



کاربرد داده کاوی بر روی داده های آموزش عالی دانشجویان دانشکده دندانپزشکی شهر رشت با استفاده از تکنیکهای طبقه بندی و خوش بندی

علی زواره^۱، امیر کوچکی^۲ و علی رهنما رودپشتی^۳

^۱ کارشناس ارشد مدیریت فناوری اطلاعات، دانشکده دندانپزشکی رشت ، a_zavareh@gums.ac.ir

^۲ مدرس بخش فناوری اطلاعات، سازمان فنی و حرفه ای شهر رشت، SecondAuthor@Email

^۳ مدرس، دانشگاه علمی کاربردی رشت، ThirdAuthor@Email

چکیده - هدف بررسی اطلاعات آموزشی دانش آموختگان رشته دندانپزشکی بین سالهای ۱۳۷۵ تا ۱۳۹۰ و بدست آوردن مدلی جهت پیش بینی دانشجویان ممتاز(قوی) و دانشجویان تحت ریسک مشروطی (ضعیف) و عوامل مؤثر بر نمرات و رفتار آموزشی آنها است.

داده کاوی داده ها با استفاده از یک فرآیند چند مرحله ای بنام *CRIPS* صورت گرفت در این تحقیق از مدل *CHAID* برای پیش بینی از طریق طبقه بندی و مدل های *K-Means* و *Two-Step* برای خوش بندی استفاده گردید.

نتایج این مطالعه نشان داد فیلدهای جنس، سن، سهمیه، معدل کل، تعداد ترم مشروطی و فیلد *pass* در خوش بندی بعنوان فیلدهای با اهمیت بدست آمدند. مدل *Two-Step* بهتر از *K-Means* شناخته شد و چهار خوش شامل: خوشه ۱(بهداشتکاران)، خوشه ۲(دانشجویان متوسط با٪۱۰۰ جمعیت مرد)، خوشه ۳(دانشجویان ضعیف) و خوشه ۴(دانشجویان قوی) تولید نمود. هر دو مدل در رابطه با خوش بندی دانشجویان متوسط، ضعیف عمل نمودند. مدل *CHAID* در تکنیک طبقه بندی برای پیش بینی دانشجویان قوی(متاز)، دانشجویان متوسط و دانشجویان ضعیف قوانینی با درصد صحت بالا تولید نمود.

داده کاوی آموزشی نشان داد که از مدل *CHAID* در تکنیک طبقه بندی، جهت پیش بینی دانشجویان ضعیف و قوی و دانشجویانی که در خطر مشروطی هستند و از مدل *Two-Step* می توان در تکنیک خوش بندی جهت بررسی داده ها و رفتار دانشجویان، برای بهبود کیفیت و برنامه ریزی استراتژیک آموزشی استفاده نمود.

کلید واژه- تکنیک خوش بندی، تکنیک طبقه بندی، داده کاوی آموزشی، Data Mining

انسان علم تجزیه و تحلیل داده ها یا داده کاوی پا به عرصه گذاشت.

۱- مقدمه

داده کاوی^۱ (DM) فرآیندی است که در آغاز دهه ۹۰ مطرح شد و با نگرشی نو، به مسئله استخراج اطلاعات از پایگاه داده ها پرداخت. از سال ۱۹۹۵ داده کاوی به صورت جدی وارد مباحث آمار شد و در سال ۱۹۹۶، اولین شماره مجله کشف دانش و معرفت از پایگاه داده ها^۲ (KDD) منتشر شد. محققانی نظری

با رشد فناوری اطلاعات و روش های تولید و جمع آوری داده ها، داده های مربوط به تبادلات تجاری، کشاورزی، اینترنت، جزئیات مکالمات تلفنی، داده های پژوهشی و غیره سریعتر از هر روز جمع آوری و انبارش می شوند. لذا از اواخر دهه ۸۰ میلادی پسر به فکر دست یابی به اطلاعات نهفته در این پایگاه داده های حجمی افتاد، زیرا سیستمهای سنتی قادر به این کار نبودند. به دلیل رقابت در عرصه های سیاسی، نظامی، اقتصادی، علمی و اهمیت دست یابی به اطلاعات در کمترین زمان و بدون دخالت

¹ Data Mining

² Knowledge Discovery in Database



استفاده از داده کاوی در سیستم های آموزشی را با توجه به گرایش آن به سه گروه می توان تقسیم نمود:

۱- گرایش به سمت دانشجویان: هدف یادگیری فعال، بهبود یادگیری و پیشنهاد تجارب یادگیری خوب و پیشنهاد کوتاه کردن مسیر یادگیری به یادگیرنده و موفقیت آنها است.

۲- گرایش به سمت آموزش دهنگان: هدف بازخورد نتیجه آموزش، مدرسین. ارزیابی ساختار مطالب دوره و اثربخشی آن در روند آموزش، طبقه بندی یادگیرنده ها، کشف اطلاعات مفید که سبب برنامه ریزی در سازماندهی دروس میشود.

۳- گرایش به سمت مسئولان و مدیران دانشگاهی: هدف داشتن اقداماتی در مورد چگونگی سازماندهی بهتر منابع رسمی و سازمانی و پیشنهاد آموزشی می باشد.

هدف ما در این تحقیق گرایش نوع اول یعنی گرایش به سمت دانشجویان می باشد. کشف دانش از داده های ذخیره شده توسط داده کاوی، در واقع یک فرآیند پشت سرهم هست که دارای مراحل زیر می باشد:

۲-۱- مراحل داده کاوی

تعیین اطلاعات گذشته

پاک سازی داده ها و پردازش اولیه- در این مرحله خطاهای داده ها تصحیح می شوند و داده های اشتباہ جایگزین میشوند. این مرحله ممکن است تا ۶۰ درصد از زمان داده کاوی را در برگیرد.

یکپارچه سازی داده ها- معمولاً داده ها از منابع متفاوتی جمع آوری می شوند باید به صورتی در آیند که یک مخزن از داده های مناسب ایجاد شود تا بتوان عملیات داده کاوی را بهتر انجام داد.

انتخاب مجموعه داده های هدف.

یافتن ویژگیهای مورد استفاده و تعیین ویژگی های جدید نمایش داده ها به صورتی که بتوان برای داده کاوی استفاده نمود.

انتخاب عملیات داده کاوی (طبقه بندی، خوشه بندی، پیش بینی و غیره)

انتخاب روش داده کاوی (شبکه های عصبی، درخت تصمیم و نظایر آن).

داده کاوی و جستجو برای یافتن الگوی مناسب.
ارزیابی و تحلیل الگوی به دست آمده و حذف الگوهای نامناسب.

تفسیر نتایج داده ها و استنتاج از اطلاعات با ارزش.

