

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده فنی و مهندسی

بخش مهندسی کامپیوتر

پایان نامه تحصیلی برای دریافت درجه کارشناسی ارشد
رشته مهندسی کامپیوتر گرایش هوش مصنوعی

کاربرد محاسبات دی.ان.ای در حوزه داده کاوی

مؤلف:

رامین معاذالهی

استاد راهنما:

دکتر علی اکبر نیک نفس

استاد مشاور:

دکتر مجید محمدی

بهمن ماه ۱۳۹۱



این پایان‌نامه به عنوان یکی از شرایط درجه کارشناسی ارشد به

بخش مهندسی کامپیوتر

دانشکده فنی و مهندسی

دانشگاه شهید باهنر کرمان

تسلیم شده است و هیچگونه مدرکی به عنوان فراغت از تحصیل دوره مزبور شناخته نمی‌شود.

دانشجو: رامین معاذالهی

استاد راهنما: دکتر علی اکبر نیک‌نفس

استاد مشاور: دکتر مجید محمدی

دور ۱: دکتر فرامرز صادقی

دور ۲: دکتر حمید میروزی

معاونت پژوهشی و تحصیلات تکمیلی دانشکده: دکتر مریم احتشام زاده

حق چاپ محفوظ و مخصوص به دانشگاه شهید باهنر کرمان است.

تقدیم به خدایی که آفرید جهان را، انسان را، عقل را، علم را، معرفت را، عشق را و به کسانی که عشقشان را در وجودم دمید.

با سپاس از سه وجود مقدس:

آنان که ناتوان شدند تا ما به توانایی برسیم...

موهایشان سپید شد تا ما روسفید شویم...

و عاشقانه سوختند تا گرمابخش وجود ما و روشنگر راهمان باشند...

پدرانمان

مادرانمان

استادانمان

چکیده

در سال‌های اخیر قدرت محاسبات دی.ان.ای در حل مسائل پیچیده کامپیوتری به اثبات رسیده است. به دلیل قدرت پردازش موازی بسیار زیاد مولکول‌های دی.ان.ای، امکان جستجوی تمام فضای راه‌حل با استفاده از عملگرهای بیولوژیکی میسر می‌شود. در این تحقیق دو روش بر پایه محاسبات دی.ان.ای برای دسته‌بندی داده‌ها و استخراج قوانین همبستگی که در حوزه داده‌کاوی مطرح می‌شوند پیشنهاد می‌کنیم. در روش اول نشان می‌دهیم که چگونه می‌توان مدل دسته‌بندی داده‌ها را با استفاده از رشته‌های دی.ان.ای تولید کرده و با اعمال توالی از عملگرهای بیولوژیکی روی مدل ایجادشده، دسته‌بندی داده‌های جدید را انجام داد. نتایج شبیه‌سازی روش پیشنهادی روی پنج مجموعه داده از UCI و مقایسه با روش C5.0، برتری روش پیشنهادی را از لحاظ دقت دسته‌بندی نشان می‌دهد. در روش دوم نشان می‌دهیم که چگونه می‌توان به طریق مشابه کلیه قوانین همبستگی را با استفاده از رشته‌های دی.ان.ای تولید کرده و با اعمال توالی از عملگرهای بیولوژیکی، قوانینی که حداقل پشتیبان و اطمینان تعیین شده را ارضا می‌کنند در زمان چندجمله‌ای استخراج کرد. لازم به ذکر است که استخراج کلیه قوانین همبستگی در کامپیوترهای سیلیکونی در بدترین حالت دارای پیچیدگی زمانی نمایی است، بنابراین مزیت روش پیشنهادی آشکار می‌شود.

کلمات کلیدی:

محاسبات دی.ان.ای - داده‌کاوی - مدل آدلما لپتون - دسته‌بندی - استخراج قوانین همبستگی.

