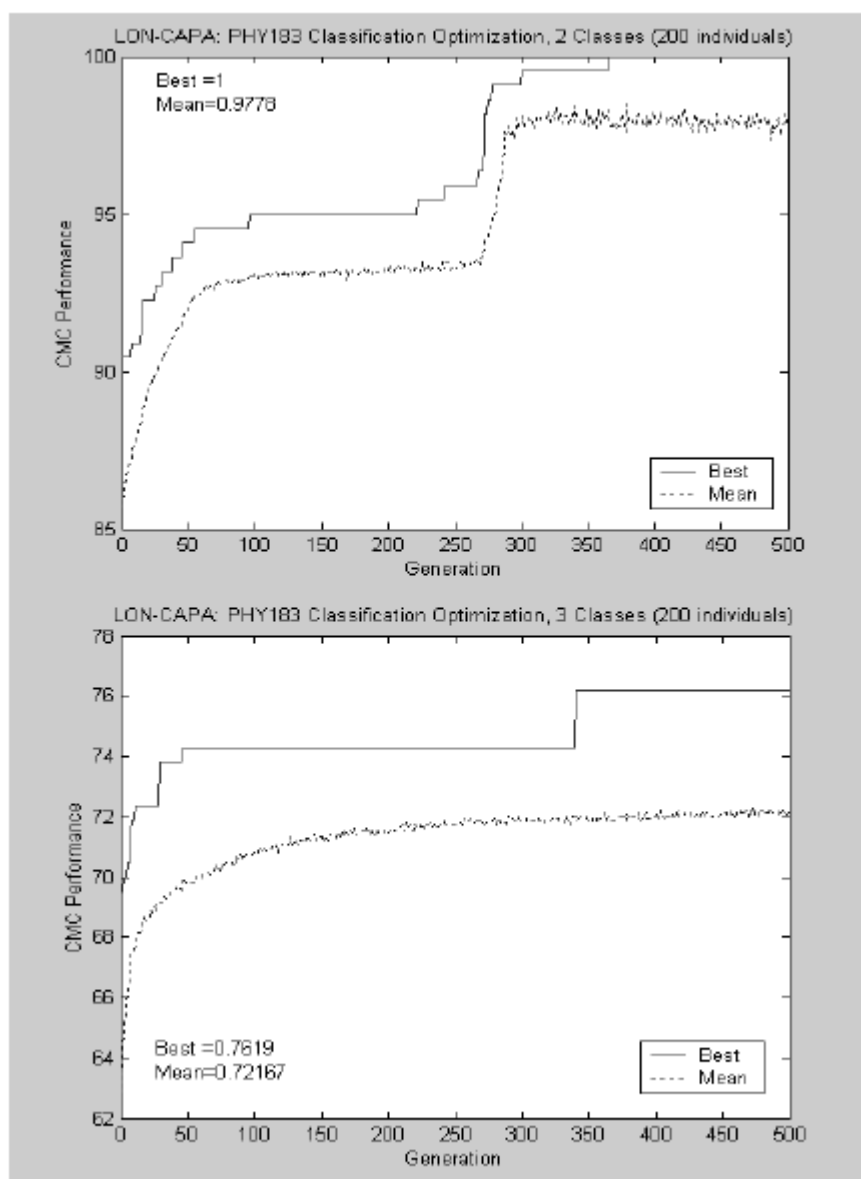


Fig. 2. Graph of GA Optimized CMC performance in the case of 2, and 3-Classes

در نهایت، ما می‌توانیم افراد (اوزان) را برای ویژگی‌ها به وسیله‌ی به دست آوردن نتایج بهبود یافته، بررسی کنیم. این وزندهی به ویژگی‌ها، اهمیت آن‌ها را مشخص می‌سازد، و برای ساختن دسته بندی مورد نیاز به کار می‌آید. در بیش‌تر موارد نتایج مشابه رگرسیون خطی چندگانه یا نرم‌افزار مبتنی بر درخت که در روش‌های آماری، برای اندازه‌گیری اهمیت، مورد استفاده قرار می‌گیرد، می‌باشد. جدول ۶ اهمیت ۶ ویژگی در حالت ۳-دسته‌ای، با استفاده از معیار تفکیک آنتروپی، مشخص شده است. بر مبنای آنتروپی، یک خصوصیت آماری بهره‌ی اطلاعاتی^۱ نامیده می‌شود، که نشان می‌دهد، چگونه ویژگی داده شده، نمونه‌های آموزشی را در رابطه با کلاس‌های مقصدشان، تفکیک می‌کند.

