

مطالعاتی

محسن فرهادی^۱، مریم چوپانی^۲

چکیده

رشد روز افزون اطلاعات در سازمان‌ها و شرکت‌ها موجب افزایش پیچیدگی روند نگهداری، بازیابی و تحلیل اطلاعات شده است. داده کاوی یا کشف دانش در پایگاه داده‌ها، راهبرد جدیدی است در راستای کشف محتوی اطلاعات معتبر پنهان در پایگاه‌های داده. مدیریت، بازیابی، نگهداری و تجزیه و تحلیل داده‌ها برای پیدا کردن الگوهای تازه و کاربردی در پایگاه داده‌های بزرگ از این طریق میسر می‌شود. اغلب ارگانها به تصمیم گیری‌های استراتژیک و اتخاذ خط مشی‌های جدید برای عملکرد بهتر نیاز دارند. در واقع داده کاوی اطلاعاتی را در اختیار ما قرار می‌دهد، که برای گرفتن تصمیم هوشمندانه در انبوهی از اطلاعات به آن نیاز است. ابزار داده کاوی، داده را می‌گیرد و یک دید کلی از واقعیت را به شکل مدل می‌سازد، این مدل روابط موجود در داده‌ها را شرح می‌دهد. داده کاوی در امر آموزش یکی از راهبردهای مطرح در راستای بهبود و توسعه اهداف آموزشی است. مطالعه و تحلیل اطلاعات در بخش آموزش عالی، عملکرد و کارایی روال‌های اجرایی را مورد سنجش قرار می‌دهد، در صورت لزوم امر پیش بینی را میسر می‌سازد و زمینه و ایده لازم برای اصلاح و بهبود بهره وری را فراهم می‌آورد. یکی از موارد کاربردی در این زمینه، پیش بینی عملکرد دانشجویان بر اساس داده‌های موجود در پایگاه اطلاعاتی سیستم آموزشی، از قبیل زمان آزمون، فرصت مطالعاتی و .. می‌باشد. ما در این مقاله تأثیر جنبه‌های مختلف آزمون را در نمره‌ی اکتسابی دانشجویان در دروس مختلف را مورد بررسی قرار داده‌ایم.

واژه‌های کلیدی

داده کاوی، آموزش عالی، آزمون، پیش‌بینی نمرات

کنفرانس داده کاوی ایران

۱. عضو هیئت علمی دانشکده کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی شاهرود narfarhadi@yahoo.com

۲. فارغ التحصیل دانشکده کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی شاهرود maryamchoopany@yahoo.com

Students' scores prediction based on the exam date and the preparation time

Mohsen Farhadi¹Maryam Choopany²

Abstract

Daily increasing amount of data in organizations and companies, the process of data storage, analysis and recovery is becoming more complex than ever. Data-mining or knowledge exploration in database is a new procedure in finding hidden valid knowledge and it also makes management, recovery, storage, and analysis of data to find new and practical patterns in database possible. Almost all companies and organizations need to make strategic decisions and determine new procedures to promote their actions and that is when data-mining comes in to provide us with required data to do so, among tons of data.

As a matter of fact, by means of data-mining, data is chosen and analyzed to give out a general perspective of facts, in the form of a model, explaining the relationships among the data. The study and analysis of data in higher education can evaluate the practicality of the procedures in use, make the necessary predictions possible, and can also provide us with ideas to reform the procedures and increase the efficiency.

One of the implements of data-mining is the prediction of students' performance at universities based on the existing data in the database of educational system such as the date of the exam, preparation time of the exam, etc. In this article, the relationship between different aspects of an exam and the students' scores has been studied.

Keywords

Data-mining, higher education, exams, score prediction

پیش‌بینی آینده در زمینه‌های مختلف همواره برای انسان جالب و جذاب بوده است. با اطمینان می‌توان گفت که پیش‌بینی آینده و روند تغییرات در همه‌ی حوزه‌ها از دغدغه‌های اصلی و همیشگی مدیران سطح بالا و میانی می‌باشد. اما همواره مشکلات فراوانی در برابر آن وجود داشته است که انجام پیش‌بینی‌های دقیق و قابل اعتماد را تقریباً غیر ممکن نموده است. وجود پارامترهای زیاد و در بسیاری موارد پنهان اینگونه موارد را به مسائلی بسیار پیچیده تبدیل نموده است که الگوریتم‌های غول‌پیکر ریاضی نیز از ارائه راهکاری مناسب برای ساخت یک مدل پیش‌بینی کارآمد عاجز مانده‌اند. [۲]