فهرست مطالب

صفحه	عنوان
۱.....	فصل اول: مقدمه
۲.....	۱-۱ مقدمه
۳.....	۲-۱ داده کاوی
۴.....	۳-۱ محاسبات دی.ان.ای
۵.....	۴-۱ مروری بر کارهای پیشین
۱۶.....	فصل دوم: داده کاوی
۱۷.....	۱-۲ مقدمه
۱۷.....	۲-۲ داده کاوی
۱۸.....	۳-۲ فرآیند داده کاوی
۲۳.....	۴-۲ دسته بندی
۲۵.....	۱-۴-۲ روش C5.0
۲۶.....	۲-۴-۲ ارزیابی مدل دسته بندی
۲۷.....	۵-۲ خوشه بندی
۲۸.....	۶-۲ قوانین همبستگی
۲۸.....	۱-۶-۲ تحلیل سبد خرید
۲۹.....	۲-۶-۲ ارزیابی قوانین
۳۱.....	۳-۶-۲ استخراج قوانین همبستگی
۳۲.....	۷-۲ پیچیدگی زمانی
۳۳.....	فصل سوم: محاسبات دی.ان.ای
۳۴.....	۱-۳ مقدمه
۳۵.....	۲-۳ دی.ان.ای

- ۳-۳ محاسبات با مولکول‌های دی.ان.ای ۳۷
- ۴-۳ کاربرد محاسبات دی.ان.ای ۳۷
- ۵-۳ مزیت کامپیوترهای بیولوژیکی نسبت به الکترونیکی ۳۸
- ۶-۳ عملگرهای پایه ۳۹
- ۱-۶-۳ ترکیب ۳۹
- ۲-۶-۳ واسرشته سازی ۳۹
- ۳-۶-۳ بازپخت ۴۰
- ۴-۶-۳ انعقاد ۴۰
- ۵-۶-۳ جداسازی ۴۱
- ۶-۶-۳ الکتروفورز ژلی ۴۲
- ۷-۶-۳ واکنش زنجیره‌ای پلیمرز (پی.سی.آر) ۴۲
- ۸-۶-۳ برش ۴۴
- ۷-۳ بررسی اولین مساله حل شده با محاسبات دی.ان.ای ۴۴
- ۸-۳ انواع مدل‌های محاسبات دی.ان.ای ۵۰
- ۱-۸-۳ مدل آدلما-لیپتون ۵۰
- ۲-۸-۳ مدل استیکر ۵۳
- ۳-۸-۳ مدل فیلتر کردن موازی ۵۶
- ۴-۸-۳ مدل فیلتر کردن بوسیله انسداد ۵۸
- ۵-۸-۳ مدل سطحی ۶۰
- ۹-۳ ساخت رشته‌های دی.ان.ای مناسب ۶۲
- ۱-۹-۳ محدودیت‌ها ۶۳
- ۲-۹-۳ روش‌ها ۶۵
- فصل چهارم: شرح مساله و روش پیشنهادی ۶۶**
- ۱-۴ مقدمه ۶۷
- ۲-۴ روش پیشنهادی ۱: دسته‌بندی با استفاده از محاسبات دی.ان.ای ۶۷
- ۱-۲-۴ ساخت مدل ۶۸

۷۲.....	۲-۲-۴ دسته‌بندی داده‌های جدید
۷۳.....	۳-۲-۴ شبیه‌سازی مثال نمونه
۷۷.....	۳-۴ روش پیشنهادی ۲: استخراج قوانین همبستگی با استفاده از محاسبات دی.ان.ای
۸۴.....	فصل پنجم: ارزیابی، نتیجه‌گیری و پیشنهادات
۸۵.....	۱-۵ مقدمه
۸۵.....	۲-۵ ارزیابی روش پیشنهادی ۱
۸۷.....	۳-۵ مقایسه روش پیشنهادی ۱ با روش C5.0
۸۸.....	۴-۵ ارزیابی روش پیشنهادی ۲
۸۹.....	۵-۵ نتیجه‌گیری و پیشنهادات
۹۲.....	مراجع

