

کاربرد الگوریتم‌های تکاملی در داده‌کاوی

تحقیق درس: سیستم‌های خبره و مهندسی دانش

استاد گرامی: دکتر ناصر قاسم‌آقایی

نویسنده: محسن مؤمنی (mohsenmomeni@yahoo.com)

(۱) چکیده

گسترش روز افزون کاربردهای داده‌کاوی، و به‌کارگیری آن در مورد داده‌هایی که از حجم نمونه‌ی بالایی برخوردارند، باعث گردیده، مشکلات کار با ابزارهای متداول داده‌کاوی بیش از پیش نمایان‌گردد. الگوریتم‌های تکاملی راهبرد پر قدرتی است، که توانایی چیرگی بر این مشکلات را در ترکیب با ابزارهای دیگر، در مراحل مختلف داده‌کاوی، به ما می‌دهد. مقاله‌ی حاضر سعی دارد، به مروری کلی بر چگونگی به‌کارگیری الگوریتم‌های تکاملی در مراحل گوناگون داده‌کاوی بپردازد.

(۲) مقدمه

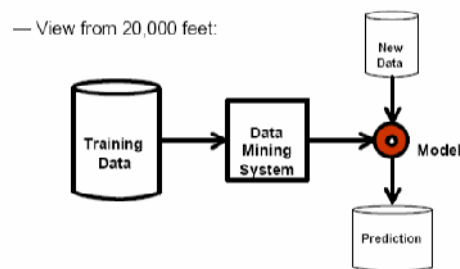
امروزه با وجود حجم بالا و روبه‌رشد پایگاه‌های داده و پدیدارشدن مجموعه‌های داده‌ای ترابایتی و حتا گاهی پتابایتی، یافتن اطلاعات پنهان (اطلاعاتی که بعضاً بسیار مهم و استراتژیک نیز هستند)، در میان داده‌های چنین مجموعه‌هایی، دیگر به آسانی امکان‌پذیر نیست.

داده‌کاوی مجموعه‌ای از ابزارها را در اختیار ما می‌گذارد، تا قادر باشیم، به کشف این اطلاعات پنهان بپردازیم. با گسترش روزافزون به‌کارگیری داده‌کاوی در کاربردهای مختلف، مشکلات آن نیز به مرور بیش‌تر جلوه‌می‌کند. ابزارهای مختلفی را می‌توان برای غلبه بر این مشکلات جدید به‌کارگرفت. الگوریتم‌های تکاملی ابزارهای پر قدرتی هستند که می‌توانند در حل این مشکلات به‌کارآیند و باعث بهبود کارایی ابزارهایی دیگر گردند.

در این مقاله سعی بر آن داریم که ابتدا شناختی کلی از داده کاوی به دست دهیم و بر آن مروری داشته باشیم و سپس در ادامه به بررسی کارکرد الگوریتم‌های تکاملی در مراحل مختلف داده کاوی بپردازیم. اما باید دانست که بررسی داده کاوی، بی نظر کردن به فرآیند بزرگ‌تری که داده کاوی جزئی از آن است، میسر نخواهد شد و ما بدون دانستن کلیات داده کاوی در درک مفهوم آن مشکل خواهیم داشت. لذا این مقاله پس از مرور داده کاوی بررسی گذرایی نیز به کلیات فرآیند استخراج دانش دارد و به بررسی آن می پردازد.

۲-۱) مروری کوتاه بر داده کاوی

داده کاوری فرآیندی است که از ابزارهای مختلف تحلیل داده بهره می گیرد و به طور خودکار، به کشف الگوها و روابط میان داده‌ها می پردازد. سپس با استفاده از این تحلیل‌ها، یک پیش‌بینی معتبر انجام می گیرد. در داده کاوی ابتدا به ساختن یک مدل پیش‌بینی کننده (Predictive Model) مبادرت می کنیم، که برپایه الگوهای مشخص شده از نتایج شناخته شده به دست آمده است. سپس، به وسیله نمونه‌هایی که از مجموعه‌ی منبع داده‌های نمونه، کنار گذاشته شده بودند، مدل خود را آزمایش نموده و در مرحله‌ی نهایی به بازبینی (verify) عملکرد مدل می پردازیم.



شکل ۱- شمای کلی آنچه در داده کاوی صورت می گیرد. بدون در نظر گرفتن جزئیات واقعی، و به عبارتی از ارتفاع ۲۰۰۰۰ پایی

۲-۲) ۴ تعریف شاخص برای داده کاوی

در زیر برخی تعاریف متداول داده کاوی را مرور می کنیم:

- ۱) «کلیفتون»، داده کاوی را، استخراج الگوهای مفید یا دانش (و نه جزئی، ضمنی و ناآشکار، ناشناخته و بالقوه مفید) از مبالغ عظیم داده، می داند.
- ۲) فرآیندی است برای کشف دانش، از پایگاه داده‌ی یکپارچه‌ی بزرگ، در مورد تصمیم‌گیری‌های استراتژیک و تاکتیکی، به وسیله‌ی استخراج اطلاعاتی که تاکنون ناشناخته بوده و قابل تعقیب هستند.
- ۳) داده کاوی یک استخراج داده‌ی خودکار اطلاعات پیش‌بینی کننده‌ی پنهان از پایگاه داده است. «تیرلینگ»
- ۴) داده کاوی از الگوریتم‌هایی با قابلیت مطالعه‌ی کمیت‌های عظیم داده‌ی پیچیده و رسیدن به تفاسیر و مقایسه‌های شرطی است.

در رابطه با تعاریف بالا روی دو چیز باید تأکید گردد: نخست قدرت پیش‌بینی کنندگی که در داده کاوری در پی دستیابی به آن هستیم و دوم، انجام این عمل به طور خودکار.

اما در این میان دو نکته حائز اهمیت است: نخست آن که داده‌کاوی شما را از متخصصان فن و مدیران مجرب بی‌نیاز نمی‌کند، تنها ابزاری برای سهولت و بهبود کار آن‌ها فراهم می‌کند. و نکته‌ی دوم این است که برای قوانین و روابطی که میان داده‌ها در داده‌کاوی کشف می‌شود، لزوماً دلیل طبیعی و عقلانی وجود ندارد.

۲-۳) چه چیزهایی داده‌کاوی نیست؟

داده‌کاوی را نباید با بعضی فرآیندها و روش‌های دیگر اشتباه گرفت. اغلب مفهوم داده‌کاوی به علت نزدیکی با بعضی مفاهیم دیگر، همانند انبار داده، OLAP، بصری‌سازی داده و عامل‌های نرم‌افزاری، با آن‌ها یکی گرفته می‌شود. حال آن‌که ابداً این‌گونه نیست. در این‌جا به بیان تفاوت‌های داده‌کاوی با دو مفهوم دیگر می‌پردازیم: انبار داده و OLAP.

۲-۳-۱) تفاوت‌های داده‌کاوی و انبار داده

غالباً آنچه می‌خواهد مورد داده‌کاوی قرار گیرد، از انبار داده استخراج شده و به یک مرکز داده‌ی داده‌کاوی منتقل می‌گردد. من انبار داده را قبرستان داده‌های عملیاتی می‌نامم. در انبار داده معمولاً داده‌های عملیاتی مرتباً ذخیره می‌گردد و از این داده‌های مجتمع در پیش‌بینی‌ها و تخمین‌ها استفاده می‌گردد. ویژگی خوب انبار داده آن است که داده‌های آن سازگار شده‌اند و جامعیت دارد. فلذا یکی از شرایط لازم برای اطمینان از نتایج داده‌کاوی یعنی جامعیت داده در داده‌های استخراج شده از انبار داده، فراهم است. حُسن دیگر انبار داده آن است که داده‌ها در آن تمیز شده ذخیره می‌شوند. یکی از مراحل پیش‌پردازشی در داده‌کاوی تمیز نمودن داده‌ها است. این عمل برای داده‌هایی که از انبار داده استخراج می‌گردند، مورد نیاز نیست، زیرا قبلاً داده‌ها در ورود به انبار داده تمیز شده‌اند. تفاوت دیگر این است که لزومی به وجود پایگاه‌داده‌ی فیزیکی جداگانه برای داده‌کاوی نیست. پایگاه‌داده‌ی داده‌کاوی می‌تواند به صورت منطقی روی یک پایگاه‌داده‌ی دیگر وجود داشته باشد. این درحالی است که انبار داده لزوماً به انبارفیزیکی داده‌ها نیاز دارد. البته چنانچه امکان حضور پایگاه‌داده‌ی داده‌کاوی بر روی پایگاه‌داده‌ی دیگر نباشد، باید یک پایگاه‌داده‌ی جداگانه داشته باشیم. ضمناً باید دانست که برای انجام داده‌کاوی لزوماً به وجود انبار داده نیاز نیست. اما وجود آن به انجام به‌تر داده‌کاوی کمک می‌کند.

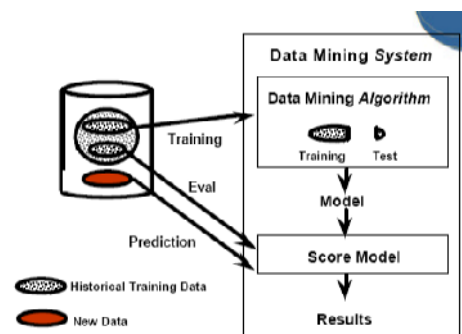
۲-۳-۲) تفاوت‌های داده‌کاوی و OLAP

Olap در واقع به طیفی گسترده از ابزارهای پشتیبانی داده گفته می‌شود. در olap یک سری فرضیه مطرح می‌گردد، و سپس با استفاده از یک سری پرس‌وجو روی داده‌ها به آن فرضیه‌ها پاسخ داده می‌شود. یک تفاوت مهم میان olap و داده‌کاوی آن است که olap یک فرآیند استنتاجی است و داده‌کاوی یک عمل استقرائی. اما باید دانست که olap و داده‌کاوی یک وجه افتراق اساسی با یکدیگر دارند؛ و آن این است که در olap الگویی فرضی و با توجه به اطلاعات پیشین خبره‌ای که با داده‌ها کار می‌کند، مورد ارزیابی قرار می‌گیرد و این در حالی است که در داده‌کاوی الگوها در خود داده‌ها جست‌وجو می‌شوند.