امروزه در اکثر دانشگاه‌های ایران اسلامی بانک‌های اطلاعاتی وسیعی از ویژگی‌های دانشجویان موجود است که حجم بالایی از اطلاعات مربوط به سوابق آموزشی، تحصیلی و ... را شامل می‌شود. در سال‌های اخیر با توجه به مشکلات موجود در زمینه استخراج دانش از چنین پایگاه‌داده‌های عظیمی و استفاده از آنها به منظور تصمیم‌گیری‌های مناسب، تکنیک‌های داده کاوی و استفاده از آنها در امر آموزش مورد توجه قرار گرفته است. [۱] هدف پژوهش‌هایی این چنین نیز استخراج قوانین و الگوهایی قابل استنتاج و مرتبط از روی پایگاه داده دانشگاه به منظور بهبود کیفیت امور در آموزش عالی می‌باشد.

از جمله تحقیقات مشابه انجام شده در این زمینه می‌توان مواردی همچون به‌کارگیری ابزارهای داده‌کاوی جهت پیش‌بینی موفقیت دانشجویان در فارغ‌التحصیلی توسط فرهادی و عطفایی به صورت ترمی [۱]، پیش‌بینی میزان موفقیت و عدم موفقیت دانشجویان رشته فناوری اطلاعات در فارغ‌التحصیلی، بر اساس معدل، توسط توحیدی مقدم و فرهادی [۲] که هر دو روش از ابزار درخت تصمیم استفاده کرده‌اند، را نام برد.

در یک نظرسنجی از داده کاوان در مورد اینکه در ۱۲ ماه اخیر از چه زبان برنامه نویسی برای تحلیل داده و داده کاوی استفاده کرده‌اید، زبان sql با ۳۲٪ در رتبه دوم قرار گرفت [۵]. در نظرسنجی سالیانه دیگر از داده‌کاوان در مورد موضوع کاری داده کاوی، در سال ۲۰۱۲ نیز، کاربرد داده کاوی در امر آموزش با اختصاص ۱۶٫۲٪ در رتبه پنجم قرار گرفته است [۵] که نشان دهنده آن است که بحث آموزش زمینه‌ای است پرتعداد و نیازمند تحقیق. بدین ترتیب این تحقیق با انتخاب زبان پرکاربرد SQL و موضوع آموزش عالی با بررسی تاثیر عوامل جانبی آزمون بر نمره اکتسابی آغاز شد، در ادامه اطلاعات گرد آوری شده و در پایگاه داده یکپارچه و اصلاح شد، سپس دسته‌هایی از اطلاعات برای تست و تحلیل به صورت تصادفی انتخاب شد و به منظور ایجاد یک الگو و سنجش آن تعیین شد، پس از آن نتایج تحلیل در حالت‌های مختلف بررسی شد، که نشان دهنده عوامل موثر در نتیجه مطلوب دانشجویان است.

۲. آماده سازی داده‌ها و پیش‌بینی نمرات

پس از انتخاب زبان اطلاعات مورد نظر از پایگاه داده دانشگاه دریافت شد، اطلاعات در مورد دانشجو، نمرات، گروه‌های درسی و درس‌ها در جداول جداگانه از پایگاه داده استخراج شد. در مرحله بعد با کمک نرم افزار SQL Server Management Studio ابتدا داده‌ها برای تحلیل آماده شدند. محاسباتی روی اطلاعات انجام شد و صفات مورد نیاز برای آنالیز برآورد شد و اطلاعات اضافی حذف شد. صفاتی همچون معدل متوالی دانشجو (معدل کل دانشجو از ترم اول تا ترم قبل از اخذ درس مورد پیش‌بینی)، زمان امتحان (صبح یا ظهر بودن) برای شرکت در تحلیل، آماده شدند. صفاتی دیگر همچون معدل امتحان قبل در پایگاه داده (معیاری برای اندازه‌گیری درجه سختی امتحان قبلی دانشجو)، فرصت مطالعاتی (فاصله زمانی امتحان تا امتحان قبلی)، معدل نمره دهی استاد (میانگین نمرات ثبت شده در سیستم از یک استاد) و معدل درس (میانگین نمرات ثبت شده در سیستم از یک درس) محاسبه شدند، که تمامی صفات ورودی الگوریتم در بخش ضمایم در **Error! Unknown switch argument.** آمده است. زمانی که داده‌ها برای داده کاوی آماده شدند می‌توان وارد مرحله دوم از انجام داده کاوی شد.

