

پیش‌بینی افزایش ثبت‌نام در مؤسسات آموزش عالی

با استفاده از مدل MLP

بهزاد نخکوب^۱، مریم خادمی^۲ و علی برومندنی^۳^۱ دانشجوی کارشناسی‌ارشد مهندسی کامپیوتر- نرم‌افزار، دانشگاه آزاد اسلامی واحد تهران جنوب، bnakh@yahoo.com^۲ استادیار گروه ریاضی کاربردی، دانشگاه آزاد اسلامی واحد تهران جنوب، khademi@azad.ac.ir^۳ استادیار گروه مهندسی کامپیوتر، دانشگاه آزاد اسلامی واحد تهران جنوب، broumandnia@azad.ac.ir

چکیده - در این مقاله، با استفاده از تکنیک داده‌کاوی بر روی مجموعه اطلاعات دواطلبان مقطع ارشد آزمون سراسری دانشگاه آزاد اسلامی، الگوریتمی جهت پیش‌بینی میزان جذب و افزایش ثبت‌نام دانشجویان در سال‌های آتی ارائه دادیم. در این راستا ابتدا به مدل‌سازی داده‌ها در ۱۵ مدل شبکه عصبی پرداختیم و به منظور افزایش دقت از روش جمعی بوستینگ استفاده کردیم. سپس دو مدل: شبکه‌های عصبی و درخت تصمیم را بر روی مجموعه داده‌ها دوباره اعمال نمودیم و با شاخص ارزیابی دقت کاپا، نتایج حاصله را با هم مقایسه کردیم. در نهایت، دقت روش بوستینگ به منظور پیش‌بینی دانشجویانی که در رشته - محل قبول شده‌اند و ثبت‌نام می‌نمایند، حدوداً ۹۶٪ برآورد شد.

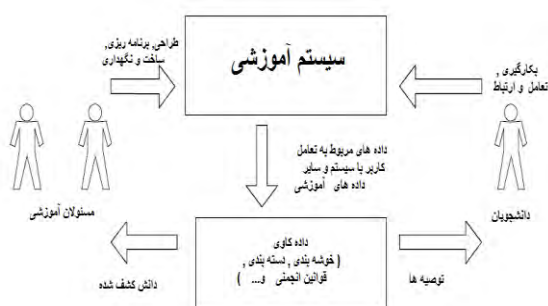
کلید واژه- مدل MLP، شبکه‌های عصبی، درخت تصمیم، شاخص دقت کاپا.

مالی دانشگاه‌ها و موارد دیگر کاربرد داشته باشد. همان‌طور که در شکل (۱) می‌بینید آموزش دهندگان و مسئولان آکادمیک، مسئول طراحی، برنامه‌ریزی و حفظ سیستم آموزشی هستند که فراگیران از آن استفاده می‌کنند و با آن تعامل دارند. با استفاده از تمام داده‌های موجود، تکنیک‌های متفاوتی از کاوش اطلاعات به-کارگرفته می‌شود تا اطلاعات مفیدی برای بهبود فرآیند آموزشی به‌دست آید. بنابراین کاربرد کاوش داده‌ها در سیستم آموزشی می‌تواند گرایش به سمت افرادی متفاوت و با طرز نگرش متفاوت داشته باشد (Romero, 1996)

۱ - مقدمه

در سال‌های اخیر تکنیک‌های داده‌کاوی از اهمیت به‌سزای برخوردار شده است. زیرا داده‌کاوی به بررسی و تجزیه و تحلیل مقادیر عظیمی از داده‌ها به منظور کشف الگوها و قوانین پنهان و معنی‌دار درون داده‌ها می‌پردازد. به عبارتی داده‌کاوی از منظری وسیع‌تر دیدگاهی است که مطابق آن کلیه فعالیت‌های تجاری باید براساس یادگیری باشد.