اغلب استفاده از **olap** در مراحل اولیه‌ی استخراج دانش صورت می‌گیرد. انجام عمل **olap** گرچه برای استخراج دانش الزامی نیست، ولی در مواردی که فهم داده‌ها مشکل است، ضرورت می‌یابد. در کل **olap** در فهم داده که نقش به‌سزایی در فرآیند کشف دانش دارد، بسیار اهمیت دارد.

۲-۴) داده‌کاوی چگونه انجام می‌شود؟

در شکل ۲ می‌توان شمایی کلی از آنچه در داده‌کاوی انجام می‌شود، را دید. ان‌گونه که می‌بینید، در داده‌کاوی ما داده را به دو قسمت **training data** و **test data** تقسیم می‌کنیم. میزان داده‌ی اختصاص یافته برای قسمت آزمون بسته به مسائل گوناگون متفاوت است، اما معمولاً $\frac{1}{3}$ داده‌ها برای این کار اختصاص می‌یابد. مدل طراحی شده به‌وسیله‌ی داده‌های آموزشی با داده‌های قسمت آزمون مورد ارزیابی قرار می‌گیرد. در صورت کامیابی مدل ساخته شده برای استفاده و پیش‌بینی در مورد داده‌های جدید به‌کار گرفته می‌شود و در صورت عدم موفقیت، برای آموزش دوباره و ایجاد تغییرات به قسمت یادگیری بازمی‌گردد.



شکل ۲ - تقسیم‌بندی داده‌های آموزشی و چگونگی به‌کارگیری هر یک در فرآیند داده‌کاوی

۲-۶) داده‌های آموزشی

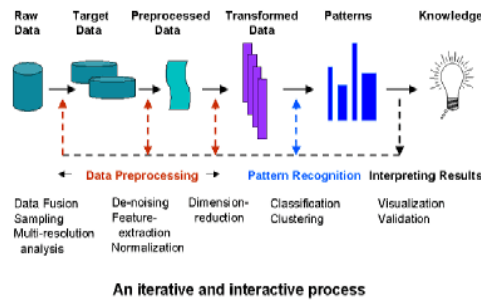
آن‌چنان که پیش از این توضیح دادیم، داده‌کاوی از دسته‌ای داده به عنوان داده‌ی آموزشی استفاده می‌کند. این داده‌ها ویژگی خاصی ندارند و در واقع همان داده‌های خام پایگاه داده هستند، که برای ارزیابی به داده‌کاوی محول شده‌اند. اما نکته‌ی لازم به ذکر این است که پیش از انجام داده‌کاوی باید عملیاتی روی این داده‌ها صورت گیرد. یک‌سری صفات به‌عنوان صفات ورودی داده‌کاوی مشخص‌گردند و در صورتی که یادگیری بانظارت است، یعنی از ابتدا ویژگی داده‌های هر دسته را می‌دانیم، باید صفاتی را به‌عنوان برچسب کلاس و صفات خروجی مشخص‌سازیم. داده‌ها برحسب این صفات خروجی در دسته‌های مختلف قرار می‌گیرند. در یادگیری بدون نظارت هیچ صفتی به‌عنوان خروجی نداریم و لذا تعداد دسته‌ها از قبل نامشخص است.

به طور کلی داده‌ها را می‌توان به دو نوع categorical (real-value) و numerical تقسیم کرد. باید دانست که صفات خروجی در یادگیری نظارت‌شده حتماً بایستی از میان داده‌های categorical انتخاب شوند.

Row-id	A	B	C	D	Class label
1	a_1	b_1	c_1	d_1	A
2	a_1	b_2	c_1	d_2	B
3	a_2	b_3	c_2	d_3	A
4	a_1	b_2	c_3	d_3	C
5	a_1	b_2	c_1	d_3	C

Table 1. A training data set.

شکل ۳ - یک نمونه از داده‌های آموزشی



شکل ۴ - مراحل مختلف داده‌کاوی از منظر کاربرد الگوریتم‌های تکاملی

۲-۷) مروری بر چند روش شاخص داده‌کاوی

پیش از آن که به مرور چند روش داده‌کاوی بپردازیم، لازم است توجه کنیم، که داده‌کاوی شامل چندین مرحله است. در بین مراحل مختلف آن، استخراج ویژگی‌ها از میان داده‌ها، انتخاب چند ویژگی برای داده‌کاوی و بالأخره ایجاد مدل پیش‌بینی‌کننده، با استفاده از الگوریتم‌های دسته‌بندی یا خوشه‌بندی، شاخص‌ترین قسمت‌ها هستند. در این مقاله، پس از بررسی کوتاه الگوریتم‌های داده‌کاوی و مروری بر مراحل استخراج دانش، به بررسی عمیق‌تر کارکرد الگوریتم‌های تکاملی در این مراحل شاخص (که در شکل ۴ دیده می‌شوند) خواهیم پرداخت.

۲-۷-۱) دسته‌بندی با درخت‌های تصمیم‌گیری

در میان ابزارهای مختلف داده‌کاوی نخست نگاهی گذرا به درخت‌های تصمیم‌گیری خواهیم داشت. درخت‌های تصمیم‌گیری شماتیکی ساده از چگونگی ساخت یک مدل شرطی برای طبقه‌بندی داده‌ها هستند. هر گره در درخت تصمیم‌گیری بیانگر شرطی است، که روی یک صفت داده اعمال شده و شاخه‌ی مناسب را برای هر نمونه تعیین می‌کند. در هر سطح یکی از صفات که مناسب‌ترین تقسیم‌کننده‌ی ممکن است (یعنی به بهترین وجه بتواند نمونه‌ها را از یک‌دیگر جدا کند و شاخه‌های کم‌تری نیز داشته‌باشد). انتخاب‌شده و بر اساس تقسیم‌بندی آن شاخه‌های زیرین‌اش مشخص می‌گردند. چنانچه تمام نمونه‌های آموزشی یک شاخه، در یک دسته (از میان دسته‌هایی که از پیش مشخص گردیده) قرار گیرند، کار در این شاخه به پایان می‌رسد. در همه‌ی شاخه‌ها کار تا این سطح به پیش می‌رود و در نهایت مدلی خواهیم داشت که داده‌های آموزشی را تماماً به درستی دسته‌بندی می‌کند. حال این درخت به دست آمده را با داده‌های آزمون، می‌آزماییم و در صورت کسب درجه‌ی تعیین شده (که معمولاً با ماتریس

هم‌جوشی تعیین می‌شود). برای پیش‌بینی به‌کار می‌رود. البته در میان همه‌ی درختان حائز شرایط، درختی که کم‌ترین سطوح و در هر سطح کم‌ترین شاخه‌ها را داشته‌باشد برگزیده می‌شود.

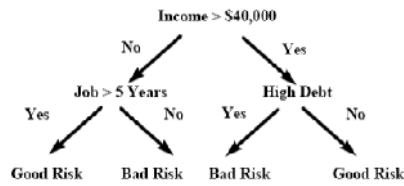


Figure 7. A simple classification tree.

شکل ۵- شمای یک درخت تصمیم‌گیری ساده

۲-۷-۲) قوانین هم‌باش

قوانین هم‌باش یکی از زیرگروه‌های، نوعی داده‌کاوی به نام تحلیل لینک هستند. تحلیل لینک در کل به معنی تحلیل روابط موجود بین داده‌هایی است که به‌همراه یک‌دیگر می‌آیند. در کل دو نوع تحلیل لینک داریم:

۱) کشف هم‌باشی (تطابق)

که به کشف قانون در مورد چیزهایی که با هم به‌نظر می‌رسند، می‌پردازد.

۲) کشف روند

که همان کشف تطابق در طول زمان و با توجه به عنصر زمان است.

یک تطابق (هم‌باشی) به صورت کلی $A \rightarrow B$ نشان داده می‌شود. بدین مفهوم که وجود A ، لزوم وجود B را تعیین می‌کند. به بیان دیگر، چنین قانونی نشان می‌دهد که هرگاه A در تاپلی از پایگاه داده دیده‌شد، احتمال B نیز در آن تاپل وجود خواهد داشت. این چنین قانونی قدرت پیش‌بینی را در مورد نمونه‌هایی که بعضی ویژگی‌هاشان برای ما مبهم است، فراهم می‌کند.

در رابطه با یک هم‌باشی دو چیز مطرح می‌شود:

۱) پشتیبانی

که به معنای میزان اتکایی است که در کل پایگاه داده بر این قانون وجود دارد و از فرمولی ساده و مشخص به دست می‌آید:

تعداد تاپل‌هایی که در طرف چپ قانون آمده / تعداد کل تاپل‌های پایگاه داده

Rule-id	Rule	Support	Confidence
1	$abc \rightarrow A$	80	80%
2	$abcd \rightarrow A$	63	90%
3	$abe \rightarrow B$	36	60%
4	$bcd \rightarrow C$	210	70%

Table 2. Rules found in a training data set.