با کمک نرم افزار Microsoft decision trees از SQL Server BusinessIntelligence Development Studio برای پیش‌بینی استفاده شد. در مجموع ۶۵۵۲۵ داده از پایگاه داده دانشگاه صنعتی شاهرود به فرایند تحلیل داده شد. ۹۰٪ از داده‌ها برای آموزش و تحلیل به صورت تصادفی انتخاب شد و ۱۰٪ دیگر برای تست الگوی به دست آمده اختصاص داده شد.

۱-۲ درخت تصمیم

درخت تصمیم یکی از رایج‌ترین تکنیک‌های داده‌کاوی است. معمول‌ترین زمینه‌ی کاری درخت تصمیم، دسته‌بندی می‌باشد. الگوریتم درخت تصمیم مایکروسافت که توسط سرویس‌های خدمات تحلیلی ارائه شده است، تکنیک‌های دسته‌بندی و رگرسیون را پشتیبانی می‌نماید و برای مدل‌های پیشگویانه بسیار خوب عمل می‌کند. با استفاده از این الگوریتم مشخصه‌های گسسته و پیوسته قابل پیش‌بینی

هستند. مقصود اصلی در درخت تصمیم، تقسیم داده‌ها به صورت بازگشتی به زیرمجموعه‌هایی است؛ به گونه‌ای که هر زیرمجموعه در برگزیده وضعیت همگنی از متغیر هدف می‌باشد. زمانیکه پردازش بازگشتی کامل شد، درخت تصمیم شکل گرفته است [۴].

- ۱-۱-۲ ویژگی‌ها و مزایای درخت تصمیم میکروسافت
- از مزایای استفاده از درخت تصمیم نسبت به سایر الگوریتم‌های داده‌کاوی این است که مدل آن سریع ساخته و آسان‌تر تفسیر می‌شود [۴].
- پیش‌بینی‌ها بر اساس درخت تصمیم مؤثرتر است [۴].
- روش‌های مختلفی برای رشد یک درخت وجود دارد و می‌توان از فرمول‌های متنوعی برای تعیین چگونگی تقسیم درخت استفاده نمود [۴].
- بر پایه تنظیم پارامتری می‌باشد یعنی اینکه رشد درخت بر اساس تعداد متغیرهای ورودی و با تنظیم پارامترها قابل کنترل است [۴].
- یکی از ویژگی‌های منحصر به فرد درخت تصمیم میکروسافت این است که می‌تواند برای آنالیز وابستگی‌ها نیز به کار گرفته شود [۴].

پس از اجرای الگوریتم درخت تصمیم میکروسافت و شبکه وابستگی‌ها به دست آمد. برای مثال یکی از برگ‌های درخت تصمیم برای دانشجویانی که معدل آنها بزرگتر از ۱۲,۰۶۷ و کوچکتر از ۱۴,۲۶۷ باشد به صورت تصویر ۱ بدست آمد که ۶۰۵۶ مورد در این شرط صدق می‌کرد، به این معنی که پیش‌بینی می‌شود مواردی که در این بازه معدل کل قرار دارند نمره درس آنها با فرمول بدست آمده برای grade در تصویر ۱ پیش‌بینی می‌شود.

جدول ۱

1. Gpa معدل متوالی دانشجو از ترم اول تا ترم قبل از اخذ درس مورد پیش‌بینی
2. Ostadavg معدل نمره دهی استاد
3. Crsnameavg معدل درس مورد پیش‌بینی در کل پایگاه داده
4. Grade نمره مورد پیش‌بینی یا نمره مورد انتظار

Existing Cases: 6056
Missing Cases: 0
 $GRADE = 13.362 + 0.961 * (GPA - 14.521) + 0.665 * (Ostadavg - 13.927) + 0.649 * (Crsnameavg - 13.367)$