براچمن^۳ و آناند^۴ کلیه مراحل واقع گرایانه و رو به جلو کشف دانش از پایگاه داده ها را تشخیص دادند.

در حال حاضر، داده کاوی مهمترین فناوری جهت بهره برداری موثر از داده های حجمی است و اهمیت آن رو به فزونی است. به طوریکه تخمین زده شده است که مقدار داده ها در جهان هر ۲۰ ماه به حدود دو برابر می رسد. در یک تحقیق که بر روی گروه های تجاری بسیار بزرگ در جمع آوری داده ها صورت گرفت مشخص گردید که ۱۹ درصد از این گروه ها دارای پایگاه داده هایی با سطح بیشتر از ۵۰ گیگابایت می باشند و ۵۹ درصد از آنها انتظار دارند که در آینده ای نزدیک در چنین سطحی قرار گیرند[۱]. در سالهای اخیر تحقیقات زیادی در زمینه بکارگیری فرآیند داده کاوی در امر آموزش صورت گرفته است. این زمینه تحقیقاتی جدید، داده کاوی آموزشی^۵ (EDM) نامیده می شود که به امر توسعه روشهای کشف دانش از داده های محیط های آموزشی خصوصاً دانشجویان می پردازد.

۱-۱- داده کاوی آموزشی

داده کاوی آموزشی به عنوان ناحیه پژوهش علمی ای تعریف شده است که در حول روش های توسعه برای ایجاد کشفیاتی که در نوع خود منحصر به فرد است، متمرکز گردیده است. استفاده از آن روشی برای درک بهتر دانشجویان و مجموعه ای که آنها در آنجا آموزش می بینند، می باشد.

به عنوان مثال، در داده کاوی بررسی دانشجویان در استفاده کردن از نرم افزارهای آموزشی، ممکن است به طور همزمان داده ها را در نظر گرفت و می تواند سطح زمان فشردن صفحه کلید، سطح پاسخ، سطح جلسه، سطح دانشجویان، سطح کلاس درس، و سطح مدرسه، داده های ارزشمندی باشند. بحث زمان، ترتیب و مفاد نیز نقش مهمی در مطالعه داده های آموزشی ایفا می کند.

داده کاوی آموزشی به عنوان یک زمینه تحقیقاتی مستقل در سال های اخیر پدید آمده است، اوج آن در سال ۲۰۰۸ با ایجاد کنفرانس بین المللی سالانه در مورد داده کاوی آموزشی، و مجله های آموزشی داده کاوی می باشد[۲].

³ Brachman

⁴ Anand

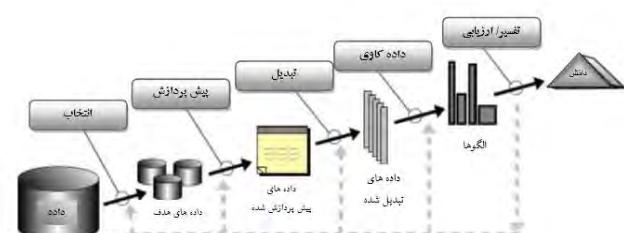
⁵ Education Data Mining

ونوس شکورنیاز، آرش حاجی علی اکبری (۱۳۸۷) دومین کنفرانس انجمن داده کاوی ایران) در مقاله‌ای با عنوان «خوشه بندی داده‌های آماری دانشجویان دانشگاه علم و صنعت و استخراج نمایه ساز توصیفی برای دانشجویان موفق» از اطلاعات آماری ۳۱۱ دانشجوی ورودی ۸۲ تا ۸۵ برای سالهای ۸۲-۸۵ استفاده نمودند و از الگوریتم‌های آماری و میانگین کا توسعه نرم افزار SQL با توجه به معیارهای خاص، دانشجویان را از لحاظ موفقیت، پیشرفت، آینده شغلی و ... خوشه بندی نمودند و بر روی تعداد ترم مشروطی تمرکز نمودند و ویژگی‌های دانشجویانیکه مشروط می‌شوند را مشخص کردند، برای مشخص نمودن ارتباطات میان مشخصه‌های ورودی تاثیرگذار به منظور تعیین قدرت و استقامت مشخصه‌ها در پیش‌بینی هدف، از تکنیک درخت تصمیم به منظور پیش‌بینی و تحلیل معدل کل دانشجویان و همچنین تعداد ترم‌های مشروطی آنها استفاده گردید. نتایج بدست آمده از تکنیک‌های درخت تصمیم و خوشه بندی عبارت است:

دو خوشه شماره ۱ و ۲ در بردارنده ممتازترین دانشجویان با بهترین معدل‌ها می‌باشند. معدل ترم اول دانشجویان موفق دانشگاه علم و صنعت، با تفاوت بسیار زیاد، بالاتر از سایر دانشجویان می‌باشد. این موارد به پیگیری و تلاش آنها از سالهای اولیه تحصیل اشاره دارد. این خوشه‌ها پایین‌ترین درصد دانشجویان مربوط به مهندسی سازه‌های ریلی را دارا می‌باشند و این مطلب احتمال دارد به سختگیری بیش از حد استادان این رشته مربوط گردد. با بررسی دانشجویان خوشه‌های ۵ و ۶ نیز که عنوان خوشه دانشجویان ضعیف نام برده می‌شوند، می‌توان به صحت این موضوع پی برد. میانگین تعداد واحدهای اخذ شده در هر ترم بین ۱۶ و ۱۷ واحد درسی می‌باشد. این می‌تواند بعنوان یک راهکار مناسب جهت راهنمایی دانشجویان بکار گرفته شود. شاید بهترین نحوه انتخاب واحد، داشتن تعداد واحد واحد در هر ترم را پیشنهاد نماید. این دو خوشه بیشترین تعداد دختران را نسبت به سایر خوشه‌ها در خود جای داده اند و این مطلب، علاقمندی بیشتر دانشجویان دختر به مطالعه نتایج مدل پسران را نشان میدهد.^[۵]

ابراهیم صحافی زاده (۱۳۸۸)، در سومین کنفرانس انجمن داده کاوی ایران) در مقاله خود با عنوان «تحلیل عوامل مؤثر بر نمرات دانشجویان دانشگاه پیام نور بوشهر با استفاده از داده کاوی» بر روی بیش از ۲۵۰ هزار نمره دانشجویان دانشگاه پیام نور بوشهر در بین سالهای ۸۱ تا ۸۷ با استفاده از خوشه بندی به

در شکل ۱ میتوان مراحل داده کاوی را به اختصار نشان داد.^[۳]



شکل ۱- مراحل داده کاوی

Jan M. Mykola Pechenizkiy، Gerben W. Dekker و Vleeshouwers (سال ۲۰۰۹) در دومین کنفرانس داده کاوی آموزشی در کشور اسپانیا مقاله‌ای با عنوان «پیش‌بینی دانشجویانی که ترک تحصیل می‌نمایند: مطالعه موردی» را ارائه نمودند. آنها بر روی دانشجویان مهندسی برق این مطالعه را انجام دادند و بر روی کسانی که پس از ترم اول ترک تحصیل می‌نمایند.