فهرست جداول

شماره و عنوان	صفحه
جدول ۱-۲ اطلاعات بیماران نمونه به همراه داروی تجویز شده	۲۵.....
جدول ۲-۲ مشخصات بیمار جدید	۲۵.....
جدول ۳-۲ مجموعه داده شامل شش تراکنش خرید مشتریان	۳۰.....
جدول ۱-۳ کد کردن گراف شکل ۳-۱۲ با رشته‌های دی.ان.ای مناسب	۴۶.....
جدول ۲-۳ مسیرهای ممکن از گره A_1 به A_3 در گراف شکل ۳-۱۴	۵۲.....
جدول ۳-۳ رشته‌های دی.ان.ای معادل گره‌های گراف شکل ۳-۱۴	۵۲.....
جدول ۱-۵ پیچیدگی زمانی روش پیشنهادی ۱	۸۶.....
جدول ۲-۵ مجموعه داده‌های انتخاب شده برای ارزیابی روش پیشنهادی ۱ از لحاظ دقت دسته‌بندی	۸۷.....
جدول ۳-۵ نتایج دسته‌بندی داده‌ها با روش پیشنهادی ۱ و روش C5.0	۸۸.....
جدول ۴-۵ پیچیدگی زمانی روش پیشنهادی ۲	۸۹.....

فهرست شکل ها

شماره و عنوان	صفحه
شکل ۱-۲ فرآیند کلی داده کاوی	۱۹.....
شکل ۲-۲ ماتریس گيجی نمونه مربوط به دسته بندی بیماران	۲۶.....
شکل ۳-۲ خوشه بندی داده ها در فضای دوبعدی. (الف) در دو گروه (ب) در سه گروه (ج) در چهار گروه	۲۷.....
شکل ۱-۳ ساختار نوکلئوتید	۳۵.....
شکل ۲-۳ ساختار رشته دی.ان.ای با ترتیب ۳'-T-A-C-G-۵'	۳۶.....
شکل ۳-۳ پیوند رشته ۳'-T-A-C-G-۵' با مکمل واتسون خود و تشکیل مارپیچ دوتایی	۳۶.....
شکل ۴-۳ ماشین ترکیب کننده	۳۹.....
شکل ۵-۳ فرآیند واسرشته سازی (ذوب)	۴۰.....
شکل ۶-۳ فرآیند بازپخت	۴۰.....
شکل ۷-۳ فرآیند انعقاد	۴۱.....
شکل ۸-۳ جداسازی رشته های دی.ان.ای حاوی زیر رشته GCTA	۴۱.....
شکل ۹-۳ فرآیند الکتروفورز ژلی برای مرتب سازی رشته های دی.ان.ای بر اساس طولشان	۴۲.....
شکل ۱۰-۳ نحوه برش یک رشته دی.ان.ای توسط آنزیم RSAAI	۴۴.....
شکل ۱۱-۳ گراف جهت دار استفاده شده در آزمایش آدلما	۴۵.....
شکل ۱۲-۳ گراف جهت دار بین چهار شهر	۴۵.....
شکل ۱۳-۳ توپ های فلزی مورد استفاده برای استخراج مسیرهایی که از شهر بوستون میگذرند	۴۹.....
شکل ۱۴-۳ گراف مورد استفاده برای تولید ۲ ^۲ حالت مختلف مقداردهی متغیرهای بولی X و Y	۵۱.....

- شکل ۳-۱۵ ساختار یک رشته حافظه K بیتی ۵۳
- شکل ۳-۱۶ یک کتابخانه (۲,۵) از رشته‌های حافظه ۵۴
- شکل ۴-۱ ساختار مجموعه داده آموزشی ۷۰
- شکل ۴-۲ مجموعه داده آموزشی نمونه ۷۳
- شکل ۴-۳ محتویات تست تیوب T_0 پس از اجرای گام ۱ ۷۴
- شکل ۴-۴ محتویات تست تیوب T_0 پس از اجرای گام ۲ ۷۵
- شکل ۴-۵ محتویات تست تیوب T_1 پس از اجرای گام ۳ ۷۵
- شکل ۴-۶ محتویات تست تیوب T_2 پس از اجرای گام ۳ ۷۶
- شکل ۴-۷ محتویات تست تیوب T_1 پس از اجرای گام ۴ ۷۶
- شکل ۴-۸ محتویات تست تیوب T_2 پس از اجرای گام ۴ ۷۶
- شکل ۴-۹ قوانین استخراج شده برای دسته‌بندی داده $X=(1,3,2)$ ۷۶
- شکل ۴-۱۰ محتویات تست تیوب‌های T_1 تا T_4 پس از اجرای گام ۱ برای $M = 4$ ۷۹
- شکل ۴-۱۱ تراکنشی نمونه از یک فروشگاه با ۴ کالا ۸۰
- شکل ۴-۱۲ ساختار رشته‌های دی.ان.ای MINSUPPTAG ۸۱