تصویر ۱

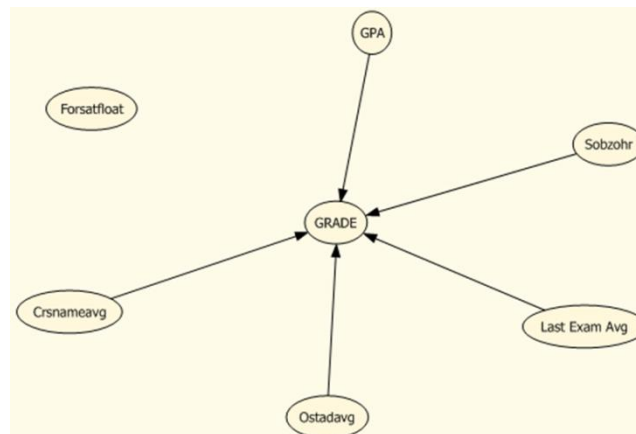
کنفرانس داده‌کاوی ایران

۲-۲ شبکه وابستگی‌ها

شبکه وابستگی‌ها معرف قدرت مشخصه‌های ورودی و پیش‌بینی کننده در پیش‌بینی متغیر هدف می‌باشد. هر نود در شبکه وابستگی‌ها بیانگر یک مشخصه و هر بردار نشانگر وجود یک رابطه بین دو نود می‌باشد. هر بردار دارای یک جهت مشخصه‌ی ورودی به سمت مشخصه‌ی قابل پیش‌بینی می‌باشد. استفاده از این نمودار زمانی مفید است که تعداد زیادی مشخصه‌ی قابل پیش‌بینی وجود داشته باشد. این نمودار در تحلیل‌های اکتشافی داده‌ها بسیار سودمند است [۴].

۳. تحلیل پیش‌بینی‌ها

پس از بدست آمدن درخت تصمیم از روی ۹۰٪ از داده‌ها که به صورت تصادفی انتخاب شدند، پیش‌بینی با ۱۰٪ داده‌های باقی‌مانده انجام شد تا دقت عمل درخت تصمیم بررسی شود، در نتیجه **Error! Unknown switch argument.** با score قابل قبول 1.44 بدست آمد که نشان می‌دهد پیش‌بینی درخت تصمیم به خوبی انجام شده است، سپس شبکه وابستگی‌ها به صورت **Error! Unknown switch argument.** حاصل شد.



نمودار ۱

ترتیب صفات موثر بر نمره نهایی به صورت جدول 2 بدست آمد که همانطور که آمده است، معدل نمره دهی استاد اولین و موثرترین صفت در تعیین نمره بدست آمده است.

جدول 2

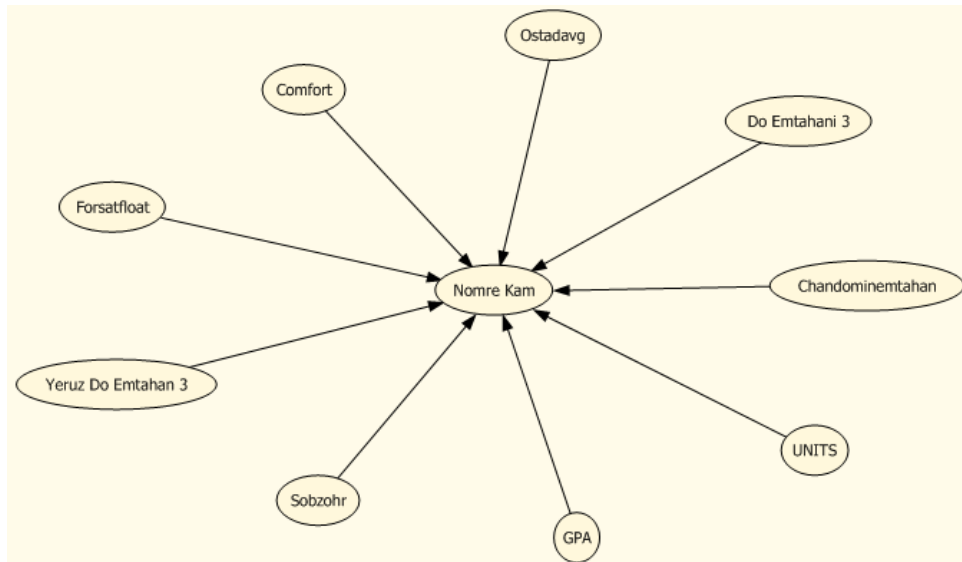
نمره مورد پیش‌بینی یا نمره مورد انتظار Grade
1. معدل نمره دهی استاد Ostadavg
2. معدل متوالی دانشجو از ترم اول تا ترم قبل از اخذ درس مورد پیش‌بینی Gpa
3. معدل درس مورد پیش‌بینی در کل پایگاه داده Crsnameavg
4. Sobzohr صبح یا ظهر بودن زمان امتحان
5. Last exam avg معدل امتحان قبل
6. Forsat float یا فرصت مطالعاتی فاصله زمانی امتحان تا امتحان قبل

باقی صفات ورودی که در **Error! Unknown switch argument.** آمده‌اند اما در شبکه وابستگی **Error! Unknown switch argument.** پدیدار نشدند حاکی از این مطلب است که الگوریتم این صفات را در پیش‌بینی بی‌تأثیر دانسته است. در طی اجرای پروژه بر آن بودیم تا رابطه نمره و دو امتحانی بودن را بیابیم اما به این هدف نرسیدیم لذا موضوع پیش‌بینی را به کسری نمره یعنی تفاضل نمره از ۲۰ تغییر دادیم تا ببینیم چه عواملی موجب نتیجه نامطلوب دانشجویان می‌شود.