¹ information gain

آنتروپی ناخالصی^۱ یک مجموعه‌ی اختیاری را از نمونه‌ی S در یک گره N خاص، تعیین می‌کند. در [5] ناخالصی گره‌ی N با $i(N)$ مشخص شده‌است، مثلاً به این صورت:

$$Entropy(S) = i(N) = -\sum_j P(\omega_j) \log_2 P(\omega_j) \quad (5)$$

where $P(\omega_j)$ is the fraction of examples at node N that go to category ω_j .

Table 6. Feature Importance in 3-Classes Using Entropy Criterion

Feature	Importance %
Total Correct Answers	100.00
Total Number of Tries	58.61
First Got Correct	27.70
Time Spent to Solve	24.60
Total Time Spent	24.47
Communication	9.21

نتایج الگوریتم ژنتیک نیز نشان می‌دهد که، «تعداد کل پاسخ‌های صحیح» و «تعداد کل تلاش‌ها» مهم‌ترین ویژگی‌ها برای دسته‌بندی بوده‌اند. ستون دوم در جدول ۶ درصد اهمیت ویژگی‌ها را نشان می‌دهد.

۵ نتایج و کوشش‌های پیش رو

چهار دسته بند در تفکیک دانشجویان مورد استفاده قرار گرفت. ترکیبی از دسته‌بندهای چندتایی به بهبود قابل ملاحظه‌ی دقت، در هر سه مورد، منجر گردید. وزندهی ویژگی‌ها و استفاده از الگوریتم‌های ژنتیک برای کمینه‌سازی نرخ خطا، دقت پیش‌بینی را حداقل ۱۰٪، در هر سه حالت ۲، ۳ و ۹ دسته‌ای، بهبود داد. در حالت‌هایی که تعداد ویژگی‌ها کم است، وزندهی ویژگی‌ها بسیار بهتر از انتخاب ویژگی کار می‌کند. بهینه‌سازی موفق دسته‌بندی دانشجویان، در هر سه حالت شایستگی استفاده از داده‌های گن-کاپا برای پیش‌بینی نمرات نهایی دانشجویان را برپایه‌ی ویژگی‌های آنها، که از داده‌های مربوط به تکالیف شان استخراج شده بود، نشان می‌داد. ما درصدد آن هستیم که از برنامه‌سازی ژنتیکی برای تولید ترکیب‌های به‌مراتب متفاوت تری از ویژگی‌ها بهره‌گیریم، تا ویژگی‌های جدیدی استخراج کنیم و پیش‌بینی را بهبود بخشیم. ما قصد داریم با استفاده از الگوریتم‌های تکاملی برای دسته‌بندی دانشجویان و مسائل به صورت مستقیم نیز، استفاده کنیم. همچنین قصد داریم، الگوریتم‌های ژنتیکی را برای یافتن قوانین هم‌باش و وابستگی میان گروه‌های مسائل تکالیف گن-کاپا به کاربریم.^۲

سپاس‌گزاری

این کوشش تاحدودی توسط بنیاد دانش ملی^۳ مورد حمایت قرار گرفته‌است. تحت ITR ۰۰۸۵۹۲۱.