با توجه به این که آموزش عالی همواره با داده‌ها و اطلاعات بسیار زیادی در مورد دانشگاه‌ها، دانشجویان، اعضای هیأت علمی، کارکنان، منابع مادی و ... روبرو است و در اکثر مواقع این داده‌ها می‌تواند حامل اطلاعات و الگوهای با ارزشی باشند، لذا به نظر می‌رسد یکی از مهم‌ترین کاربردهای داده‌کاوی در آموزش عالی است. استفاده از تکنیک‌های پیشرفته داده‌کاوی مانند خوشه‌بندی، طبقه‌بندی، و ... می‌تواند در طبقه‌بندی دانشگاه‌ها، یافتن الگوهای خاص و با ارزش در مورد دانشجویان موفق، یافتن یک برنامه یا روش موفق تدریس، یافتن نقاط بحرانی در مدیریت



شکل ۱: چرخه به کارگیری داده کاوی در سه سیستم های آموزشی

شکل ۲: کریسپ (مدل اجرایی)

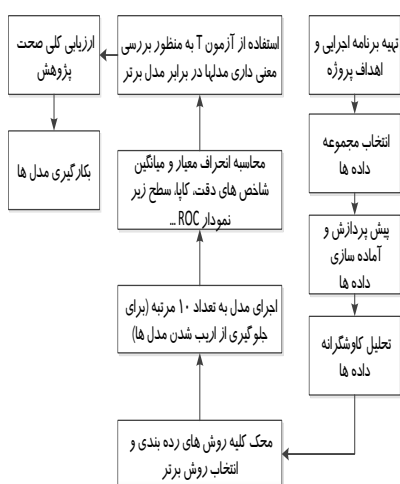
۲- روش اجرای تحقیق

برای هر تحقیقی با توجه به ابزار مورد استفاده در آن تحقیق روش اجرا وجود دارد. در تحقیق حاضر با توجه به این که داده کاوی ابزار اصلی آن است باید از استاندارد اجرای پروژه ها و تحقیقات داده کاوی استفاده نمود. یکی از معروف ترین و کاربردی ترین روش های اجرای پروژه های داده کاوی روش کریسپ است.

مدل اجرایی کریسپ یک مدل تعاملی، تکراری و چرخه ای است (Chapman, 2000; Larose, 2005). تعاملی از این جهت که مراحل مختلف در این مدل به هم وابستگی زیادی دارند و خروجی ها و ورودی های آنها کامل بر روی یکدیگر تأثیرگذار بوده و متغیرها و شاخص های خارج از مدل، عموماً تأثیری بر روی مدل اجرایی ندارند. این مدل تکراری است یعنی همان طور که در شکل (۲) نیز مشخص شده است برخی از مراحل مانند آماده سازی و مدل سازی در این فرآیند شاید بارها و بارها به منظور رسیدن به مدل هایی با دقت بیشتر تکرار شوند. همچنین در این مدل پس از این که یک بار مدل سازی انجام شد و پروژه به اتمام رسید دوباره باید با ورود داده های جدید مدل سازی را دوباره انجام داد تا در صورت به کارگیری نتایج، مدل ها با دقت و صحت بیشتری مدل سازی را در فضای واقعی انجام دهند. در ادامه مراحل این فرآیند با توجه به مشخصات، اهداف و داده های تحقیق حاضر تشریح شده است.

۳- نرم افزارهای تحقیق

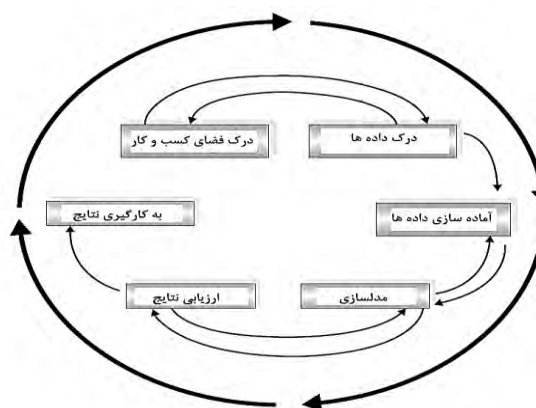
در این تحقیق از دو نرم افزار وکا و کلمنتاین استفاده شده است. روش کار این نرم افزار به صورت تولید جریان داده ای با استفاده از گره ها و ارتباط دادن بین آنها است. مدل کلی تحقیق