در این مطالعه اطلاعات طول سالهای ۲۰۰۰-۲۰۰۹ که شامل اطلاعات مربوط به تمام دانشجویان که در برنامه رشته مهندسی برق درگیر جمع آوری شده بودند. در انتخاب مجموعه داده‌های هدف از اطلاعات ۶۴۸ دانشجو استفاده نمودند.

در مطالعه، از تکنیک‌های اثربخش طبقه بندی توسعه نرم افزار WEKA استفاده شد. الگوریتم درخت تصمیم به نظر می‌رسد کمی بهتر از الگوریتم‌های دیگر: الگوریتم‌های CART و C4.5 تنها مواردی بودند که قادر به بهبود مدل خط پایه بصورت قابل توجهی بودند.

هنگامی که به مقایسه بخش‌های مختلف مجموعه داده پرداخته شد، مشخص شد که پیش‌بینی نمرات دانشگاه نسبت به داده‌های پیش‌دانشگاهی بهتر است. به نظر می‌رسد ویژگی‌های داده‌های پیش‌دانشگاهی یا خیلی بی‌ربط و یا به شدت در ارتباط است، و در نتیجه دقت پیش‌بینی با استفاده از میانگین نمرات دروس علوم حدود ۷۰ درصد است. مقایسه نتایج مدل‌های مختلف نشان میدهد که پیش‌بینی نمره درس جبر خطی بدنبال پیش‌بینی نمره درس حساب دیفرانسیل و انتگرال امکان پذیر است. نتیجه دیگری که گرفته شد مشخص گردید اگر در مدل مورد بررسی مجموعه داده‌ها را بطور جداگانه بررسی کنیم مفید نمی‌باشد بلکه بهتر است مجموعه داده‌ها را بصورت ترکیبی بکار برد.^[۴]



میانگین نمرات دانشجویان، از ورودیهای ۷۷ به بعد رشد صعودی داشته برای دانشجویان ورودی ۸۴ به اوج رسیده اما برای دانشجویان ورودی بعد از ۸۴ به صورت نزولی کاهش می یابد. هر چند میانگین سنی دانشجویان دانشگاه پیام نور به شدت کاهش یافته است اما سن دانشجویان تاثیر معناداری در نمره ندارد. جنسیت تاثیری در نتایج ندارد و نمرات هر دو گروه جنسی در طبقات مشابه قرار دارند. دروس عملی بالاترین نمره را در بین دروس دارند. با توجه به اینکه نمره عملی دست استاد می باشد، تاثیر وجود استاد و برگزاری امتحان به صورت غیر متمرکز کاملا مشخص است. هر چند استاد تصحیح اوراق تشريحی را بر عهده دارند اما تفاوتی در نمرات تمام تستی و تستی- تشریحی مشاهده شود که این نشان دهنده تصحیح دقیق بخش تشریحی نمی شود. توسط استاد و مطابق با کلید سازمان مرکزی باشد. نمرات امتحانات تشريحی که بدون منبع هستند و آزمون آنها توسط استاد طرح و تصحیح می شود جزو بالاترین نمرات هستند و نهایتاً اینکه میانگین نمرات نیمسال بهار بهتر از نیمسال پاییز است [۶].

۲- روش تحقیق

ساختار اجرایی این تحقیق بر اساس دو رویکرد دانش محور^۶ و داده محور^۷ بنا گردیده است که در ادامه به آن اشاره گردیده است. ارزیابی قوانین حاصله از مدلسازی داده کاوی (تکنیک خوشه بندی و طبقه بندی) ماهیت دانش محور دارد، زیرا براساس دانش خبرگان دانشگاهی شامل مدیران، استادی و کارشناسان آموزشی صورت می پذیرد و با استفاده از تایید صحت قوانین اخذ شده از طرف ایشان ارزیابی مدلهای تولید شده انجام می پذیرد.

ماهیت اصلی تحقیق داده محور بوده و پایه اصلی تحقیق حاضر بر کشف دانش پایگاه داده های بانک مورد مطالعه نهاده شده است.

از فرآیند^۸ CRIPS و از نرم افزار 12 SPSS Clementine به منظور داده کاوی استفاده شده است.

روش k-mean، نمرات خوشه بندی شده و تاثیر عواملی از قبیل نوع آزمون، نوع ورود به آموزش عالی، نوع رائیه درس (خود خوان یا حضوری)، واحدهای درس (عملی - نظری)، تعداد جلسات ارائه، دوره، نیمسال ارائه درس، جنسیت و سن بر روی هر خوشه بررسی شده و قوانین تداعی موجود بین این فیلدها استخراج شده است.

در تحلیل بر اساس معدل تعداد گروهی از دانشجویان ورودی بین سالهای ۸۱ تا ۸۶ که معدل آنها خوشه بندی شد ۸۵۰۱ دانشجو بود. این دانشجویان براساس معدل کل آنها تا نیمسال اول سال تحصیلی ۸۷-۸۸ در ۵ خوشه طبقه بندی شدند. خوشه ۱ بیشترین فراوانی را دارد و خوشه ۳ که بالاترین معدل ها را دارد دارای کمترین فراوانی است و فقط حدود ۸ درصد از دانشجویان را شامل می شود. با توجه به اینکه خوشه ۴ بعد از خوشه ۱ بیشترین فراوانی را دارد می توان گفت که نزدیک به ۶۰ درصد از دانشجویان دانشگاه پیام نور استان بوشهر به لحاظ معدل در خوشه های ۱ و ۴ قرار می گیرند. بنابراین معدل کل ۶۰ درصد دانشجویان بین ۱۲ تا ۱۴ می باشد. با استفاده از ابزارهای تولید قانون در نرم افزار Clementine اولین شاخص اثر گذار بر خوشه بندی معدل کل، پارامتر نحوه ورود به آموزش عالی بود. دانشجویانی که با سهمیه خانواده شهدا وارد دانشگاه می شوند از سال ۸۵ به بعد دچار افت تحصیلی شده و در خوشه شماره ۲ قرار گرفته اند که این خوشه دارای معدل پایینی می باشد. آموزگاران زیر ۳۰ سال معدل بهتری نسبت به آموزگاران بالای ۳۰ سال دارند. همچنین در تحلیلی که بر روی دانشجویانیکه نوع ورود به آموزش عالی آنها عادی بود انجام گرفت مشخص شد که دانشجویان دوره فراغیر از نظر معدل وضعیت بهتری نسبت به دانشجویان دوره رسمی دارند. دانشجویان آزمون محور رسمی و فراغیر و الکترونیکی نیز در خوشه دو قرار گرفتند این خوشه وضعیت ضعیفی دارد. بالاترین میانگین معدل متعلق به دانشجویان ورودی سال ۸۴ و پایین ترین میانگین معدل متعلق به دانشجویان ورودی ۸۱ می باشد.