۴. اجرای دوباره آنالیز با کسری نمره

همانطور که پیش‌تر گفته شد در این بخش متغیری به نام کسری نمره تعریف شد تا تأثیر عناصری چون دو امتحان در یک روز بر نمره سنجیده شود. با حذف معدل امتحان قبلی و معدل درس (last exam avg, crsnameavg) و اکتفا به دسته بندی چهار سطحی

انجام شده برای همین صفات (در متغیر last examcomfort و comfort) الگوریتم اجرا شد. ورودی‌های الگوریتم در این بخش در **Error! Unknown switch argument.** آمده است و دو صفت حذف شده، با علامت ضربدر در جلوی آنها، مشخص شده اند. نتیجه این تغییر، یعنی حذف دو متغیر پیوسته‌ی avg و اکتفا به بازبندی در ۵ سطح، کاهش تعداد سطوح و تعداد برگها به معنی دسته بندی بهتر و افزایش score را به همراه داشت و مشاهده شد که عوامل تاثیر گزار تغییر کرد.



نمودار ۲

ترتیب تاثیرات در dependency network بر روی کسری نمره (nomrekam) در جدول ۳ آمده است

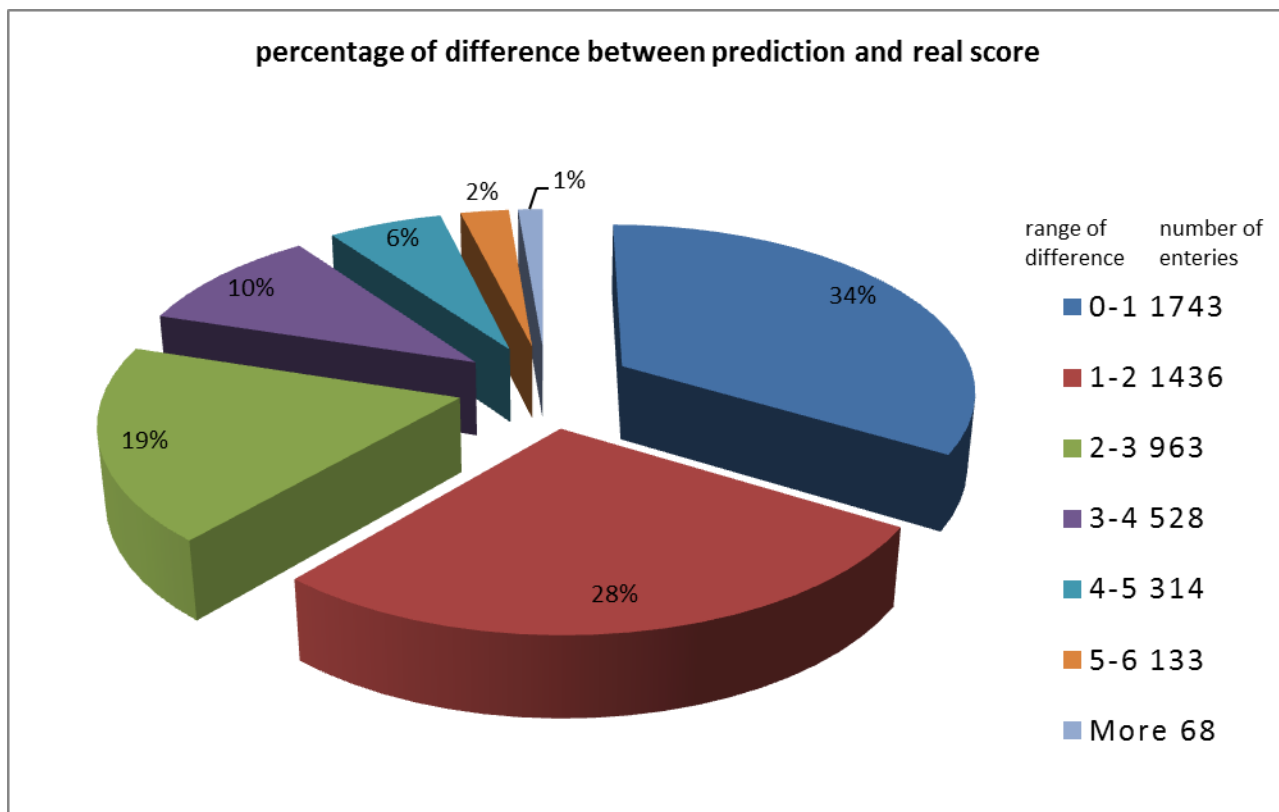
جدول ۳

1. Ostadavg معدل نمرات ثبت شده توسط هر استاد
2. Gpa معدل متوالی دانشجو
3. Comfort درجه راحتی یک درس
4. Units تعداد واحدهای یک درس
5. Chandominemtahan تعداد امتحانهایی که یک دانشجو تا یک روز داشته است
6. Yeruz do emtahan دانشجو در یک روز دو امتحان در دو ساعت متفاوت دارد
7. Forsatfloat فرصت مطالعاتی یا فاصله زمانی با امتحان قبلی
8. Do emtahani نشان می‌دهد دانشجو در یک روز، امتحانی دیگر، در ساعت یکسان با درس مورد پیش بینی دارد
9. Sobzohr صبح یا ظهر بودن امتحان

۵. گزارش گیری و سنجش نتایج

پس از آنالیز گزارش گرفته شد و بدست آمدن الگوها، برای ۱۰٪ از داده‌ها (۵۱۸۵ داده) که از قبل برای تست جدا شده بودند پیش بینی دوباره انجام شد. لازم به ذکر است که این داده‌ها از ابتدا جدا شده بودند و در جریان بدست آمدن الگو نقشی نداشتند، پس از آن نمرات پیش‌بینی شده با نمرات واقعی در مقابل هم قرار گرفتند و تفاوت آنها محاسبه شد. برای مثال در تصویر ۲ خطای پیش بینی و