شکل ۶- نمونه‌ای از قوانین هم‌باش

۲) قابلیت اعتماد

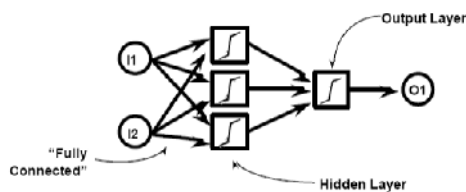
میزان قابلیت اعتمادی که به این قانون می‌توان داشت را نشان می‌دهد. به عبارت دیگر نشان می‌دهد که این قانون در مورد چند درصد از داده‌ها درست کار می‌کند. این ویژگی نیز فرمولی ساده دارد:

تعداد تکرار طرف چپ و راست قانون باهم دیگر | تعداد تکرار طرف چپ قانون در کل (به همراه و بدون طرف راست)

۲-۷-۳) شبکه‌های عصبی

شبکه‌های عصبی مصنوعی با بهره‌گیری از طبیعت کار شبکه‌ی عصبی حیوانی، ایجاد گردیده‌اند. آن‌گونه که می‌دانید، کل ساختار عصبی، از نورون‌هایی تشکیل شده که سه قسمت مهم دارند، شاخه‌های ورودی، هسته و شاخه‌ی خروجی. با به هم پیوستن این عصب‌های ساده، شبکه‌ی پیچیده‌ای به وجود می‌آید که به وسیله‌ی آموزش، قادر به یادگیری و فهم می‌گردد. شبکه‌های عصبی مصنوعی با الهام از این ایده، سعی بر آن دارند، تا ساختار اعصاب را روی ماشین پیاده‌سازی کرده و آن را قادر به فهم سازند.

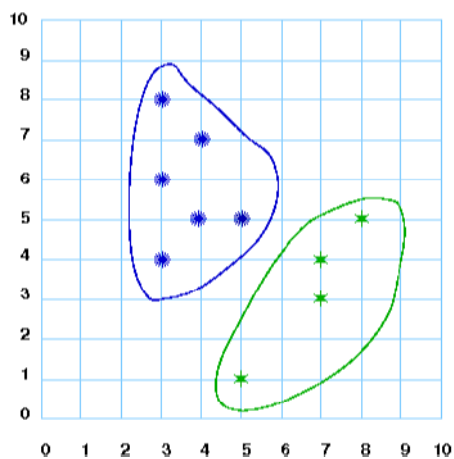
هر شبکه‌ی عصبی یک لایه‌ی ورودی و یک لایه‌ی خروجی و چندین لایه‌ی پنهان میانی دارد. این لایه‌ها که در هر یک تعداد مشخصی گره‌ی عصبی وجود دارد، به وسیله‌ی اتصالاتی به یکدیگر متصل‌اند. اتصالات، خروجی یک گره را دریافت کرده، آن را به شکل تقویت شده و یا تضعیف شده (بر مبنای وزنی که هر اتصال دارد). به ورودی گره‌ی بعدی انتقال می‌دهند. هر گره بر اساس ورودی‌های دریافت شده و تابعی که دارد، خروجی مشخصی خواهد داشت. بنابراین می‌توان صفات ورودی یک دسته را به عنوان ورودی شبکه و دسته‌ی آن را به عنوان خروجی شبکه محسوب نمود و شبکه‌ی عصبی را نیز همانند دیگر مدل‌های پیش‌بینی از بیرون دید. شبکه‌ی عصبی با تغییر وزن اتصالات خود، کاری می‌کند، که میزان اشتباه در تعیین دسته‌ی نمونه‌ها کمینه‌گردد.



شکل ۷- شمای کلی یک شبکه‌ی عصبی

۲-۷-۴) خوشه‌بندی

در خوشه‌بندی آیت‌ها (تاپل‌های) پایگاه داده را به گروه‌های متفاوت تقسیم می‌کنیم. هدف خوشه‌بندی یافتن گروه‌هایی است که اعضای آن‌ها بسیار مشابه یکدیگر و ضمناً این گروه با گروه دیگر بسیار متفاوت باشد. برخلاف دسته‌بندی در آغاز عملیات خوشه‌بندی از ویژگی‌های هر خوشه آگاه نیستیم. پس از خوشه‌بندی یک فرد آگاه در مورد آن داده‌ها و عرصه‌ی مورد داده‌کاوی، باید به تفسیر نتایج پردازش. پس از یافتن خوشه‌ها، آن‌ها می‌توانند برای دسته‌بندی داده‌های جدید مورد استفاده قرار گیرند.



شکل ۸- خوشه بندی

در این جا لازم است، به نکته‌ی مهمی اشاره کنیم و آن تفاوت خوشه بندی با بخش بندی است. بخش بندی یک مسئله‌ی کلی است. مسئله‌ی اصلی در آن تعیین گروه‌هایی است که خصوصیات مشترکی دارند. خواه این خصوصیت مشترک از قبل تعیین شده باشد و خواه نشده باشد.

بخش بندی درحالتی که خصوصیات گروه‌ها از قبل تعیین شده است، دسته بندی و درحالتی که این خصوصیات از قبل تعیین نشده باشند، خوشه بندی نامیده می شود.

در پایان به نکات مهمی که در مورد داده کاوی بایستی همیشه آن‌ها را به خاطر سپرد اشاره می کنیم:

نخست آن که داده کاوی تنها قسمتی کوچک از یک فرایند بسیار بزرگ به نام کشف دانش است. دوم آن که در بسیاری از موارد ضریب قابلیت اعتماد مهم ترین معیار درستی و کارایی داده کاوی نیست. و سوم آن که در داده کاوی این خود داده است که اهمیت دارد. بسیاری گمان می کنند که در داده کاوی الگوریتم‌ها مهم هستند، درحالی که این چنین نیست. در داده کاوی الگوریتم مناسب با توجه به نوع داده و تشخیص خبرگان فن در مورد خصوصیات آن تعیین می شود و در مورد داده های مختلف، الگوریتم‌های متفاوتی به ترین انتخاب هستند. به ترین الگوریتم با توجه به داده‌ای که مورد داده کاوی قرار می گیرد، تعیین می گردد.

در ادامه ابتدا به بررسی فرایند استخراج دانش می پردازیم و فرآیندی کلی که داده کاوی تنها مرحله‌ای کوچک از آن را در بر می گیرد را بر می رسم. سپس به بیان بخش اصلی مقاله، یعنی کارکردهای الگوریتم‌های تکاملی در مراحل مختلف داده کاوی می پردازیم.

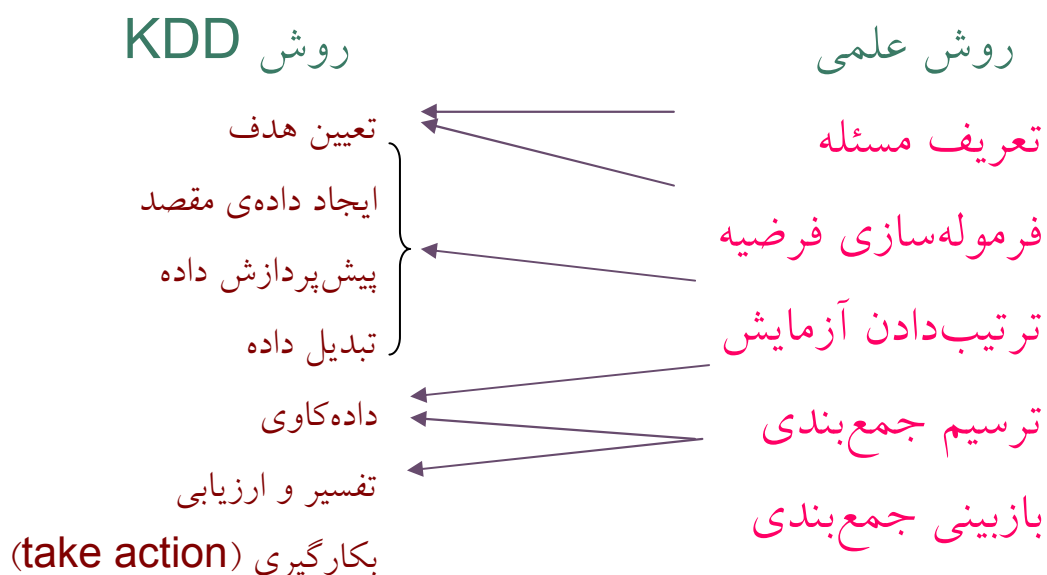
۳) KDD و فرآیند استخراج دانش

۳-۱) پیشینه‌ی KDD

در سال ۱۹۳۰ سر فرانسیس بیکن، برای نخستین بار روش علم تجربی را در کتابش «ارغنون جدید» توضیح داد. او آن را یک فرآیند چهار مرحله‌ای توصیف نمود، مراحلی که آن‌ها را می توان به صورت زیر خلاصه کرد:

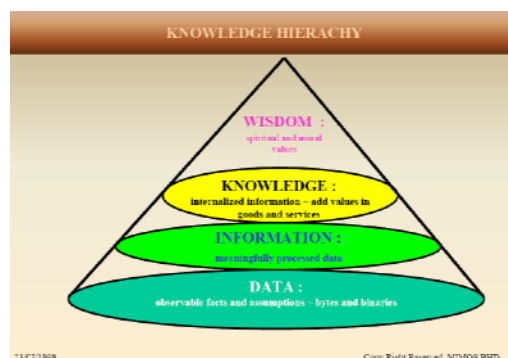
- ۱) تعریف مسئله‌ای که باید حل شود.
- ۲) فرموله سازی یک فرضیه در مورد آن مسئله
- ۳) ترتیب دادن یک یا چند آزمایش برای تأیید یا تکذیب فرضیه
- ۴) ترسیم نتایج و بازبینی

ویژگی خاص این روش آن است که برخلاف روش‌های علمی قدیمی که همه روش‌هایی استنتاجی بودند، یک روش استقرائی است و اجازه می‌دهد که دانشی را که از طریق ارتباط با طبیعت و دانسته‌های محدود از یک فضای کلی به دست می‌آوریم را در مقوله‌ی علم بگنجانیم و قوانین به دست آمده از این روش را معتبر بدانیم. در مدل استخراج دانش کامپیوتری، از این مدل معرفی شده در حدود ۴۰۰ سال پیش اقتباس شده و مدلی جدید برای استخراج چنین دانشی این بار از داده‌هایی که در پایگاه‌های داده هستند، پیشنهاد شده است. در روش جدید ۴ مرحله‌ی روش علم تجربی به ۷ مرحله‌ی استخراج دانش تصویر شده است.