مدل سازی تحقیق

در این بخش ابتدا داده های تحقیق که از واحد اداری یک دانشگاه استخراج شده اند، توصیف گردید. در تشریح این داده ها از روش های آماری و همچنین روش های مصورسازی داده ها استفاده شده است و تحلیل نتایج به دست آمده نیز موجود است. سپس داده های تحقیق پاک سازی و آماده سازی شده اند و وارد مرحله مدل سازی شده اند. مدل سازی در این تحقیق با استفاده از نرم افزار وکا صورت پذیرفته است. در این راستا ابتدا مدل های مختلف شبکه های عصبی با تغییر پارامترهای اصلی بر روی مجموعه داده اعمال شده اند و بهترین مدل از میان مدل ها انتخاب شده است. در مرحله بعد نتیجه مدل برتر شبکه های عصبی با مدل رده بندی گروهی بگینگ مقایسه گردیده و در نهایت دقت آن مدل محاسبه شده است.

درک داده ها



معدل (GPA): این متغیر نشان دهنده معدل کل دوره کارشناسی داوطلب است.

کد واحد کارشناسی (Bs-Ucode): این متغیر یک متغیر عددی سه رقمی است که بیانگر کد واحدی است که داوطلب دوره کارشناسی خود را در آن گذرانده است.

کد رشته انتخابی (Ms-Course Code): متغیر کد رشته انتخابی یک متغیر عددی است که شامل ۵ رقم است و مختص رشته ای است که داوطلب آن را انتخاب کرده است.

کد واحد انتخابی (Ms-UCode): این متغیر نیز عددی است و دارای سه مقدار است که نمایانگر کد واحدی است که داوطلب تمایل به ادامه آن در مقطع کارشناسی ارشد را دارد.

نام رشته انتخابی (Ms-Course): رشته انتخابی نام متغیری است که نشان دهنده نام رشته انتخابی کارشناسی ارشد توسط داوطلب است.

نام واحد انتخابی (Ms-UName): واحد انتخابی نیز متغیری است که نام واحدی که داوطلب به عنوان واحد کارشناسی ارشد انتخاب کرده است را دارا است.

نام واحد کارشناسی (Bs-Uname): این متغیر مقداری دسته ای است که بیانگر نام واحدی است که داوطلب دوره کارشناسی خود را در آن واحد گذرانده است.

پذیرش (Admission): متغیر پذیرش نشان دهنده قبولی داوطلب در کنکور کارشناسی ارشد واحد مربوطه است. این متغیر دارای سه مقدار صفر، ۱ و ۲ است که مقدار اول به معنای عدم پذیرش و مقدار دوم پذیرش بدون سهمیه و مقدار سوم به معنای پذیرش با سهمیه است.

ثبت نام (Registration): متغیر ثبت نام متغیری است که نمایانگر ثبت نام در صورت پذیرش در رشته ای مورد نظر است.

تشریح داده ها

درک داده ها عبارت است از جمع آوری داده های اولیه، توصیف داده ها، بازرسی و بررسی داده ها و اعتبارسنجی کی فیت داده ها. کارایی داده کاوی مستقیماً با داده های مورد استفاده مرتبط است. هر اندازه داده ها دقیق تر، جامع تر و با کیفیت بالاتری باشند خروجی داده کاوی کارا تر خواهد بود. داده ای با کیفیت خوانده می شود که صحیح، کامل، سازگار، به روز، قابل قبول، با ارزش، قابل تفسیر و در دسترس باشد. از این رو انتخاب و جمع آوری داده های درست و یکپارچه ساز قالب آن ها برای استفاده در داده کاوی از اهمیت بالایی برخوردار است.

متغیرهای موجود در مجموعه داده در زیر شرح داده شده اند:

شماره پرونده (File-Number): این متغیر نمایانگر شماره پرونده دانشجویانی است که دانشگاه مورد نظر را برای ادامه تحصیل در رشته کارشناسی ارشد انتخاب کرده اند.