منابع

1. Baker, J. E.: Reducing bias and inefficiency in the selection algorithm, Proceeding ICGA 2, Lawrence Erlbaum Associates, Publishers, (1987) 14-21

¹ impurity

² Mathematical, Optional Response, Numerical, Java Applet, and so forth

³ National Science Foundation

2. Bala J., De Jong K., Huang J., Vafaie H.: and Wechsler H. Using learning to facilitate the evolution of features for recognizing visual concepts. *Evolutionary Computation* 4(3) – Special Issue on Evolution, Learning, and Instinct: 100 years of the Baldwin Effect. (1997)
3. Bandyopadhyay, S., and Muthy, C.A.: *Pattern Classification Using Genetic Algorithms*, *Pattern Recognition Letters*, Vol. 16, (1995) 801-808
4. De Jong K.A., Spears W.M. and Gordon D.F.: Using genetic algorithms for concept learning. *Machine Learning* 13, (1993) 161-188
5. Duda, R.O., Hart, P.E., and Stork, D.G.: *Pattern Classification*. 2nd Edition, John Wiley & Sons, Inc., New York NY. (2001)
6. Falkenauer E.: *Genetic Algorithms and Grouping Problems*. John Wiley & Sons, (1998).
7. Freitas, A.A.: A survey of Evolutionary Algorithms for Data Mining and Knowledge Discovery, See: www.pgia.pucpr.br/~alex/papers. A chapter of: A. Ghosh and S. Tsutsui. (Eds.) *Advances in Evolutionary Computation*. Springer-Verlag,(2002)
8. Goldberg, D.E.: *Genetic Algorithms in Search, Optimization, and Machine Learning*, MA, Addison-Wesley, (1989)
9. Guerra-Salcedo C. and Whitley D.: Feature Selection mechanisms for ensemble creation: a genetic search perspective. In: Freitas AA (Ed.) *Data Mining with Evolutionary Algorithms: Research Directions*, Technical Report WS-99-06. AAAI Press, (1999)
10. Jain, A. K.; Zongker, D.: Feature Selection: Evaluation, Application, and Small Sample Performance, *IEEE Transaction on Pattern Analysis and Machine Intelligence*, Vol. 19, No. 2, February (1997)
11. Kuncheva , L.I., and Jain, L.C.: Designing Classifier Fusion Systems by Genetic Algorithms, *IEEE Transaction on Evolutionary Computation*, Vol. 33 (2000) 351-373
12. Martin-Bautista MJ and Vila MA. A survey of genetic feature selection in mining issues. *Proceeding Congress on Evolutionary Computation (CEC-99)*, Washington D.C., July(1999) 1314-1321
13. Michalewicz Z.: *Genetic Algorithms + Data Structures = Evolution Programs*. 3rd Ed. Springer-Verlag, (1996)
14. Muhlenbein and Schlierkamp-Voosen D.: Predictive Models for the Breeder Genetic Algorithm: I. Continuous Parameter Optimization, *Evolutionary Computation*, Vol. 1, No. 1, (1993) 25-49
15. Park Y and Song M.: A genetic algorithm for clustering problems. *Genetic Programming 1998: Proceeding of 3rd Annual Conference*, Morgan Kaufmann, (1998), 568-575.
16. Pei, M., Goodman, E.D. and Punch, W.F.: Pattern Discovery from Data Using Genetic Algorithms, *Proceeding of 1st Pacific-Asia Conference Knowledge Discovery & Data Mining (PAKDD-97)* Feb. (1997)
17. Pei, M., Punch, W.F., and Goodman, E.D.: Feature Extraction Using Genetic Algorithms, *Proceeding of International Symposium on Intelligent Data Engineering and Learning '98 (IDEAL'98)*, Hong Kong, Oct. (1998)

18. Punch, W.F., Pei, M., Chia-Shun, L., Goodman, E.D.: Hovland, P., and Enbody R. Further research on Feature Selection and Classification Using Genetic Algorithms, In 5th International Conference on Genetic Algorithm, Champaign IL, (1993) 557-564
19. Siedlecki, W., Sklansky J., A note on genetic algorithms for large-scale feature selection, Pattern Recognition Letters, Vol. 10, (1989) 335-347
20. Skalak D. B.: Using a Genetic Algorithm to Learn Prototypes for Case Retrieval and Classification. Proceeding of the AAAI-93 Case-Based Reasoning Workshop, Washington, D.C., American Association for Artificial Intelligence, Menlo Park, CA, (1994) 64-69