شکل ۹- چگونگی تصویر شدن علم تجربی بر فرآیند استخراج دانش

این مدل با گردشی **bottom-up** در سلسله مراتب دانش (که در شکل ۱۰ نشان داده شده است). به پی‌جویی دانش و به دست آوردن دانش جدید می‌پردازد.



شکل ۱۰- سلسله مراتب دانش

۲-۳) مراحل استخراج دانش

مراحل مختلف KDD به این شرح است:

۱) تعیین هدف

که خود شامل مراحل زیر می‌باشد:

- نخست باید یک بیان واضح از مسئله فراهم آوریم.
- در مرحله‌ی بعد باید تخمین هزینه‌ی پروژه انجام شود.
- سپس نوبت تخمین زمان پروژه می‌رسد.
- در نهایت ارائه‌ی طرحی برای نگهداری سیستم را در این مرحله داریم.

۲) ایجاد یک مجموعه‌ی داده‌ای نمونه

از میان داده‌های انبار داده، داده‌های عملیاتی و گاهی فایل‌های ساده دسته‌ای از داده‌ها به عنوان داده‌ی نمونه برای آموزش و تست مدل برگزیده می‌شوند.

۳) پیش‌پردازش داده

۱) متناسب‌ساختن داده‌های مغشوش

این عمل در چند مورد باید انجام گیرد:

- ۱) جایی که داده‌ی تکرار شده داریم.
- ۲) جایی که صفات غیر صحیح داریم.
- ۳) در مورد داده‌های ناهم‌وار
- ۴) داده‌های جدا افتاده

۲) تصحیح داده‌های گم‌شده

چند راهکار در قبال داده‌های گم‌شده وجود دارد:

- ۱) از بین بردن رکوردی که داده‌ی گم‌شده در آن وجود دارد.
- ۲) جایگزینی داده‌ی گم‌شده با میانگین آن صفت در تمام نمونه‌ها
- ۳) استفاده از مقدار آن صفت در نزد شبیه‌ترین نمونه به نمونه‌ای که دارای داده‌ی گم‌شده است.

۴) تبدیل داده

- ۱) نرمال‌سازی
- ۲) تبدیل نوع
- ۳) انتخاب صفات و نمونه‌ها

۵) داده‌کاوی

- ۱) انتخاب مجموعه‌ی آموزشی و داده‌ی تست از استخر نمونه‌های در دسترس
- ۲) طراحی یک مجموعه از صفات ورودی
- ۳) اگر آموزش نظارت‌شده است، انتخاب یک یا چند صفت خروجی
- ۴) انتخاب مقادیری برای پارامترهای یادگیری
- ۵) استفاده از ابزار داده‌کاوی برای ساخت یک مدل عمومی در مورد آن داده

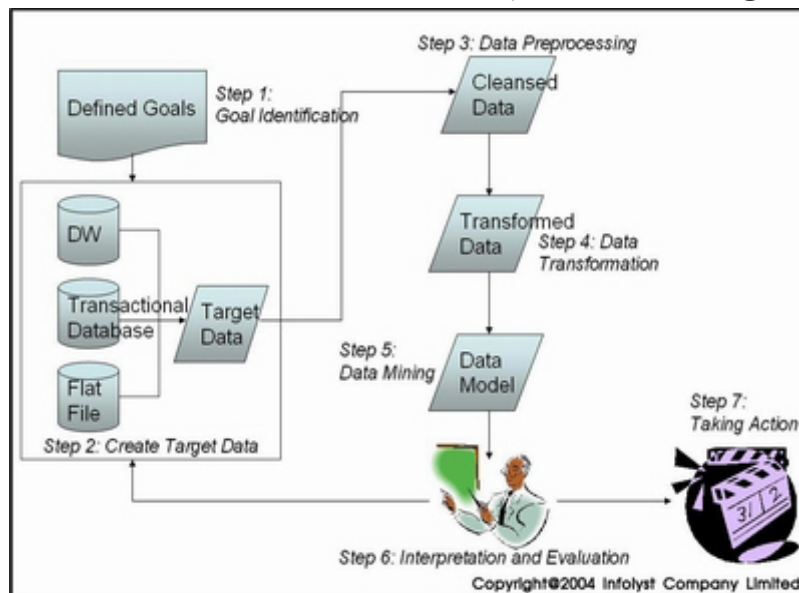
۶) تفسیر و ارزیابی

در این مرحله چهار نوع تحلیل روی نتایج به دست آمده صورت می‌گیرد: (۱) تحلیل آماری، (۲) تحلیل اکتشافی، (۳) تحلیل تجربی و (۴) تحلیل انسانی.

۷) به کارگیری

در مرحله‌ی به کارگیری فعالیت‌هایی برای ایجاد امکان استفاده‌ی خوب از مدل به دست آمده صورت می‌گیرد و مثلاً مستنداتی برای استفاده تهیه می‌شوند.

آن گونه که در شکل ۱۱ می‌بینید، چنانچه در مرحله‌ی ششم، یعنی ارزیابی، مدل موفق نشود کارایی لازم را تدمین کند، مراحل داده‌کاوی از مرحله‌ی دوم تا ششم دوباره باید تکرار شوند. برای پرداختن به جزئیات هر یک از مراحل استخراج دانش مجالی دیگر لازم است. در اینجا تنها کلیات مسئله را مورد بررسی قرار دادیم.



شکل ۱۱ - شمای کلی مراحل KDD و ترتیب کلی به کارگیری مراحل مختلف آن

۳-۳) مدل فرآیند CRISP-DM

کنسرسیومی از چند شرکت تجارتي یک فرآیند را برای ساخت مدل داده‌کاوی استاندارد نمودند. این فرآیند با نام Cross Industry Standard Process for Data Mining خوانده می‌شود، که معمولاً به شکل مخفف از آن به عنوان CRISP یاد می‌شود. در زیر تنها فازهای کلی آن را نام می‌بریم.

شامل ۶ فاز است:

- ۱) فهم تجارتي
- ۲) فهم داده‌ای
- ۳) آماده‌سازی داده
- ۴) مدل‌سازی
- ۵) ارزیابی
- ۶) Deployment و به کارگیری

طی مراحل **crisp** معرفی می‌کند و با رعایت جزئیات آن‌ها، می‌توانیم داده‌کاوی را در یک فرآیند از پیش مشخص و مهندسی شده به‌انجام‌رسانیم و از نتایج آن اطمینان حاصل کنیم. پس از مرور کوتاهی که در این بخش بر کلیات استخراج دانش داشتیم، نوبت آن رسیده که به بحث اصلی این مقاله وارد شویم. اما پیش از آن لازم به نظر می‌رسد که مروری بسیار کوتاه بر چگونگی کار الگوریتم‌های تکاملی و انواع آن‌ها داشته باشیم.

۴) مروری بر الگوریتم‌های تکاملی

الگوریتم‌های تکاملی رویه‌های جست‌وجوی تصادفی‌ای هستند، که به‌وسیله‌ی مکانیسم ژنتیک و انتخاب طبیعی کار بهینه‌یابی را به‌انجام می‌رسانند. الگوریتم‌های تکاملی اغلب در بهینه‌سازی مورد استفاده قرار می‌گیرند. در کاربردهای داده‌کاوی نیز غالباً کاربرد آن‌ها در همین مورد است.

EAها روی یک جمعیت از افراد کار می‌کنند. هر فرد حل ممکن برای یک مسئله است و با یک کروموزوم ارائه می‌شود. جمعیت ابتدایی افراد، می‌تواند تصادفی تولید گردد و یا آنکه با استفاده از دانش به‌دست‌آمده در حل‌ها یا شناخت‌های قبلی به‌وجود آید. الگوریتم به ارزیابی کروموزوم‌ها می‌پردازد تا مشخص کند که هر یک برای حل مسئله چه قدر مناسب هستند. این کار با تابع برازندگی انجام می‌گیرد. با استفاده از نسل فعلی، افرادی از این نسل برگزیده شده و نسل جدید با استفاده از اعمال عملگرهای ژنتیکی روی این افراد، ایجاد می‌شود. این عمل تا یافتن حل قانع‌کننده دنبال می‌شود.

به طور کلی الگوریتم‌های تکاملی را می‌توان در چهار دسته ی کلی قرارداد:

- ۱) الگوریتم‌های ژنتیکی
- ۲) برنامه‌سازی ژنتیکی
- ۳) استراتژی تکاملی
- ۴) برنامه‌سازی تکاملی

در همه‌ی این‌ها مفاهیم پایه‌ای ثابت است و تنها تفاوت در شکل تصویرکردن مسئله و عملگرهایی است که برای ایجاد نسل بعدی به‌کار می‌رود.

Layout of Genetic algorithm

Input: A problem instance

Output: A (sub-optimal) solution

- 1) $t = 0$, initialize $P(t)$, and evaluate the fitness of the individuals in $P(t)$
- 2) While (termination condition is not satisfied) do
 - (a) $t = t+1$
 - (b) select $P(t)$, recombine ($P(t)$) and evaluate($P(t)$)
- 3) Output the best solution among all the population as the (sub-optimal) solution

شکل ۱۲- شمای کلی یک الگوریتم ژنتیک

الگوریتم‌های تکاملی با چندین ورودی کنترل می‌شوند:

(۱) اندازه‌ی جمعیت

(۲) نرخ اعمال جهش و آمیزش

عموماً تضمینی وجود ندارد که الگوریتم ژنتیک جواب بهینه را در یک مسئله بیابد، اما راه‌هایی وجود دارد، که شانس موفقیت را افزایش می‌دهد. راه‌های فراوانی برای تصویر کردن یک حل مسئله بر کروموزوم وجود دارد. هر مسئله، انتخاب مناسب خود را دارد و یک انتخاب، که در همه‌ی مسئله‌ها به‌ترین باشد، وجود ندارد. در الگوریتم‌های ژنتیک، عموماً کروموزوم‌ها رشته‌ای از صفر و یک را ارائه می‌دهند، اما تصویرهای دیگری نیز از مسئله می‌توان ارائه کرد. در برنامه‌سازی ژنتیکی، حل‌ها به‌صورت برنامه‌های کامپیوتری تصویر می‌شوند. استراتژی تکاملی و برنامه‌سازی تکاملی از اعداد اعشاری استفاده می‌کنند، که برای پارامترهای عملگرهای شان، مناسب‌تر است. غالباً شیوه‌ی تصویر کردن^۱ مسئله، بستگی به ورودی عملگرها دارد. استراتژی تکاملی^۲ و برنامه‌سازی تکاملی^۳ بیش‌تر از عملگرهای جهش^۴ تصادفی بهره‌می‌گیرند.