در تحلیل نمرات دانشجویان دانشگاه پیام نور بوشهر مشخص شد که حدود ۶۰ درصد از نمرات دانشجویان بین ۱۰ تا ۱۶ می باشد. وضعیت نمرات دانشجویان دوره فراغیر بهتر از دانشجویان رسمی است. رائیه درس به صورت خودخوان و حضوری در نتیجه آزمونها تاثیر متفاوتی ندارند. میانگین نمرات دانشجویان تا ۵ ترم اول به صورت صعودی بوده و از ترم ۵ به بعد روند کاهش نمرات دانشجویان آغاز و تا آخرین ترمها به صورت نزولی افت می کند.

⁶ Knowledge Oriented

⁷ Data Oriented

⁸ CRoss-Industry Standard Process for Data Mining

های مشروطی بیشتری هستند؟»، «چرا میزان مردودی دانشجویان در بعضی از دورس بیشتر از دروس دیگر است؟»، «چرا دانشجویان ورودی یک سال خاص نسبت به سالهای ورودی دیگر قویتر می‌باشند»، از اهداف استخراج داده‌ها، مجزا نمودن داده‌هایی که محقق را در شناسایی عوامل مشروطی و مردودی افراد در ترم و درس و همچنین رفتار آموزشی دانشجویان کمک نماید. معیار موفقیت پرورزه داده کاوی کشف قوانینی که به اهداف پرورزه نزدیک باشد و توسط خبرگان آموزش دانشکده مورد تایید باشد.

۳-۲- مرحله درک داده‌ها

درک داده‌ها به بررسی داده‌های مورد نیاز می‌پردازد. این مرحله می‌تواند شامل جمع آوری داده‌های اولیه، توصیف داده‌ها، اکتشاف داده‌ها، و تایید کیفیت داده‌ها باشد. اکتشاف داده‌ها از قبیل مشاهده اطلاعات آماری خلاصه (که شامل نمایش تصویری از طبقه بندي متغیرها است) می‌تواند در پایان این مرحله رخ دهد.

جامعه آماری تحقیق شامل ۵۷۷ دانشجو که از این تعداد ۳۲۸ نفر دانشجو شاغل به تحصیل (فعال)، و ۲۳۹ دانشجو از نظر تحصیلی، فارغ التحصیل، اخراج یا جابجایی به دانشگاه‌های دیگر داشته‌اند (غیرفعال). منتهی‌پس از بررسی نهایی و حذف رکوردهایی که انتقال، اخراج یا مهمان از دانشگاه‌های دیگر بودند تعداد رکوردها به ۵۱۹ عدد رسید. البته جهت فرآیند داده کاوی با توجه به طرح در نظر گرفته برای الگوهای پیش‌بینی، تعداد ۳۵۲ رکورد که بیش از ۱۰ ترم تحصیلی را گذرانده بودند به عنوان یک منبع جداگانه برای تکنیکهای پیش‌بینی طبقه بندي و خوش‌بندی تفکیک گردیدند.

فیلدهای دانشجویان شامل دو قسمت اطلاعات شخصی و اطلاعات آموزشی بودند. این فیلدها ۲۱ عدد و شامل:

شماره دانشجویی، سال شروع به تحصیل، جنسیت، سهیمه ثبت نامی، تاریخ تولد، معدل کل، تعداد واحدهای گذرانده، تعداد واحدهای مردود شده، تعداد ترم مشروطی، معدل ترم اول، تعداد واحدهای پاس شده در ترم اول، تعداد واحدهای پاس نشده در ترم اول، معدل ترم دوم، تعداد واحدهای پاس شده در ترم دوم، تعداد واحدهای پاس نشده در ترم دوم، معدل ترم سوم، تعداد واحدهای پاس شده در ترم سوم، تعداد واحدهای پاس شده در ترم چهارم، تعداد واحدهای پاس شده در ترم چهارم، تعداد واحدهای پاس نشده در چهارم بودند.

۱-۲- فرآیند CRIPS-DM

محصول کنسرسیومی متشكل از سه شرکت (سال ۱۹۹۶) دایملر-کرایسلر، SPSS و NCR است تاکید بر چرخشی بودن فرآیند دارد.

شامل فازهای اصلی:

درک کسب و کار

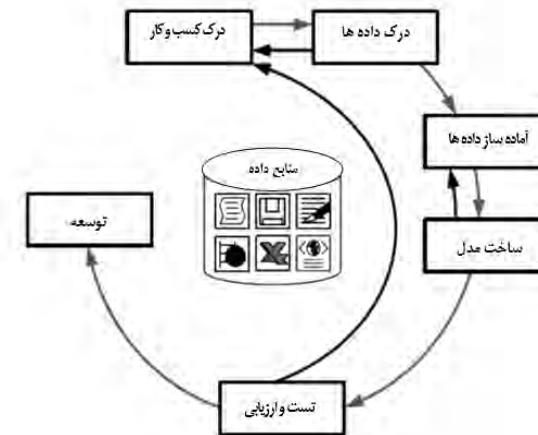
شناسایی و درک داده‌ها

آماده سازی داده‌ها

مدلسازی

تست و ارزیابی

توسعه



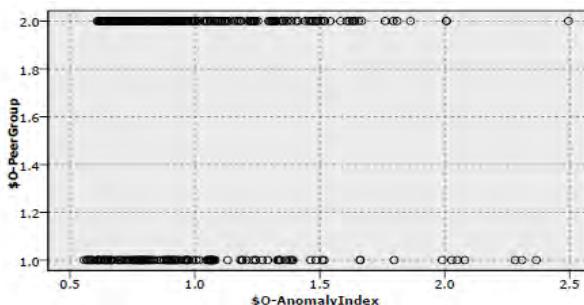
شکل ۲- فرآیند CRIPS-DM

۲-۲- مرحله درک کسب و کار

در این مرحله شناسایی پس زمینه فعالیت، اهداف و معیارهای موفقیت فعالیت کسب و کار انجام گرفت. دانشکده دندانپزشکی گیلان در سال ۱۳۷۵ اولین گروه دانشجویان خود را از طریق کنکور سراسری پذیرش نمود. از ابتدای شروع بکار دارای مشکلات متعددی بوده است. یکی از مهمترین مشکلات آن نداشتن فضای فیزیکی مناسب است. هدف کسب و کار در دانشکده دندانپزشکی گیلان عبارت است از « تربیت دانش پژوهی دکتری دندانپزشکی است که پس از گذراندن ۶ سال تحصیلی بعنوان فردی متعدد و مجرب بتواند در جامعه دندانپزشکی کشور قادر به خدمت رسانی به ملت عزیز ایران باشد».