فصل اول:

مقدمه

۱-۱ مقدمه

امروزه حجم عظیمی از داده‌های خام توسط شرکت‌ها و سازمان‌ها در پایگاه داده‌ها ذخیره شده و روز به روز به میزان این داده‌ها افزوده می‌شود. اگر این داده‌ها به روش مناسبی پردازش شوند، اطلاعات مفیدی از آنها استخراج می‌شود که در نهایت منجر به ارتقا سازمان و ارائه بهتر خدمات توسط آنها می‌شود. روش‌های زیادی برای تحلیل داده‌ها و استخراج اطلاعات سودمند از آنها وجود دارد که یکی از این روش‌ها داده کاوی^۱ است.

بسیاری از تکنیک‌های مورد استفاده در داده کاوی نیازمند پردازش‌های زیاد توسط سیستم کامپیوتری است که در برخی موارد خارج از توان سیستم‌های موجود می‌باشد. از سوی دیگر در سال‌های اخیر، توانایی سیستم‌های بیولوژیکی در حل مسائل ان.پی.کامل^۲ به اثبات رسیده است. در این سیستم‌ها با بهره‌گیری از قابلیت پردازش موازی بسیار زیاد مولکول‌های دی.ان.ای^۳، کلیه محاسبات در لوله‌های آزمایشی انجام می‌پذیرد که منجر به ظهور علم میان‌رشته‌ای جدیدی بنام محاسبات دی.ان.ای^۴ شده

¹ data mining

⁴ DNA computing

² NP-complete

³ DNA (deoxyribonucleic acid)

است. به دلیل قدرت پردازش موازی بسیار زیاد مولکول‌های دی.ان.ای می‌توان از آن‌ها برای حل مسائل پیچیده کامپیوتری استفاده کرد.

علی‌رغم قدرت زیاد محاسبات دی.ان.ای هنوز پیوند مناسبی میان این حوزه و حوزه داده‌کاوی ایجاد نشده است. از اینرو در این تحقیق می‌کوشیم تا کاربرد محاسبات دی.ان.ای برای دسته‌بندی داده‌ها و استخراج قوانین همبستگی از حوزه داده‌کاوی را بررسی کنیم. در ادامه فصل اول ضمن معرفی اجمالی داده‌کاوی و محاسبات دی.ان.ای به برخی از روش‌هایی که تاکنون در مقالات برای دسته‌بندی و استخراج قوانین همبستگی پیشنهاد شده است اشاره ای داریم. همچنین نگاهی اجمالی به برخی مقالات مهم در زمینه محاسبات دی.ان.ای داریم. در فصل دوم به داده‌کاوی و برخی تکنیک‌های آن می‌پردازیم. محاسبات دی.ان.ای به همراه کاربردها، مزایا، انواع مدل‌ها و مسائل نمونه به صورت کامل در فصل سوم شرح داده می‌شود. در فصل چهارم، دو روش، یکی برای حل مساله دسته‌بندی و دیگری برای استخراج قوانین همبستگی در مساله تحلیل سبد خرید، با استفاده از محاسبات دی.ان.ای پیشنهاد می‌کنیم. در نهایت فصل پنجم به ارزیابی روش‌های پیشنهادی و نتایج مهمی که از این تحقیق حاصل می‌شود اختصاص دارد. در فصل پنجم همچنین پیشنهاداتی برای ادامه تحقیقات در این زمینه ارائه می‌شود.

۱-۲ داده‌کاوی

داده‌کاوی یا کشف دانش فرآیندی است که طی آن داده‌ها از جنبه‌های گوناگون تجزیه و تحلیل شده و بصورت اطلاعات مفید-اطلاعاتی که می‌توانند باعث افزایش درآمد و یا کاهش هزینه شوند- خلاصه می‌شوند. داده‌ها می‌توانند واقعیات، اعداد و یا متن‌هایی باشند که بوسیله کامپیوتر قابل پردازش باشند. امروزه سازمان‌ها حجم وسیعی از داده‌هایشان را در قالب شکل‌های متفاوت و در پایگاه داده‌های گوناگون ذخیره می‌کنند. به عنوان مثال داده‌ها می‌توانند لیست فروش، هزینه و یا سرمایه‌گذاری در یک سازمان باشند.