سال تولد (YBirth): متغیر سال تولد بیانگر دو رقم سال تولد دانشجویان است.

جنسیت (Sex): این متغیر جنس شرکت کنندگان در کنکور را نمایش می دهد.

وضعیت داوطلب (Candidate-Status): وضعیت داوطلب متغیری است که دارای ۹ مقدار است و این مقادیر عبارت اند از: (۱) رزمنده، (۲) بسیجی فعال، (۳) جهادگر، (۴) آزاده، (۵) جانباز، (۶) همسر و فرزند جانباز بالای ۵۰٪، (۷) همسر و فرزند جانباز بین ۲۵٪ - ۴۰٪، (۸) فرزند شهید و (۰) عادی.

کد پستی (PostCode): این متغیر مقدار عددی ۱۰ رقمی است که مربوط به کد پستی محل زندگی داوطلب است.

آدرس (Address): این متغیر شامل آدرس محل زندگی داوطلب است.

کد استان (City-Code): کد استان متغیری عددی و ۲ رقمی است که نشان دهنده کد استانی است که داوطلب در آن زندگی می کند.

در این قسمت مدل سازی شامل دو بخش اصلی است. با توجه به عنوان این پژوهش در بخش اول با استفاده از شبکه های عصبی مصنوعی ابتدا مدل سازی بر روی داده های تحقیق انجام گرفت. در این بخش پارامترهای مختلف شبکه های عصبی به منظور دستیابی به بهینه ترین پارامترها تغییر داده شده اند. ملاک یافتن بهترین پارامترها در این قسمت افزایش شاخص های کارایی مدل ها است. در بخش دوم مدل بهینه ای که در قسمت اول به دست آمد، با دیگر مدل رده بندی کننده مانند مدل یادگیری جمعی بوستینگ مقایسه شده و مدل های رده بندی معمول مانند بیز ساده و درخت تصمیم مقایسه شده و مدل برتر انتخاب شده است.

انتخاب پارامترهای بهینه شبکه

توجه به این که شبکه های عصبی مصنوعی دارای پارامترهای متعددی به منظور تنظیم نمودن می باشند، لازم است بهترین پارامترها به منظور رسیدن به شبکه بهینه پیدا شود. با توجه به تحقیق منتشر شده توسط Swingler در سال ۱۹۹۶، تعداد گره های موجود در لایه پنهان نباید از دو برابر تعداد گره های ورودی یا همان متغیرهای ورودی بیش تر باشند. بدین ترتیب

در این مرحله از پردازش داده ها هدف استفاده از روش های آماری تک متغیره و چند متغیره و همچنین روش های مصورسازی به منظور به تصویر کشیدن مشخصات هر متغیر و همچنین روابط بین متغیرها است. گزارشات استخراج شده در این مرحله می تواند به منظور تشریح هرچه بهتر نتایج حاصل از مدلسازی کمک نماید. تمامی گزارشات این قسمت با استفاده از نرم افزار داده کاوی کلمنتاین انجام شده است.

پیش پردازش داده ها

در این قسمت از تحقیق باید از روش های مختلف به منظور پاک سازی داده های اشتباه یا دور افتاده استفاده نمود. ابتدا در این قسمت باید داده های دور افتاده را شناسایی نمود و آن ها را حذف کرد تا مدل را تحت تأثیر خود قرار ندهند. بدین ترتیب مقادیر از دست رفته در داده ها ایجاد خواهند شد که باید این مقادیر با مقادیر مناسبی جایگزین گردند. پس از پاک سازی داده ها باید متغیرهای تأثیرگذار بر روی متغیر هدف را شناسایی کرد و فقط آنها را وارد مدل سازی نمود. سپس به منظور شناسایی نقاط دور افتاده از روش نمودارهای کنترلی شوهارت براساس فرمول زیر استفاده شده است:

$$\mu \pm a \cdot \sigma = \text{محدوده کنترلی}$$

در این معادله μ برابر میانگین و σ برابر با انحراف معیار متغیر مورد نظر است. مقدار σ در آن برابر با عدد ۳ به منظور کشف تعداد نقاط دور افتاده و عدد ۵ به منظور کشف تعداد نقاط بسیار دور افتاده تنظیم شده است. نتیجه حاصل از استفاده این روش بر روی داده های تحقیق در جدول زیر قابل مشاهده است.