تعداد آن در بازه های متفاوت آمده است که چه تعداد از ۵۱۸۵ پیش‌بینی، خطایی بین صفر تا یک داشته اند و چند درصد از کل را تشکیل می‌دهند. همانطور که مشخص است بیشترین نسبت مربوط به خطاها با تفاضل کمتر از یک است.



تصویر ۲ خطای پیش‌بینی و تعداد آن در بازه های متفاوت

۶. نتیجه‌گیری

با توجه به تاثیر قابل توجه اساتید در نمره دانشجویان دریافت می‌شود که نحوه تدریس و نمره دهی استاد به طور مستقیم بر دو عنصر نمره‌ی دانشجو و کسری نمره از بیست موثر است، لذا توجه به این موضوع از جهات مختلف، موضوعی است که باید بیش از اینها مورد بررسی قرار گیرد. پیشنهاد می‌گردد برای هر استاد معیاری به نام معدل استاد ایجاد و نهادینه شود تا مقایسه و بررسی عملکرد اساتید تسهیل شود چراکه این معیار اولین ملاک موثر چه در کسب نمره و چه در نتیجه منفی دانشجویان در آزمونی باشد. فرصت مطالعاتی، داشتن دو امتحان در یک روز و یک ساعت عوامل دیگر در نتیجه منفی دانشجو در امتحانات بود که ناشی از چینش نامناسب برنامه امتحانی است، لذا برای بازدهی بیشتر و برنامه ریزی صحیح امتحانات بایستی تلاش کرد تا این امور با بهره گیری از هوش مصنوعی مکانیزه شود.

Mining model structure:

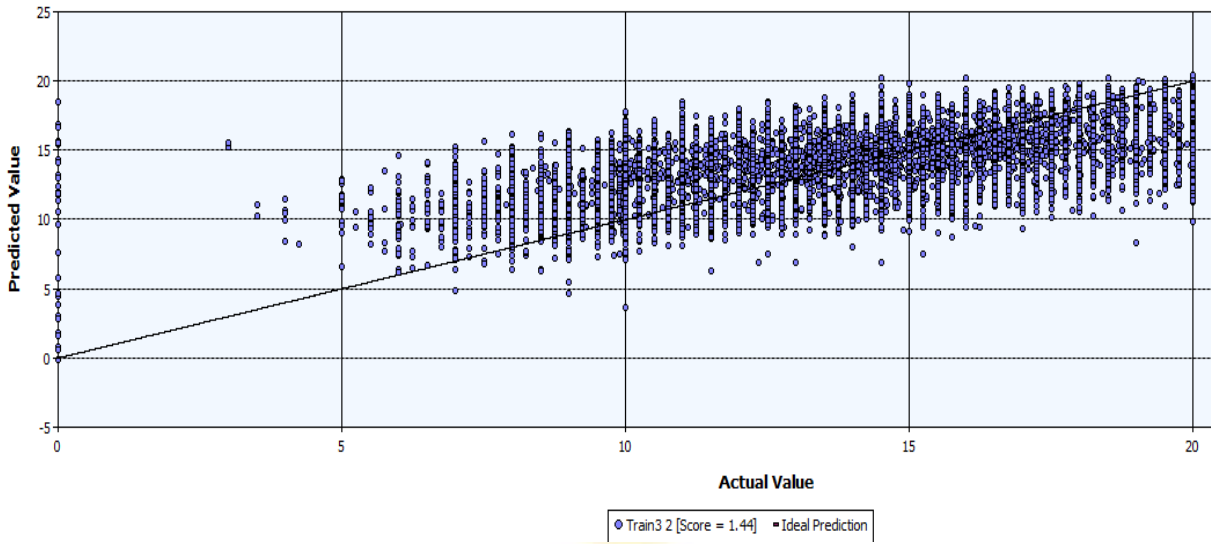
	Content Type	Data Type
Chandominemtahan	Ordered	Long
Comfort	Ordered	Long
Crsnameavg	Continuous	Double
Do Emtahani 3	Discrete	Text
Forsatfloat	Continuous	Double
GPA	Continuous	Double
GRADE	Continuous	Double
Key Id Trm Crs	Key	Text
Last Exam Avg	Continuous	Double
Last Exams Comfort	Discretized	Long
Ostadavg	Continuous	Double
Sobzohr	Discrete	Text
Unit	Discrete	Text
Yeruz Do Emtahan 3	Discrete	Text