۵) نقش الگوریتم‌های تکاملی در داده‌کاوی

الگوریتم‌های تکاملی می‌توانند در بهبود دقت^۵ و استحکام^۶ تکنیک‌های قدیمی‌ای که در مراحل مختلف داده‌کاوی، یعنی «استخراج ویژگی»، «انتخاب ویژگی»، «دسته‌بندی» و یا «خوشه‌بندی» مفید و مؤثر باشد.

۵-۱) الگوریتم‌های تکاملی و استخراج ویژگی

فرآیند استخراج ویژگی بسیار وابسته به داده و مسئله است. در بعضی انواع داده، نسبتاً به آسانی می‌توان ویژگی‌ها را مشخص کرد. مثلاً در داده‌ی متنی، کلمات ویژگی اصلی هستند و یا در تحلیل سبدبازار، ویژگی مشخص و مهم، اقلام جنسی هستند که در یک دادوستد، خریداری شده‌اند. در هر مورد، پردازش خاصی برای استخراج ویژگی از میان مجموعه‌ی داده‌ای خام لازم است.

در مورد داده‌ی تصویری یکی از مشکل‌ترین قسمت‌ها استخراج ویژگی است. و ضمناً با گسترش فراوان به کارگیری تصویر و داده‌کاوی‌های تصویری، ایجاد یک تکنیک عمومی، مستحکم و دقیق در مورد این مسئله مهم به‌نظر می‌رسد. الگوریتم‌های تکاملی به‌صورت بسیار مؤثری، برای حل این مشکل، کارایی دارند. موارد عمده‌ی استفاده‌ی از الگوریتم‌های تکاملی در داده‌کاوی تصویری را می‌توان در بخش‌بندی تصاویر، تصاویر طبی، تشخیص لبه‌ها در تصویر و ثبت تصویر دانست.

¹ encode

² ES

³ EP

⁴ Mutation

⁵ Accuracy

⁶ Robustness

۵-۲) الگوریتم‌های تکاملی در انتخاب ویژگی

چندین دلیل برای دست‌یازیدن به انتخاب ویژگی وجود دارد، که وجود آن‌را در غالب مواردی که از داده‌کاوی استفاده می‌شود، امری ضروری می‌سازد:

- ۱) در بسیاری موقیت‌ها، دانستن این‌که کدام‌یک از ویژگی‌های استخراج شده با مسئله‌ای که در دست است، بیش‌تر مرتبط‌اند، کاری ممکن و میسر نیست.
- ۲) وجود ویژگی‌های نامربوط نه تنها پیچیدگی زمانی اجرای بسیاری از الگوریتم‌ها را افزایش می‌دهد، بلکه زمان مورد نیاز برای استخراج این ویژگی‌ها از میان داده‌های خام نیز تبدیل به معضلی خاص و زمان‌گیر می‌شود.
- ۳) علاوه‌براین، با توجه به این‌که ابعاد فضای ویژگی انتخاب می‌شود، در هنگام رورویی با ویژگی‌های کاهش‌یافته، نمونه‌های آموزشی کم‌تری برای یادگیری مورد نیاز هستند.
- ۴) ضمن آن‌که امکان دارد برخی از ویژگی‌ها با توجه به تطابق زمانی و مکانی، برای استخراج هزینه‌ی تطبیق زمانی و مکانی نیز داشته باشند. نتیجتاً این موارد باید در طول پروسه‌ی داده‌کاوی وزن‌دهی شوند. همه‌ی این موارد ما را به سوی یک مسئله می‌کشاند و آن این است که، از میان مجموعه‌ی خصوصیات داده‌ها، زیرمجموعه‌ای که برای ارائه‌ی الگوهای متناسب با مسئله، مناسب‌ترند را، از میان مجموعه‌ی بزرگ ویژگی‌ها، که احتمالاً بسیاری از آن‌ها زائد و بی‌مصرف هستند، انتخاب کنیم. ساده‌ترین راه برای حذف خصوصیات نامرتبط و زائد، به‌کارگیری دانش آن عرصه، توسط خبرگان آن فن می‌باشد. مثلاً در خوشه‌بندی مستندات متنی، ما جزئیاتی مثل "a", "an", "the" و... را متغیرهای بی‌ربط قلمداد نموده و آن‌ها را حذف می‌کنیم. اما این کار تنها زمانی امکان‌پذیر است که یک خبره‌ی دامنه به‌آسانی بتواند صفات نامرتبط را تشخیص دهد، که این امر در موارد نادری امکان‌پذیر است. تکنیک‌های پیچیده‌تر، مثل تحلیل مؤلفه‌های عمده نیز می‌توانند در تعیین ترکیبات خطی ویژگی‌ها به وسیله‌ی انتخاب آن‌ها از میان موارد مستقیمی که بیش‌ترین واریانس را دارند، به‌کارگرفته شوند. الگوریتم تکاملی برای انتخاب خصوصیت، غالباً ترکیب انتخاب با یادگیری را مورد استفاده قرار می‌دهد. (چیزی مانند راهبرد wrapper) در این راهبرد، قابلیت زیرمجموعه‌های ویژگی‌ها، که در طول محاسبات تکاملی مشخص می‌گردد، با استفاده از خود الگوریتم یادگیری ارزیابی می‌شود. نخستین تلاش در این زمینه را [اسکلانسکی و سیدلکی، ۸۹] انجام دادند. یک کروموزوم، با مجموعه‌ای از اوزان صفر و یک، که در آن یک نشان‌دهنده‌ی حضور ویژگی متناظر آن سلول در دسته بندی و صفر نشان‌گر عدم حضور آن ویژگی بود، از مجموعه‌ی ویژگی‌ها به وجود می‌آمد. سپس از الگوریتم k-NN برای ارزیابی میزان مناسبیت هر فرد، براساس دقت دسته‌بندی آن و تعداد ویژگی‌هایی که مورد استفاده قرار داده‌بود به‌دست می‌آمد. در بقیه‌ی تحقیقاتی که در این زمینه پس از آن انجام شد نیز، از همین روش تصویر مسئله (صفر و یک، برای حضور و عدم حضور هر ویژگی) مورد استفاده قرار گرفت. مثلاً [بریل، برون و مارتین، ۹۰] که دسته بندی را با استفاده از شبکه‌های عصبی انجام دادند و یا [برادرتون و سیمپسون، ۹۵].

[پانچ و دیگران، ۹۳] ایده‌ی انتخاب ویژگی دودویی را با ارئه‌ی یک کروموزوم که با سری اوزان از صفر تا ۱۰ برای هر ویژگی، مشخص می‌شد، توسعه دادند. در این راهبرد بعضی از ویژگی‌ها مهم‌تر از دیگر ویژگی‌ها به حساب می‌آمدند و در داده‌کاوی تعیین‌کنندگی بیش‌تری داشتند. در نهایت پانچ و دیگران به این نتیجه رسیدند که در

مواردی که داده‌ها مغتشش هستند، روش وزن‌دهی (صفر تا ۱۰) بسیار بهتر از روش صفر و یک کار می‌کند. از آن‌جا که در دنیای واقعی غالباً با داده‌های مغتشش روبه‌رو هستیم، این روش در بیش‌تر موارد می‌تواند کارایی به‌تری داشته‌باشد.

[وفائی و دژانگ، ۹۸] نیز از یک راهبرد مشابه برای انتخاب خصوصیت در استفاده از درخت تصمیم‌گیری برای دسته‌بندی، سود جستند. در کارهای آن‌ها به جای آن‌که وزن‌دهی تنها به هر ویژگی صورت‌یرد، اجازه‌ی وزن‌دهی به ترکیبی از ویژگی‌های موجود نیز داده می‌شد. ترکیب ویژگی‌ها که با استفاده از اعمال ساده‌ای چون جمع و ضرب و تقسیم صورت‌می‌گرفت، می‌توانست ویژگی جدیدی را از میان ویژگی‌های موجود فراهم‌آورد. این تبدیل فضای ویژگی متغیر، کاهش قابل ملاحظه‌ای را در تعداد ویژگی‌ها فراهم می‌آورد و دقت دسته‌بندی را به نحو چشم‌گیری بهبود می‌داد. از کارهای دیگری که در این رابطه انجام‌گرفته می‌توان به [یانگ و هاناوار، ۹۷] که برای ایجاد یک سیستم دسته‌بند، از یک شبکه‌ی عصبی بهره‌می‌جست و استراتژی ساده‌ی صفر و یک را برای وزن‌دهی به ویژگی‌ها اختیار کرده بود، اشاره کرد.

۵-۳) الگوریتم‌های تکاملی و دسته‌بندی

الگوریتم‌های تکاملی را می‌توان برای بهبود دسته‌بندی انجام‌گرفته، در ترکیب با ابزارهای متداول دسته‌بندی به‌کاربرد. این عمل باعث قابل اعتمادتر شدن مدل دسته‌بندی ساخته‌شده می‌گردد. در زیر کاربرد الگوریتم تکاملی را در سه دسته‌ی شاخص از ابزارهای دسته‌بندی، یعنی سیستم‌های مبتنی بر قانون، شبکه‌های عصبی و درخت‌های تصمیم‌گیری می‌پردازیم. اولی و سومی را به‌طور کلی برمی‌رسیم و دومی را با جزئیات بیش‌تری مورد بررسی قرار می‌دهیم.