تعداد نقاط خیلی دور افتاده	تعداد نقاط دور افتاده	نوع متغیر	نام متغیر
0	905	Continuous	GPA
13	250	Continuous	Age

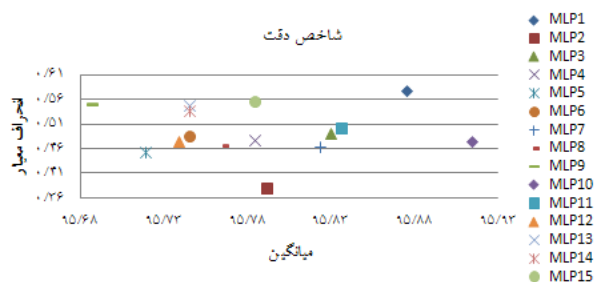
جدول ۱: شاخص های آماری متغیرهای پیوسته مورد استفاده در مدلسازی

مدل سازی:

مدل	نرخ یادگیری	تعداد لایه های پنهان	میانگین		انحراف معیار	
			میانگین	انحراف معیار	میانگین	انحراف معیار
MLP1	0.3	5	0.675	0.077	0.633	0.07
MLP2	0.35	5	0.751	0.370	0.623	0.051
MLP3	0.4	5	0.83	0.49	0.623	0.064
MLP4	0.3	5	0.784	0.477	0.626	0.069
MLP5	0.35	5	0.719	0.452	0.619	0.055
MLP6	0.4	5	0.745	0.487	0.62	0.062
MLP7	0.3	10	0.823	0.163	0.627	0.062
MLP8	0.35	10	0.765	0.364	0.627	0.067
MLP9	0.4	10	0.867	0.55	0.614	0.064
MLP10	0.3	12	0.914	0.474	0.636	0.045
MLP11	0.35	12	0.886	0.503	0.631	0.042
MLP12	0.4	12	0.729	0.472	0.612	0.077
MLP13	0.3	14	0.745	0.449	0.612	0.067
MLP14	0.35	14	0.745	0.357	0.617	0.061
MLP15	0.4	14	0.794	0.377	0.62	0.065

تصمیم گرفته شد با توجه به این که تعداد متغیرهای ورودی ۸ عدد است، تعداد گره های لایه پنهان از مقداری کم تر از متغیرهای ورودی یعنی ۶ عدد شروع شده و با افزایش دو واحدی تا عدد ۱۴ تغییر کرده و نتایج حاصل از شبکه ثبت شده است. از طرفی دیگر در هر مرحله افزایش، سه مقدار برای یادگیری در نظر گرفته شده است که این مقادیر عبارتند از : ۰/۳ ، ۰/۳۵ و ۰/۴. در جدول زیر نتایج اجرای آموزش در شبکه های عصبی وجود دارد. همان طور که گفته شد، با توجه به پارامترهای تنظیم شده در این مرحله ۱۵ مدل تولید شده است. همچنین به منظور اعتبارسنجی داده ها در این قسمت از روش ۱۰ دسته اعتبارسنجی استفاده شده است. در این روش ابتدا ۹۰٪ از داده ها به عنوان داده های آموزشی و ۱۰٪ به عنوان داده های آزمایشی تعیین می گردد. سپس در مرحله بعد ۱۰٪ قبلی که به منظور آموزش به کار گرفته شده، وارد مجموعه داده آموزشی شده و ۱۰٪ از داده های آموزشی که در مرحله قبل استفاده شده است به منظور آزمایش به کار می رود.