حاصل پردازش داده‌ها، اطلاعات می‌باشد که بیانگر الگوها، وابستگی‌ها و یا ارتباط بین داده‌هاست. به عنوان مثال، تحلیل داده‌های فروش یک فروشگاه می‌تواند اطلاعات سودمندی از قبیل نحوه چیدمان کالاها به صورتیکه احتمال فروش آن‌ها با هم بیشتر باشد را در اختیار فروشنده قرار دهد.

از جمله کارهایی که در حوزه داده‌کاوی قابل انجام است می‌توان به دسته‌بندی^۱، خوشه‌بندی^۲، استخراج قوانین همبستگی^۳، تخمین^۴ و پیش‌بینی^۵ اشاره کرد. در دسته‌بندی مجموعه‌ای از داده‌ها را در اختیار داریم که در دسته‌های متفاوت قرار داده شده‌اند. هدف دسته‌بندی، پیش‌بینی دسته صحیح برای هر داده است که این کار از طریق پیدا کردن قوانین دسته‌بندی^۶ صورت می‌پذیرد. هر قانون بیانگر رابطه منطقی بین داده‌ها و دسته‌ای است که داده‌ها در آن قرار دارند. پس از استخراج قوانین دسته‌بندی، می‌توان از آن‌ها برای دسته‌بندی داده‌های جدید که دسته آن‌ها مشخص نیست استفاده کرد. خوشه‌بندی نیز مانند دسته‌بندی است با این تفاوت که فقط داده‌ها را در اختیار داریم و دسته‌بندی داده‌ها بایستی با توجه به میزان شباهتشان به یکدیگر صورت پذیرد. در بحث استخراج قوانین همبستگی، هدف پیدا کردن روابط همبستگی و الگوهای تکرار شونده بین مجموعه‌ای از کالاها در پایگاه داده‌های تراکنشی است. با استفاده از قوانین استخراج شده، می‌توان سبدهای خرید کالا را بگونه‌ای طراحی کرد که بیشترین شانس فروش را داشته باشند. در تخمین و پیش‌بینی که شبیه دسته‌بندی می‌باشند، هدف تعیین یک مقدار عددی به عنوان خروجی برای هر داده می‌باشد. تفاوت پیش‌بینی با تخمین این است که در پیش‌بینی نتیجه مربوط به آینده می‌باشد، به عنوان مثال پیش‌بینی نرخ طلا در سه ماهه آینده سال جاری.

۳-۱ محاسبات دی.ان.ای

محاسبات دی.ان.ای شیوه‌ای بدیع از محاسبات است که در آن از مولکول‌های دی.ان.ای بجای مدارات منطقی دیجیتال استفاده می‌شود. در این شیوه محاسبات، هر سلول بیولوژیکی مانند یک کامپیوتر پیچیده با توان پردازشی بسیار بالاست. الفبای محاسبات در این روش، بازهای آمینواسیدی بوده که با حروف A، T، C و G نمایش داده شده و اجزای اصلی تشکیل دهنده مولکول‌های دی.ان.ای می‌باشند. برای کد کردن اطلاعات از حروف ذکر شده بجای صفر و یک رایج در کامپیوترهای الکترونیکی استفاده می‌شود. مزیت اصلی محاسبات دی.ان.ای توان پردازش موازی بسیار زیاد مولکول‌های دی.ان.ای است که امکان حل مسائل پیچیده و بخصوص مسائل ان.پی.کامل را که در برخی موارد از توان پردازشی کامپیوترهای الکترونیکی خارج است، فراهم می‌کند.