۵-۳-۱) الگوریتم‌های تکاملی و کاربرد آن‌ها در سیستم‌های Rule Based

ارائه‌ی مفاهیم به شکل مجموعه‌ای از قوانین، مدت‌هاست که یکی از محبوب‌ترین روش‌ها در یادگیری ماشین می‌باشد. علت این است که قوانین به آسانی می‌توانند به‌صورت دستورات کامپیوتری ارائه شوند و انسان‌ها به آسانی قادر به تفسیر آن‌ها هستند.

برای استفاده از الگوریتم‌های تکاملی در سیستم‌های Rule Based دو راهبرد عمده وجود دارد:

۱) راهبرد میشیگان [holland,75] که در [بوکر، هلند و گلدبرگ، ۸۹] مورد بررسی قرار گرفته است. در این راهبرد هر فرد جمعیت یک قانون با طول ثابتی را ارائه می‌کند. در این‌جا همه‌ی جمعیت بر روی هم مفهوم را ارائه می‌دهند.

۲) راهبرد پیتزبورگ [smith,80,83] که در [دژانگ، اسپیرس و گوردون، ۹۳] مورد شرح و بررسی قرار گرفته است. هر فرد در این راهبرد، ارائه‌کننده‌ی یک مجموعه‌ی کلی از قوانین است و لذا طولی متغیر دارد.

هر دو راهبرد، محاسن و عایب خود را دارند و هر یک با موفقیت در برخی انواع سیستم‌های دسته‌بند، به‌کار گرفته شده‌اند.

سیستم‌های دسته‌بند، سیستم‌های rule-based هستند که الگوریتم‌های تکاملی را با یادگیری تقویتی ترکیب می‌کنند. حلقه‌ی اصلی در یک سیستم دسته‌بند، آن است که سیستم با ورودی‌هایی از محیط، مهیا می‌شود،

ورودی‌ها به پیام تبدیل می‌شوند. پیام‌هایی که به یک لیست پیام اضافه می‌شوند و قوی‌ترین قانونی که با هر پیام منطبق شود، انگیزته می‌شود. قوانین یک مقدار برآزندگی مبتنی بر پاداش برگشته از محیط دارند. یک الگوریتم ژنتیک به عنوان مؤلفه‌ی کشف سیستم به‌کار می‌رود. که قوانین جدید را بر مبنای به‌ترین قوانین فعلی شکل می‌دهد. در اینجا فرصت پرداختن به دسته‌بندها [cs] نیست. خواننده‌ی علاقه‌مند می‌تواند به کتاب [گلدبرگ، ۸۹] به عنوان مقدمه‌ای بر cs، یا مقالات [ویلسون، ۹۵ و 2000a] که توسعه‌های مختلف را شرح می‌دهند مراجعه کند. [ویلسون و گلدبرگ، ۸۹] و [ویلسون، 2000b] نیز به مرور چشم‌اندازی بر تحقیقات انجام گرفته بر روی XCS می‌پردازد. سیستم‌های دسته‌بند، معمولاً به‌عنوان سیستم کنترل در محیط‌های متغیر و غیرمطمئن به‌کار می‌روند. جایی که احتمال دانش خیره‌ی روشن و کافی وجود ندارد، تا بتوانیم از روش‌های متداول استفاده کنیم. [گلدبرگ، 1983].

در مورد مبحث مورد بحث ما، داده‌کاوی، سیستم‌های دسته‌بند، برای آموختن توابع بولی مورد استفاده‌اند. [ویلسون، ۹۵] که استفاده‌ی آن‌ها بسیار با اهمیت است، زیرا توانایی یادگیری مفاهیم پیچیده‌ی غیرخطی را برآورده می‌سازند. کاربردهای دیگری شامل دسته‌بندی نامه‌ها [فری، ۹۱] و تشخیص سرطان پستان [ویلسون، ۲۰۰۰] نیز برای آن‌ها ارائه شده است.

طرف چپ قوانین در این سیستم‌ها یک اجتماع از عبارات است. این امر قدرت توصیفی این قوانین را در قیاس با (مثلاً) منطق درجه‌ی اول، بسیار بالایی برد. مشکل منطق درجه‌ی اول چه در راهبرد معین و چه در راهبرد اکتشافی، ماندن در بهینه‌های محلی است. برای حل این مشکل می‌توان از الگوریتم تکاملی بهره گرفت. در اینجا مسئله‌ی عمده پیدا کردن فرمی برای قوانین است، که عملگر تکاملی بتواند روی آن‌ها به صورت مؤثری کارگرفتند. و بدین سان پیشرفت قوانین محسوس باشد.

۵-۳-۲) الگوریتم‌های تکاملی و شبکه‌های عصبی

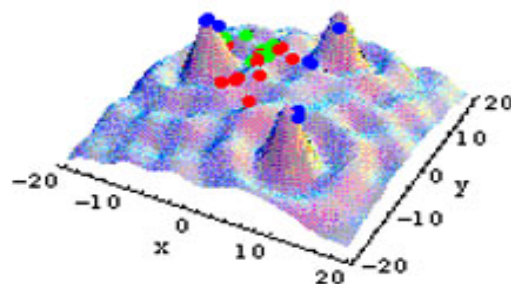
دو روش عمده برای استفاده از الگوریتم‌های تکاملی و شبکه‌های عصبی به‌همراه یک‌دیگر وجود دارد:

- ۱) الگوریتم تکاملی در آموزش شبکه‌ی عصبی مورد استفاده قرار گیرد.
- به‌خصوص، از الگوریتم‌های تکاملی در یافتن اوزان یال‌های شبکه‌ی عصبی و جست‌وجو برای مقدار مناسب پارامترهای یادگیری استفاده می‌شود.
- ۲) نوع عمده‌ی دیگر همکاری الگوریتم‌های تکاملی و شبکه‌های عصبی، استفاده از الگوریتم‌های تکاملی در طراحی شبکه‌ی عصبی است.
- ساختار شبکه‌ی عصبی تأثیر فراوانی بر کارایی آن و توانایی‌اش در حل مسائل دارد. این امر مشهور است که در مسائل جداسازی غیرخطی شبکه حداقل به یک لایه‌ی میانی نیاز دارد. اما تشخیص تعداد و اندازه‌ی لایه‌های پنهان، بیش‌تر یک امر آزمون و خطاست. از الگوریتم تکاملی می‌توان در جست‌وجوی این ساختار و پارامترهای آن استفاده کرد.

۵-۳-۱) استفاده از الگوریتم تکاملی در یادگیری شبکه‌ی عصبی

نخست باید بدانیم که یادگیری شبکه‌ی عصبی مسئله‌ی بهینه‌سازی است، که هدف در آن یافتن مجموعه‌ای از اوزان است، که اندازه‌ی خطا را کمینه می‌کند. فضای جست‌وجو ابعاد فراوانی دارد و دربردارنده‌ی چندین بهینه‌ی محلی است. در روش‌های سنتی آموزش شبکه، مانند انتشار عقب‌گرد، امکان درج‌آزاد الگوریتم در یک بهینه‌ی

محلی و گیر افتادن در آن وجود دارد. البته این امر را می‌توان با اصلاح پویای پارامتر لاندا بهبود داد، اما در کل این مشکل برای اغلب الگوریتم‌های انتشار عقب‌گرد وجود دارد. در مقابل الگوریتم ژنتیک با نمونه‌سازی هم‌زمان چندناحیه از فضای مسئله از به دام افتادن در بهینه‌های محلی اجتناب می‌کند. [شکل ۵] این امر بستگی مستقیمی به احتمال جهش در الگوریتم تکاملی دارد. با بالا رفتن احتمال جهش، احتمال گیر افتادن در بهینه‌ی محلی کاسته می‌شود.



شکل ۱۳- شیوه‌ی کار الگوریتم تکاملی در یافتن جواب بهینه موجب اجتناب از افتادن در بهینه‌ی محلی می‌شود.

یک ترکیب ساده از الگوریتم‌های تکاملی و شبکه‌های عصبی، استفاده از الگوریتم‌های تکاملی در جست‌وجوی اوزانی است که باعث می‌شوند، شبکه آن‌گونه که مورد انتظار است، عمل نماید. در این‌جا برازندگی هر شبکه با اندازه‌گیری دقت دسته‌بندی یا رگرسیون روی مجموعه‌ی آموزشی به دست می‌آید. بنابراین برای هر ارزیابی باید کل مجموعه‌ی آموزشی مورد ارزیابی قرار گیرد. اگر مجموعه‌ی آموزشی بزرگ باشد، این کار می‌تواند بسیار وقت‌گیر باشد. اما قابلیت را می‌توان تنها با نمونه‌ی کوچکی از مجموعه‌ی آموزشی نیز تخمین زد. گرچه در کاربردهای مختلف، قابلیت‌ها متفاوت است، الگوریتم‌های تکاملی به جست‌وجوی خوب در داده‌های دچار اغتشاش، معروف‌اند.

سه روش مختلف اصلی برای آموزش شبکه‌ی عصبی به کمک الگوریتم تکاملی وجود دارد:

(۱) شروع از یک جمعیت تصادفی و به کارگیری اوزانی که توسط الگوریتم تکاملی یافته شده‌اند، بدون هیچ‌گونه پالودگی (تقویت) بیش‌تر. از این روش در [کودل و دلان، ۸۹]، [مانتانا و دیویس، ۸۹]، [واتیلی و هانسن، ۸۹] استفاده شده است.

(۲) استفاده از روش انتشار عقب‌گرد یا روش دیگری برای پالایش اوزان یافته شده به وسیله‌ی الگوریتم تکاملی. [کیتانو، ۹۰]

انگیزه‌ی استفاده از این روش آن است که الگوریتم تکاملی به سرعت می‌تواند نواحی امیدبخش را جست‌وجو کند، اما در این کار عمل تنظیم پارامترها را به‌کندی انجام می‌دهد. بنابراین برای یافتن یک مجموعه‌ی نواحی امیدبخش اولیه از یک روش مبتنی بر گرادیان (مثلاً روش انتشار عقب‌گرد) بهره می‌برد. بنابراین مثلاً در هر دوره‌ی آموزش، در میان داده‌های آموزشی مورد قبول واقع شده، از یک انتشار عقب‌گرد استفاده می‌شود. و این امر باعث افزایش زمان پردازش می‌گردد.