در جدول فوق دو پارامتر تنظیم شده نرخ یادگیری و تعداد گره های لایه پنهان برای هر یک از مدل های ساخته شده شبکه های عصبی ذکر شده است. با توجه به این که هر مدل به تعداد ۱۰ بار اجرا شده است، در هر بار اجرا نتایج هر مدل تغییراتی هر چند ناچیز خواهند داشت. در بین این مدل های ساخته شده مدلی به عنوان مدل برتر با پارامترهای بهینه انتخاب خواهد شد که دارای بیش ترین میانگین و کم ترین انحراف معیار برای هر شاخص باشد. به منظور بررسی هر یک از مدل های ساخته شده از نمودارهای پراکنندگی زیر برای میانگین و انحراف معیار با شاخص دقت استفاده شده است.

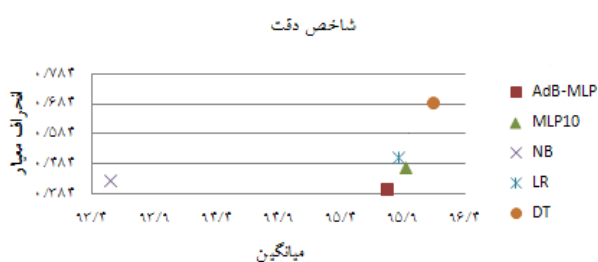


شکل ۳: نمودار پراکنندگی میانگین و انحراف معیار برای شاخص دقت

با توجه به شکل فوق مشخص می گردد که دسته کننده MLP10 به عنوان بهترین دسته کننده از نظر شاخص دقت با مقدار میانگین ۹۵/۹۱۴ و انحراف معیار ۰/۴۷۴ انتخاب

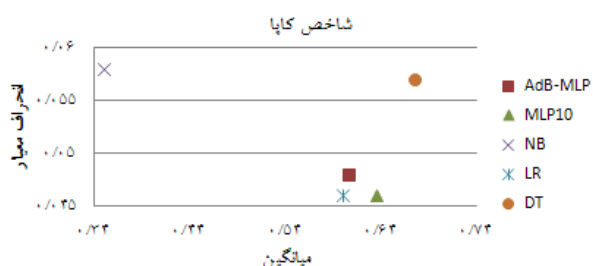
¹ 10 Cross-fold-validation

عصبی، بیز ساده و درختان تصمیم استفاده شده است. لازم به ذکر است که روش بوستینگ جزء روش های رده بندی جمعی می باشند که انتظار می رود در این قسمت عملکرد بهتری را نسبت به مدل پایه شبکه های عصبی ارائه دهند. در نمودار زیر شاخص دقت تمامی مدل ها با توجه به میانگین و انحراف معیار آن ها ترسیم شده اند.



شکل ۵: نمودار پراکندگی میانگین و انحراف معیار برای شاخص دقت مقایسه ۵ مدل با یکدیگر

همان طور که در نمودار نیز مشخص است روش MLP10 و درخت تصمیم در مقایسه با سایر روش ها دارای میانگین بیشتری می باشند. بین این دو روش نیز روش MLP10 با مقدار کمتر انحراف معیار به عنوان روش بهتر قابل شناسایی است. از طرفی دیگر مشخص است که روش بیز ساده از نظر شاخص دقت دارای عملکرد کاملاً ضعیف تری نسبت به ۴ مدل دیگر بوده است. در ادامه نمودار کاپا برای مقایسه دقیق تر نتایج مدل ها با یکدیگر ترسیم شده است.



شکل ۶: نمودار پراکندگی میانگین و انحراف معیار برای شاخص دقت مقایسه ۵ مدل با یکدیگر

با توجه به شکل مدل درخت تصمیم دارای بیشترین میانگین است در حالی که مدل MLP10 با انحراف معیار کم تر و میانگین نسبتاً مشابه در رده بهتری قرار گرفته است.

می گردد. به دنبال آن دسته کننده MLP1 با میانگین ۸۷۵/۹۵ و انحراف معیار ۰/۵۷۷ بهترین دقت را از نظر این شاخص کسب کرده است. از طرفی دسته کننده MLP9 بدترین نتیجه را در بین مدل های ساخته شده به دست آورده است.