در این مرحله اهداف استخراج داده‌ها نیز مشخص گردید، مشکلاتی که در حین مطالعه سبب مشخص شدن اهداف تحقیق گردید عبارت بودند «چرا عده‌ای از دانشجویان دارای تعداد ترم

دانشجویان که در دوره فوق دپیلم خود آنها را پاس نموده اند، می باشد و اشتباه در ورود اطلاعات داشتیم که بر طرف گردید. Quality در خصوص مقادیر پرت و گمشده به سربرگ مراجعه گردید در فیلدهای تعداد ترم مشروطی و تعداد واحد مردودی چون دارای مقادیر خالی بودند، مقادیر گمشده مشاهده گردید که مشکل با دادن مقادیر صفر در داخل مقادیر خالی SPSS برطرف شد. با استفاده از نود Anomaly نرم افزار Clementine رکوردهای پرت شناسایی شدند همانطور که در شکل ۳ مشاهده می شود رکوردهایی که مقدار AnomalyIndex آنها از ۰/۲ بیشتر بود عنوان رکورد پرت شناسایی شدند و مجرما گردیدند پس از بررسی این رکوردها، این نتیجه حاصل شد که با توجه به اینکه از لحاظ آماری نقاط پرت مشاهده شدند ولی از لحاظ آموزشی این مقادیر پرت منطقی و توجیه پذیر بودند و با توجه به اینکه شناسایی رفتار آنها جزء اهداف این تحقیق می باشد بنابراین آنها حذف نگردیدند.

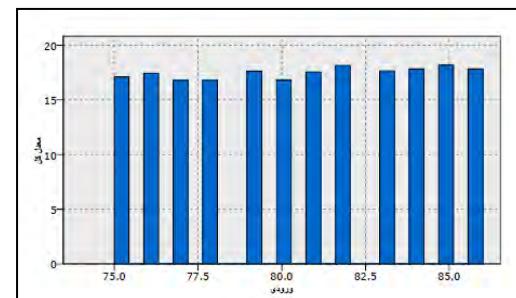


شکل ۳- نمایش نقاط پرت با استفاده از نود Anomaly

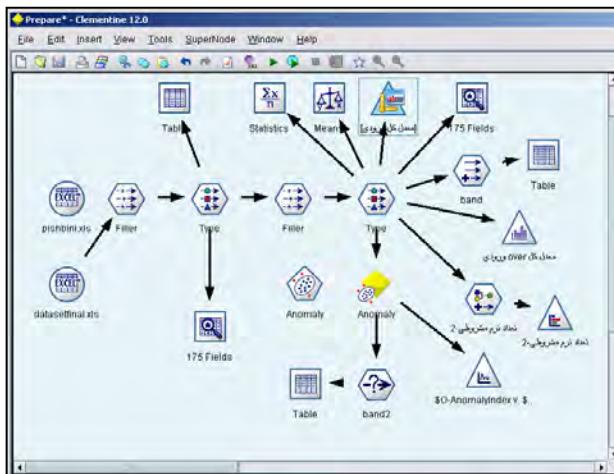
برای بی بردن رابطه بین معدل کل با معدل ترم ۱، ورودی، سن، تعداد واحد قبولی، تعداد ترم مشروطی و از طریق ضریب پیرسون متوجه شدیم که بجز ورودی که دارای رابطه ضعیفی است بقیه دارای رابطه قوی با معدل کل می باشند.

برای بی بردن رابطه بین تعداد ترم مشروطی با معدل کل، معدل ترم ۱، ورودی، سن، تعداد واحد قبولی از طریق ضریب پیرسون متوجه شدیم که تعداد واحد قبولی که دارای رابطه ضعیفی است بقیه دارای رابطه قوی می باشند.

همانطور که در ابتدای این بخش گفته شد، اکتشاف داده ها از طریق مشاهده اطلاعات آماری خلاصه (که شامل نمایش تصویری از طبقه بندی متغیرها است) نیز صورت می گیرد. این اکتشاف اولیه جالب است که قبل از عمل مدل سازی انجام می پذیرد. عنوان مثال نمودار ۱ نشان میدهد اختلاف فاحشی بین کسانی که دارای بالاترین معدل کل در هر ورودی بودند مشاهده نمی شود و ورودیهای ۸۲ و ۸۵ تقریبا نسبت به سایر ورودیها بالاترین معدل کل را دارند.



نمودار ۱- نمودار میله ای ارتباط بین ورودیها و حداقل معدل کل هر ورودی



شکل ۴- نمایی از کارهای انجام گرفته برای درک داده ها در نرم افزار Clementine

در این مرحله همچنین کیفیت داده ها مورد بررسی قرار گرفت که شامل، بررسی عدم انطباق داده ها، مقادیر پرت و نقاط گمشده بود. برای اینکار از نود Data Audit نرم افزار Clementine استفاده گردید. ابتدا از سربرگ Audit کیفیت داده ها از لحاظ صحت مقادیر کمینه و بیشینه و همچنین مراجعه به ساختار نمودار آمده شده آنها بررسی شد و در صورت مشاهده مورد اختلاف رکورد مربوطه را یافته و اختلاف را بررسی و در صورت لزوم اصلاح گردید. به عنوان مثال در بررسی تعداد واحدهای پاس شده ترم دوم مواردی عدم انطباق مشاهده گردید بطوریکه رکوردهایی دارای مقادیر بالای ۲۰ واحد پاس شده داشتند، رکوردها مربوط به دانشجویان بهداشتکار می باشد و مقادیر بالای ۲۰ واحد مربوط به واحدهای تطبیقی این



هایی از داده‌های مشابه می‌باشیم بدون اینکه از قبل، پیش‌بینی در مورد شباهت‌های موجود داشته باشیم. جهت انجام تکنیک طبقه‌بندی از الگوریتم‌های معروف CRT، CHAID و C5 استفاده گردید. جهت انجام تکنیک خوش‌بندی از الگوریتم‌های معروف TwoStep و K-Means استفاده گردید.

۶-۲- تست و ارزیابی

پس از مدل سازی مدل‌ها مورد تست و ارزیابی قرار گرفت و چندین بار اصلاح گردید تا مدل مورد قبول با اهداف ما بدست آمد.

۷-۲- توسعه

این مرحله بعد از تولید مدل و تایید دانش بدبست آمده، شروع می‌شود. نیاز به تغییر مدل‌های تولید شده در شرایط عملیاتی باید بررسی شود، زیرا ممکن است قوانین تولید شده توسط مدل‌ها در سالهای بعد درست نباشند. اگر تغییرات پیش‌آمده در محیط کسب و کار قابل توجه باشند مدل‌ها باید دوباره بازسازی شوند.