(۳) استفاده از الگوریتم تکاملی در پالودن نتایج یافت شده به وسیله‌ی الگوریتم‌های یادگیری سنتی شبکه‌ی عصبی [کابادا و نایگارد، ۱۹۹۰]. جمعیت اولیه در این روش توسط الگوریتم انتشار عقب‌گرد، تعیین می‌شود.

اما به راحتی چرا چنین تلاش‌هایی در ترکیب الگوریتم تکاملی با دیگر روش‌ها صورت می‌گیرد؟ نخست آن که، طول کروموزوم‌ها به سرعت با اندازه‌ی شبکه رشد می‌کند. از آن‌جا که غالباً لایه‌های هم‌جوار در شبکه‌ی عصبی کاملاً متصل هستند و همه‌ی اوزان در این روش باید ارائه‌گردند، تعداد ایت‌هایی که باید در هر کروموزوم ارائه‌شوند، با $O(n^2)$ متناسب است. (که در آن n تعداد نورون‌هاست).

علاوه‌براین، عموماً کروموزوم‌های این چنین نیاز به جمعیت بزرگ‌تری برای آموزش دارند، که هزینه‌ی محاسباتی بدین ترتیب فوق‌العاده افزایش می‌یابد. بنابراین الگوریتم تکاملی در این حالت فاقد *scalability* است. و در شبکه‌های عصبی بزرگ کارایی خود را از دست می‌دهد.

عیب دیگری که وجود دارد و اصطلاحاً مشکل جایگشت نامیده می‌شود، آن است که اگر ترتیب گره‌های پنهانی جایگشتی باشد، ارائه‌ی وزن‌ها متفاوت خواهد بود. زیرا ترتیب گره‌ها تناسب‌هایی را به وجود می‌آورد و الگوریتم تکاملی می‌تواند این ترتیب‌های مناسب گره‌های پنهانی را در هم‌ریزد. برای حل این مسئله نیز راه‌هایی مانند امکان تبدیل ورودی و خروجی گره‌ها به یک‌دیگر، پیشنهاد شده است. هنکک، در سال ۹۲ برای اعتقاد بود که مشکل جایگشت‌ها به این سختی که نشان داده می‌شود، نباید باشد. و تیرنس در سال ۹۵، شیوه‌ای از ارائه‌ی الگوریتم‌های تکاملی به وجود آورد، که از مشکل جایگشت کاملاً اجتناب می‌کرد.

۵-۳-۲) استفاده از الگوریتم تکاملی در طراحی شبکه‌ی عصبی

دو راهبرد کلی برای استفاده از الگوریتم تکاملی در طراحی توپولوژی شبکه‌ی عصبی وجود دارد:

۱) استفاده از ارائه‌ی مستقیم، برای مشخص نمودن اتصالات شبکه

۲) فراهم آوردن یک توصیف غیرمستقیم از اتصالات گره‌های شبکه

شبکه‌ای که در نتیجه‌ی این طراحی به وجود می‌آید، می‌تواند با یکی از راه‌های یادگیری شبکه‌ی عصبی، مورد آموزش قرار گیرد. می‌توان از الگوریتم تکاملی در یافتن پیکربندی و اوزان، به صورت هم‌زمان بهره‌برد.

معمولاً شبکه‌ی عصبی به صورت یک گراف در نظر گرفته می‌شود. که در آن هر گره، نشان‌گر یک نورون و هر یال نمایاننده‌ی یک اتصال در شبکه‌ی عصبی است. روش معمول آن است که با استفاده از یک ماتریس دودویی اتصال (که در آن عنصر (i,j) برابر یک است، اگر اتصالی بین گره i و j وجود داشته باشد و در غیر این صورت صفر است)، گراف را به صورت مستقیم تصویر می‌کنند. به آسانی می‌توان با *concat* سطرها یا ستون‌های این ماتریس به یک‌دیگر آن‌را به صورت یک رشته‌ی دودویی درآورد.

چنانچه توپولوژی شبکه به این صورت (یعنی با استفاده از الگوریتم تکاملی) تعیین شود، فرآیند یادگیری سریع‌تر از حالت عادی طی خواهد شد. اما در این روش نیز طول کروموزوم $O(n^2)$ است و هم‌چنان *scalability* نخواهیم داشت.

یک روش برای بهبود آن است که نیروی خود را بر مسئله‌ی خاصی که روی آن کار می‌کنیم متمرکز کنیم و الگوریتم یادگیری خاصی برای آن تدوین نماییم. و سپس از الگوریتم‌های تکاملی برای تعیین پارامترهایی که توصیف شبکه را کامل می‌کنند، بهره‌گیریم. مثلاً در توپولوژی *fully connected feedforward* از الگوریتم تکاملی می‌توانیم برای جست‌وجوی تعداد لایه‌ها و گره‌های هر لایه برای به‌ترین کارایی بهره‌گیریم.

یک راهبرد پیچیده‌تر ارائه‌ی غیرمستقیم، استفاده از گرامر برای کدگذاری قوانینی است، که می‌توانند یک شبکه را به وجود آورند. کیتانو در سال ۱۹۹۰ نخستین بار راهبرد مبتنی بر گرامر را معرفی نمود. بدین سان در کروموزوم‌ها

قوانین نگهداری می‌شوند و در موقع ارزیابی از روی این قوانین ماتریس گراف تشکیل می‌شود و سپس ماتریس اتصال برای ساختن شبکه مورد تفسیر قرار می‌گیرد.

۳-۳-۵) الگوریتم تکاملی و درخت تصمیم‌گیری

درخت تصمیم‌گیری یکی از راه‌های محبوب دسته‌بندی است. زیرا به آسانی می‌توان آنرا ساخت و نخبگان به آسانی آنرا تفسیر می‌کنند. گره‌های داخلی درخت ارائه‌کننده‌ی تستی برای ویژگی‌ها هستند و برگ‌ها ارائه‌کننده‌ی برچسب دسته‌اند. یک مسیر از ریشه‌ی درخت تا یکی از برگ‌ها ارائه‌کننده‌ی یک اجتماع از تست‌هاست. از آنجا که برنامه‌سازی ژنتیکی به‌طور سنتی از درخت‌ها برای نمایش بهره‌می‌گیرد، استفاده از آن برای یافتن درخت تصمیم‌گیری مناسب است.

۴-۵) استفاده از الگوریتم‌های تکاملی در خوشه‌بندی

دو راهبرد کلی در این زمینه وجود دارد:

۱) هر جایگاه در کروموزوم، ارائه‌کننده‌ی یک آیت‌م از مجموعه‌ی آموزشی است. (این راهبرد را راهبرد «مبتنی بر بخش‌بندی» می‌نامیم).

۲) تنها مراکز خوشه‌ها در هر کروموزوم، مورد بررسی قرارگیرند. (راهبرد مبتنی بر مرکز)

۱-۴-۵) راهبرد مبتنی بر بخش‌بندی

در این راهبرد هر جایگاه در کروموزوم ارائه‌کننده‌ی یک آیت‌م از مجموعه‌ی آموزشی است. در این‌جا وظیفه‌ی الگوریتم تکاملی، یافتن خوشه‌ی مناسب برای هر آیت‌م است. اگر تعداد خوشه‌ها k (ازپیش مشخص) باشد، مقدار هر جایگاه در بازه‌ی $[1, k]$ خواهد بود. مشکل این روش عدم $scalability$ آن است. مثال بارز آن هم [مورتی و چاوداری، ۱۹۹۶] است.

صورت تغییر یافته‌ای از این روش تصویر مسئله‌ی خوشه‌بندی بر مسئله‌ی بخش‌بندی گراف (Graph Partitioning) است. [پارک و سونگ، ۹۸] این تغییر را در نمایش مستقیم دادند. در این راهبرد آیت‌م‌های مجموعه‌ی داده گره‌های گراف هستند. اندیس نزدیک‌ترین همسایگان ذخیره‌می‌شود. که در این‌جا تعداد همسایگان تابع ورودی مسئله است. اگر آیت‌م i دارای مقدار j باشد، یعنی یک اتصال در گراف بین i و j وجود دارد. مقدارهای یک آیت‌م، نشان‌دهنده‌ی اندیس نزدیک‌ترین همسایگان آن است. کار پارک و سونگ روی یک فرهنگ لغت جامع بود و نتایج الگوریتم خود را با الگوریتم‌های خوشه‌بندی دیگر مقایسه کردند. آن‌ها به این نتیجه رسیدند که لزومی به تعیین تعداد خوشه‌ها در این روش وجود ندارد.

در این مسئله نیز هم‌چنان مشکل $scalability$ وجود دارد. زیرا هم‌چنان طول هر کروموزوم به تعداد آیت‌م‌های مجموعه‌ی داده بستگی دارد و محاسبه‌ی نزدیک‌ترین همسایگان یک آیت‌م، برای ابعاد زیاد عملی دشوار است.

۲-۴-۵) راهبرد مبتنی بر مرکز

در [هال، اوزیوت و بازدک، ۱۹۹۹]، راهبردی تکاملی برای تعدیل مرکز خوشه‌ها معرفی گردید. مطالعه‌ی آن‌ها دربرگیرنده‌ی هر دو نوع خوشه‌بندی: $fuzzy, hard$ بود. مسئله‌ی کارایی این الگوریتم مشکل‌ساز است و بسیار

زمان برتر از الگوریتم‌های معمولی مثل FCM/HCM است. در برابر این عیب، یک حُسن نیز دارد و آن امکان استفاده از آن در جاهایی است، که از هیچ روش بهینه‌سازی متداول دیگری نمی‌توان بهره‌گرفت. [بابو و مورتی، ۹۳] راهبردی برای استفاده از الگوریتم‌های تکاملی در جست‌وجوی اولیه‌ی مراکز خوشه‌ها و سپس استفاده از الگوریتم‌های خوشه‌بندی متداول، کشف‌کردند.