¹ classification

² clustering

³ association rule mining

⁴ estimation

⁵ prediction

⁶ classification rules

مسائل ان.پی کامل دسته‌ای از مسائل تصمیم‌گیری هستند که هیچ راه‌حل سریعی برای آن‌ها وجود ندارد. به عبارت دیگر با افزایش اندازه ورودی این مسائل، زمان مورد نیاز برای حل آن‌ها با استفاده از الگوریتم‌های فعلی بصورت نمایی رشد می‌کند. این مسائل نمونه‌های مناسبی برای حل با استفاده از محاسبات دی.ان.ای می‌باشند، زیرا با استفاده از قدرت پردازش موازی مولکول‌های دی.ان.ای امکان بررسی تمام جواب‌های ممکن و غیرممکن برای مساله در زمان قابل قبول امکان‌پذیر است. مساله مسیر همیلتونی^۱ [۱]، فروشنده دوره‌گرد^۲ [۲,۳]، رنگ‌آمیزی گراف^۳ [۴,۵]، مسئله گروه‌گ^۴ [۶]، مجموعه مستقل^۵ [۷]، افزار مجموعه^۶ [۸]، مساله کوله‌پشتی^۷ [۹]، مساله ارضا قیده‌های منطقی^۸ [۱۰] و مساله هم‌ریختی گراف^۹ [۱۱] از جمله مسائل ان.پی می‌باشند که راه‌حل آن‌ها با استفاده از محاسبات دی.ان.ای ارائه شده است.

۱-۴ مروری بر کارهای پیشین

در بحث دسته‌بندی داده‌ها هدف ساخت مدل مناسبی است که بتوان از آن برای دسته‌بندی داده‌های جدید استفاده کرد. کوینلان [۱۲] در سال ۱۹۸۶ الگوریتم ID3 را برای ساخت مدل دسته‌بندی بر اساس درخت‌های تصمیم ارائه داد. در درخت ساخته‌شده هر مسیر از ریشه به برگ معرف یک قانون دسته‌بندی است. برای ساخت درخت تصمیم، از جستجوی حریصانه بالا به پایین در فضای درخت‌های موجود استفاده می‌شود. بدین صورت که در هر مرحله آن ویژگی که به بهترین شکل داده‌ها را از هم جدا می‌کند به عنوان گره بعدی درخت انتخاب می‌شود. معیار انتخاب بهترین ویژگی نیز بر اساس مقداری آماری به نام بهره‌اطلاعات صورت می‌پذیرد. الگوریتم‌های C4.5 و C5.0 نیز نسخه‌های تکامل‌یافته الگوریتم ID3 می‌باشند. از مزیت‌های درخت‌های تصمیم می‌توان به عدم حساسیت به داده‌های نویزی و کارآمد بودن آن برای حجم زیاد داده‌ها اشاره کرد اما از آنجاییکه ساخت درخت به روش حریصانه انجام می‌شود، بنابراین تضمینی برای یافتن درخت بهینه وجود ندارد.

در سال ۱۹۹۶ لو و همکاران [۱۳] روشی برای استخراج قوانین دسته‌بندی با استفاده از شبکه عصبی مصنوعی^{۱۰} پیشنهاد دادند. در روش پیشنهادی آن‌ها، ابتدا یک شبکه عصبی سه لایه ساخته شده و با

¹ hamiltonian path problem

² traveling salesman problem

³ vertex coloring

⁴ clique problem

⁵ independent set problem

⁶ set splitting problem

⁷ knapsack problem

⁸ satisfiability problem

⁹ graph isomorphism

¹⁰ artificial neural network

استفاده از داده‌های آموزشی، آموزش داده می‌شود تا به دقت مورد نظر برسد. سپس با استفاده از یک الگوریتم هرس‌بندی، اتصالات و گره‌های اضافی به گونه‌ای هرس می‌شوند که دقت دسته‌بندی کاهش نیابد که منجر به تولید شبکه عصبی ساده تری می‌شود. در مرحله بعد قوانین دسته‌بندی از این شبکه استخراج می‌شوند. یکی از معایب این روش زمان نسبتاً زیادی است که برای آموزش شبکه عصبی صرف می‌شود. نویسندگان مقاله برای رفع این مشکل دو روش پیشنهاد می‌دهند که یکی استفاده از آموزش تدریجی و دیگری کاهش تعداد ورودی‌های شبکه عصبی با استفاده از الگوریتم‌های انتخاب ویژگی^۱ است.