در جدول زیر میانگین و انحراف معیار شاخص کاپا برای هر کدام از مدل ها در یک نمودار پراکندگی رسم شده است. این شاخص نسبت احتمال رده بندی توسط مدل را به احتمال رده بندی تصادفی محاسبه می کند. مقدار بیشتر این شاخص از عدد صفر نمایان گر عملکرد بهتر مدل رده بندی نسبت به حالت تصادفی است.



شکل ۴: نمودار پراکندگی میانگین و انحراف معیار برای شاخص کاپا

با توجه به شکل بالا مدل MLP10 به عنوان مدل برتر شناخته می شود. هم چنین مدل های MLP1 و MLP11 نیز به عنوان مدل های دوم و سوم انتخاب شده اند. مدل های MLP14 و نیز MLP5 نیز با کمترین مقدار میانگین به عنوان مدل های ضعیف از دید این شاخص انتخاب می شوند. با دقت در نمودار رسم شده فوق مشخص است که باز هم مدل های MLP1 و MLP10 با اختلاف نسبتاً زیادی از سایر مدل ها نتیجه بهتری را کسب کرده اند. از طرفی مدل های MLP12 و MLP13 نیز به عنوان ضعیف ترین مدل ها انتخاب شده اند.

مقایسه مدل بهینه با سایر مدل ها

با انتخاب مدل بهینه MLP10 به عنوان بهترین مدل از مرحله قبل، در این مرحله از برخی از مدل های رده بندی به منظور مقایسه کارایی این مدل ها با مدل شبکه عصبی انتخاب شده، استفاده شده است. بدین ترتیب از مدل های بوستینگ شبکه های

4. Romero, C., & Ventura, S. (2007). Educational data mining: A survey from 1995 to 2005. *Expert Systems with Applications*, 33(1), 135-146.

5. Arruabarrena, R., Pérez, T. A., López-Cuadrado, J., Gutiérrez, J., & Vadillo, J. A. (2006, January). On evaluating adaptive systems for education. In *Adaptive Hypermedia and Adaptive Web-Based Systems* (pp. 363-367). Springer Berlin Heidelberg.

با توجه به نتایج در شکل بالا مشخص است که MLP10 دارای عملکرد بهتری نسبت به سایر مدلها میباشد. نکته قابل توجه این نمودار مقدار شاخص برای مدل بیز ساده است که با اختلاف زیادی ضعیفتر از این مدلها عمل کرده است.

لبتوجه به عملکرد بهتر روش MLP10 در شاخص دقت و قرار گرفتن این مدل در گوشه ایترین قسمت در سمت راست نمودارها می توان نتیجه گرفت که مدل فوق برتری نسبی نسبت به سایر مدلها در بحث پیش بینی را داشته است و در مقایسه با سایر روشها از عملکرد قابل قبول تری برخوردار بوده است.

نتیجه گیری

نتایج حاصل از این تحقیق در پیش بینی ثبت نام در دانشگاه ها بسیار مؤثر خواهد بود و تاکنون هیچ تحقیقی در این زمینه با استفاده از تکنیک های داده کاوی در ایران انجام نشده است. هم چنین توجه به این که داده ها مستقل از نوع آموزشگاه و بر پایه اطلاعات اصلی متقاضیان انتخاب شده است قابل تعمیم برای هر یک از آموزشگاهها و مؤسسات آموزش عالی است.

1. Rad, A., Naderi, B., & Soltani, M. (2011). Clustering and ranking university majors using data mining and AHP algorithms: A case study in Iran. *Expert Systems with Applications*, 38(1), 755-763.

2. Ramaswami, M., & Bhaskaran, R. (2009). A study on feature selection techniques in educational data mining. arXiv preprint arXiv:0912.3924.

3. Thomas, J., Chongwatpol, J., Pengnate, F., & Hass, M. (2011). *Data Mining in Higher Education: University Student Declaration of Major*.