همانند دیگر عرصه‌ها، در اینجا نیز می‌توانیم از دانش عرصه، برای بهبود کارایی بهره‌گیریم. این استفاده در این‌جا به چندین طریق می‌تواند صورت‌گیرد. یک طریق، طراحی عملگرهای خاص برای آن عرصه به شکلی است که نسبت به عملگرهای معمول کارا تر باشند. طریق دیگر ترکیب الگوریتم تکاملی با الگوریتم‌های خوشه‌بندی معمول است. [فرانتی و دیگران، ۱۹۹۷] هر دوی این روش‌ها را به‌کارگرفتند. آن‌ها ۵ روش متفاوت آمیزش را استفاده‌نمودند که ۳ تای آن‌ها اختراع خودشان بود. پس از هر آمیزش، عضو جدید به تکراری از الگوریتم k -mean نارد می‌گردید. آن‌ها بعدها الگوریتم‌شان را برای تغییر خودکار پارامترها توسعه‌دادند. [کیویجاروی، ۲۰۰۰] در [فرانتیو دیگران، ۹۷] ضرورت استفاده از تکرارهای k -mean بررسی و ثابت‌گردیده‌بود. این امر باعث افت سرعت فوق‌العاده‌ی این الگوریتم می‌گردید. بعدها نشان‌داده‌شد که تنها به تعداد اندکی تکرار k -mean برای کارکرد درست و کارایی الگوریتم نیاز است. [کریشنا و مورتی، ۱۹۹۹] این الگوریتم را با یک و تنها یک تکرار k -mean اجرا کردند و نتیجه‌ی مطلوب را به دست‌آوردند. یک راهبرد برای بهبود این الگوریتم استفاده از نمونه‌های خوب در شروع کار است. مثلاً می‌توان از نتایج چند تکرار k -mean روی داده‌ها به‌عنوان مراکز اولیه‌ی خوشه‌ها استفاده‌کرد. [مورتی و چاداری، ۱۹۹۶] از این راهبرد استفاده‌کردند.

حُسن کلّی این روش در این است که $scalable$ است، زیرا طول کروموزوم، بستگی به ابعاد مسئله و تعداد مراکز خوشه‌ها دارد و نه به تعداد نمونه‌های داده‌ی آموزشی. بنابراین روش مبتنی بر مرکز از این لحاظ بر روش مبتنی بر بخش‌بندی ارجحیت دارد. حُسن دیگر این روش آن است که امکان فرم‌پذیری بیش‌تری برای خوشه‌ها فراهم‌می‌کند و مثلاً قسمت‌های ناهم‌جوار فضای داده می‌توانند به یک خوشه‌تعلق‌بگیرند.

۵-۵) کارایی الگوریتم‌های تکاملی

الگوریتم‌های تکاملی کارایی خود را در حل مسائل واقعی داده‌کاوی اثبات‌کرده‌اند. خصوصاً در مواردی که داده‌ها دچار اغتشاش است و یا به حلی نیاز داریم، که مسئله‌ی بهینه‌سازی چندم‌اموریتی را انجام‌دهد، این الگوریتم‌ها کارایی فراوانی دارند. اما در هر حال خالی از عیب نیز نیستند. مهم‌ترین عیبی که در چندین جا و توسط چندین نویسنده در عرصه‌های مختلف برای آن‌ها ذکر شده، آن است که می‌توانند بسیار وقت‌گیر باشند. و تبدیل به گلوگاه زمانی برنامه شوند.

چندین روش برای رفع این مشکل پیشنهاد شده‌است:

- ۱) استفاده از قسمتی از نمونه‌ی آموزشی به‌جای استفاده از همه‌ی آن
- ۲) نگه‌داری از یک جمعیت متشکل از $global\ fit\ individuals$ و به کارگیری آن‌ها در اکثر موارد. (این مورد تنها در کاربردهای کاوش تصویر مطرح است.)
- ۳) استفاده از پردازش موازی و موازی‌سازی الگوریتم‌های ژنتیک
- ۴) استفاده از ارائه‌ی به تری از مسئله با توجه به ویژگی‌های خاص آن و یا طراحی عملگرهای مقتضی همان مسئله.

۶) مرور مطالب و نتیجه گیری

- ✓ الگوریتم های تکاملی می توانند در تکمیل بسیاری از الگوریتم های داده کاوی موجود مفید و مؤثر باشند.
- ✓ در استخراج ویژگی، گزینش ویژگی، وزن دهی ویژگی، یادگیری شبکه ی عصبی، یافتن قوانین دسته بندی، ساختن درخت های تصمیم گیری و خوشه بندی به کار آیند.
- ✓ الگوریتم های تکاملی به خصوص در هنگامی که مسئله شامل بهینه سازی یک تابع ناهموار و چند بُعدی (not smooth and differentiable) یا تابعی که در آن مقدار تابع هدف در طول زمان تغییر می کند، بسیار مفید فایده اند.
- ✓ با این وجود هزینه ای که این الگوریتم ها برای اجرا دارند بسیار زیاد است، که البته با توجه به پیشرفت های پردازشی و توانایی های موازی سازی، امیدواریم این نقیصه در آینده ی نزدیک برطرف شود.
- ✓ در این مقاله به مطالعه ی کلی کاربردهای الگوریتم های تکاملی در قسمت های مختلف داده کاوی پرداختیم. سعی بر آن است که در آینده ی کاربردهای جزئی مورد مطالعه قرار گرفته و راهکارهای جدیدی در عرصه ی مسائل مختلف داده کاوی ارائه گردد. هم چنین نگارنده قصد دارد به بررسی امکان موازی سازی این راهبردهای مختلف معرفی شده در این جا، بپردازد.

سپاس گذاری و حدیث آخر

در این جا لازم می دانم از استاد عزیز جناب آقای دکتر ناصر قاسم آقائی که مسبب انجام ایت تحقیق و ارائه ی آن در درس سیستم های خبره و مهندسی دانش بودند، تشکر نمایم. همچنین از آقای بهروز شاهقلی که به مطالعه ی بخشی از مقاله پرداخته و راهنمایی های لازم را جهت تصحیح آن به عمل آوردند، نهایت سپاس گذاری را دارم. این مقاله برای درس سیستم های خبره و مهندسی دانش تهیه گردیده و بیش تر مروری بر مراجع معرفی شده در قسمت مرجع داشته است. ضمناً به عنوان ضمیمه ی این مقاله ترجمه ی دو مقاله از مقاله های مرجع (شماره ی ۱ و ۴) انجام گرفته است.

روزی که این موضوع را به عنوان موضوع سمینار خود برای درس سیستم های خبره تعیین نمودم، با توجه به سابقه ی ذهنی ای که از داده کاوی و الگوریتم های ژنتیک داشتم، می پنداشتم که کاربردهای الگوریتم های ژنتیک تنها به بخشی محدود از داده کاوی منحصر است و راه کوتاهی در این بررسی پیش روی دارم. اما بررسی بیش تر نشان داد، که این موضوع دامنه ی وسیعی را دربرمی گیرد، که برای بررسی کامل آن باید ده ها هزار application و مقاله ی معرفی شده در ارتباط با آن ها مورد مطالعه قرار گیرد. و ناچار تنها به بررسی مباحث عمده ی مطرح شده در این موضوع پرداخته و فرصت نظر در جزئیات هر کاربرد را پیدانکردم. اما گستردگی این بحث حُسنی نیز دارد و آن شاخه های مختلف تحقیق خیزی است که برای بهبود کارکردها و توسعه ی کاربردها، در این موضوع وجود دارد. نگارنده ی این سطور امیدوار است که این مقاله ابتدای به کارگیری این راهکارهای کلی در جهت بهبود کاربردهای مختلف باشد.

Main References:

- 1) Cantú-Paz, E. and Kamath, C., "ON THE USE OF EVOLUTIONARY ALGORITHMS IN DATA MINING"(2001), Center for Applied Scientific Computing, Lawrence Livermore National Laboratory.
- 2) Two Crows Corporation, "Introduction to data mining and knowledge discovery", third edition(1999), www.twocrows.com.
- 3) Data Mining Techniques,(2001) by Michael Berry and Gordon Linoff.
- 4) Minaei-Bidgoli B., Punch III W. F. , Using Genetic Algorithms for Data Mining Optimization in an Educational Web-based System (2003) , GARAGE , Michigan State University.
- 5) Management Information System for Information Age,Third edition(2002), Haag, Cummings, McCubbery,McGraw-Hill Publishing,P 109,154,418.
- 6) Roiger. R. J., Geatz M. W. , "Data Mining (A tutorial based primer),(2001), Addison-Wesley published.

Other References:

- 1) Paradolas, P. M. and Pitsoulis, L. and Marridon, T. and Resede, M. G. C., "parallel Search for combinational optimization: genetic algorithm, simulated annealing, grasp", (1996), P.1~5.
- 2) Cantú-Paz. Eric, "A Survey of Parallel Genetic Algorithms", 1996, Department of computer science and Illinois Genetic Algorithms Laboratory, University of Illinois at Urbana-Champaign.
- 3) Alex A. Freitas , A Survey of Evolutionary Algorithms for Data Mining and Knowledge Discovery (2002) ,Postgraduate Program in Computer Science, Pontificia Universidade Catolica do Parana , Brazil.
- 4) Goldberg D. E., "Genetic Algorithms in Search, Optimization, and Machine Learning", 1989, University of Alabama, 26printing,2004, addison-wesley published.
- 5) افسانه فاطمی، "ایجاد یک درخت جست و جوی دودویی نزدیک بهینه با استفاده از الگوریتم ژنتیک، ۱۳۸۴، دانشگاه اصفهان.
- 6) Hsu, H. W. and Cheng, Y. and Gue, H. and Gustafson, S. M., "Genetic Algorithms for Reformulation of Large-Scale KDD Problems with Many Irrelevant Attributes", (2001), research proposal.