در سال ۲۰۰۲ پارینلی و همکاران [۱۴] روشی بنام آنت‌ماینر برای استخراج قوانین دسته‌بندی با استفاده از الگوریتم اجتماع مورچه ارائه دادند. در روش پیشنهادی آن‌ها از معیار بی‌نظمی^۲ به عنوان معیار شهودتی^۳ در تصمیم‌گیری مورچه‌ها برای انتخاب میان صفات استفاده می‌شود. هر مسیری که توسط یک مورچه طی شود منجر به تولید یک قانون دسته‌بندی می‌شود. پس از ساخت قانون توسط مورچه، آن قانون هرس شده که باعث افزایش دقت، سادگی و قابل فهم‌تر شدن قانون می‌شود. لیو و همکاران [۱۵] با ارائه معیار شهودتی ساده‌تری که بر پایه تخمین چگالی کار می‌کرد، نسخه دوم آنت‌ماینر را ارائه کردند. ایده اصلی آن‌ها در استفاده از معیار ساده‌تر این بود که لازم نیست معیار شهودتی خیلی دقیق باشد. نسخه دوم آنت‌ماینر نسبت به اولین نسخه پیچیدگی محاسباتی کم‌تر و دقت دسته‌بندی تقریباً یکسانی را ارائه می‌داد. در سال ۲۰۰۴ سومین نسخه آنت‌ماینر توسط لیو و همکاران [۱۶] ارائه شد که در آن استراتژی بروزرسانی فرمون‌ها و قانون انتقال حالت بهبود یافته بودند که باعث افزایش دقت دسته‌بندی نسبت به نسخه‌های پیشین شد.

در سال ۲۰۰۳ روشی برای استخراج قوانین دسته‌بندی با استفاده از برنامه‌نویسی عبارات ژنی^۴ توسط ژو و همکاران [۱۷] ارائه شد. برنامه‌نویسی عبارات ژنی که ترکیب الگوریتم ژنتیک^۵ و برنامه‌نویسی ژنتیک^۶ است، ارتباط بین صفات را به صورت روابط ریاضی استخراج می‌کند. در روش پیشنهادی آن‌ها برای حل مساله دسته‌بندی با n دسته از تکنیک یادگیری یکی در مقابل همه^۷ استفاده می‌شود که مساله را به n مساله دسته‌بندی باینری تبدیل کرده که با استفاده از برنامه‌نویسی عبارات ژنی قابل حل

¹ feature selection

² entropy

³ heuristic

⁴ gene expression programming

⁵ genetic algorithm

⁶ genetic programming

⁷ one-against-all learning

می‌باشند. همچنین تابع ارزیابی بگونه‌ای تعریف می‌شود که برای هر قانون، هم کامل بودن و هم سازگار بودن در نظر گرفته می‌شوند. از مزایای روش ارائه شده عدم حساسیت آن به داده‌های نویزی به دلیل استفاده از اصل حداقل طول توصیفی^۱ و هرس کردن مجموعه قوانین استخراج شده است. هرس کردن قوانین در دو مرحله انجام می‌شود: یکی در حین یادگیری قوانین که از ایجاد قوانین پیچیده (که معمولاً بدلیل داده‌های نویزی ایجاد می‌شوند) جلوگیری کرده و دیگری پس از فرآیند یادگیری است که برای حذف تضاد بین قوانین لازم است. مشکل عمده روش‌هایی که از برنامه‌نویسی عبارات ژنی استفاده می‌کنند این است که در برخی موارد قوانین تولید شده بسیار پیچیده بوده و فهم آن‌ها برای ما دشوار است.

چیو [۱۸] در سال ۲۰۰۵ از الگوریتم ژنتیک مبتنی بر قیود برای استخراج قوانین دسته‌بندی استفاده کرد. از آنجا که هر داده مجموعه‌ای از صفات است، در روش پیشنهادی می‌توان قیود را بصورت رابطه بین این صفات بیان کرد. همچنین برای تولید کروموزوم‌هایی که همه قیود را ارضا کنند از استدلال مبتنی بر قیود استفاده شده است. از آنجا که با افزایش تعداد صفات، تعداد قوانین دسته‌بندی ممکن بصورت نمایی افزایش می‌یابد، استفاده از الگوریتم ژنتیک برای جستجوی جهت‌یافته این فضای بزرگ مناسب به نظر می‌رسد.