۳- نتایج

۱-۳- نتایج خوش‌بندی

کشف دانش موجود در اطلاعات ورودی و مشخص کردن ورودی‌های با اهمیت از نتایجی است که از طریق استفاده از این تکنیک بدبست می‌آید. ورودی‌های با اهمیت الگوریتم‌های Two-Step و K-Means یکسان بوده و عبارت بودند از فیلد‌های pass و Y-Variable، سن، سهمیه و معدل کل. این امر تعداد ترم مشروطی، جنس، سن، سهمیه و معدل کل. این امر اشاره به آن دارد متغیرهایی که ما جهت فرآیند داده کاوی در نظر گرفته ایم با متغیرهای تشخیص داده شده مدل داده کاوی مطابقت داشته و انتخاب صحیحی نموده ایم. نتایج بدست آمده در هر دو تکنیک در جدول ۱ نشان داده شده است.

۴-۲- آماده سازی داده‌ها

هنگامی که منابع داده‌ای در دسترس شناسایی شدند، نیاز به آن دارند که انتخاب شوند، پاکسازی شوند، و به شکل مورد نظر ساخته، و ساختار بندی شوند. پاکسازی داده‌ها و تبدیل داده‌ها در هنگام آماده سازی مدل داده‌ها، باید در این مرحله رخ دهد. هدف از آماده سازی داده‌ها تمیز کردن داده‌ها انتخاب شده برای کیفیت بهتر است.

فیلد سال ورود به تحصیل که حاوی اطلاعات سال، ماه و روز ورود به تحصیل دانشجو بود به فیلد ورودی تبدیل شد و فقط سال ورود به این فیلد وارد گردید. از فیلد تاریخ تولد دانشجو سال تولد دانشجو گرفته شد. همچنین در فیلد‌های تعداد واحد‌های پاس شده و پاس نشده ترم‌های ۱ تا ۴، بجای اینکه در نتیجه گیری تعداد واحد‌های پاس شده و پاس نشده چهار ترم اول در قانون‌های تولید شده در تصمیم گیری دخیل باشد با استفاده از فیلد‌های جدید Pass1 تا Pass4 نتیجه گیری بهتری خواهیم داشت. که بعنوان مثال فرمول (۱) زیر روش بدست آمده فیلد Pass1 را نشان می‌دهد که به همین صورت برای سه فیلد دیگر محاسبه می‌شود.

(۱)

$$\text{Pass1} = \frac{\text{تعداد واحد‌های پاس شده ۱}}{\text{تعداد واحد‌های پاس شده ۱} + \text{تعداد واحد‌های پاس نشده ۱}}$$

فرمت تمام داده‌های کمی بجز معدل از حالت اعشاری به عدد صحیح تبدیل شد.

۵-۲- مدل سازی

پس از آماده سازی داده‌ها نوبت به تشخیص این است که برای رسیدن به اهداف داده کاوی خود از چه تکنیک‌ها و الگوریتم‌هایی باید استفاده کنیم. تکنیک‌های داده کاوی بر اساس اهداف داده کاوی به دو گروه تقسیم می‌شوند (الف) توضیح گذشته (ب) پیش‌بینی آینده

با توجه به هدف، در این تحقیق ما از تکنیک‌های پیش‌بینی آینده استفاده کردیم، تکنیک‌های پیش‌بینی آینده که به چهار دسته (۱) طبقه‌بندی (۲) رگرسیون (۳) خوش‌بندی (۴) قوانین انجمانی تقسیم می‌شوند. که ما در تحقیق خود از تکنیک طبقه‌بندی و خوش‌بندی استفاده نمودیم. که طبقه‌بندی جهت پیشگویی مقادیر گسسته و اسمی مورد استفاده قرار می‌گیرد، خوش‌بندی در واقع یک عملیات غیر نظارتی می‌باشد. این عملیات هنگامی استفاده می‌شود که ما به دنبال یافتن گروه



گرفته اند. در این خوشه $\frac{۳۱}{۲۵}\%$ دانشجویان بهداشتکار جزء دانشجویان ضعیف بوده و $\frac{۶۸}{۷۵}\%$ جزء دانشجویان متوسط می باشند. میانگین معدل کل آنها $\frac{۱۴}{۴۰}$ با ضریب انحراف معیار $\frac{۰}{۸۱۸}$ می باشد.

در تقسیم بندی دانشجویان قوی (خوشه ۴) با میانگین معدل کل $\frac{۱۶}{۱۲}$ ، با توجه به میانگین معدل کل $\frac{۱۶}{۴۵}$ در تحقیق که برای تفیکیک دانشجویان قوی تعیین گردیده بود، بهتر عمل نموده است و درصد فراوانی جنسیتی و سهمیه ای را با تطبیق با مجموعه داده دقیقتر محاسبه نموده است ($\frac{۹۲}{۹۱}\%$ زن و $\frac{۷}{۰}\%$ مرد).

در مورد دانشجویان ضعیف (خوشه ۳) هر چند که میانگین معدل کل آنها $\frac{۱۴}{۰۶}$ با ضریب انحراف معیار $\frac{۱}{۰۲۱}$ با استانداردی که ما در این تحقیق در نظر گرفته ایم فاصله ای مشاهده میشود، اما از لحاظ جنسیت، سهمیه و تعداد ترم مشروطی و بدون مشروطی تفیکیک صحیح تری را انجام داده است. همچنین در مدل Two-Step برخلاف مدل K-Means دانشجویان ضعیفی که در طول تحصیل بدون مشروطی می باشند، مشاهده می شوند.

در مورد دانشجویان متوسط هر دو مدل خوب عمل ننموده اند. زیرا ۱۰۰% جنسیت تشکیل دهنده آنها مرد بوده اند و دانشجویان زن متوسط در هر دو مدل بیشتر در خوشه دانشجویان ضعیف قرار گرفتند. با توجه به اینکه هدف اصلی تحقیق ما دانشجویان قوی (ممتاز) و ضعیف (تحت ریسک) بود وضعیت دانشجویان متوسط مورد بررسی ما نمی باشد.

۲-۳ نتایج طبقه بندی

هدف اصلی از تکنیک طبقه بندی پیش بینی دانشجویان ممتاز (قوی) و تحت ریسک مشروطی (ضعیف) می باشد خروجی مدل را براساس دو فیلد قراردادیم یکی معدل کل و دیگری تعداد ترم مشروطی برای انجام آن از مدلهای CRT، CHAID و C5 استفاده نمودیم. بطورکلی در قوانین بوجود آمده، با توجه به تعداد رکوردهای شرکت کننده، قوانینی دارای شرط لازم و کافی برای انتخاب شدند که درصد پشتیبانی (support) آن دارای حداقل ۵% تعداد رکوردهای مورد استفاده در آموزش (Training) باشند و درصد اطمینان (confidence) آن حداقل ۸۰% باشد.