تصویر ۲ صفات ورودی الگوریتم و نوع داده ای آنها

Structure	Train4 2
	Microsoft_Decision_Trees
Chandominemtahan	Input
Comfort	Input
Crsnameavg	Ignore X
Do Emtahani 3	Input
Forsatfloat	Input
GPA	Input
Key Id Trm Crs	Key
Last Exam Avg	Ignore X
Last Exams Comfort	Input
Nomre Kam	PredictOnly
Ostadavg	Input
Sobzohr	Input
UNITS	Input
Yeruz Do Emtahan 3	Input

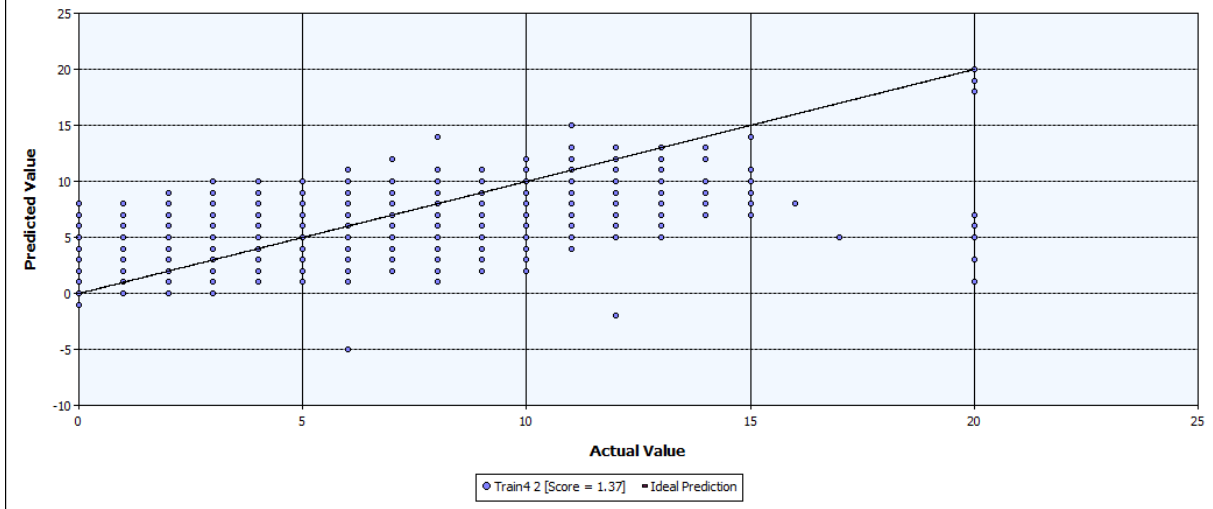
تصویر ۲

Data Mining Scatter Plot for Mining Structure: Train3 2

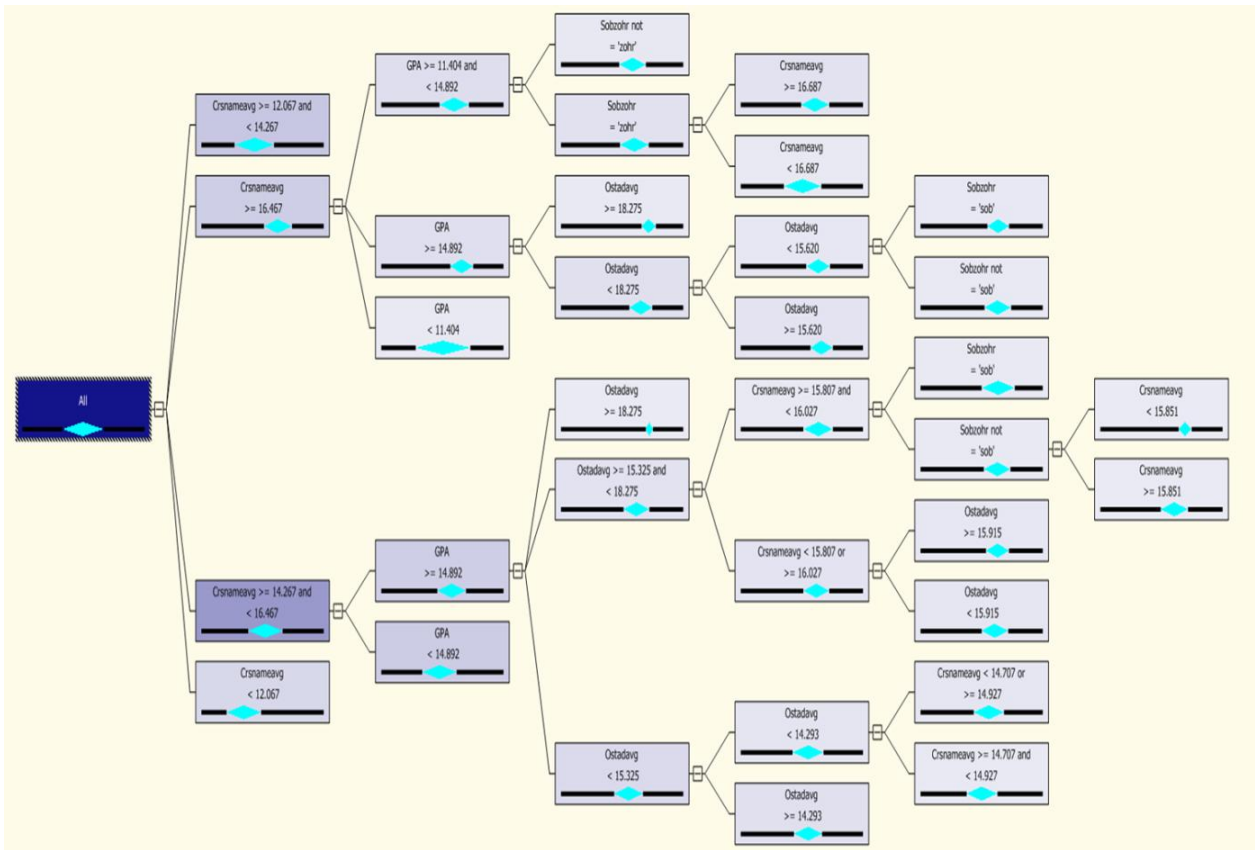


نمودار ۳

Data Mining Scatter Plot for Mining Structure: Train4 2



نمودار ۴



تصویر ۴ درخت تصمیم

۸. مراجع

- [۱] فرهادی، محسن؛ عطفانی، مهسا؛ "به کارگیری ابزارهای داده کاوی جهت پیش بینی موفقیت یا عدم موفقیت دانشجویان رشته کامپیوتر در فارغ التحصیلی"، کنفرانس داده کاوی ایران، دوره ششم، دانشگاه صنعتی امیرکبیر، ۱۳۹۱.
- [۲] توحیدی مقدم، مریم؛ فرهادی، محسن؛ "پیش بینی میزان موفقیت و عدم موفقیت دانشجویان رشته فناوری اطلاعات در فارغ التحصیلی"، کنفرانس داده کاوی ایران، دوره ششم، دانشگاه صنعتی امیرکبیر، ۱۳۹۱.
- [۳] نوری، بهاره؛ مقصودی، بهروز؛ شیخ احمدی، سیدامیر؛ "پیدا کردن دروس مرتبط از طریق الگوریتم سید خرید و تاثیر آن در پیش بینی نمره از طریق الگوریتم درخت تصمیم" چهارمین کنفرانس داده کاوی، دانشگاه صنعتی شریف، ۱۳۸۹.
- [۴] شهرابی، جمال؛ شکورنیا، ونوس؛ داده کاوی در SQL Server (اولین مرجع داده کاوی کاربردی به زبان فارسی)، انتشارات جهاد دانشگاهی واحد صنعتی امیرکبیر، تهران، اول، ۱۳۸۸.

[۵] <http://www.kdnuggets.com/>



آکادمی داده 
dataacademy.ir