از دیگر مطالعات انجام شده در زمینه داده کاوی می‌توان به استخراج قوانین همبستگی اشاره کرد. یکی از کاربردهای مهم این قوانین استفاده از آن‌ها در سیستم‌های توصیه گر است. فرآیند استخراج قوانین همبستگی از دو مرحله تشکیل می‌شود: [۱۹] در مرحله اول سبد کالاهای مکرر-سبد کالاهایی که تعداد تکرارشان در کل پایگاه داده از مقدار مشخصی بیشتر است- شناسایی شده و سپس در مرحله دوم قوانین همبستگی با استفاده از این سبد کالاها تولید می‌شوند. از آنجا که مرحله دوم ساده و سراسر است، بیشتر تحقیقات انجام شده در زمینه استخراج قوانین همبستگی در رابطه با مرحله اول یعنی شناسایی سبد کالاهای مکرر بوده است.

آگراوال و همکاران [۱۹] در سال ۱۹۹۳ اولین الگوریتم را برای استخراج قوانین همبستگی بنام الگوریتم ایز^۲ ارائه دادند. در این الگوریتم برای شناسایی سبد کالاهای مکرر^۳، کل پایگاه داده

¹ minimum description length

² AIS (Agrawal, Imielinski, Swami)

³ frequent itemset

چندین بار پیمایش می‌شود که یکی از معایب این الگوریتم هم به شمار می‌رود. در اولین پیمایش پایگاه داده، کل سبد کالاهای مکرر تک‌عضوی شناسایی شده و با اضافه کردن کالاهای باقیمانده به هر کدام از این سبدها، سبد کالاهای دو عضوی منتخب ساخته می‌شوند که با پیمایش مجدد پایگاه داده، سبد کالاهای مکرر دو عضوی از میان کل سبد کالاهای منتخب شناسایی می‌شوند. مشکل اصلی الگوریتم ایز این است که بسیاری از سبد کالاهایی که امکان تبدیل شدن به سبد کالاهای مکرر را ندارند بیهوده تولید شده و زمان و فضای زیادی برای پردازش آنها تلف می‌شود؛ هر چند در این الگوریتم روش‌هایی برای هرس کردن سبد کالاهای منتخب و همچنین مدیریت حافظه پیشنهاد شده است.

در سال ۱۹۹۴ آگراوال و سریکانث [۲۰] الگوریتم ایپریوری^۱ را ارائه دادند تا مشکلات موجود در الگوریتم ایز را برطرف سازند. پیمایش بسیار زیاد پایگاه داده و تولید سبد کالاهای منتخب غیرضروری از مشکلات الگوریتم ایز می‌باشند که در الگوریتم ایپریوری به مقدار قابل توجهی کاهش یافته‌اند. در این الگوریتم تکنیک تولید سبد کالاهای منتخب و هرس کردن بهبود داده شده‌اند که همین امر باعث افزایش کارایی آن در مقایسه با الگوریتم ایز شده است. روش تولید سبد کالاهای منتخب بدین صورت است که ابتدا با یکبار پیمایش پایگاه داده، مجموعه تمام سبد کالاهای مکرر تک‌عضوی شناسایی شده و از طریق پیوند زدن مجموعه حاصل با خودش، سبد کالاهای منتخب دو عضوی ایجاد می‌شوند. در هر مرحله مجموعه‌ی سبد کالاهای منتخب k عضوی با خودش پیوند زده شده و مجموعه‌ی سبد کالاهای منتخب $k+1$ عضوی ساخته می‌شود. همچنین با استفاده از این ویژگی که تمام زیرمجموعه‌های k عضوی از یک سبد کالای مکرر $k+1$ عضوی، خودشان سبد کالای مکرر می‌باشند، تعداد سبد کالاهای منتخب k عضوی کاهش می‌یابد. بدین ترتیب تعداد پیمایش‌های پایگاه داده به نسبت زیادی کاهش می‌یابد. با این وجود هنوز هم پیمایش کل پایگاه داده به تعداد زیاد و تکنیک پیچیده تولید مجموعه سبد کالاهای منتخب از مشکلات اصلی الگوریتم ایپریوری است. انتخاب بخشی از پایگاه داده بجای کل آن برای پیمایش، نمونه‌برداری از داده‌ها، استفاده از تکنیک‌های هرس کردن متفاوت مجموعه سبد کالاهای منتخب از جمله اصلاحاتی است که در مطالعات مختلف برای افزایش کارایی الگوریتم ایپریوری پیشنهاد شده است. از جمله این الگوریتم‌ها می‌توان به [۲۰، ۲۱، ۲۲] اشاره کرد.

¹ apriori