- قانون بدست آمده مدل CHAID بر اساس معدل کل برای دانشجویان قوی:

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۴ آنها بزرگتر از $\frac{۱۶}{۷۵}$ و معدل ترم ۳ آنها بزرگتر از $\frac{۱۶}{۵۰}$ بود،

جدول ۱- مقایسه دو مدل خوشه بندی K-Means و TwoStep

K-Means		
دانشجویان متوسط	خوشه ۱	
معدل کل $\frac{۱۵}{۱۴}\%$ - $\frac{۱۰۰}{۱۰۰}\%$ مرد - تعداد ترم مشروطی صفر - همه سهمیه وجود دارند بالاترین درصد فراوانی منطقه ۲ - کمترین شاهد		
معدل کل $\frac{۱۳}{۷۷}\%$ - $\frac{۶۸}{۶۸}\%$ زن و $\frac{۳۲}{۳۲}\%$ مرد هستند - ۱ $\frac{۵۲}{۵۲}\%$ یا ۲ مشروط شده اند و $\frac{۳}{۴۸}\%$ بار یا بیشتر - بجز بهداشتکار بقیه سهمیه میشوند، بالاترین درصد فراوانی شاهد و کمترین بهداشتکار ($\frac{۰}{۰}\%$)	دانشجویان ضعیف	خوشه ۲
معدل کل $\frac{۱۳}{۷۶}\%$ - $\frac{۱۰۰}{۹۷}\%$ مرد - $\frac{۶}{۶}\%$ تعداد ترم مشروطی ۱ یا ۲ بار و $\frac{۴۴}{۴۴}\%$ بار یا بیشتر - همه سهمیه ها بجز منطقه ۳ وجود دارند بالاترین درصد بهداشتکار و کمترین منطقه ۳ ($\frac{۰}{۰}\%$)	دانشجویان ضعیف	خوشه ۳
معدل کل $\frac{۱۵}{۹۷}\%$ - $\frac{۱۰۰}{۱۰۰}\%$ زن - $\frac{۹۷}{۹۷}\%$ مرد - تعداد ترم مشروطی صفر - همه سهمیه های وجود دارند بالاترین درصد مربوط به منطقه ۱ و کمترین رزمندگان و جانبازان	دانشجویان قوی	خوشه ۴
Two-Step		
دانشجویان بهداشتکار	خوشه ۱	
معدل کل $\frac{۱۴}{۴}\%$ - $\frac{۳۱}{۲۵}\%$ داشجویان بهداشتکار که ۱ یا ۲ ترم مشروط شدند جز ضعیف محسوب میشوند و $\frac{۶}{۶}\%$ داشجویان بهداشتکار که تعداد ترم مشروطی آنها صفر است که جزو دانشجویان متوسط محسوب می شوند	دانشجویان بهداشتکار	خوشه ۱
معدل کل $\frac{۱۵}{۲۳}\%$ - $\frac{۱۰۰}{۱۰۰}\%$ مرد - تعداد ترم مشروطی صفر - همه سهمیه ها بجز بهداشتکار هستند - بالاترین درصد منطقه ۲ و کمترین شاهد، فرزند هیئت علمی	دانشجویان متوسط	خوشه ۲
معدل کل $\frac{۱۴}{۰۶}\%$ - $\frac{۵۲}{۱۷}\%$ زن و $\frac{۴۷}{۸۳}\%$ مرد می باشند، $\frac{۱۸}{۸۴}\%$ ترم یا بیشتر مشروط شده اند و $\frac{۷۴}{۸۳}\%$ یک یا دو ترم مشروط شده اند و $\frac{۳۳}{۳۳}\%$ مشروط نشده اند همه سهمیه ها در آن شرکت دارند بجز بهداشتکاران، بالاترین درصد مربوط به شاهد در این سهمیه است و کمترین بهداشتکار	دانشجویان ضعیف	خوشه ۳
معدل کل $\frac{۱۶}{۱۲}\%$ - $\frac{۹۲}{۹۱}\%$ زن - $\frac{۷۹}{۷۹}\%$ مردمه یا جانباز بیشترین مقدار مربوط به منطقه ۲ و کمترین رزمند و جانباز	دانشجویان قوی	خوشه ۴

مدل Two-Step را می توان بهتر از K-Means بدلایل زیر در نظر گرفت:

خوشه ۱ مختص دانشجویان بهداشتکار بوده و هیچ سهمیه دیگری در آن جا نگرفته است و این نشان می دهد رفتار آموزشی دانشجویان بهداشتکار جدای سهمیه های دیگر بوده است. البته ۱۰۰% دانشجویان بهداشتکار در این خوشه قرار



* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها بزرگتر از ۱۵/۶۰ و مقدار فیلد pass3 آنها بزرگتر از ۰/۷۱۴ باشد، ۹۹/۱۳٪ آنها ترم مشروطی ندارند.

۹۴/۴ از آنها معدل بالای ۱۶/۴۴ داشتند که طبق تقسیم بندی تحقیق جزء دانشجویان قوی محسوب میشوند.
قانون بدست آمده مدل CHAID بر اساس معدل کل برای دانشجویان متوسط:

۴- نتیجه گیری

در مدل Two-Step، ۴ خوش بوجود آمد، خوش ۱ مربوط به دانشجویان بهداشتکار می باشد که با توجه به تعداد آنها و نوع رفتار آموزشی آنها یک خوش به آنها اختصاص یافت و در سایر خوش ها دیده نشدند. و در بررسی داده های آموزشی، خوش ۲ بعنوان خوش دانشجویان متوسط نام گرفت با توجه به اینکه ۱۰۰٪ جمعیت آن مرد بودند دانشجویان زن متوسط در خوشه های ۳ و ۴ قرار گرفتند و این از نقاط ضعف این مدل می باشد، خوش ۳، خوش دانشجویان ضعیف نام گرفت، و خوش ۴، خوش دانشجویان قوی که تعداد ترم مشروطی آنها صفر و ۹۲/۹۱٪ آنها زن بودند. که نشان دهنده تلاش و پشتکار و اهمیت دادن به درس خواندن در نزد خانم ها می باشد.

اطلاعات دیگری که از خوش بندی جهت تحلیل رفتار دانشجویان حاصل شد، شامل:

- افزایش درصد پذیرفته شدگان زن نسبت به مرد از سال ۸۸ به بعد

- بالا بودن میانگین معدل کل خانم ها در تمامی ورودیها نسبت به آقایان که تایید کننده بالا بودن تعداد زن ها در دانشجویان قوی نسبت به مرد ها می باشد.

- با توجه به وجود دانشجویان با تعداد ترم مشروطی ۳، ۴ و ۵ فقط در خوش ۳ با صحت بیشتری می توان این خوش را بعنوان خوش دانشجویان ضعیف نامگذاری کرد.

- با توجه به مشاهده شدن دانشجویان با تعداد ترم مشروطی ۴، ۳ و ۵ در سهمیه های شاهد، رزمندگان و فرزندان هیئت علمی، نیاز به بازنگری آموزشی در این سه گروه می باشد، توصیه می شود نقاط ضعف آموزشی این دانشجویان شناسایی و مسئولین جهت بهبود آنها برنامه ریزی نمایند.

- دانشجویان ورودی ۷۸ و ۸۴ دارای کمترین تعداد مشروطی در بین ورودیها بوده اند(ورودی ۷۸ یک نفر و ورودی ۸۴ دو نفر) که نسبت به دانشجویان ورودی سالهای دیگر دارای راندمان آموزشی خوبی بودند.

- با توجه به بالا بودن میانگین معدل خانم ها نسبت به آقایان و قرار گرفتن اکثریت خانم ها در خوش دانشجویان قوی

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۴ آنها بزرگتر از ۱۳/۶۲ و کمتر از ۱۵/۰۳ باشد و معدل ترم ۳ آنها بزرگتر از ۱۲/۱۰ و کمتر از ۱۷/۳۳ باشد، ۹۰٪ از آنها معدل بین ۱۳/۹۵ و ۱۶/۴۴ را داشتند که طبق تقسیم بندی تحقیق جزء دانشجویان متوسط محسوب میشوند.

قانون بدست آمده مدل CHAID بر اساس معدل کل برای دانشجویان ضعیف:

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۴ آنها کوچکتر یا مساوی ۱۳/۶ باشد و معدل ترم ۱ کوچکتر مساوی ۱۳/۴۷ باشد و مقدار pass1 آنها کوچکتر مساوی ۰/۹۵ باشد، ۷۵٪ از آنها معدل کمتر از ۱۳/۹۵ داشتند که طبق تقسیم بندی تحقیق جزء دانشجویان ضعیف محسوب میشوند.

- قانون بدست آمده مدل CHAID بر اساس تعداد ترم مشروطی:

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها کوچکتر مساوی ۹/۳۵ باشد، ۱۰۰٪ از آنها تعداد ترم مشروطی مساوی یا بیشتر از ۳ ترم دارند.

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها بزرگتر از ۹/۳۵ و کوچکتر مساوی ۱۱/۳۶ باشد، ۹۰٪ از آنها تعداد ترم مشروطی ۱ یا ۲ ترم دارند.

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها بزرگتر از ۱۲/۴۲ و کوچکتر مساوی ۱۳/۳۸ باشد، ۷۵٪ از آنها تعداد ترم مشروطی ۱ یا ۲ ترم دارند.

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها کوچکتر مساوی ۱۳/۳۸ و بزرگتر از ۱۲/۴۲ باشد و معدل ترم ۴ آنها کوچکتر مساوی ۱۳/۳۸ و بزرگتر از ۱۲/۵۷ باشد، ۶۶٪ از آنها تعداد ترم مشروطی ۱ یا ۲ ترم دارند.

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها کوچکتر مساوی ۱۴/۱۰ و بزرگتر از ۱۳/۳۸ و سهمیه آنها فرزند هیئت علمی یا انتقالی از خارج باشد، ۶۰٪ از آنها تعداد ترم مشروطی ۱ یا ۲ ترم دارند.

* در مجموعه داده ها، از دانشجویانی که معدل ترم ۲ آنها بزرگتر از ۱۴/۱۰۳۵ و کوچکتر مساوی ۱۵/۶۰ و معدل ترم ۴ آنها بزرگتر از ۱۴/۷۹ باشد، ۱۰۰٪ آنها ترم مشروطی ندارند.



مراجع

می توان ، جنسیت را جزء عوامل موثر بر نمرات دانشجویان برشمرد.

-با وجود قرارگرفتن تمامی سهمیه ها در خوشه دانشجویان ضعیف و متوسط می توان گفت متغیر سن تاثیری در این دو گروه ندارد، ولی با قرارگرفتن دانشجویان سهمیه مناطق در خوشه قوی و داشتن میانگین سنی زیر ۲۰ سال در این سهمیه ، می توان گفت متغیر سن در کسب نمرات بهتر در این خوشه تأثیر گذار است.

-نتایج حاصل از تحقیق نشان می دهد نوع سهمیه ورودی دانشجویان نیز بر نمرات کسب شده توسط دانشجویان مؤثر است. بطوریکه دانشجویان سهمیه های شاهد، رزمنه و فرزند هیئت علمی دارای نمرات پایین تری بودند و نیاز به توجه آموزشی دارند.

در مدل CHAID تکنیک طبقه بندی معدل ترم ۴ جزء پراهمیت ترین فیلد در درخت تصمیم می باشد. قوانین برگرفته از آن دارای ضریب اطمینان ۷۵ تا ۱۰۰ درصد بود، قوانین بدست آمده از نظر پیش بینی دانشجویان تحت ریسک و ممتاز بسیار ارزشمند بوده و می توان از آنها بعنوان قوانینی در تحلیل رفتار دانشجویان استفاده نمود.

- [1] Thearling, K., 2007. An Introduction to Data Mining [Online]. Available at: <http://www.thearling.com>
- [2] Ryan S.J.d. Baker., 2008. Data Mining for Education. [pdf] Available at: <<http://users.wpi.edu/~rsbaker/Encyclopedia%20Chapter%20Draft%20v10%20-fw.pdf>>
- [3] Fayyad, U., Piatetsky-Shapiro G., and Smyth P., 1996. From Data Mining to Knowledge Discovery in Databases, American Association for Artificial Intelligence, [pdf] Available at: www.aaai.org/aitopics/assets/PDF/AIMag17-03-2-article.pdf
- [4] Dekker, G., Pechenizkiy, M. and Vleeshouwers, J., 2009. Predicting Students Drop Out: A Case Study. p.41-50. Conference Proceedings of the 2nd International Conference on Educational Data Mining, Cordoba, Spain [pdf]. Available at: <http://www.win.tue.nl/~mpechen/projects/edm/internshipreport_090409.pdf>
- [5] شکورنیاز. ونوس و حاجی علی اکبری. آرش، سال ۱۳۸۷. «خوشه بندی داده های آماری دانشجویان دانشگاه علم و صنعت و استخراج نمایه ساز توصیفی برای دانشجویان موفق»، مقاله دومین کنفرانس انجمن داده کاوی ایران.
- [6] - صحافی زاده. ابراهیم، سال ۱۳۸۸. «تحلیل عوامل مؤثر بر نمرات دانشجویان دانشگاه پیام نور بوشهر با استفاده از داده کاوی»، مقاله سومین کنفرانس انجمن داده کاوی ایران.